

# A SUBSPACE ACCELERATION METHOD FOR MINIMIZATION INVOLVING A GROUP SPARSITY-INDUCING REGULARIZER\*

FRANK E. CURTIS , YUTONG DAI , AND DANIEL P. ROBINSON†

**Abstract.** We consider the problem of minimizing an objective function that is the sum of a convex function and a group sparsity-inducing regularizer. Problems that integrate such regularizers arise in modern machine learning applications, often for the purpose of obtaining models that are easier to interpret and that have higher predictive accuracy. We present a new method for solving such problems that utilize subspace acceleration, domain decomposition, and support identification. Our analysis shows, under common assumptions, that the iterate sequence generated by our framework is globally convergent, converges to an  $\epsilon$ -approximate solution in at most  $O(\epsilon^{-(1+p)})$  (respectively,  $O(\epsilon^{-(2+p)})$ ) iterations for all  $\epsilon$  bounded above and large enough (respectively, all  $\epsilon$  bounded above) where  $p > 0$  is an algorithm parameter, and exhibits superlinear local convergence. Preliminary numerical results for the task of binary classification based on regularized logistic regression show that our approach is efficient and robust, with the ability to outperform a state-of-the-art method.

**Key words.** nonlinear optimization, convex optimization, worst-case iteration complexity, regularization methods, group regularizer, sparsity, logistic regression, subspace acceleration

**AMS subject classifications.** 49M37, 65K05, 65K10, 65Y20, 68Q25, 90C30, 90C60

**1. Introduction.** We consider the minimization of a function that may be written as the sum of a convex function and a nonoverlapping group sparsity-inducing regularizer. Specifically, given a convex and twice continuously differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , a collection of  $n_{\mathcal{G}} > 0$  nonoverlapping groups  $\mathcal{G} := \{\mathcal{G}_i\}_{i=1}^{n_{\mathcal{G}}}$  that forms a partition of  $\{1, 2, \dots, n\}$  (i.e.,  $\mathcal{G}_i \cap \mathcal{G}_j = \emptyset$  for all  $i \neq j$  and  $\cup_{i=1}^{n_{\mathcal{G}}} \mathcal{G}_i = \{1, 2, \dots, n\}$ ), and group-wise weighting parameters  $\{\lambda_i\}_{i=1}^{n_{\mathcal{G}}} > 0$ , our algorithm solves the problem

$$(1.1) \quad \min_{x \in \mathbb{R}^n} \{f(x) + r(x)\}, \text{ where } r(x) := \sum_{i=1}^{n_{\mathcal{G}}} \lambda_i \|[x]_{\mathcal{G}_i}\|_2$$

and  $[x]_{\mathcal{G}_i}$  is the subvector of  $x$  corresponding to elements in  $\mathcal{G}_i$ . The regularizer  $r$  generalizes the  $\ell_1$ -norm, which is recovered by choosing  $\mathcal{G}_i = \{i\}$  for all  $i \in \{1, 2, \dots, n\}$ .

Despite the successes of  $\ell_1$ -norm regularization, its inadequacy in the context of many modern machine learning applications has been noticed by researchers, and is one motivation for the use of group regularization. In some machine learning applications the covariates come in groups (e.g., genes that regulate hormone levels in microarray data [23]), in which case one may wish to select them jointly. Also, integrating group information into the modeling process can improve both the interpretability and accuracy [35] of the resulting model. Yuan and Lin [34] observed that in the multi-factor analysis-of-variance problem, where each factor is expressed through a set of dummy variables, deleting an irrelevant factor is equivalent to deleting a *group* of dummy variables; the  $\ell_1$ -norm regularizer fails to achieve this goal.

**1.1. State-of-the art methods.** There is a long history of algorithms for solving regularized problems of the form (1.1) (see [1] and the references therein). Here, we review some of the state-of-the-art approaches for solving sparsity-promoting problems that are most closely related to our proposed approach.

---

\*This material is based upon work supported by the U.S. National Science Foundation under the Division of Computing and Communication Foundations (Award Number CCF-1740796).

†Department of Industrial and Systems Engineering, Lehigh University, Bethlehem, PA, USA; E-mail: [frank.e.curtis@gmail.com](mailto:frank.e.curtis@gmail.com), [yud319@lehigh.edu](mailto:yud319@lehigh.edu), [daniel.p.robinson@gmail.com](mailto:daniel.p.robinson@gmail.com)

41 **First-order methods.** Proximal methods are designed to solve problems of the  
 42 form (1.1) and have received attention in the machine learning community [3, 7, 31].  
 43 A well-known example for  $\ell_1$ -norm regularized problems is the iterative shrinkage-  
 44 thresholding algorithm (ISTA), which is obtained by applying a proximal gradient  
 45 (PG) iteration to minimize a smooth function plus the  $\ell_1$ -norm regularizer [10, 12].  
 46 Under certain assumptions, one can prove a worst-case complexity bound on the num-  
 47 ber of iterations required by the PG method before it correctly identifies the support  
 48 of the optimal solution [28]. Combined with the acceleration technique proposed  
 49 by Nesterov [26, 27], one obtains the algorithm FISTA [3]. One obtains a related,  
 50 but distinct approach from ISTA by posing an equivalent smooth reformulation of  
 51 the problem—separating the positive and negative parts of the variables—and apply-  
 52 ing a gradient projection method to the resulting formulation [13, 15]. All of these  
 53 approaches have been shown to work well in practice, at least compared to other  
 54 first-order methods such as the subgradient algorithm. However, these algorithms  
 55 are often inferior in practice compared to alternative approaches that employ space  
 56 decomposition techniques and/or second-order derivatives [5, 6, 18].

57 As an alternative to PG and gradient projection techniques, researchers have con-  
 58 sidered (block) coordinate descent for solving  $\ell_1$ -norm regularized problems. Such a  
 59 strategy is appealing, since when minimizing an  $\ell_1$ -norm regularized objective along  
 60 coordinate directions, it is common that the objective is minimized with variables be-  
 61 ing zero. These approaches are also easy to implement to exploit parallel computing;  
 62 see, e.g., the accelerated randomized proximal coordinate gradient method in [20], the  
 63 parallel coordinate descent methods in [29], and the asynchronous coordinate descent  
 64 technique in [22]. A downside of these approaches is that the space decomposition is  
 65 performed in a prescribed manner, rather than in an adaptive way that can benefit  
 66 from information acquired during the solution process. Also, these approaches do not  
 67 effectively exploit second-order derivative information and require exact minimization  
 68 along coordinate directions. An exception to this latter criticism is the inexact coord-  
 69 inate descent algorithm from [30], although this approach does not effectively exploit  
 70 second-order derivatives and uses a prescribed space decomposition strategy.

71 Various other approaches have been proposed for solving problems involving spe-  
 72 cific regularizers. In [21], the authors discuss various methods for sparse learning  
 73 that make use of projection techniques. A well-known package is GLMNET [16], which  
 74 is designed for solving problems with the elastic-net regularization. Finally, let us  
 75 mention the work in [32], which proposes and tests a groupwise-majorization-descent  
 76 algorithm (called `gglasso`) for solving problems involving the group- $\ell_1$ -norm regular-  
 77 izer. A potential downside of this approach is that it updates variables by groups in  
 78 a cycle, rather than by using an adaptive space decomposition technique.

79 **Second-order methods.** Relatively few second-order methods have been pro-  
 80 posed for minimizing sparsity-promoting objective functions. In [17], an acceler-  
 81 ated regularized Newton scheme is proposed. A similar proximal-Newton method is  
 82 proposed in [19], which under some assumptions can be shown to converge locally  
 83 superlinearly. These approaches can be effective in practice, although they appear  
 84 to lack good worst-case guarantees in terms of identification of the optimal solution  
 85 support. Other approaches, such as the orthant-based method in [18], can predict the  
 86 solution support, but in practice are often outperformed by a closely related method  
 87 called `FaRSA` [5, 6]. As for publicly available solvers based on second-order methods,  
 88 most have been designed for specific loss functions and regularizers. For example,  
 89 `newGLMNET` in [33] is designed for  $\ell_1$ -regularized logistic regression and the method in  
 90 [14] is designed for regularized logistic regression and support vector machines.

91 **1.2. Notation and assumptions.** Let  $\mathbb{R}$  denote the set of real numbers,  $\mathbb{R}^n$   
92 denote the set of  $n$ -dimensional real vectors, and  $\mathbb{R}^{m \times n}$  denote the set of  $m$ -by- $n$ -  
93 dimensional real matrices. The set of natural numbers is denoted as  $\mathbb{N} := \{0, 1, 2, \dots\}$ .  
94 For any set  $\mathcal{I} \subseteq \{1, 2, \dots, n\}$ , we define the projection of  $x \in \mathbb{R}^n$  onto the subspace  
95 spanned by the coordinate vectors indexed by the entries of  $\mathcal{I}$  as  $P_{\mathcal{I}}(x)$ , so that

$$96 \quad (1.2) \quad [P_{\mathcal{I}}(x)]_i := \begin{cases} x_i & \text{if } i \in \mathcal{I}, \\ 0 & \text{if } i \notin \mathcal{I}. \end{cases}$$

For a function  $h : \mathbb{R}^n \rightarrow \mathbb{R}$ , vector  $x \in \mathbb{R}^n$ , and direction  $d \in \mathbb{R}^n$ , the directional derivative of  $h$  at  $x$  in the direction  $d$  is defined as the following limit:

$$D_h(x; d) := \lim_{t \searrow 0} \frac{h(x + td) - h(x)}{t}.$$

97 The following assumption is assumed to hold throughout the paper.

98 **ASSUMPTION 1.1.** *The function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  used in the definition of the objective*  
99 *function of problem (1.1) is convex and continuously differentiable. It follows that*  
100 *there exists a constant  $L_f$  such that  $\|\nabla f(x)\|_2 \leq L_f$  for all  $x \in \mathcal{L} := \{x \in \mathbb{R}^n :$   
101  $f(x) + r(x) \leq f(x_0) + r(x_0)\}$  for any initial estimate  $x_0$  of a solution to problem (1.1).*  
102 *The objective function  $f + r$  is bounded below and the gradient function  $\nabla f$  is Lipschitz*  
103 *continuous on  $\mathcal{L}$  with Lipschitz constant  $L_g$ .*

104 **1.3. Organization.** In Section 2, we present preliminary results related to PG  
105 calculations. In Section 3, by using PG-calculations as a starting point, we pro-  
106 pose a reduced-space second-order domain decomposition algorithm for solving prob-  
107 lem (1.1). The algorithm is analyzed in Section 4 and numerical results are presented  
108 in Section 5. Finally, in Section 6, we provide concluding remarks.

109 **2. Preliminaries.** In this section, we discuss preliminary material related to the  
110 objective function  $f + r$  and its associated PG calculations. (All proofs may be found  
111 in Appendix A.) For any  $\bar{x} \in \mathbb{R}^n$  and  $\bar{\alpha} > 0$ , we define the PG *update* as

$$112 \quad (2.1) \quad T(\bar{x}, \bar{\alpha}) := \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ \frac{1}{2\bar{\alpha}} \|x - (\bar{x} - \bar{\alpha} \nabla f(\bar{x}))\|_2^2 + r(x) \right\}$$

113 and the associated PG *step* as

$$114 \quad (2.2) \quad s(\bar{x}, \bar{\alpha}) := T(\bar{x}, \bar{\alpha}) - \bar{x}.$$

115 The next result shows that the directional derivative of  $f + r$  along the PG step is  
116 negative with magnitude proportional to the squared norm of the PG direction.

**LEMMA 2.1.** *For any  $\bar{x} \in \mathbb{R}^n$  and  $\bar{\alpha} > 0$ , the PG step  $s(\bar{x}, \bar{\alpha})$  in (2.2) satisfies*

$$D_{f+r}(\bar{x}; s(\bar{x}, \bar{\alpha})) \leq -\frac{1}{\bar{\alpha}} \|s(\bar{x}, \bar{\alpha})\|_2^2.$$

117 The PG update defined in (2.1) can be computed group-wise for each  $\mathcal{G}_i \in \mathcal{G}$  by

$$118 \quad (2.3) \quad \begin{aligned} [T(\bar{x}, \bar{\alpha})]_{\mathcal{G}_i} &= \left[ \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ \frac{1}{2\bar{\alpha}} \|x - (\bar{x} - \bar{\alpha} \nabla f(\bar{x}))\|_2^2 + \sum_{i=1}^{n_{\mathcal{G}}} \lambda_i \|x_{\mathcal{G}_i}\|_2 \right\} \right]_{\mathcal{G}_i} \\ &= \max \left\{ 1 - \frac{\bar{\alpha} \lambda_i}{\|[\bar{x}]_{\mathcal{G}_i} - \bar{\alpha} \nabla_{\mathcal{G}_i} f(\bar{x})\|_2}, 0 \right\} \left( [\bar{x}]_{\mathcal{G}_i} - \bar{\alpha} \nabla_{\mathcal{G}_i} f(\bar{x}) \right). \end{aligned}$$

119 Combining this observation with Lemma 2.1 leads to the following corollary, which will  
120 be relevant to the manner in which we design the algorithm we propose in Section 3.

121 LEMMA 2.2. For any  $\bar{x} \in \mathbb{R}^n$ ,  $\bar{\alpha} > 0$ , and set  $\mathcal{I}$  equal to the union of a subset of  
 122  $\{\mathcal{G}_i\}_{i=1}^{n_g}$ , the PG step  $s(\bar{x}, \bar{\alpha})$  defined in (2.2) satisfies

$$123 \quad (2.4) \quad D_{f+r}(\bar{x}; P_{\mathcal{I}}(s(\bar{x}, \bar{\alpha}))) \leq -\frac{1}{\bar{\alpha}} \|P_{\mathcal{I}}(s(\bar{x}, \bar{\alpha}))\|_2^2$$

124 where the projection operator  $P_{\mathcal{I}}$  is defined through (1.2).

125 Our next result quantifies the decrease in  $f + r$  that one can expect to obtain by  
 126 taking a PG step  $s(\bar{x}, \bar{\alpha})$ , provided the PG parameter  $\bar{\alpha}$  is sufficiently small.

LEMMA 2.3. For any  $\bar{x} \in \mathbb{R}^n$ ,  $\bar{\alpha} \in (0, 2/L)$ , and  $\mathcal{I}$  equal to the union of a subset  
 of  $\{\mathcal{G}_i\}_{i=1}^{n_g}$ , the objective function decrease satisfies

$$f(\bar{x} + P_{\mathcal{I}}(\bar{x}, \bar{s})) + r(\bar{x} + P_{\mathcal{I}}(\bar{x}, \bar{s})) \leq f(\bar{x}) + r(\bar{x}) - \left(\frac{1}{\bar{\alpha}} - \frac{L}{2}\right) \|P_{\mathcal{I}}(s(\bar{x}, \bar{\alpha}))\|_2^2.$$

127 The next result shows that, when restricted to certain groups, the size of the PG  
 128 step is bounded above by the gradient of the objective function.

LEMMA 2.4. If the pair  $(\bar{x}, \bar{\alpha})$  and group  $\mathcal{G}_i$  satisfy  $\bar{\alpha} \in (0, 1]$ ,  $[\bar{x}]_{\mathcal{G}_i} \neq 0$ , and  
 $[\bar{x} + s(\bar{x}, \bar{\alpha})]_{\mathcal{G}_i} \neq 0$ , where  $s(\bar{x}, \bar{\alpha})$  is defined in (2.2), then

$$\|\nabla_{\mathcal{G}_i}(f + r)(\bar{x})\|_2 \geq \|[s(\bar{x}, \bar{\alpha})]_{\mathcal{G}_i}\|_2.$$

129 With the preliminaries now completed, we can propose our new algorithm.

130 **3. Proposed Algorithm Framework.** We propose Algorithm 3.1, which we  
 131 call **FaRSA-Group** (Fast Reduced-Space Algorithm for Group sparsity-inducing regular-  
 132 ization) for solving problem (1.1) that uses ideas related to domain decomposition,  
 133 subspace acceleration, and support identification. An overview of the algorithm is  
 134 described in Section 3.1. During each iteration of our method, at least one of three  
 135 subroutines is called. The three subroutines are described in Sections 3.2–3.4.

136 **3.1. Main algorithm (Algorithm 3.1).** Our main algorithm is formally stated  
 137 as Algorithm 3.1. At the beginning of the  $k$ th iteration,  $x_k$  and  $\alpha_k > 0$  denote the  
 138 current solution estimate for problem (1.1) and the PG parameter, respectively. We  
 139 then compute  $s_k$  in Line 5 as the PG step associated with problem (1.1), namely,

$$140 \quad (3.1) \quad s_k := s(x_k, \alpha_k) \quad \text{with } s(x_k, \alpha_k) \text{ defined in (2.2).}$$

141 Although the repeated computation of PG steps is the basis for a first-order method,  
 142 here we primarily use it to *predict* the zero/nonzero structure of a solution and to  
 143 formulate optimality measures. Specifically, in Line 6 we compute the index set

$$144 \quad (3.2) \quad \bar{\mathcal{I}}_k^{\text{cg}} := \{j \in \mathcal{G}_i : [x_k]_{\mathcal{G}_i} \neq 0, [x_k + s_k]_{\mathcal{G}_i} \neq 0, \text{ and} \\ \|[x_k]_{\mathcal{G}_i}\|_2 \geq \kappa_1 \|\nabla_{\mathcal{G}_i}(f + r)(x_k)\|_2\}$$

145 for some  $\kappa_1 \in (0, \infty)$ . The groups of variables that compose  $\bar{\mathcal{I}}_k^{\text{cg}}$  are *candidates* for use  
 146 in a Newton-type calculation aimed to accelerated convergence. Before using them,  
 147 however, we first check to see if each candidate block is sufficiently far from zero, and  
 148 those that are not are removed. Specifically, we first define

$$149 \quad (3.3) \quad \mathcal{I}_k^{\text{small}} := \{j \in \mathcal{G}_i : \mathcal{G}_i \subseteq \bar{\mathcal{I}}_k^{\text{cg}} \text{ and } \|[x_k]_{\mathcal{G}_i}\|_2 < \kappa_2 \|\nabla_{\bar{\mathcal{I}}_k^{\text{cg}}}(f + r)(x_k)\|_2^p\}$$

150 for some  $\{\kappa_2, p\} \subset (0, \infty)$ , and then define in Line 7 the sets and optimality measures

$$151 \quad (3.4) \quad \left\{ \begin{array}{l} \mathcal{I}_k^{\text{cg}} := \bar{\mathcal{I}}_k^{\text{cg}} \setminus \mathcal{I}_k^{\text{small}} \\ \mathcal{I}_k^{\text{pg}} := \{1, 2, \dots, n\} \setminus \mathcal{I}_k^{\text{cg}} \end{array} \right\} \quad \text{and} \quad \left\{ \begin{array}{l} \chi_k^{\text{cg}} := \|[s_k]_{\mathcal{I}_k^{\text{cg}}}\|_2 \\ \chi_k^{\text{pg}} := \|[s_k]_{\mathcal{I}_k^{\text{pg}}}\|_2 \end{array} \right\}$$

152 where by convention  $\|[\cdot]_{\emptyset}\|_2 = 0$ . (See Lemma 4.1 for a justification that these sets  
 153 together represent a measure of optimality.) This construction of sets also ensures  
 154 that the subvector of  $x_k$  that corresponds to  $\mathcal{G}_i$  for each  $\mathcal{G}_i \subseteq \mathcal{I}_k^{\text{CG}}$  is at least a distance

$$155 \quad (3.5) \quad \rho_{k,i} := \max\{\kappa_1 \|\nabla_{\mathcal{G}_i}(f+r)(x_k)\|_2, \kappa_2 \|\nabla_{\mathcal{I}_k^{\text{CG}}}(f+r)(x_k)\|_2^p\}$$

156 away from zero (see Lemma 4.5(i)), which is crucial in our analysis.

157 Armed with  $\chi_k^{\text{PG}}$  and  $\chi_k^{\text{CG}}$ , Algorithm 3.1 seeks decrease in the objective function  
 158 in a subspace that is likely to allow for significant progress. We consider two cases.

159 **Case 1: the condition  $\chi_k^{\text{PG}} \leq \chi_k^{\text{CG}}$  checked in Line 8 holds.** In this case, the  
 160 inequality  $\chi_k^{\text{PG}} \leq \chi_k^{\text{CG}}$  indicates that significant reduction in the objective function can  
 161 be achieved by focusing on variables in the set  $\mathcal{I}_k^{\text{CG}}$ . Therefore, in Line 9 we choose any  
 162 index set  $\mathcal{I}_k$  that is (i) a subset of  $\mathcal{I}_k^{\text{CG}}$ , (ii) equal to the union of some subset of groups  
 163 from  $\mathcal{G}$ , and (iii) the size of the PG step restricted to the index set  $\mathcal{I}_k$  is at least a  
 164 fraction of the size of the PG step when restricted to the index set  $\mathcal{I}_k^{\text{CG}}$ . The easiest  
 165 choice that satisfies these conditions is  $\mathcal{I}_k \equiv \mathcal{I}_k^{\text{CG}}$ , but for large-scale problems it may  
 166 be beneficial to restrict  $|\mathcal{I}_k|$ . The opposite extreme choice is selecting  $\mathcal{I}_k$  as the group  
 167  $\mathcal{G}_i$  contained in  $\mathcal{I}_k^{\text{CG}}$  with largest associated PG step, in which case one would choose  
 168  $\varphi = 1/\sqrt{n_{\mathcal{G}}}$  for the user-defined parameter in Line 9. Once  $\mathcal{I}_k$  has been selected, a  
 169 *reduced-space* gradient  $g_k$  and *reduced-space* positive-definite matrix  $H_k$  is computed  
 170 in Line 10, where the derivatives are taken with respect to variables in  $\mathcal{I}_k$ . (In practice,  
 171  $H_k$  could be selected based on  $\nabla_{\mathcal{I}_k}^2(f+r)(x_k)$  to ensure a fast local convergence  
 172 rate.) Note that such derivatives exist since by construction  $\mathcal{I}_k \subseteq \mathcal{I}_k^{\text{CG}} \subseteq \bar{\mathcal{I}}_k^{\text{CG}}$ , and  
 173 from (3.2) the objective function  $f+r$  is differentiable with respect to groups of  
 174 variables in  $\bar{\mathcal{I}}_k^{\text{CG}}$ . Next,  $g_k$  and  $H_k$  are used to compute a direction  $\bar{d}_k$  of sufficient  
 175 descent for  $f+r$  by calling the subroutine CG\_DIRECTION (see Section 3.2). Once a  
 176 full-space vector  $d_k$  is obtained by padding  $\bar{d}_k$  with zeros in Line 12, a *projected* line  
 177 search is performed by calling subroutine UPDATE\_CG in Line 13 (see Section 3.3).

178 **Case 2: the condition  $\chi_k^{\text{PG}} \leq \chi_k^{\text{CG}}$  checked in Line 8 does not hold.** In this case,  
 179 the inequality  $\chi_k^{\text{PG}} > \chi_k^{\text{CG}}$  indicates that significant reduction in the objective function  
 180 can be achieved by focusing on variables in the set  $\mathcal{I}_k^{\text{PG}}$ . Therefore, in Line 16, we  
 181 choose any index set  $\mathcal{I}_k$  that is (i) a subset of  $\mathcal{I}_k^{\text{PG}}$ , (ii) equal to the union of some  
 182 subset of groups from  $\mathcal{G}$ , and (iii) the size of the PG step restricted to the index set  
 183  $\mathcal{I}_k$  is at least a fraction of the size of the PG step restricted to the index set  $\mathcal{I}_k^{\text{PG}}$ . The  
 184 easiest choice that satisfies these conditions is  $\mathcal{I}_k \equiv \mathcal{I}_k^{\text{PG}}$ . Once  $\mathcal{I}_k$  has been chosen, the  
 185 next iterate is obtained by performing a line search along the PG direction in Line 17  
 186 by calling the subroutine UPDATE\_PG (for details, see Section 3.4). If the subroutine  
 187 returns  $\text{flag}_k^{\text{PG}} = \text{decrease}_{\alpha}$ , the PG parameter is decreased for the next iteration.

188 **3.2. Computing a CG direction (Algorithm 3.2).** This subroutine returns  
 189 a reduced-space direction  $\bar{d}_k$  that satisfies conditions (3.7)–(3.9). We call it a reduced-  
 190 space vector because the inputs  $g_k$  and  $H_k$  are elements in  $\mathbb{R}^{|\mathcal{I}_k|}$  and  $\mathbb{R}^{|\mathcal{I}_k| \times |\mathcal{I}_k|}$ , re-  
 191 spectively, where  $\mathcal{I}_k$  is computed in Line 9 of Algorithm 3.1. Condition (3.7) ensures  
 192 that  $\bar{d}_k$  is a descent direction for the objective function as a consequence of how the  
 193 reference direction  $d_k^R$  is computed in Line 25. Condition (3.8) ensures that  $\bar{d}_k$  reduces  
 194 the model  $m_k$  at least as much as a zero step. Finally, condition (3.9) promotes fast  
 195 *local* convergence of the iterate sequence  $\{x_k\}$  (see Section 4.2), but its enforcement  
 196 (or lack of enforcement) is irrelevant with respect to the complexity result that we  
 197 prove in Section 4.1. The subroutine name CG\_DIRECTION indicates our intent to  
 198 use the linear CG algorithm in our implementation, although other possible options

199 include a block-wise coordinate descent method applied to the model  $m_k$  in (3.6). In  
 200 particular, the direction associated with every iteration of the CG algorithm satis-  
 201 fies conditions (3.7)–(3.8), and condition (3.9) is satisfied by all sufficiently large CG  
 202 iterations. Thus, the requirements of this subroutine can always be met.

### 203 3.3. Reduced-space search using the CG direction (Algorithm 3.3).

204 This subroutine performs a search using the direction  $d_k$  returned by the subroutine  
 205 CG\_DIRECTION in Line 11 of Algorithm 3.1. For an illustration of this search, which  
 206 incorporates projections, see Figure 3.1. The approach uses the direction  $d_k$ , without  
 207 modification, for each block of variables  $\mathcal{G}_i$  such that the ray  $\{[x_k + \tau d_k]_{\mathcal{G}_i} : \tau \geq 0\}$   
 208 does not intersect the ball centered at zero of radius  $\bar{\rho}_{k,i} = \min\{\rho_{k,i}, \sin(\theta)\|[x_k]_{\mathcal{G}_i}\|_2\}$ ,  
 209 where  $\rho_{k,i}$  is defined in (3.5) and  $\theta \in (0, \pi/2)$  is a user-defined parameter. When they  
 210 do intersect, we first compute  $\tau_{k,i}$  as the smallest step along the Newton direction  
 211 (restricted to block  $\mathcal{G}_i$ ) that intersects the ball. Then, during the search that follows,  
 212 anytime the trial step size  $\xi^j$  is larger than  $\tau_{k,i}$ , the trial step for block  $\mathcal{G}_i$  is set to  
 213 zero; otherwise, the Newton direction is used so that the trial step (with respect to  
 214 block  $\mathcal{G}_i$ ) is  $[x_k + \xi^j d_k]_{\mathcal{G}_i}$  (see Line 41). If termination occurs in Line 42, then a new  
 215 block of variables will become zero, in which case we require the objective function  
 216 not to increase (see Line 43). On the other hand, if termination occurs in Line 48,  
 217 then it indicates that the objective function has been sufficiently reduced (see Line 47)  
 218 and no new groups of zeros have been formed.

### 219 3.4. Reduced-space line search along a PG step (Algorithm 3.4).

220 This subroutine performs a line search along the PG direction  $P_{\mathcal{I}}(s_k)$ . The search ensures  
 221 that the next iterate yields decrease in the objective of size at least  $(\eta \xi^j / \alpha_k) \|P_{\mathcal{I}_k}(s_k)\|_2^2$   
 222 for some positive integer  $j$  computed within the while loop in Line 53. Once the  
 223 while loop terminates, the update  $\text{flag}_k^{\text{PG}} \leftarrow \text{same\_}\alpha$  is made if  $j = 0$ , and set as  
 224  $\text{flag}_k^{\text{PG}} \leftarrow \text{decrease\_}\alpha$  otherwise. The motivation for this update is Lemma 2.3, which  
 225 shows that the while loop in Line 53 will terminate with  $j = 0$  if the PG parameter  
 226  $\alpha_k$  is sufficiently small. Therefore, anytime  $j > 0$ , Algorithm 3.4 returns  $\text{flag}_k^{\text{PG}} \leftarrow$   
 227  $\text{decrease\_}\alpha$  to Algorithm 3.1 in Line 17 so that the PG parameter value for the next  
 228 iteration is reduced by a factor of  $\xi \in (0, 1)$  in Line 19.

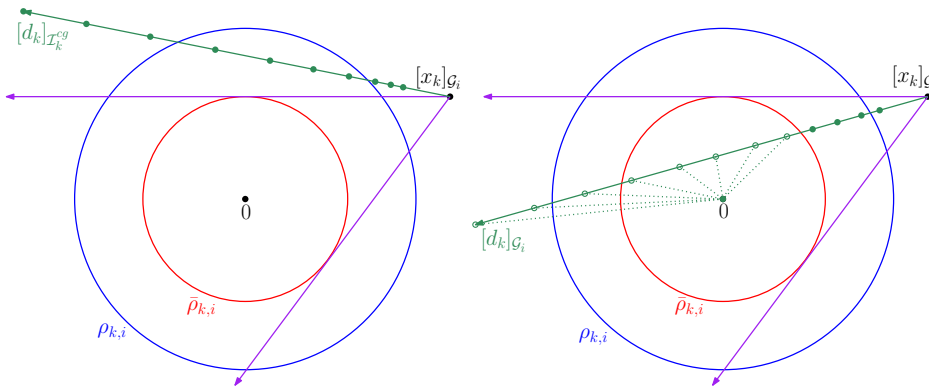


FIG. 3.1. The reduced-space projected search based on the Newton-CG direction  $d_k$  described in Section 3.3. In the figure on the left, the direction  $d_k$  does not intersect the ball of radius  $\bar{\rho}_{k,i}$ . In this case, standard backtracking is used, as indicated by the solid green dots. In the figure on the right, the direction  $d_k$  does intersect the ball of radius  $\bar{\rho}_{k,i}$ . In this case, all points after the first point of intersection (indicated by hollow green circles) are projected to zero. Once the backtracking points leave the ball of radius  $\bar{\rho}_{k,i}$  (indicated as solid green dots), standard backtracking is resumed.

---

**Algorithm 3.1** FaRSA-Group for solving problem (1.1).

---

- 1: **Input:**  $x_0$
  - 2: **Constants:**  $\{\varphi, \xi, \eta, \zeta\} \subset (0, 1)$ ,  $\{\kappa_1, \kappa_2, p\} \subset (0, \infty)$ ,  $\theta \in (0, \pi/2)$ , and  $q \in [1, 2]$ .
  - 3: Choose any initial PG parameter  $\alpha_0 \in (0, 1]$ .
  - 4: **for**  $k = 0, 1, 2, \dots$  **do**
  - 5:     Compute the step  $s_k$  from (3.1).
  - 6:     Compute the set  $\bar{\mathcal{I}}_k^{\text{cg}}$  from (3.2).
  - 7:     Compute  $\mathcal{I}_k^{\text{cg}}$  and  $\bar{\mathcal{I}}_k^{\text{pg}}$  and their optimality measures  $\chi_k^{\text{cg}}$  and  $\chi_k^{\text{pg}}$  from (3.4).
  - 8:     **if**  $\chi_k^{\text{pg}} \leq \chi_k^{\text{cg}}$  **then**
  - 9:         Choose any  $\mathcal{I}_k \subseteq \mathcal{I}_k^{\text{cg}}$  such that
 
$$\|[s_k]_{\mathcal{I}_k}\|_2 \geq \varphi \|[s_k]_{\mathcal{I}_k^{\text{cg}}}\|_2 \equiv \varphi \chi_k^{\text{cg}}$$
 and  $\mathcal{I}_k$  is the union of some  $\{\mathcal{G}_j\}$ .
  - 10:     Set  $g_k \leftarrow \nabla_{\mathcal{I}_k}(f + r)(x_k)$  and pick a positive-definite  $H_k \in \mathbb{R}^{|\mathcal{I}_k| \times |\mathcal{I}_k|}$ .
  - 11:     Call Algorithm 3.2 to obtain  $\bar{d}_k \leftarrow \text{CG\_DIRECTION}(g_k, H_k)$ .
  - 12:     Set  $[d_k]_{\mathcal{I}_k} \leftarrow \bar{d}_k$  and  $[d_k]_{\mathcal{I}_k^c} \leftarrow 0$ .
  - 13:     Call Algorithm 3.3 to obtain  $(x_{k+1}, \text{flag}_k^{\text{cg}}) \leftarrow \text{UPDATE\_CG}(x_k, d_k, \mathcal{I}_k)$ .
  - 14:     Set  $\alpha_{k+1} \leftarrow \alpha_k$ .
  - 15:     **else**
  - 16:         Choose any  $\mathcal{I}_k \subseteq \mathcal{I}_k^{\text{pg}}$  such that
 
$$\|[s_k]_{\mathcal{I}_k}\|_2 \geq \varphi \|[s_k]_{\mathcal{I}_k^{\text{pg}}}\|_2 \equiv \varphi \chi_k^{\text{pg}}$$
 and  $\mathcal{I}_k$  is the union of some  $\{\mathcal{G}_j\}$ .
  - 17:     Call Algorithm 3.4 to obtain  $(x_{k+1}, \text{flag}_k^{\text{pg}}) \leftarrow \text{UPDATE\_PG}(x_k, s_k, \alpha_k, \mathcal{I}_k)$ .
  - 18:     **if**  $\text{flag}_k^{\text{pg}} = \text{decrease\_}\alpha$  **then**
  - 19:          $\alpha_{k+1} \leftarrow \zeta \alpha_k$ .
  - 20:     **else**
  - 21:          $\alpha_{k+1} \leftarrow \alpha_k$ .
- 

---

**Algorithm 3.2** Computing  $\bar{d}_k$  in Line 11 of Algorithm 3.1.

---

- 22: **procedure**  $\bar{d}_k = \text{CG\_DIRECTION}(g_k, H_k)$
- 23:     **Constant:**  $q$  is provided by Algorithm 3.1.
- 24:     Define the model

$$(3.6) \quad m_k(d) := g_k^T d + \frac{1}{2} d^T H_k d.$$

- 25:     Compute the reference direction (an approximate minimizer of  $m_k$ ) as

$$d_k^R \leftarrow -\beta_k g_k, \text{ where } \beta_k \leftarrow \|g_k\|_2^2 / (g_k^T H_k g_k).$$

- 26:     Choose  $\mu_k \in (0, 1]$  and then compute any  $\bar{d}_k \approx \underset{d}{\text{argmin}} m_k(d)$  that satisfies

$$(3.7) \quad g_k^T \bar{d}_k \leq g_k^T d_k^R,$$

$$(3.8) \quad m_k(\bar{d}_k) \leq m_k(0), \text{ and}$$

$$(3.9) \quad \|H_k \bar{d}_k + g_k\|_2 \leq \mu_k \|g_k\|_2^q.$$

- 27:     **return**  $\bar{d}_k$
-



---

**Algorithm 3.3** Computing  $x_{k+1}$  in Line 13 of Algorithm 3.1.

---

```

28: procedure  $(x_{k+1}, \text{flag}_k^{\text{cg}}) = \text{UPDATE\_CG}(x_k, d_k, \mathcal{I}_k)$ 
29:   Constants:  $\eta, \xi,$  and  $\theta$  provided by Algorithm 3.1.
30:   for each  $i$  such that  $\mathcal{G}_i \subseteq \mathcal{I}_k$  do
31:     Compute  $\rho_{k,i}$  as defined in (3.5).
32:     Set  $\bar{\rho}_{k,i} \leftarrow \min\{\rho_{k,i}, \sin(\theta)\|[x_k]_{\mathcal{G}_i}\|_2\}$ .
33:     if  $\{[x_k + \tau d_k]_{\mathcal{G}_i} : \tau \geq 0\} \cap \{x \in \mathbb{R}^{|\mathcal{G}_i|} : \|x\|_2 \leq \bar{\rho}_{k,i}\} = \emptyset$  then
34:       Set  $\tau_{k,i} \leftarrow \infty$ .
35:     else
36:       Set  $\tau_{k,i}$  as the smallest positive root of  $\|[x_k + \tau d_k]_{\mathcal{G}_i}\|_2 = \bar{\rho}_{k,i}$ .
37:   Set  $j \leftarrow 0$  and  $\tau_k := \min_i\{\tau_{k,i} : \mathcal{G}_i \subseteq \mathcal{I}_k\}$ .
38:   while  $\xi^j \geq \tau_k$  do
39:     Set  $[y_j]_{\mathcal{I}_k^c} \leftarrow [x_k]_{\mathcal{I}_k^c}$ .
40:     for each  $i$  such that  $\mathcal{G}_i \in \mathcal{I}_k$  do
41:       Set  $[y_j]_{\mathcal{G}_i} \leftarrow \begin{cases} [x_k]_{\mathcal{G}_i} + \xi^j [d_k]_{\mathcal{G}_i} & \text{if } \xi^j < \tau_{k,i}, \\ 0 & \text{if } \xi^j \geq \tau_{k,i}. \end{cases}$ 
42:       if  $f(y_j) + r(y_j) \leq f(x_k) + r(x_k)$  then
43:         return  $x_{k+1} \leftarrow y_j$  and  $\text{flag}_k^{\text{cg}} \leftarrow \text{new\_zero}$ 
44:       Set  $j \leftarrow j + 1$ .
45:   loop
46:     Set  $y_j \leftarrow x_k + \xi^j d_k$ .
47:     if  $f(y_j) + r(y_j) \leq f(x_k) + r(x_k) + \eta \xi^j \nabla_{\mathcal{I}_k}(f + r)(x_k)^T [d_k]_{\mathcal{I}_k}$  then
48:       return  $x_{k+1} \leftarrow y_j$  and  $\text{flag}_k^{\text{cg}} \leftarrow \text{suff\_descent}$ 
49:     Set  $j \leftarrow j + 1$ .

```

---

**Algorithm 3.4** Computing  $x_{k+1}$  in Line 17 of Algorithm 3.1.

---

```

50: procedure  $(x_{k+1}, \text{flag}_k^{\text{pg}}) = \text{UPDATE\_PG}(x_k, s_k, \alpha_k, \mathcal{I}_k)$ 
51:   Constants:  $\eta$  and  $\xi$  provided by Algorithm 3.1.
52:   Set  $j \leftarrow 0$  and  $y_0 \leftarrow x_k + P_{\mathcal{I}_k}(s_k)$ .
53:   while  $f(y_j) + r(y_j) > f(x_k) + r(x_k) - \eta \xi^j \frac{1}{\alpha_k} \|P_{\mathcal{I}_k}(s_k)\|_2^2$  do
54:     Set  $j \leftarrow j + 1$  and then  $y_j \leftarrow x_k + \xi^j P_{\mathcal{I}_k}(s_k)$ .
55:   if  $j = 0$  then
56:     return  $x_{k+1} \leftarrow y_j$  and  $\text{flag}_k^{\text{pg}} \leftarrow \text{same\_}\alpha$ 
57:   else
58:     return  $x_{k+1} \leftarrow y_j$  and  $\text{flag}_k^{\text{pg}} \leftarrow \text{decrease\_}\alpha$ 

```

---

229 **4. Analysis.** Our analysis considers worst-case complexity (Section 4.1) and local  
230 convergence (Section 4.2) properties of Algorithm 3.1. To identify an approximate  
231 solution to problem (1.1), we use the measure  $\max\{\chi_k^{\text{pg}}, \chi_k^{\text{cg}}\}$ , as we now justify.

232 **LEMMA 4.1.** *Let  $\mathcal{K} \subseteq \mathbb{N}$  be such that  $\lim_{k \in \mathcal{K}} x_k = x_*$  and  $\lim_{k \in \mathcal{K}} \alpha_k = \alpha_* > 0$ .*  
233 *Then,  $x_*$  is a solution to problem (1.1) if and only if  $\lim_{k \in \mathcal{K}} \max\{\chi_k^{\text{pg}}, \chi_k^{\text{cg}}\} = 0$ .*

234 *Proof.* First, we may apply [8, Theorem 3.2.8], with the choice  $y = (\bar{x}, \bar{\alpha})$  and  
235 the set map  $\mathcal{C}(y) = \mathbb{R}^n$ , to the objective function appearing in (2.1) to conclude that  
236  $T(\bar{x}, \bar{\alpha})$  is continuous on  $\mathbb{R}^n \times (0, \infty)$ . Combining this property with the definition  
237 of  $T$  in (2.1) and the assumption that  $\lim_{k \in \mathcal{K}} (x_k, \alpha_k) = (x_*, \alpha_*)$  with  $\alpha_* > 0$  shows



238 that  $\lim_{k \in \mathcal{K}} s_k = \lim_{k \in \mathcal{K}} (T(x_k, \alpha_k) - x_k) = T(x_*, \alpha_*) - x_*$ . It follows from this limit  
 239 and the fact that Assumption 1.1 and [2, Theorem 10.7] together show that  $x_*$  is a  
 240 solution to problem (1.1) if and only if  $T(x_*, \alpha_*) = x_*$ .  $\square$

241 If  $\max\{\chi_k^{\text{cg}}, \chi_k^{\text{pg}}\} = 0$  for some  $k \in \mathbb{N}$ , then Lemma 4.1 implies that  $x_k$  is a solution  
 242 to problem (1.1). Hence, all that remains is to consider the behavior of Algorithm 3.1  
 243 when an infinite number of iterations is performed. To focus on this case, we make the  
 244 following assumption, which is assumed to hold throughout the rest of this section.

245 **ASSUMPTION 4.1.** *For all iterations  $k \in \mathbb{N}$ , it holds that  $\max\{\chi_k^{\text{cg}}, \chi_k^{\text{pg}}\} > 0$ .*

246 Since our analysis considers the properties of the sequence of iterates, it is con-  
 247 venient to define the following partition of iterations performed by Algorithm 3.1:

248  $\mathcal{K}^{\text{cg}} := \{k \in \mathbb{N} : \text{Line 13 is reached during the } k\text{th iteration}\},$   
 249  $\mathcal{K}_0^{\text{cg}} := \{k \in \mathcal{K}^{\text{cg}} : \text{subroutine UPDATE\_CG returns flag}_k^{\text{cg}} = \text{new\_zero in Line 13}\},$   
 250  $\mathcal{K}_{\text{sd}}^{\text{cg}} := \{k \in \mathcal{K}^{\text{cg}} : \text{subroutine UPDATE\_CG returns flag}_k^{\text{cg}} = \text{suff\_descent in Line 13}\},$   
 251  $\mathcal{K}^{\text{pg}} := \{k \in \mathbb{N} : \text{Line 17 is reached during the } k\text{th iteration}\},$   
 252  $\mathcal{K}_{\rightarrow}^{\text{pg}} := \{k \in \mathcal{K}^{\text{pg}} : \text{subroutine UPDATE\_PG returns flag}_k^{\text{pg}} = \text{same\_}\alpha \text{ in Line 17}\}, \text{ and}$   
 253  $\mathcal{K}_{\downarrow}^{\text{pg}} := \{k \in \mathcal{K}^{\text{pg}} : \text{subroutine UPDATE\_PG returns flag}_k^{\text{pg}} = \text{decrease\_}\alpha \text{ in Line 17}\},$   
 254

255 so that  $\mathcal{K}^{\text{cg}} = \mathcal{K}_0^{\text{cg}} \cup \mathcal{K}_{\text{sd}}^{\text{cg}}$ ,  $\mathcal{K}^{\text{pg}} = \mathcal{K}_{\rightarrow}^{\text{pg}} \cup \mathcal{K}_{\downarrow}^{\text{pg}}$ , and  $\mathbb{N} = \mathcal{K}^{\text{cg}} \cup \mathcal{K}^{\text{pg}}$ .

256 Finally, we assume that the symmetric and positive-definite matrices required in  
 257 Line 10 are chosen to be bounded and uniformly positive definite.

258 **ASSUMPTION 4.2.** *The matrix sequence  $\{H_k\}_{k \in \mathcal{K}^{\text{cg}}}$  chosen in Line 10 is bounded  
 259 and uniformly positive definite. That is, there exist constants  $0 < \mu_{\min} \leq \mu_{\max} < \infty$   
 260 such that  $\mu_{\min} \|v\|_2^2 \leq v^T H_k v \leq \mu_{\max} \|v\|_2^2$  for all  $k \in \mathcal{K}^{\text{cg}}$  and  $v \in \mathbb{R}^{|\mathcal{I}_k|}$ .*

261 **4.1. Complexity result.** We first focus our attention on iterations in  $\mathcal{K}^{\text{pg}}$ . The  
 262 next result shows that Algorithm 3.4 is well posed and that the new iterate that it  
 263 produces satisfies a decrease property that will be useful for our complexity analysis.

264 **LEMMA 4.2.** *For each  $k \in \mathcal{K}^{\text{pg}}$ , Algorithm 3.4 is called in Line 17 and successfully  
 265 returns  $x_{k+1}$  and  $\text{flag}_k^{\text{pg}}$ . Moreover, the value of  $\text{flag}_k^{\text{pg}}$  indicates whether  $k \in \mathcal{K}_{\downarrow}^{\text{pg}}$  or  
 266  $k \in \mathcal{K}_{\rightarrow}^{\text{pg}}$ , and for these respective cases the following properties hold:*

267 (i) *If  $k \in \mathcal{K}_{\rightarrow}^{\text{pg}}$ , then  $\alpha_{k+1} = \alpha_k$  and*

$$268 \quad (4.1) \quad f(x_{k+1}) + r(x_{k+1}) \leq f(x_k) + r(x_k) - \frac{\eta \varphi^2}{\alpha_k} (\chi_k^{\text{pg}})^2.$$

269 (ii) *If  $k \in \mathcal{K}_{\downarrow}^{\text{pg}}$ , then  $\alpha_{k+1} = \xi \alpha_k$  and  $f(x_{k+1}) + r(x_{k+1}) < f(x_k) + r(x_k)$ .*

270 *Proof.* Since  $k \in \mathcal{K}^{\text{pg}}$ , we know that the condition tested in Line 8 of Algorithm 3.1  
 271 must not hold, meaning that  $\chi_k^{\text{pg}} > \chi_k^{\text{cg}}$ . Combining this observation with Line 16 of  
 272 Algorithm 3.1 shows that the set  $\mathcal{I}_k$  defined in Line 16 satisfies

$$273 \quad (4.2) \quad \|P_{\mathcal{I}_k}(s_k)\|_2 = \|[s_k]_{\mathcal{I}_k}\|_2 \geq \varphi \chi_k^{\text{pg}} > 0.$$

274 Combining this result with Lemma 2.2 (using  $\mathcal{I} = \mathcal{I}_k$ ,  $\bar{x} = x_k$ , and  $\bar{\alpha} = \alpha_k$ ) yields

$$275 \quad (4.3) \quad D_{f+r}(x_k; P_{\mathcal{I}_k}(s_k)) \leq -\frac{1}{\alpha_k} \|P_{\mathcal{I}_k}(s_k)\|_2^2 < 0.$$

276 It is possible that Algorithm 3.4 terminates in Line 56 because the inequality  
 277 in Line 53 does not hold for  $j = 0$ . In this case, Algorithm 3.4 successfully returns

278  $x_{k+1} = y_0 = x_k + P_{\mathcal{I}_k}(s_k)$  and  $\text{flag}_k^{\text{pg}} = \text{same\_}\alpha$ , also indicating that  $k \in \mathcal{K}_{\downarrow}^{\text{pg}}$ . Since  
 279 the while-loop in Line 53 terminates with  $j = 0$ , we can conclude that

$$280 \quad (4.4) \quad f(x_{k+1}) + r(x_{k+1}) \equiv f(y_0) + r(y_0) \leq f(x_k) + r(x_k) - \frac{\eta}{\alpha_k} \|P_{\mathcal{I}_k}(s_k)\|_2^2.$$

281 Combining this inequality with (4.2) shows that (4.1) holds. Finally, since  $\text{flag}_k^{\text{pg}} =$   
 282  $\text{same\_}\alpha$ , it follows from Line 21 that  $\alpha_{k+1} = \alpha_k$ , completing the proof in this case.

283 It remains to consider the case when Algorithm 3.4 is unable to terminate in  
 284 Line 56 because the inequality in Line 53 holds for  $j = 0$ . In this case, it follows from  
 285 (4.3) and standard results for a backtracking Armijo line search that, for all sufficiently  
 286 large  $j$ , the vector  $y_j \leftarrow x_k + \xi^j P_{\mathcal{I}_k}(s_k)$  defined in Line 54 of Algorithm 3.4 satisfies

$$287 \quad (4.5) \quad \begin{aligned} f(y_j) + r(y_j) &\leq f(x_k) + r(x_k) + \eta \xi^j D_{f+r}(x_k; P_{\mathcal{I}_k}(s_k)) \\ &\leq f(x_k) + r(x_k) - \eta \xi^j \frac{1}{\alpha_k} \|P_{\mathcal{I}_k}(s_k)\|_2^2. \end{aligned}$$

288 This inequality shows that the while loop starting in Line 53 of Algorithm 3.4 will  
 289 terminate finitely, and thus Algorithm 3.4 successfully returns  $x_{k+1} = y_j = x_k +$   
 290  $\xi^j P_{\mathcal{I}_k}(s_k)$  for some  $j > 0$  and  $\text{flag}_k^{\text{pg}} = \text{decrease\_}\alpha$ , also indicating that  $k \in \mathcal{K}_{\downarrow}^{\text{pg}}$ .  
 291 Combining (4.5),  $y_j = x_{k+1}$ , and (4.3) proves that  $f(x_{k+1}) + r(x_{k+1}) < f(x_k) + r(x_k)$ ,  
 292 as claimed. Finally, since  $\text{flag}_k^{\text{pg}} = \text{decrease\_}\alpha$ , we see in Line 19 that  $\alpha_{k+1} = \xi \alpha_k$ .  $\square$

293 Next, we prove that the PG parameter remains bounded away from zero.

294 LEMMA 4.3. *The PG parameter sequence generated by Algorithm 3.1 satisfies*

$$295 \quad (4.6) \quad 1 \geq \alpha_k \geq \alpha_{\min} := \min \left\{ \alpha_0, \frac{2\xi(1-\eta)}{L} \right\} > 0 \quad \text{for all } k \in \mathbb{N}.$$

296 Moreover, a bound on the number of times the PG parameter is decreased is given by

$$297 \quad (4.7) \quad |\mathcal{K}_{\downarrow}^{\text{pg}}| \leq c_{\downarrow}^{\alpha} := \max \left\{ 0, \left\lceil \log \left( \frac{\alpha_0 L}{2(1-\eta)} \right) / \log(\xi^{-1}) \right\rceil \right\}.$$

298 *Proof.* We first prove (4.6). Since  $\alpha_0 \in (0, 1]$  in Line 3 and  $\alpha_{k+1} \leq \alpha_k$  for all  
 299  $k \in \mathbb{N}$ , we need only prove the lower bound on  $\alpha_k$  in (4.6). With that goal in mind,  
 300 for the purpose of obtaining a contradiction, suppose that there exists an iteration  $k$   
 301 satisfying  $\alpha_k \leq 2(1-\eta)/L < 2/L$ , with the latter inequality holding since  $\eta \in (0, 1)$ .

302 First suppose that  $k \in \mathcal{K}^{\text{pg}}$ . With  $y_0 = x_k + P_{\mathcal{I}_k}(s_k)$  as defined in Line 52 of  
 303 Algorithm 3.4, it follows from Lemma 2.3 with  $\bar{x} = x_k$ ,  $\bar{\alpha} = \alpha_k$ , and  $s(\bar{x}, \bar{\alpha}) = s_k$  that

$$304 \quad \begin{aligned} f(y_0) + r(y_0) &\leq f(x_k) + r(x_k) - \left( \frac{1}{\alpha_k} - \frac{L}{2} \right) \|P_{\mathcal{I}}(s_k)\|_2^2 \\ 305 \quad &\leq f(x_k) + r(x_k) - \left( \frac{1}{\alpha_k} - \frac{2(1-\eta)}{2\alpha_k} \right) \|P_{\mathcal{I}}(s_k)\|_2^2 \\ 306 \quad &= f(x_k) + r(x_k) - \frac{\eta}{\alpha_k} \|P_{\mathcal{I}}(s_k)\|_2^2. \end{aligned}$$

308 This inequality implies that the condition checked in Line 53 for  $j = 0$  will not hold,  
 309 meaning that  $j = 0$  when Line 55 is reached so that  $\text{flag}_k^{\text{pg}} \leftarrow \text{same\_}\alpha$  in Line 56. Thus,  
 310 when Line 18 in Algorithm 3.1 is reached, the update  $\alpha_{k+1} \leftarrow \alpha_k$  will take place.  
 311 Second, if  $k \in \mathcal{K}^{\text{cg}}$ , then Algorithm 3.1 sets  $\alpha_{k+1} \leftarrow \alpha_k$ . To summarize, anytime  
 312  $\alpha_k \leq 2(1-\eta)/L$ , the update  $\alpha_{k+1} \leftarrow \alpha_k$  takes place. Combining this property with  
 313 the fact that when the PG parameter is decreased the update  $\alpha_{k+1} \leftarrow \xi \alpha_k$  is used  
 314 (see Line 19 in Algorithm 3.1), shows that (4.6) holds.

315 We now prove (4.7). Let us observe from the first paragraph in this proof that  
 316 if  $\alpha_0 \leq 2(1-\eta)/L$  then  $|\mathcal{K}_{\downarrow}^{\text{pg}}| = 0$ , which verifies that (4.7) holds. Therefore, for

317 the remainder of the proof, suppose that  $\alpha_0 > 2(1 - \eta)/L$ . Combining this bound  
 318 with the fact that when the PG parameter is decreased the update  $\alpha_{k+1} \leftarrow \xi\alpha_k$  is  
 319 used, we can see that an upper bound on  $|\mathcal{K}_\downarrow^{\text{PG}}|$  is the smallest integer  $\ell$  such that  
 320  $\alpha_0\xi^\ell \leq 2(1 - \eta)/L$ . Solving this inequality for  $\ell$  shows that the result in (4.7) holds.  $\square$

321 We now switch our attention to iterations in  $\mathcal{K}^{\text{cg}}$ . The next result establishes  
 322 that Algorithm 3.2 is well posed, and that the direction  $d_k$  that results from it when  
 323 called by Algorithm 3.1 satisfies a certain descent property.

324 LEMMA 4.4. *For each  $k \in \mathcal{K}^{\text{cg}}$ , Algorithm 3.2 is well posed. Moreover, the result-*  
 325 *ing direction  $\bar{d}_k$ , which is used to compute  $d_k$  in Line 12, guarantees that  $d_k$  satisfies*

- 326 (i)  $\nabla_{\mathcal{I}_k}(f + r)(x_k)^T[d_k]_{\mathcal{I}_k} \leq -\frac{1}{\mu_{\max}}\|\nabla_{\mathcal{I}_k}(f + r)(x_k)\|_2^2 < 0$ , and  
 327 (ii)  $\|d_k\|_2 \leq (2/\mu_{\min})\|\nabla_{\mathcal{I}_k}(f + r)(x_k)\|_2$

328 where  $\mathcal{I}_k \subseteq \mathcal{I}_k^{\text{cg}}$  is the set in Line 9 used as an input to Algorithm 3.2 in Line 13.

329 *Proof.* Since  $k \in \mathcal{K}^{\text{cg}}$ , Algorithm 3.2 is called in Line 11 with input  $\mathcal{I}_k$  defined in  
 330 Line 9. We first prove that  $g_k = \nabla_{\mathcal{I}_k}(f + r)(x_k)$ , as defined in Line 10, is nonzero. For  
 331 a proof by contradiction, suppose that  $g_k = 0$  so that  $\nabla_{\mathcal{G}_i}(f + r)(x_k) = 0$  for all  $i$  such  
 332 that  $\mathcal{G}_i \subseteq \mathcal{I}_k$ . Consider arbitrary such  $i$ . Note that  $[x_k]_{\mathcal{G}_i} \neq 0$  and  $[x_k + s_k]_{\mathcal{G}_i} \neq 0$  since  
 333  $\mathcal{G}_i \subseteq \mathcal{I}_k \subseteq \mathcal{I}_k^{\text{cg}}$  (see Line 9) and by how  $\mathcal{I}_k^{\text{cg}}$  is defined. This allows us to conclude from  
 334 Lemma 2.4 that  $[s_k]_{\mathcal{G}_i} = 0$ , i.e., that  $[s_k]_{\mathcal{I}_k} = 0$  since  $i$  with  $\mathcal{G}_i \subseteq \mathcal{I}_k$  was arbitrary.  
 335 This fact and Line 9 yields  $\chi_k^{\text{cg}} = 0$ , but since the inequality in Line 8 must hold, we  
 336 also have  $\chi_k^{\text{PG}} = 0$ . This contradicts Assumption 4.1, thus establishing that  $g_k \neq 0$ .  
 337 Now, it follows from Lines 10, 12, 26, and 25,  $g_k \neq 0$ , and Assumption 4.2 that

$$\begin{aligned} \nabla_{\mathcal{I}_k}(f + r)(x_k)^T[d_k]_{\mathcal{I}_k} &\equiv g_k^T \bar{d}_k \leq g_k^T d_k^R = -\beta_k \|g_k\|_2^2 \\ &= -\|g_k\|_2^4 / (g_k^T H_k g_k) \leq -\frac{1}{\mu_{\max}} \|g_k\|_2^2. \end{aligned}$$

338 The result in (i) follows from this inequality and  $g_k = \nabla_{\mathcal{I}_k}(f + r)(x_k) \neq 0$ .

340 Part (ii) is precisely [5, Lemma 3.8] under our Assumption 4.2 since our conditions  
 341 placed upon the step  $d_k$  are exactly the same as those used in [5].  $\square$

342 The next lemma shows that, for  $k \in \mathcal{K}^{\text{cg}}$ , a local Lipschitz property holds along a  
 343 certain portion of the search path defined by the reduced-space Newton-CG direction.

344 LEMMA 4.5. *Let  $k \in \mathcal{K}^{\text{cg}}$  so that  $\mathcal{I}_k$  is computed in Line 9. The following hold:*

- 345 (i) *The constant  $\theta \in (0, \pi/2)$  and index set  $\mathcal{I}_k$  passed into Algorithm 3.3 satisfy,*  
 346 *for each  $i$  such that  $\mathcal{G}_i \subseteq \mathcal{I}_k$  with  $\rho_{k,i}$  computed in (3.5) and  $\bar{\rho}_{k,i}$  computed in*  
 347 *Line 32, the following conditions:*

- 348 (a)  $\|[x_k + s_k]_{\mathcal{G}_i}\|_2 \neq 0$ ,  
 349 (b)  $\|[x_k]_{\mathcal{G}_i}\|_2 \geq \rho_{k,i} \geq \bar{\rho}_{k,i} \geq \sin(\theta)\rho_{k,i} > 0$ , and  
 350 (c)  $\|[x_k]_{\mathcal{G}_i}\|_2 - \bar{\rho}_{k,i} \geq \kappa_2(1 - \sin(\theta))\|\nabla_{\mathcal{I}_k}(f + r)(x_k)\|_2^p$ .

- 351 (ii) *For all step sizes  $\beta \in [0, \tau_k)$  with  $\tau_k$  computed in Line 37, it holds, with*

$$(4.8) \quad \lambda_{\max} := \max\{\lambda_1, \lambda_2, \dots, \lambda_{n_{\mathcal{G}}}\} \quad \text{and} \quad \rho_{k,\min} := \min_i \{\rho_{k,i} : \mathcal{G}_i \subseteq \mathcal{I}_k\}$$

353 that  $\|\nabla_{\mathcal{I}_k}(f + r)(x_k) - \nabla_{\mathcal{I}_k}(f + r)(x_k + \beta d_k)\|_2 \leq \beta(L + \frac{\lambda_{\max}}{\rho_{k,\min}})\|d_k\|_2$ .

354 *Proof.* We first prove part (i). Consider arbitrary  $i$  with  $\mathcal{G}_i \subseteq \mathcal{I}_k$ , where  $\mathcal{I}_k \subseteq \mathcal{I}_k^{\text{cg}}$   
 355 is passed into Algorithm 3.3 and constructed to satisfy the condition in Line 9. Part  
 356 (a) follows from  $\mathcal{I}_k^{\text{cg}} \subseteq \bar{\mathcal{I}}_k^{\text{cg}}$  and the definition of  $\bar{\mathcal{I}}_k^{\text{cg}}$  in (3.2). The first inequality in  
 357 part (b) follows from  $\mathcal{I}_k^{\text{cg}} \subseteq \bar{\mathcal{I}}_k^{\text{cg}}$ , and how  $\mathcal{I}_k^{\text{cg}}$ ,  $\mathcal{I}_k^{\text{small}}$ , and  $\bar{\mathcal{I}}_k^{\text{cg}}$  are defined. The second  
 358 inequality in (b) follows from how  $\bar{\rho}_{k,i}$  is defined in Line 32. The third inequality in

359 (b) follows from Line 32 and the first inequality in (b). To complete the proof for  
360 part (b), we must prove that  $\rho_{k,i} > 0$ . For a proof by contradiction, assume that  
361  $\rho_{k,i} = 0$ , which by (3.5) means that  $\|\nabla_{\mathcal{I}_k^{\text{cg}}}(f+r)(x_k)\|_2 = 0$ . It follows from this  
362 fact that each  $i$  with  $\mathcal{G}_i \subseteq \mathcal{I}_k \subseteq \mathcal{I}_k^{\text{cg}}$  satisfies  $\|\nabla_{\mathcal{G}_i}(f+r)(x_k)\|_2 = 0$ , which in light  
363 of Lemma 2.4 (using  $\bar{x} = x_k$ ,  $\bar{\alpha} = \alpha_k$ , and  $s(\bar{x}, \bar{\alpha}) = s_k$ ) and the definition of  $\mathcal{I}_k^{\text{cg}}$   
364 implies that  $\|[s_k]_{\mathcal{G}_i}\|_2 = 0$  for each  $\mathcal{G}_i \subseteq \mathcal{I}_k$ , i.e., that  $\|[s_k]_{\mathcal{I}_k}\|_2 = 0$ . It now follows  
365 from Line 9 that  $\chi_k^{\text{cg}} = 0$ , which combined with the inequality in Line 8 shows that  
366  $\chi_k^{\text{pg}} = 0$ . Since we have reached a contradiction to Assumption 4.1, we must conclude  
367 that  $\rho_{k,i} > 0$ , as claimed. Finally, we aim to prove part (c). It follows from Line 32,  
368  $\theta \in (0, \pi/2)$ , part (b), (3.5), and the fact that  $\mathcal{I}_k \subseteq \mathcal{I}_k^{\text{cg}}$  that

$$\begin{aligned} 369 \quad \|[x_k]_{\mathcal{G}_i}\|_2 - \bar{\rho}_{k,i} &\geq \|[x_k]_{\mathcal{G}_i}\|_2 - \sin(\theta)\|[x_k]_{\mathcal{G}_i}\|_2 = (1 - \sin(\theta))\|[x_k]_{\mathcal{G}_i}\|_2 \\ 370 &\geq (1 - \sin(\theta))\rho_{k,i} \geq \kappa_2(1 - \sin(\theta))\|\nabla_{\mathcal{I}_k^{\text{cg}}}(f+r)(x_k)\|_2^p \\ 371 &\geq \kappa_2(1 - \sin(\theta))\|\nabla_{\mathcal{I}_k}(f+r)(x_k)\|_2^p, \end{aligned}$$

373 which completes the proof of part (c).

374 To prove part (ii), let  $\beta \in [0, \tau_k)$ . It follows from part (i) and the definition of  
375  $\tau_k$  in Line 37 that every point on the segment that connects  $[x_k]_{\mathcal{G}_i}$  to  $[x_k + \beta d_k]_{\mathcal{G}_i}$  is  
376 outside of the ball in  $\mathbb{R}^{|\mathcal{G}_i|}$  centered at zero of radius  $\bar{\rho}_{k,i} > 0$ . This means that both  
377  $\|[x_k]_{\mathcal{G}_i}\| \geq \bar{\rho}_{k,i}$  and  $\|[x_k + \beta d_k]_{\mathcal{G}_i}\| \geq \bar{\rho}_{k,i}$ . It now follows that

$$\begin{aligned} &\|\nabla_{\mathcal{G}_i} r(x_k) - \nabla_{\mathcal{G}_i} r(x_k + \beta d_k)\|_2 \\ 378 \quad (4.9) &= \lambda_i \left\| \frac{[x_k]_{\mathcal{G}_i}}{\|[x_k]_{\mathcal{G}_i}\|_2} - \frac{[x_k + \beta d_k]_{\mathcal{G}_i}}{\|[x_k + \beta d_k]_{\mathcal{G}_i}\|_2} \right\|_2 = \frac{\lambda_i}{\bar{\rho}_{k,i}} \left\| \frac{\bar{\rho}_{k,i}[x_k]_{\mathcal{G}_i}}{\|[x_k]_{\mathcal{G}_i}\|_2} - \frac{\bar{\rho}_{k,i}[x_k + \beta d_k]_{\mathcal{G}_i}}{\|[x_k + \beta d_k]_{\mathcal{G}_i}\|_2} \right\|_2 \\ &\leq \frac{\lambda_i}{\bar{\rho}_{k,i}} \|[x_k]_{\mathcal{G}_i} - [x_k + \beta d_k]_{\mathcal{G}_i}\|_2 = \frac{\lambda_i \beta}{\bar{\rho}_{k,i}} \|[d_k]_{\mathcal{G}_i}\|_2, \end{aligned}$$

379 where the (only) inequality follows from the nonexpansive property of the projection  
380 (of  $[x_k]_{\mathcal{G}_i}$  and  $[x_k + \beta d_k]_{\mathcal{G}_i}$ ) onto the ball of radius  $\bar{\rho}_{k,i}$ . From (4.9) we have

$$\begin{aligned} 381 &\|\nabla_{\mathcal{I}_k} r(x_k) - \nabla_{\mathcal{I}_k} r(x_k + \beta d_k)\|_2^2 \\ 382 &= \sum_{i: \mathcal{G}_i \subseteq \mathcal{I}_k} \|\nabla_{\mathcal{G}_i} r(x_k) - \nabla_{\mathcal{G}_i} r(x_k + \beta d_k)\|_2^2 \leq \beta^2 \sum_{i: \mathcal{G}_i \subseteq \mathcal{I}_k} \frac{\lambda_i^2}{\bar{\rho}_{k,i}^2} \|[d_k]_{\mathcal{G}_i}\|_2^2 \\ 383 \quad (4.10) &\leq \frac{\beta^2 \lambda_{\max}^2}{\rho_{k,\min}^2} \sum_{i: \mathcal{G}_i \subseteq \mathcal{I}_k} \|[d_k]_{\mathcal{G}_i}\|_2^2 = \frac{\beta^2 \lambda_{\max}^2}{\rho_{k,\min}^2} \|[d_k]_{\mathcal{I}_k}\|_2^2. \end{aligned}$$

385 It follows from Assumption 1.1,  $[d_k]_{\mathcal{I}_k^c} = 0$ , the triangle inequality, and (4.10) that

$$\begin{aligned} 386 &\|\nabla_{\mathcal{I}_k}(f+r)(x_k) - \nabla_{\mathcal{I}_k}(f+r)(x_k + \beta d_k)\|_2 \\ 387 &\leq \|\nabla_{\mathcal{I}_k} f(x_k) - \nabla_{\mathcal{I}_k} f(x_k + \beta d_k)\|_2 + \|\nabla_{\mathcal{I}_k} r(x_k) - \nabla_{\mathcal{I}_k} r(x_k + \beta d_k)\|_2 \\ 388 &\leq L\beta\|d_k\|_2 + \left(\beta \frac{\lambda_{\max}}{\rho_{k,\min}}\right) \|[d_k]_{\mathcal{I}_k}\|_2 = \beta \left(L + \frac{\lambda_{\max}}{\rho_{k,\min}}\right) \|[d_k]_{\mathcal{I}_k}\|_2, \end{aligned}$$

390 which completes the proof.  $\square$

391 We now show that Algorithm 3.4 is well posed and that the new iterate it produces  
392 satisfies a decrease property that will be used in the final complexity result.

393 LEMMA 4.6. *For each  $k \in \mathcal{K}^{\text{cg}}$ , Algorithm 3.3 is called in Line 13 and successfully*  
394 *returns  $x_{k+1}$  and  $\text{flag}_k^{\text{cg}}$ . Moreover, the value of  $\text{flag}_k^{\text{cg}}$  indicates whether  $k \in \mathcal{K}_0^{\text{cg}}$  or*  
395  *$k \in \mathcal{K}_{\text{sd}}^{\text{cg}}$ , and for these respective cases the following properties hold:*

396 (i) If  $k \in \mathcal{K}_0^{\text{cg}}$ , then  $f(x_{k+1}) + r(x_{k+1}) \leq f(x_k) + r(x_k)$  and  $x_{k+1}$  has at least one  
 397 additional block of zeros compared to  $x_k$ .

398 (ii) If  $k \in \mathcal{K}_{\text{sd}}^{\text{cg}}$ , then

$$399 \quad (4.11) \quad f(x_{k+1}) + r(x_{k+1}) \leq f(x_k) + r(x_k) - \min\{c_1(\chi_k^{\text{cg}})^{1+p}, c_2(\chi_k^{\text{cg}})^{2+p}\}$$

400 where

$$401 \quad (4.12) \quad c_1 := \frac{\eta \xi \mu_{\min} \kappa_2 (1 - \sin(\theta)) \varphi^{1+p}}{2\mu_{\max}} > 0 \quad \text{and}$$

$$c_2 := \frac{\kappa_2 \mu_{\min}^2 \xi \eta (1 - \eta) \varphi^{2+p}}{2\mu_{\max}^2 (L\kappa_2(L_f + \lambda_{\max}\sqrt{n\bar{c}})^p + \lambda_{\max})} > 0.$$

402 *Proof.* Throughout, we use  $F := f + r$ . It is possible that Algorithm 3.3 success-  
 403 fully terminates in Line 43, in which case it follows from Line 43 and Line 42 that the  
 404 returned  $x_{k+1}$  and  $\text{flag}_k^{\text{cg}}$  satisfy  $F(x_{k+1}) \leq F(x_k)$  and  $\text{flag}_k^{\text{cg}} = \text{new\_zero}$ , indicating  
 405 that  $k \in \mathcal{K}_0^{\text{cg}}$ . Moreover, upon termination, the value  $j$  satisfies  $\xi^j \geq \tau_k$  (see Line 38),  
 406 which combined with Line 41 shows that at least one additional group of variables  
 407 has become zero at  $x_{k+1}$ . This proves that part (i) holds.

408 Next, suppose that Algorithm 3.3 does not terminate in Line 43. Observe from  
 409 the definition of  $\tau_k$  in Line 37 that  $\tau_k > 0$  (this follows from Lemma 4.5(i) and the  
 410 definition of  $\bar{\rho}_{k,i}$ ). Therefore, it follows that the **while** loop starting in Line 38 will  
 411 terminate with the smallest nonnegative integer  $\bar{j}$  such that  $\xi^{\bar{j}} < \tau_k$ , and the **loop** in  
 412 Line 45 will begin with  $j = \bar{j}$ . We now claim that the condition in Line 47 used to  
 413 determine termination of the **loop** is satisfied for all  $j \geq \bar{j}$  such that

$$414 \quad (4.13) \quad \xi^j \in \left[ 0, \frac{2(\eta - 1) \nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k}}{(L + \lambda_{\max}/\rho_{k,\min}) \|[d_k]_{\mathcal{I}_k}\|_2^2} \right] \subset [0, \tau_k).$$

415 To see that this claim holds, we can use the integral form of Taylor's Theorem and  
 416 Lemma 4.5(ii) (using the fact that  $\gamma \xi^j \in [0, \tau_k)$  for all  $\gamma \in [0, 1]$ ) to obtain

$$417 \quad |F(x_k + \xi^j d_k) - F(x_k) - \xi^j \nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k}|$$

$$418 \quad \leq \left| \int_0^1 \xi^j [d_k]_{\mathcal{I}_k}^T (\nabla_{\mathcal{I}_k} F(x_k + \gamma \xi^j d_k) - \nabla_{\mathcal{I}_k} F(x_k)) d\gamma \right|$$

$$419 \quad \leq \xi^j \int_0^1 \|[d_k]_{\mathcal{I}_k}\|_2 \|\nabla_{\mathcal{I}_k} F(x_k + \gamma \xi^j d_k) - \nabla_{\mathcal{I}_k} F(x_k)\|_2 d\gamma$$

$$420 \quad \leq \xi^{2j} (L + \lambda_{\max}/\rho_{k,\min}) \|[d_k]_{\mathcal{I}_k}\|_2^2 \int_0^1 \gamma d\gamma = \frac{1}{2} \xi^{2j} (L + \lambda_{\max}/\rho_{k,\min}) \|[d_k]_{\mathcal{I}_k}\|_2^2.$$

422 Combining this inequality with (4.13) yields

$$423 \quad F(x_k + \xi^j d_k) \leq F(x_k) + \xi^j \nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k} + \frac{1}{2} \xi^{2j} (L + \lambda_{\max}/\rho_{k,\min}) \|[d_k]_{\mathcal{I}_k}\|_2^2$$

$$424 \quad = F(x_k) + \xi^j \nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k} + \xi^j (\eta - 1) \nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k}$$

$$425 \quad = F(x_k) + \eta \xi^j \nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k},$$

427 which establishes our claim that the inequality in Line 47 holds for all  $j \geq \bar{j}$  such  
 428 that  $\xi^j$  satisfies (4.13). This shows that the **loop** will successfully terminate with  
 429  $\text{flag}_k^{\text{cg}} = \text{suff\_descent}$  (thus indicating that  $k \in \mathcal{K}_{\text{sd}}^{\text{cg}}$ ) and  $x_{k+1}$  satisfying

$$430 \quad (4.14) \quad F(x_{k+1}) \leq F(x_k) + \eta \xi^{\bar{j}} \nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k}$$

431 for some  $\hat{j}$  satisfying

$$\begin{aligned}
432 \quad \xi^{\hat{j}} &\geq \min \left\{ \xi^{\bar{j}}, \frac{2\xi(\eta-1)\nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k}}{(L + \lambda_{\max}/\rho_{k,\min})\| [d_k]_{\mathcal{I}_k} \|_2^2} \right\} \\
433 \quad (4.15) \quad &\geq \min \left\{ \xi\tau_k, \frac{2\xi(\eta-1)\nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k}}{(L + \lambda_{\max}/\rho_{k,\min})\| [d_k]_{\mathcal{I}_k} \|_2^2} \right\} \\
434
\end{aligned}$$

435 where the second inequality follows from the fact that  $\bar{j}$  is the *smallest* nonnegative  
436 integer such that  $\xi^{\bar{j}} < \tau_k$ . We now consider two cases.

437 **Case 1:** the minimum in (4.15) is  $\xi\tau_k$ , from which we may conclude that  $\tau_k < \infty$ .

438 Using (4.14) and Lemma 4.4(i) we have that

$$439 \quad (4.16) \quad F(x_{k+1}) \leq F(x_k) + \eta\xi^{\hat{j}}\nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k} \leq F(x_k) - \frac{\eta\xi}{\mu_{\max}}\tau_k\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^2.$$

440 We now seek a lower bound on  $\tau_k$ . Consider  $i$  such that  $\tau_{k,i} < \infty$  when computed in  
441 Algorithm 3.3. The triangle inequality gives  $\bar{\rho}_{k,i} = \|[x_k + \tau_{k,i}d_k]_{\mathcal{G}_i}\|_2 \geq \|[x_k]_{\mathcal{G}_i}\|_2 -$   
442  $\tau_{k,i}\|[d_k]_{\mathcal{G}_i}\|_2$ , which together with Lemma 4.5(i)(c) and Lemma 4.4(ii) shows that

$$\begin{aligned}
443 \quad \tau_{k,i} &\geq \frac{\|[x_k]_{\mathcal{G}_i}\|_2 - \bar{\rho}_{k,i}}{\|[d_k]_{\mathcal{G}_i}\|_2} \\
444 \quad &\geq \frac{\mu_{\min}\kappa_2(1 - \sin(\theta))\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^p}{2\|\nabla_{\mathcal{I}_k} F(x_k)\|_2} = \frac{1}{2}\mu_{\min}\kappa_2(1 - \sin(\theta))\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^{p-1}. \\
445
\end{aligned}$$

446 From this, it follows that  $\tau_k \geq \frac{1}{2}\mu_{\min}\kappa_2(1 - \sin(\theta))\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^{p-1}$ . Using this in-  
447 equality with (4.16), Lemma 2.4, and the set  $\mathcal{I}_k$  from Line 9 shows that

$$\begin{aligned}
448 \quad F(x_{k+1}) &\leq F(x_k) - \frac{\eta\xi\mu_{\min}\kappa_2(1 - \sin(\theta))}{2\mu_{\max}}\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^{1+p} \\
449 \quad &\leq F(x_k) - \frac{\eta\xi\mu_{\min}\kappa_2(1 - \sin(\theta))}{2\mu_{\max}}\|[s_k]_{\mathcal{I}_k}\|_2^{1+p} \\
450 \quad &\leq F(x_k) - \frac{\eta\xi\mu_{\min}\kappa_2(1 - \sin(\theta))\varphi^{1+p}}{2\mu_{\max}}(\chi_k^{\text{cg}})^{1+p}, \\
451
\end{aligned}$$

452 thus completing the proof for this case.

453 **Case 2:** the minimum in (4.15) is  $\frac{2\xi(\eta-1)\nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k}}{(L + \lambda_{\max}/\rho_{k,\min})\| [d_k]_{\mathcal{I}_k} \|_2^2}$ . Combining this fact with

454 (4.14), (4.15), Lemma 4.4(i), and Lemma 4.4(ii) shows that

$$\begin{aligned}
455 \quad (4.17) \quad F(x_{k+1}) &\leq F(x_k) + \eta\xi^{\hat{j}}\nabla_{\mathcal{I}_k} F(x_k)^T [d_k]_{\mathcal{I}_k} \\
&\leq F(x_k) - \frac{2\xi\eta(1 - \eta)\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^4}{\mu_{\max}^2(L + \lambda_{\max}/\rho_{k,\min})\| [d_k]_{\mathcal{I}_k} \|_2^2} \\
&\leq F(x_k) - \frac{2\mu_{\min}^2\xi\eta(1 - \eta)\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^4}{4\mu_{\max}^2(L + \lambda_{\max}/\rho_{k,\min})\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^2} \\
&= F(x_k) - \frac{\mu_{\min}^2\xi\eta(1 - \eta)\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^2}{2\mu_{\max}^2(L + \lambda_{\max}/\rho_{k,\min})}.
\end{aligned}$$

456 It follows from (4.8), (3.5), and  $\mathcal{I}_k \subseteq \mathcal{I}_k^{\text{cg}}$  that  $\rho_{k,\min} \geq \kappa_2\|\nabla_{\mathcal{I}_k} F(x_k)\|_2^p$ . Combining

457 this bound with (4.17) shows that

$$\begin{aligned}
(4.18) \quad F(x_{k+1}) &\leq F(x_k) - \frac{\mu_{\min}^2 \xi \eta (1 - \eta) \|\nabla_{\mathcal{I}_k} F(x_k)\|_2^2}{2\mu_{\max}^2 (L + \lambda_{\max} / \rho_{k, \min})} \\
&\leq F(x_k) - \frac{\mu_{\min}^2 \xi \eta (1 - \eta) \|\nabla_{\mathcal{I}_k} F(x_k)\|_2^2}{2\mu_{\max}^2 (L + \lambda_{\max} / (\kappa_2 \|\nabla_{\mathcal{I}_k} F(x_k)\|_2^p))} \\
&= F(x_k) - \frac{\kappa_2 \mu_{\min}^2 \xi \eta (1 - \eta) \|\nabla_{\mathcal{I}_k} F(x_k)\|_2^{2+p}}{2\mu_{\max}^2 (L \kappa_2 \|\nabla_{\mathcal{I}_k} F(x_k)\|_2^p + \lambda_{\max})}.
\end{aligned}$$

459 Next, we know from Lemma 4.2, Lemma 4.6(i), and equations (4.16) and (4.18) that  
460  $F(x_k) \leq F(x_0)$  for all  $k \in \mathbb{N}$ , i.e.,  $x_k \in \mathcal{L}$  for all  $k \in \mathbb{N}$ . Combining this fact with the  
461 triangle inequality, Assumption 1.1, the definition of  $r$ , and (4.8) gives

$$\begin{aligned}
462 \quad \|\nabla_{\mathcal{I}_k} F(x_k)\|_2 &\leq \|\nabla_{\mathcal{I}_k} f(x_k)\|_2 + \|\nabla_{\mathcal{I}_k} r(x_k)\|_2 \\
463 \quad &= \|\nabla_{\mathcal{I}_k} f(x_k)\|_2 + \left( \sum_{i: \mathcal{G}_i \subseteq \mathcal{I}_k} \|\nabla_{\mathcal{G}_i} r(x_k)\|_2^2 \right)^{1/2} \\
464 \quad &\leq L_f + \left( \sum_{i: \mathcal{G}_i \subseteq \mathcal{I}_k} \|\lambda_i [x_k]_{\mathcal{G}_i} / \|[x_k]_{\mathcal{G}_i}\|_2\|_2^2 \right)^{1/2} \\
465 \quad &= L_f + \left( \sum_{i: \mathcal{G}_i \subseteq \mathcal{I}_k} \lambda_i^2 \right)^{1/2} \leq L_f + \left( \sum_{i: \mathcal{G}_i \subseteq \mathcal{I}_k} \lambda_{\max}^2 \right)^{1/2} \leq L_f + \lambda_{\max} \sqrt{n_{\mathcal{G}}}. \\
466
\end{aligned}$$

Combining this with (4.18) gives

$$F(x_{k+1}) \leq F(x_k) - \left( \frac{\kappa_2 \mu_{\min}^2 \xi \eta (1 - \eta)}{2\mu_{\max}^2 (L \kappa_2 (L_f + \lambda_{\max} \sqrt{n_{\mathcal{G}}})^p + \lambda_{\max})} \right) \|\nabla_{\mathcal{I}_k} F(x_k)\|_2^{2+p},$$

467 which combined with Lemma 2.4 and how the index set  $\mathcal{I}_k$  in Line 9 is defined gives

$$\begin{aligned}
468 \quad F(x_{k+1}) &\leq F(x_k) - \left( \frac{\kappa_2 \mu_{\min}^2 \xi \eta (1 - \eta)}{2\mu_{\max}^2 (L \kappa_2 (L_f + \lambda_{\max} \sqrt{n_{\mathcal{G}}})^p + \lambda_{\max})} \right) \|[s_k]_{\mathcal{I}_k}\|_2^{2+p} \\
469 \quad &\leq F(x_k) - \left( \frac{\kappa_2 \mu_{\min}^2 \xi \eta (1 - \eta) \varphi^{2+p}}{2\mu_{\max}^2 (L \kappa_2 (L_f + \lambda_{\max} \sqrt{n_{\mathcal{G}}})^p + \lambda_{\max})} \right) (\chi_k^{\text{cg}})^{2+p}, \\
470
\end{aligned}$$

471 thus completing the proof.  $\square$

472 The result in (4.11) motivates us to define the following subsets of  $\mathcal{K}_{\text{sd}}^{\text{cg}}$ :

$$(4.19) \quad \mathcal{K}_{\text{sd}, \text{big}}^{\text{cg}} := \{k \in \mathcal{K}_{\text{sd}}^{\text{cg}} : \chi_k^{\text{cg}} \geq c_1 / c_2\} \quad \text{and} \quad \mathcal{K}_{\text{sd}, \text{small}}^{\text{cg}} := \mathcal{K}_{\text{sd}}^{\text{cg}} \setminus \mathcal{K}_{\text{sd}, \text{big}}^{\text{cg}}.$$

474 This distinction plays a role in our complexity result. First, we require a lemma.

475 **LEMMA 4.7.** *The objective function  $f + r$  is monotonically decreasing over the*  
476 *sequence of iterates  $\{x_k\}$  and  $\lim_{k \rightarrow \infty} (f(x_k) + r(x_k)) =: F_{\min} > -\infty$ .*

477 *Proof.* It follows from Lemma 4.2 and Lemma 4.6 that the objective function is  
478 monotonically decreasing over the iterate sequence. The remaining conclusion of the  
479 lemma follows from the monotonicity property and Assumption 1.1.  $\square$

480 The main theorem can now be stated. It gives an upper bound on the number of  
481 iterations performed by Algorithm 3.1 before an approximate solution is obtained.



482 THEOREM 4.8. Let  $c_1$  and  $c_2$  be the constants defined in (4.12) and let us define  
 483  $c_3 := \eta\varphi^2/\alpha_0 > 0$ . For any  $\epsilon > 0$ , define  $\mathcal{K}_\epsilon := \{k \in \mathbb{N} : \max\{\chi_k^{\text{cg}}, \chi_k^{\text{pg}}\} > \epsilon\}$ . Then,

$$484 \quad (4.20) \quad \begin{aligned} |\mathcal{K}_{\rightarrow}^{\text{pg}} \cap \mathcal{K}_\epsilon| &\leq c_{\text{pg}}\epsilon^{-2} + 1, \\ |\mathcal{K}_{\text{sd},\text{big}}^{\text{cg}} \cap \mathcal{K}_\epsilon| &\leq c_{\text{big}}\epsilon^{-(1+p)} + 1, \quad \text{and} \\ |\mathcal{K}_{\text{sd},\text{small}}^{\text{cg}} \cap \mathcal{K}_\epsilon| &\leq c_{\text{small}}\epsilon^{-(2+p)} + 1 \end{aligned}$$

485 where the constants  $c_{\text{pg}}$ ,  $c_{\text{big}}$ , and  $c_{\text{small}}$  are given, respectively, by

$$486 \quad (4.21) \quad \begin{aligned} c_{\text{pg}} &:= (f(x_0) + r(x_0) - F_{\min})/c_3, \\ c_{\text{big}} &:= (f(x_0) + r(x_0) - F_{\min})/c_1, \quad \text{and} \\ c_{\text{small}} &:= (f(x_0) + r(x_0) - F_{\min})/c_2. \end{aligned}$$

487 Therefore, if  $\epsilon \geq c_1/c_2$ , then

$$488 \quad (4.22) \quad |\mathcal{K}_\epsilon| \leq (c_\downarrow^\alpha + c_{\text{pg}}\epsilon^{-2} + c_{\text{big}}\epsilon^{-(1+p)} + 2)(1 + n_{\mathcal{G}}) + n_{\mathcal{G}}$$

489 where  $c_\downarrow^\alpha$  is defined in (4.7); otherwise, i.e., if  $\epsilon < c_1/c_2$ , then

$$490 \quad (4.23) \quad |\mathcal{K}_\epsilon| \leq (c_\downarrow^\alpha + c_{\text{pg}}\epsilon^{-2} + c_{\text{big}}\epsilon^{-(1+p)} + c_{\text{small}}\epsilon^{-(2+p)} + 3)(1 + n_{\mathcal{G}}) + n_{\mathcal{G}}.$$

491 *Proof.* Note that the definitions of  $\mathcal{K}^{\text{cg}}$  and  $\mathcal{K}^{\text{pg}}$  together with Line 8 show that

$$492 \quad (4.24) \quad \chi_k^{\text{cg}} \geq \chi_k^{\text{pg}} \text{ for } k \in \mathcal{K}^{\text{cg}} \quad \text{and} \quad \chi_k^{\text{pg}} > \chi_k^{\text{cg}} \text{ for } k \in \mathcal{K}^{\text{pg}}.$$

493 Define  $\Delta_k := f(x_k) + r(x_k) - (f(x_{k+1}) + r(x_{k+1}))$  and  $m_k := \max\{\chi_k^{\text{pg}}, \chi_k^{\text{cg}}\}$ . Using  
 494 Lemma 4.2(i), Lemma 4.3, Lemma 4.6(ii), the definitions of  $c_3$  and  $\mathcal{K}_\epsilon$  in the statement  
 495 of the theorem, and (4.24) shows for arbitrary  $\bar{k} \in \mathbb{N}$  that

$$\begin{aligned} &f(x_0) + r(x_0) - (f(x_{\bar{k}+1}) + r(x_{\bar{k}+1})) = \sum_{0 \leq k \leq \bar{k}} \Delta_k \\ &\geq \sum_{\substack{k \in \mathcal{K}_{\rightarrow}^{\text{pg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} \Delta_k + \sum_{\substack{k \in \mathcal{K}_{\text{sd},\text{big}}^{\text{cg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} \Delta_k + \sum_{\substack{k \in \mathcal{K}_{\text{sd},\text{small}}^{\text{cg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} \Delta_k \\ &\geq \sum_{\substack{k \in \mathcal{K}_{\rightarrow}^{\text{pg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_3(\chi_k^{\text{pg}})^2 + \sum_{\substack{k \in \mathcal{K}_{\text{sd},\text{big}}^{\text{cg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_1(\chi_k^{\text{cg}})^{1+p} + \sum_{\substack{k \in \mathcal{K}_{\text{sd},\text{small}}^{\text{cg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_2(\chi_k^{\text{cg}})^{2+p} \\ &= \sum_{\substack{k \in \mathcal{K}_{\rightarrow}^{\text{pg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_3 m_k^2 + \sum_{\substack{k \in \mathcal{K}_{\text{sd},\text{big}}^{\text{cg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_1 m_k^{1+p} + \sum_{\substack{k \in \mathcal{K}_{\text{sd},\text{small}}^{\text{cg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_2 m_k^{2+p} \\ &\geq \sum_{\substack{k \in \mathcal{K}_{\rightarrow}^{\text{pg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_3 \epsilon^2 + \sum_{\substack{k \in \mathcal{K}_{\text{sd},\text{big}}^{\text{cg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_1 \epsilon^{1+p} + \sum_{\substack{k \in \mathcal{K}_{\text{sd},\text{small}}^{\text{cg}} \cap \mathcal{K}_\epsilon \\ 0 \leq k \leq \bar{k}}} c_2 \epsilon^{2+p}. \end{aligned}$$

497 From this inequality, Lemma 4.7, and (4.21) one finds that (4.20) follows.

498 Next, suppose that  $\epsilon \geq c_1/c_2$ . It then follows from (4.19) and (4.24) that  $\chi_k^{\text{cg}} =$   
 499  $\max\{\chi_k^{\text{pg}}, \chi_k^{\text{cg}}\} > \epsilon \geq c_1/c_2$  for all  $k \in \mathcal{K}^{\text{cg}}$ , which implies that  $\mathcal{K}_{\text{sd},\text{small}}^{\text{cg}} \cap \mathcal{K}_\epsilon = \emptyset$ . The  
 500 result in (4.22) follows from this observation, (4.20), (4.7), and since (by Lemma 4.6(i))  
 501 at most  $n_{\mathcal{G}}$  iterations in  $\mathcal{K}_0^{\text{cg}}$  can occur before the first, after the last, or between any  
 502 two iterations in  $\mathcal{K}_{\downarrow}^{\text{pg}} \cup \mathcal{K}_{\rightarrow}^{\text{pg}} \cup \mathcal{K}_{\text{sd}}^{\text{cg}}$ .

503 The final result (4.23) follows using the same argument as in the previous para-  
 504 graph, except now  $\mathcal{K}_{\text{sd},\text{small}}^{\text{cg}} \cap \mathcal{K}_\epsilon$  is no longer necessarily empty.  $\square$

505 We see from (4.23) that, for all sufficiently small  $\epsilon$ , the worst case complexity result  
 506 for Algorithm 3.1 is  $\epsilon^{-(2+p)}$ , which is worse than the  $\epsilon^{-2}$  result that holds for the PG  
 507 method. If one is concerned with such a result, the difference can be made arbitrarily  
 508 small (for a range of  $\epsilon$  values typically used in practice) by choosing  $p$  sufficiently  
 509 small. However, as is typical with well-designed second-derivative methods, although  
 510 the complexity bound is worse, it typically performs better (see Section 5).

511 **4.2. Local convergence.** We now consider the local convergence rate of the  
 512 iterates generated by Algorithm 3.1. Our analysis is performed under the following  
 513 additional assumption that will be assumed to hold throughout this section.

514 ASSUMPTION 4.3. *The function  $f$  is twice continuously differentiable and strongly*  
 515 *convex. It follows that there exists a unique solution  $x_*$  to the optimization*  
 516 *problem (1.1) with optimal support  $\mathcal{S}_* := \{i : [x_*]_{\mathcal{G}_i} \neq 0\}$ . Moreover, we assume that*  
 517  *$\nabla^2 f : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$  is Lipschitz continuous in a neighborhood of the solution  $x_*$ , and*  
 518 *that  $f + r$  is nondegenerate at  $x_*$  in the sense that  $\|[\nabla f(x_*)]_{\mathcal{G}_i}\|_2 < \lambda_i$  for all  $i \notin \mathcal{S}_*$ .*

519 Optimality conditions for problem (1.1) imply that  $\|[\nabla f(x_*)]_{\mathcal{G}_i}\|_2 \leq \lambda_i$  for all  
 520  $i \notin \mathcal{S}_*$ . Thus, the final condition in Assumption 4.3 is a strengthening of this fact.

521 ASSUMPTION 4.4. *The following algorithmic choices are made in Algorithm 3.1:*  
 522 *(i) The backtracking parameter is chosen to satisfy  $\eta \in (0, 1/2)$ .*  
 523 *(ii) For all sufficiently large  $k \in \mathbb{N}$ ,  $\mathcal{I}_k$  in Line 9/16 is chosen as*

$$524 \quad (4.25) \quad \mathcal{I}_k = \begin{cases} \mathcal{I}_k^{\text{cg}} & \text{if } k \in \mathcal{K}^{\text{cg}}, \\ \mathcal{I}_k^{\text{pg}} & \text{if } k \in \mathcal{K}^{\text{pg}}. \end{cases}$$

525 *(iii) For all sufficiently large  $k \in \mathcal{K}^{\text{cg}}$ ,  $H_k = \nabla_{\mathcal{I}_k}^2 (f+r)(x_k)$  is chosen in Line 10.*

526 The next result establishes that the iterate sequence converges to  $x_*$ .

527 THEOREM 4.9. *The iterate sequence  $\{x_k\}$  generated by Algorithm 3.1 satisfies*  
 528  *$\lim_{k \rightarrow \infty} x_k = x_*$  and  $\lim_{k \rightarrow \infty} \max\{\chi_k^{\text{pg}}, \chi_k^{\text{cg}}\} = 0$ .*

529 *Proof.* Theorem 4.8 gives  $\lim_{k \rightarrow \infty} \max\{\chi_k^{\text{pg}}, \chi_k^{\text{cg}}\} = 0$ . Since  $\{x_k\}$  is bounded due  
 530 to monotonicity of  $\{f(x_k) + r(x_k)\}$  (see Lemma 4.7) and Assumption 4.3, there exists  
 531 an infinite  $\mathcal{K} \subseteq \mathbb{N}$  and  $\hat{x}$  so that  $\lim_{k \in \mathcal{K}, k \rightarrow \infty} x_k = \hat{x}$ . It follows from Lemma 4.1 and  
 532 Lemma 4.3 that  $\hat{x}$  is a solution to problem (1.1), but with Assumption 4.3 this means  
 533 that  $\hat{x} = x_*$ , so  $\lim_{k \in \mathcal{K}, k \rightarrow \infty} x_k = x_*$ . The fact that the entire sequence  $\{x_k\}$  converges  
 534 to  $x_*$  follows from this fact, Assumption 4.3, and monotonicity of  $\{f(x_k) + r(x_k)\}$ .  $\square$

535 We now show for groups whose variables are all equal to zero at the solution that  
 536 the PG step will eventually predict them to be zero.

537 LEMMA 4.10. *For all  $i \notin \mathcal{S}_*$  and sufficiently large  $k$ , it holds that  $[x_k + s_k]_{\mathcal{G}_i} = 0$ .*

538 *Proof.* First note that Lemma 4.3 and the update strategy for  $\{\alpha_k\}$  in Algo-  
 539 rithm 3.1 ensure that there exists  $\bar{k}_1$  such that  $\alpha_k = \alpha_* > 0$  for all  $k \geq \bar{k}_1$ . Now, let  
 540  $i \notin \mathcal{S}_*$  so that  $[x_*]_{\mathcal{G}_i} = 0$ . It follows from Assumption 4.3 that

$$541 \quad \frac{\alpha_* \lambda_i}{\|[x_* - \alpha_* \nabla f(x_*)]_{\mathcal{G}_i}\|_2} = \frac{\lambda_i}{\|[\nabla f(x_*)]_{\mathcal{G}_i}\|_2} > 1.$$

542 Combining this with Theorem 4.9,  $\alpha_k = \alpha_* > 0$  for all  $k \geq \bar{k}_1$ , and Assumption 1.1  
 543 shows that there exists a  $\bar{k}_2 \geq \bar{k}_1$  such that  $1 - \alpha_k \lambda_i / \|[x_k - \alpha_k \nabla f(x_k)]_{\mathcal{G}_i}\|_2 < 0$  for all  
 544  $k \geq \bar{k}_2$ . Using this fact with (2.2) and (2.3) shows that  $[x_k + s_k]_{\mathcal{G}_i} = 0$  for all  $k \geq \bar{k}_2$ .  
 545 This completes the proof since the choice  $i \notin \mathcal{S}_*$  was arbitrary and  $n_{\mathcal{G}}$  is finite.  $\square$

546 We now show that, eventually, the set  $\mathcal{S}_*$  determines the sets  $\mathcal{I}_k^{\text{pg}}$  and  $\mathcal{I}_k^{\text{cg}}$ .

LEMMA 4.11. *For all sufficiently large  $k$ , it holds that*

$$\mathcal{I}_k^{\text{pg}} \equiv \{j \in \mathcal{G}_i : i \notin \mathcal{S}_*\} \quad \text{and} \quad \mathcal{I}_k^{\text{cg}} \equiv \{j \in \mathcal{G}_i : i \in \mathcal{S}_*\}$$

547 where the sets  $\mathcal{I}_k^{\text{pg}}$  and  $\mathcal{I}_k^{\text{cg}}$  are defined in (3.4).

548 *Proof.* Let  $\bar{k}_1$  be large enough so that the conclusion of Lemma 4.10 holds, i.e.,  
 549 if  $k \geq \bar{k}_1$  and  $i \notin \mathcal{S}_*$ , then  $[x_k + s_k]_{\mathcal{G}_i} = 0$ . Together with (3.2), this shows that  
 550  $\mathcal{G}_i \cap \bar{\mathcal{I}}_k^{\text{cg}} = \emptyset$  for all  $k \geq \bar{k}_1$  and  $i \notin \mathcal{S}_*$ , and thus  $\mathcal{G}_i \subseteq \mathcal{I}_k^{\text{pg}}$  (see (3.4)) for all  $k \geq \bar{k}_1$   
 551 and  $i \notin \mathcal{S}_*$ . In other words, it holds that  $\{j \in \mathcal{G}_i : i \notin \mathcal{S}_*\} \subseteq \mathcal{I}_k^{\text{pg}}$  for all  $k \geq \bar{k}_1$ .

552 Next, we prove that there exists  $\bar{k}_2$  such that  $\mathcal{I}_k^{\text{pg}} \subseteq \{j \in \mathcal{G}_i : i \notin \mathcal{S}_*\}$  for all  $k \geq \bar{k}_2$ .  
 553 For a proof by contradiction, suppose that there exists an infinite subsequence  $\mathcal{K} \subseteq \mathbb{N}$   
 554 and group index  $\bar{i}$  such that  $\mathcal{G}_{\bar{i}} \subseteq \mathcal{I}_k^{\text{pg}}$  and  $\bar{i} \in \mathcal{S}_*$  for all  $k \in \mathcal{K}$ . Since  $\mathcal{G}_{\bar{i}} \subseteq \mathcal{I}_k^{\text{pg}}$  for all  
 555  $k \in \mathcal{K}$ , it follows from (3.2), (3.3), and (3.4) that at least one of

$$556 \quad (4.26) \quad [x_k]_{\mathcal{G}_{\bar{i}}} = 0, \quad [x_k + s_k]_{\mathcal{G}_{\bar{i}}} = 0, \quad \|[x_k]_{\mathcal{G}_{\bar{i}}}\|_2 < \kappa_1 \|\nabla_{\mathcal{G}_{\bar{i}}}(f+r)(x_k)\|_2 \quad \text{or}$$

$$557 \quad (4.27) \quad \|[x_k]_{\mathcal{G}_{\bar{i}}}\|_2 < \kappa_2 \|\nabla_{\bar{\mathcal{I}}_k^{\text{cg}}}(f+r)(x_k)\|_2^p$$

559 holds for all  $k \in \mathcal{K}$ . However, since  $\bar{i} \in \mathcal{S}_*$ , it follows from Theorem 4.9 that the first  
 560 condition in (4.26) does not hold for all sufficiently large  $k \in \mathcal{K}$ . Also, it follows from  
 561 Theorem 4.9, the facts that  $\chi_k^{\text{pg}} \equiv \|[s_k]_{\mathcal{I}_k^{\text{pg}}}\|_2$  and  $\chi_k^{\text{cg}} \equiv \|[s_k]_{\mathcal{I}_k^{\text{cg}}}\|_2$ , and the fact that  
 562  $\mathcal{I}_k^{\text{cg}} \cup \mathcal{I}_k^{\text{pg}} = \{1, \dots, n\}$  that  $\lim_{k \rightarrow \infty} \|s_k\|_2 = 0$ , which combined with  $\bar{i} \in \mathcal{S}_*$  proves  
 563 that  $[x_k + s_k]_{\mathcal{G}_{\bar{i}}} \neq 0$  for all sufficiently large  $k$ . Hence, the second condition in (4.26)  
 564 does not hold for all sufficiently large  $k \in \mathcal{K}$ . Next, from the optimality conditions  
 565 for problem (1.1), the fact that  $\bar{i} \in \mathcal{S}_*$ , Theorem 4.9, Assumption 1.1, and the fact  
 566 that  $f+r$  is differentiable over the variables in  $\mathcal{G}_{\bar{i}}$  for sufficiently large  $k$  that we have  
 567  $\lim_{k \rightarrow \infty} \|\nabla_{\mathcal{G}_{\bar{i}}}(f+r)(x_k)\|_2 = 0$ . This limit,  $[x_k]_{\mathcal{G}_{\bar{i}}} \neq 0$ , and Theorem 4.9 show that  
 568  $\|[x_k]_{\mathcal{G}_{\bar{i}}}\|_2 \geq \kappa_1 \|\nabla_{\mathcal{G}_{\bar{i}}}(f+r)(x_k)\|_2$  for all sufficiently large  $k$ , meaning that the third  
 569 condition in (4.26) does not hold for all sufficiently large  $k \in \mathcal{K}$ . Therefore, we must  
 570 conclude that the inequality in (4.27) holds for all sufficiently large  $k \in \mathcal{K}$ . Combining  
 571 this with  $\bar{i} \in \mathcal{S}_*$  shows that there exists  $\epsilon > 0$  such that

$$572 \quad (4.28) \quad \|\nabla_{\bar{\mathcal{I}}_k^{\text{cg}}}(f+r)(x_k)\|_2 \geq \epsilon > 0 \quad \text{for all sufficiently large } k \in \mathcal{K},$$

573 which in particular shows that  $\bar{\mathcal{I}}_k^{\text{cg}} \neq \emptyset$  for all sufficiently large  $k \in \mathcal{K}$ . Since the  
 574 optimality conditions for problem (1.1) together with Theorem 4.9, Assumption 1.1,  
 575 and the fact that  $f+r$  is differentiable over the variables in  $\mathcal{G}_i$  for sufficiently large  $k$   
 576 imply that  $\lim_{k \rightarrow \infty} \|\nabla_{\mathcal{G}_i}(f+r)(x_k)\|_2 = 0$  for all  $i \in \mathcal{S}_*$ , we must conclude from (4.28)  
 577 that, for all sufficiently large  $k \in \mathcal{K}$ , there exists an  $i_k \notin \mathcal{S}_*$  such that  $\mathcal{G}_{i_k} \subseteq \bar{\mathcal{I}}_k^{\text{cg}}$ .  
 578 However, Lemma 4.10 yields  $[x_k + s_k]_{\mathcal{G}_{i_k}} = 0$  for all sufficiently large  $k \in \mathcal{K}$ , which  
 579 together with (3.2) shows that  $\mathcal{G}_{i_k} \not\subseteq \bar{\mathcal{I}}_k^{\text{cg}}$ , which is a contradiction. Therefore, there  
 580 exists  $\bar{k}_2$  such that  $\mathcal{I}_k^{\text{pg}} \subseteq \{j \in \mathcal{G}_i : i \notin \mathcal{S}_*\}$  for all  $k \geq \bar{k}_2$ .

581 The conclusions of the two previous paragraphs yields  $\mathcal{I}_k^{\text{pg}} \equiv \{j \in \mathcal{G}_i : i \notin \mathcal{S}_*\}$   
 582 for all sufficiently large  $k$ . The final assertion, namely that  $\mathcal{I}_k^{\text{cg}} \equiv \{j \in \mathcal{G}_i : i \in \mathcal{S}_*\}$ ,  
 583 follows from the fact that  $\mathcal{I}_k^{\text{pg}}$  and  $\mathcal{I}_k^{\text{cg}}$  partition  $\{1, 2, \dots, n\}$  for every iteration  $k$ .  $\square$

584 The next result shows that, for iterations  $k$  sufficiently large, the support of  $x_k$   
 585 agrees with the support of the solution  $x_*$ .

LEMMA 4.12. *For all sufficiently large  $k$ , it holds that*

$$[x_k]_{\mathcal{G}_i} \neq 0 \quad \text{for all } i \in \mathcal{S}_* \quad \text{and} \quad [x_k]_{\mathcal{G}_i} = 0 \quad \text{for all } i \notin \mathcal{S}_*.$$

586 *Proof.* Theorem 4.9 shows that  $[x_k]_{\mathcal{G}_i} \neq 0$  for all sufficiently large  $k$  and all  $i \in \mathcal{S}_*$ ,  
587 which is the first desired result. Hence, let us proceed by considering arbitrary  $i \notin \mathcal{S}_*$ .  
588 Assumption 4.4(ii), Lemma 4.10, Lemma 4.11, and Lemma 4.3 ensure the existence  
589 of an iteration  $\bar{k}$  such that, for all  $k \geq \bar{k}$ , the following hold:

$$590 \quad (4.29) \quad \mathcal{G}_i \subseteq \mathcal{I}_k^{\text{pg}}, \quad [x_k + s_k]_{\mathcal{G}_i} = 0, \quad \text{and} \quad \alpha_k = \alpha_{\bar{k}}.$$

591 We claim that the second desired result follows from (4.29) if there exists some suffi-  
592 ciently large  $\hat{k} \geq \bar{k}$  such that  $\hat{k} \in \mathcal{K}^{\text{pg}}$  and  $[x_{\hat{k}+1}]_{\mathcal{G}_i} = [x_{\hat{k}} + s_{\hat{k}}]_{\mathcal{G}_i} = 0$ . Indeed, since  $i$   
593 is an arbitrary element from  $\{1, \dots, n_{\mathcal{G}}\} \setminus \mathcal{S}_*$ ,  $n_{\mathcal{G}}$  is finite, and the second condition  
594 in (4.29) shows that values of the variables in  $\mathcal{G}_i$  can only be modified if  $k \in \mathcal{K}^{\text{pg}}$ , the  
595 existence of such  $\hat{k}$  along with (4.29) shows that iteration  $\hat{k} \in \mathcal{K}^{\text{pg}}$  sets  $[x_{\hat{k}+1}]_{\mathcal{G}_i}$  to  
596 zero, and these variables will remain zero for all future iterations.

597 Let us now show the existence of such  $\hat{k} \geq \bar{k}$ . We claim that there exists  $k \geq \bar{k}$   
598 such that  $[x_k]_{\mathcal{G}_i} = 0$ . For a proof by contradiction, suppose that  $[x_k]_{\mathcal{G}_i} \neq 0$  for all  
599  $k \geq \bar{k}$ . Combining this with Theorem 4.9,  $i \notin \mathcal{S}_*$ , and the fact that the variables  
600 in  $\mathcal{G}_i$  can have their values changed only if  $k \in \mathcal{K}^{\text{pg}}$  implies that there exists  $\hat{k} \geq \bar{k}$   
601 such that  $\hat{k} \in \mathcal{K}^{\text{pg}}$ . Now, since  $\bar{k} \in \mathcal{K}^{\text{pg}}$  and  $\alpha_k = \alpha_{\bar{k}}$  for all  $k \geq \bar{k}$ , it follows from  
602 Algorithm 3.1 that  $\text{flag}_k^{\text{pg}} = \text{same}.\alpha$  is returned in Line 17. Using this fact, the update  
603 used in Line 56, and (4.29) shows that  $[x_{\hat{k}+1}]_{\mathcal{G}_i} = [x_{\hat{k}} + s_{\hat{k}}]_{\mathcal{G}_i} = 0$ .  $\square$

604 We require one more lemma that shows that eventually all iterations are in  $\mathcal{K}_{\text{sd}}^{\text{cg}}$ .

605 **LEMMA 4.13.** *For all  $k$  sufficiently large, it holds that  $k \in \mathcal{K}_{\text{sd}}^{\text{cg}}$ .*

606 *Proof.* We first show that all sufficiently large  $k$  are in  $\mathcal{K}^{\text{cg}}$ . It follows from  
607 Lemma 4.11 that  $\mathcal{I}_k^{\text{pg}} \equiv \{j \in \mathcal{G}_i : i \notin \mathcal{S}_*\}$  for all sufficiently large  $k$ . Combining this  
608 with Lemma 4.12 and Lemma 4.10 shows that there exists an iteration  $\bar{k}$  such that  
609  $[x_k]_{\mathcal{I}_k^{\text{pg}}} = 0$  and  $[x_k + s_k]_{\mathcal{I}_k^{\text{pg}}} = 0$  for all  $k \geq \bar{k}$ , which means that  $\chi_k^{\text{pg}} = \|[s_k]_{\mathcal{I}_k^{\text{pg}}}\|_2 = 0$   
610 for all  $k \geq \bar{k}$ . It follows from this fact, Line 8, and Assumption 4.1 that  $k \in \mathcal{K}^{\text{cg}}$  for all  
611  $k \geq \bar{k}$ . Now, notice that at most  $n_{\mathcal{G}} - 1$  iterations from  $\bar{k}$  onward can be in  $\mathcal{K}_0^{\text{cg}}$  because  
612 of Lemma 4.6(i). (Every iteration  $k \in \mathcal{K}_0^{\text{cg}}$  fixes at least one new group of variables to  
613 zero and if they ever all become zero so that  $\mathcal{I}_k^{\text{cg}} = \emptyset$ , then the contradiction  $k \in \mathcal{K}^{\text{pg}}$   
614 is reached.) Therefore, it follows that all sufficiently large  $k$  must be in  $\mathcal{K}_{\text{sd}}^{\text{cg}}$ .  $\square$

615 We can now state our main local convergence result.

616 **THEOREM 4.14.** *If in Algorithm 3.2 we choose either  $q \in (1, 2]$ , or  $q = 1$  and*  
617  *$\{\mu_k\} \rightarrow 0$ , then  $\{x_k\} \rightarrow x_*$  at a superlinear rate. In particular, if we choose  $q = 2$ ,*  
618 *then the rate of convergence is quadratic.*

619 *Proof.* It follows from Lemma 4.11, Lemma 4.12, and Lemma 4.13 that, for all  
620 sufficiently large  $k$ , the iterates generated by Algorithm 3.1 satisfy the recurrence  
621  $x_{k+1} = x_k + \xi^{j_k} d_k$ , where  $j_k$  is the result of the backtracking Armijo line search  
622 in Line 47,  $\|[x_k]_{\mathcal{I}_k^{\text{pg}}}\|_2 = \|[d_k]_{\mathcal{I}_k^{\text{pg}}}\|_2 = 0$ , and  $[d_k]_{\mathcal{I}_k^{\text{cg}}} = \bar{d}_k$  with  $\bar{d}_k$  computed by  
623 Algorithm 3.2 to satisfy (3.9). In other words, for all sufficiently large  $k$ , we have  
624  $[x_k]_{\mathcal{I}_k^{\text{pg}}} = [x_*]_{\mathcal{I}_k^{\text{pg}}} = 0$  and the values of the variables in  $\mathcal{I}_k^{\text{cg}} \equiv \{j \in \mathcal{G}_i : i \in \mathcal{S}_*\}$   
625 are updated exactly as those of an inexact Newton method for computing a root of  
626  $\nabla_{\mathcal{I}_k^{\text{cg}}}(f+r)$ . Since, by Theorem 4.9, we have  $\lim_{k \rightarrow \infty} x_k = x_*$ , the desired conclusions  
627 follow under the stated conditions from [11, Theorem 3.3] and noting the well-known  
628 result that the unit step size  $\xi^{j_k} = 1$  is accepted (asymptotically) by a backtracking  
629 Armijo line search when  $\eta \in (0, 1/2)$  (see Assumption 4.4) under our assumptions.  $\square$

630 Theorem 4.14 states conditions under which Algorithm 3.1 yields a superlinear, or  
631 even quadratic, rate of local convergence. The neighborhood about  $x_*$  in which such

632 a rate will be achieved, and the explicit constants in the convergence rate that will be  
 633 achieved, depend as usual on magnitudes of a Lipschitz constant for  $\nabla_{\mathcal{I}_*}(f+r)$  and  
 634 an upper bound on a norm of the inverse of  $\nabla_{\mathcal{I}_*}^2(f+r)$ , where  $\mathcal{I}_* := \{j \in \mathcal{G}_i : i \in \mathcal{S}_*\}$ .  
 635 Due to the properties of the regularizer  $r$ , the latter of these values may be inversely  
 636 proportional to the norms of the groups of variables in the support at the solution.

637 **5. Numerical Results.** In this section, we present the results of numerical  
 638 experiments with an implementation of **FaRSA-Group** (Algorithm 3.1) applied to solve  
 639 a collection of group sparse regularized logistic regression problems of the form

$$640 \quad (5.1) \quad \min_{x \in \mathbb{R}^n} \frac{1}{N} \sum_{i=1}^N \log \left( 1 + e^{-y_i x^T d_i} \right) + \sum_{i=1}^{n_{\mathcal{G}}} \lambda_i \| [x]_{\mathcal{G}_i} \|_2,$$

641 where  $d_i \in \mathbb{R}^n$  is the  $i$ th data point,  $N$  is the number of data points in the data set,  
 642  $y_i \in \{-1, 1\}$  is the class label for the  $i$ th data point, and  $\lambda_i$  is the weight parameter  
 643 for the  $i$ th group. We first describe details of our implementation, then describe the  
 644 data sets considered in our experiments, and finally present our experimental results.

645 **5.1. Implementation details.** We have developed a Python implementation  
 646 of **FaRSA-Group** that is available upon request. The values of the input parameters  
 647 for Algorithm 3.1 and Algorithm 3.2 that we used are given in Figure 5.1 (with some  
 648 caveats that are mentioned in the following paragraph).

We initialized  $x_0$  as the zero vector and  $\alpha_0$  as an estimate of the inverse of the Lipschitz constant of  $f$  at  $x_0$ . To be precise, our software randomly generated a vector  $y_0 \in \mathbb{R}^n$  such that  $\|x_0 - y_0\|_2 = 10^{-8}$ , and then set  $\alpha_0 = \min\{1, \|x_0 - y_0\|_2 / \|\nabla f(x_0) - \nabla f(y_0)\|_2\}$ . Since  $\varphi = 1$ , it follows from Algorithm 3.1 that (4.25) holds for all  $k \in \mathbb{N}$ . (However, for data sets with  $N < n$ , we initially chose  $\varphi = 0.8$  and switched to  $\varphi = 1$  when an iteration in  $\mathcal{K}^{\text{cg}}$  satisfied  $f(x_k) - f(x_{k+1}) \leq 10^{-3}$ . When  $N < n$ , the matrix  $\nabla^2 f(x_k)$  is singular, which in practice often led to large CG directions and multiple backtracks in the line search. These ill effects were partly remedied by this scheme for updating  $\varphi$ .) When defining the set  $\mathcal{I}_k^{\text{small}}$  in (3.3), we used  $\tilde{\kappa}_{2,i} = \kappa_2 |\mathcal{G}_i| / \|\bar{\mathcal{I}}_k^{\text{cg}}\|$  in place of  $\kappa_2$  for all  $i$  such that  $\mathcal{G}_i \subseteq \bar{\mathcal{I}}_k^{\text{cg}}$ . This choice accounted for the fact that the two different norms in (3.3) are associated with vectors of different dimension. Note that since  $(1/n)\kappa_2 \leq \tilde{\kappa}_{2,i} \leq n\kappa_2$ , this choice is easily incorporated into the analysis in Section 4. The choice of  $H_k$  in Line 10 was based on a regularization of the exact second-derivatives of  $f$ . In particular, for any scalar  $\delta \geq 0$ , consider

$$\frac{1}{N} D^T \Sigma_{\delta}(x) D \approx \nabla^2 f(x)$$

where  $D^T := [d_1, d_2, \dots, d_N]$  and  $\Sigma_{\delta}(x)$  is the diagonal matrix with  $i$ th diagonal entry

$$[\Sigma_{\delta}(x)]_{ii} := \max\{\sigma_i(x)(1 - \sigma_i(x)), \delta\} \quad \text{with} \quad \sigma_i(x) := \exp(y_i d_i^T x) / (1 + \exp(y_i d_i^T x))$$

for all  $i \in \{1, 2, \dots, N\}$ . Notice that if  $\delta = 0$ , then  $(1/N)D^T \Sigma_0(x) D \equiv \nabla^2 f(x)$ . In order to use a small amount of regularization in our tests, we chose  $\delta = 10^{-8}$ . With this choice of  $\delta$ , our choice of  $H_k$  in Line 10 can now be written as

$$H_k \leftarrow \left[ \frac{1}{N} D^T \Sigma_{\delta}(x_k) D \right]_{\mathcal{I}_k \mathcal{I}_k} + \nabla_{\mathcal{I}_k \mathcal{I}_k}^2 r(x_k),$$

FIG. 5.1. Parameter values used in our tests for Algorithm 3.1 and Algorithm 3.2.

param.	value	param.	value
$\varphi$	1	$\kappa_1$	0.1
$\xi$	0.5	$\kappa_2$	$10^{-2}$
$\eta$	$10^{-3}$	$\theta$	$\pi/4$
$\zeta$	0.8	$q$	1
$p$	2	$\mu_k$	1

649 where we remind the reader that  $\nabla_{\mathcal{I}_k}^2 r(x_k)$  is well defined because the construction  
 650 of  $\mathcal{I}_k \subseteq \mathcal{I}_k^{\text{cg}}$  ensures that  $[x_k]_{\mathcal{G}_i} \neq 0$  for all  $\mathcal{G}_i \subseteq \mathcal{I}_k$ .

651 In Algorithm 3.2, we applied the CG method to the system  $H_k d = -g_k$  to ap-  
 652 proximately solve the optimization problem defined in Line 26. As pointed out in  
 653 Section 3.2, the direction associated with every iteration of the CG algorithm satisfies  
 654 condition (3.7) and condition (3.8), which were required to establish the complexity  
 655 result in Theorem 4.8. To reduce the cost of the CG computation and limit the number  
 656 of backtracking steps required by Algorithm 3.3, we terminated Algorithm 3.2 when at  
 657 least one of three conditions was satisfied. To describe these conditions checked dur-  
 658 ing the  $k$ th iteration, let  $d_{j,k}$  denote the  $j$ th CG iterate and let  $t_{j,k} := \|H_k d_{j,k} + g_k\|_2$   
 659 denote the  $j$ th CG residual. The three conditions are given by

660 (5.2a) 
$$t_{j,k} \leq \max\{\min\{0.1t_{0,k}, t_{0,k}^{1.5}\}, 10^{-10}\},$$

661 (5.2b) 
$$\|d_{j,k}\| \geq 10^3 \min\{1, \|\nabla_{\mathcal{I}_k}(f+r)(x_k)\|_2\}, \text{ and}$$

662 (5.2c) 
$$j = |\mathcal{I}_k|.$$

664 Outcome (5.2a) is the ideal termination condition since it indicates that the residual  
 665 of the linear system has been sufficiently reduced (see (3.9)). Outcome (5.2b) serves  
 666 as a trust-region constraint on the norm of the trial step  $d_k$ ; in particular, when the  
 667 inequality in (5.2b) holds, the size of the CG iterate  $d_{j,k}$  is relatively large, indicating  
 668 that  $x_k$  is not close to an optimal solution. Therefore, we restrict its size with the  
 669 intent of needing fewer backtracking steps during the subsequent line search. Out-  
 670 come (5.2c) caps the number of CG iterations to  $|\mathcal{I}_k|$  (the size of the reduced space)  
 671 since, in exact arithmetic, CG converges to an exact solution in at most  $|\mathcal{I}_k|$  iterations.

Algorithm 3.1 decreases the value of the PG parameter (see Line 19) for the next  
 iteration using a simple multiplicative factor when  $\text{flag}_k^{\text{pg}} = \text{decrease\_}\alpha$ . However,  
 in practice, we found an adaptation of the approach in [9] to be more efficient. To  
 describe this approach, let  $d_k$  and  $\xi^{j_k}$  be the search direction and step size used to  
 obtain  $x_{k+1} = x_k + \xi^{j_k} d_k$ . It is well known [2, Lemma 5.7] that if  $\alpha \in (0, 1/L_f]$ ,  
 then  $f(x_{k+1}) \leq f(x_k) + \xi^{j_k} \nabla f(x_k)^T d_k + \frac{1}{2\alpha} \|\xi^{j_k} d_k\|_2^2$ . Setting this inequality to be an  
 equality and then solving for  $\alpha$ , one obtains

$$\hat{\alpha}_k := \frac{\|\xi^{j_k} d_k\|_2^2}{2(f(x_{k+1}) - f(x_k) - \xi^{j_k} \nabla f(x_k)^T d_k)},$$

672 which can be viewed as a local Lipschitz constant estimate for  $f$  at  $x_k$ . In our tests,  
 673 we updated the PG parameter at the end of each iteration of Algorithm 3.1 as

674 (5.3) 
$$\alpha_{k+1} \leftarrow \min\{1, \hat{\alpha}_k/2\}.$$

675 Although this PG parameter update strategy worked better than the basic strategy in  
 676 Algorithm 3.1 (see Line 19 and Line 21), it is not covered by our analysis in Section 4.  
 677 However, a simple modification of our analysis would be to allow the update in (5.3)  
 678 to increase the PG parameter at most a finite number of times, say 100 times, at  
 679 which point the update  $\alpha_{k+1} \leftarrow \min\{\alpha_k, \hat{\alpha}_k/2\} \leq \alpha_k$  would be used. This strategy  
 680 is covered by our earlier analysis (with a larger constant in the complexity result).

681 We terminate our algorithm when  $\max\{\chi_k^{\text{cg}}, \chi_k^{\text{pg}}\} \leq 10^{-6} \max\{\chi_0^{\text{cg}}, \chi_0^{\text{pg}}, 1\}$ .

682 **5.2. Data sets.** We tested **FARSA-Group** on problem (5.1) using data sets from  
 683 the LIBSVM repository.<sup>1</sup> From this repository, we excluded all regression instances

<sup>1</sup><https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>



684 and multiple-class (greater than two) classification instances. We compared the per-  
685 formance of our algorithm to the well-cited package `gglasso` [32], which is a state-of-  
686 the-art group-wise majorization descent method.<sup>2</sup> Since `gglasso` does not support  
687 sparse data matrix inputs, we excluded all data sets that were too large to be stored in  
688 memory (6GB). Finally, for the adult data (a1a–a9a) and webpage data (w1a–w8a),  
689 we used only the largest instances, namely a9a and w8a. This left us with our final  
690 subset of 25 data sets that can be found in Table 5.1.

691 Scaling of the data sets can be important. If the LIBSVM website indicated that a  
692 data set was already scaled, then we used the data set without modification. However,  
693 when the website did not indicate that scaling for a data set was used, we scaled each  
694 column of the feature data (i.e., feature-wise scaling) into the range  $[-1, 1]$  by dividing  
695 each of its entries by the largest entry in absolute value. Labels for some data sets  
696 (e.g., breast-cancer, covtype, liver-disorders, mushrooms, phishing, skin-nonskin and  
697 svmguide1) do not take values in  $\{-1, 1\}$ , but rather in  $\{0, 1\}$  or  $\{1, 2\}$ . For these  
698 data sets, we mapped the smaller label to  $-1$  and the larger label to  $1$ .

TABLE 5.1

*The first column (data set) gives the name of the data set. The second column (N) and third column (n) indicate the number of data points and problem dimension, respectively. The fourth column (scale) provides the feature-wise scaling used: each feature is either scaled into the given interval or scaled to have mean zero ( $\mu = 0$ ) and variance one ( $\sigma^2 = 1$ ). The fifth column (who) indicates whether the data set came pre-scaled from the LIBSVM website (website), or it did not come pre-scaled and we scaled it (us) as described in Section 5.2. Finally, the sixth column (used) indicates the number of problem instances used in the numerical results presented in Figure 5.2.*

data set	N	n	scale	who	used
a9a	32561	123	[0,1]	website	8
australian	690	140	[-1,1]	website	2
breast-cancer	683	10	[-1,1]	website	0
cod-rna	59535	8	[-1,1]	us	8
colon-cancer	62	2000	$(\mu, \sigma^2) = (0, 1)$	website	8
covtype.binary	581012	54	[0,1]	website	8
diabetes	768	8	[-1,1]	website	0
duke breast-cancer	44	7192	$(\mu, \sigma^2) = (0, 1)$	website	8
fourclass	862	2	[-1,1]	website	0
german-numer	1000	24	[-1,1]	website	0
gisette	6000	5000	[-1,1]	website	8
heart	270	13	[-1,1]	website	2
ijcnn1	49990	22	[-1.5, 1.5]	website	8
ionosphere	351	34	[-1,1]	website	0
leukemia	38	7129	$(\mu, \sigma^2) = (0, 1)$	website	8
liver-disorders	145	5	[-1,1]	website	0
madelon	2000	500	[-1,1]	us	8
mushrooms	8124	112	[0,1]	website	6
phishing	11055	68	[0,1]	website	7
skin-nonskin	245057	3	[-1,1]	us	8
splice	1000	60	[-1,1]	website	0
sonar	208	60	[-1,1]	website	4
svmguide1	3089	4	[-1,1]	us	0
svmguide3	1243	21	[-1,1]	website	0
w8a	49749	300	[0,1]	website	8

<sup>2</sup><https://cran.r-project.org/web/packages/gglasso>



**5.3. Experimental setup and test results.** We tested **FaRSA-Group** and **gglasso** for solving problem (5.1) using the data sets in Table 5.1. All default settings for **gglasso** were used, including the same starting point  $x_0 = 0$  used by **FaRSA-Group**. We considered four group structures and two different solution sparsity levels. Specifically, we considered the four different numbers of groups

$$\text{number of groups} \in \{ \lfloor 0.25n \rfloor, \lfloor 0.50n \rfloor, \lfloor 0.75n \rfloor, n \},$$

where  $n$  is the problem dimension; notice that the last setting recovers  $\ell_1$ -norm regularization. Then, for a given number of groups, the variables were sequentially distributed (as evenly as possible) to the groups; e.g., 10 variables among 3 groups would have been distributed as  $\mathcal{G}_1 = \{1, 2, 3\}$ ,  $\mathcal{G}_2 = \{4, 5, 6\}$ , and  $\mathcal{G}_3 = \{7, 8, 9, 10\}$ . For the two different solution sparsity levels, we considered groups weights

$$\lambda_i = 0.1\lambda_{\min}\sqrt{|\mathcal{G}_i|} \quad \text{and} \quad \lambda_i = 0.01\lambda_{\min}\sqrt{|\mathcal{G}_i|}$$

699 where  $\lambda_{\min} = \min \{ \lambda \geq 0 : \text{the solution to (5.1) with } \lambda_i = \lambda\sqrt{|\mathcal{G}_i|} \text{ is } x = 0 \}$  (see [32,  
700 equation (23)]). Since there were 25 data sets, a total of 200 problem instances were  
701 tested (each data set has 8 instances). The experiments were conducted using the  
702 cluster in the Computational Optimization Research Laboratory (COR@L) at Lehigh  
703 University with an AMD Opteron Processor 6128 2.0 GHz CPU. In the following  
704 paragraphs, we compared the performance of **FaRSA-Group** with that of **gglasso**  
705 with respect to CPU time (seconds), final objective value, and solution sparsity.

706 First consider the CPU time. For  
707 each problem instance, we allowed a  
708 maximum of 1000 seconds. If the CPU  
709 time in a run went above this limit, we  
710 terminated that run and considered the  
711 algorithm to have failed. Out of the 200  
712 problem instances, **FaRSA-Group** failed 2  
713 times and **gglasso** failed 7 times. Fig-  
714 ure 5.2 illustrates a performance profile  
715 based on [24] for comparing the comput-  
716 ing times on problem instances that  
717 **FaRSA-Group** and/or **gglasso** took at  
718 least 1 second to terminate; this resulted  
719 in 109 problem instances. The last col-  
720 umn of Table 5.1 gives the number of  
721 instances for each data set used in this profile. Each bar in the plot corresponds to a  
722 problem instance, with the height of the bar given by

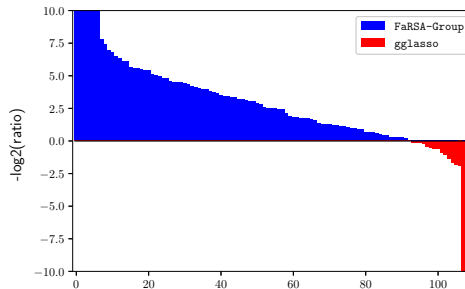


FIG. 5.2. Performance profile for CPU time (seconds). **FaRSA-Group** outperforms **gglasso** on 93 of the 109 problem instances. For each problem instance, the height of the bar is given by (5.4).

$$723 \quad (5.4) \quad -\log_2 \left( \frac{\text{time required by FaRSA-Group}}{\text{time required by gglasso}} \right).$$

724 Therefore, an upward pointing bar indicates that **FaRSA-Group** took less time to find  
725 the optimal solution for that problem instance and a downward pointing bar means  
726 that **gglasso** took less time, and in either case the size of the bar indicates the  
727 magnitude of the outperformance factor. A bar that reaches the y-axis limit of  $\pm 10$   
728 is used when indicating that an algorithm was successful when solving a problem  
729 instance while the competing algorithm was unsuccessful.

730 To compare final objective function values, let  $F_{\text{FaRSA-Group}}$  and  $F_{\text{gglasso}}$  denote (for  
731 a given problem instance) the objective values returned by **FaRSA-Group** and **gglasso**,

732 respectively. If  $F_{\text{gglasso}} - F_{\text{FaRSA-Group}} > 10^{-8}$ , then we considered **FaRSA-Group** to  
733 have obtained a lower objective function value; if  $F_{\text{FaRSA-Group}} - F_{\text{gglasso}} > 10^{-8}$ , then  
734 we considered **gglasso** to have obtained a lower objective function value; and if  
735  $|F_{\text{FaRSA-Group}} - F_{\text{gglasso}}| \leq 10^{-8}$ , then we considered them to have performed equally.  
736 From the 109 problem instances that at least one algorithm took at least one second to  
737 terminate, **FaRSA-Group** outperformed **gglasso** 95 times and **gglasso** outperformed  
738 **FaRSA-Group** 7 times. From the entire 200 instances, **FaRSA-Group** outperformed  
739 **gglasso** 153 times and **gglasso** outperformed **FaRSA-Group** 35 times.

740 In terms of solution sparsity, we considered **FaRSA-Group** to have outperformed  
741 **gglasso** if the following two conditions held: (i) all zero groups in the **gglasso** solu-  
742 tion were also zero groups in the **FaRSA-Group** solution, and (ii) the solution returned  
743 by **FaRSA-Group** had at least one zero group that was not a zero group in the **gglasso**  
744 solution. A similar criteria was used to define when **gglasso** was considered to have  
745 outperformed **FaRSA-Group**. From the 109 test instances, **FaRSA-Group** outperformed  
746 **gglasso** in 30 cases and **gglasso** outperformed **FaRSA-Group** in 7 cases. From the  
747 entire collection of 200 problem instances, **FaRSA-Group** outperformed **gglasso** in 33  
748 cases and **gglasso** outperformed **FaRSA-Group** in 8 cases.

749 **6. Conclusion.** We presented a new framework for solving optimization prob-  
750 lems that incorporate group sparsity-inducing regularization by using subspace ac-  
751 celeration, domain decomposition, and support identification. In terms of theory,  
752 we proved a complexity result on the maximum number of iterations before an  $\epsilon$ -  
753 approximate solution is computed (Theorem 4.8), and a local superlinear convergence  
754 rate (Theorem 4.14). The strong convergence theory was supported by experimental  
755 results for minimizing a group sparsity-regularized logistic function for the task of clas-  
756 sification. In terms of robustness, computational time, final objective value obtained,  
757 and solution sparsity, the numerical results showed that our proposed **FaRSA-Group**  
758 framework outperformed a state-of-the-art method.

759 **Appendix A. Proofs.** In this appendix, for completeness, we provide detailed  
760 proofs of the results from Section 2 related to the PG computations.

761 **Proof of Lemma 2.1.** Let  $x_+ = T(\bar{x}, \bar{\alpha})$  denote the PG update in (2.1) so that  
762  $x_+ = \bar{x} + s(\bar{x}, \bar{\alpha})$  with  $s(\bar{x}, \bar{\alpha})$  defined in (2.2). It follows from the optimality conditions  
763 for the problem in (2.1) that there exists  $g_+ \in \partial r(x_+)$  such that

$$764 \quad (\text{A.1}) \quad x_+ - \bar{x} + \bar{\alpha} \nabla f(\bar{x}) + \bar{\alpha} g_+ = 0.$$

766 Next, for an arbitrary  $g_{f+r} \in \partial(f+r)(\bar{x})$ , it follows from Assumption 1.1 and [4,  
767 Proposition 5.4.6] that there exists  $g_r \in \partial r(\bar{x})$  satisfying  $g_{f+r} = \nabla f(\bar{x}) + g_r$ . From the  
768 definitions of  $g_r$  and  $g_+$  and convexity of  $r$ , it follows that

$$769 \quad (\text{A.2}) \quad r(x_+) \geq r(\bar{x}) + g_r^T(x_+ - \bar{x}) \quad \text{and} \quad r(\bar{x}) \geq r(x_+) + g_+^T(\bar{x} - x_+).$$

770 Adding the two equations in (A.2) together yields  $(g_r - g_+)^T(x_+ - \bar{x}) \leq 0$ . Combining  
771 this with the definition of  $g_{f+r}$ , (A.1), and the definition of  $x_+$  that

$$772 \quad (\text{A.3}) \quad \begin{aligned} s(\bar{x}, \bar{\alpha})^T g_{f+r} &= (x_+ - \bar{x})^T (\nabla f(\bar{x}) + g_r) \\ &= \frac{1}{\bar{\alpha}} (x_+ - \bar{x})^T (\bar{x} - x_+ - \bar{\alpha} g_+ + \bar{\alpha} g_r) \\ &= -\frac{1}{\bar{\alpha}} \|x_+ - \bar{x}\|_2^2 + (x_+ - \bar{x})^T (g_r - g_+) \leq -\frac{1}{\bar{\alpha}} \|s(\bar{x}, \bar{\alpha})\|_2^2. \end{aligned}$$

Since  $g_{f+r} \in \partial(f+r)(\bar{x})$  was arbitrary, the result [25, Theorem 2.87] and (A.3) yield

$$D_{f+r}(\bar{x}; s(\bar{x}, \bar{\alpha})) = \sup_{g \in \partial(f+r)(\bar{x})} s(\bar{x}, \bar{\alpha})^T g \leq -\frac{1}{\bar{\alpha}} \|s(\bar{x}, \bar{\alpha})\|_2^2,$$

773 which is the desired result and completes the proof.

774 **Proof of Lemma 2.2.** The proof follows exactly as in the proof of Lemma 2.1 above,  
 775 but where all calculations are restricted to groups in the set  $\mathcal{I}$  (also see (2.3)).

776 **Proof of Lemma 2.3.** The result, for the case  $\mathcal{I} = \{1, 2, \dots, n\}$ , can be found in [2,  
 777 Lemma 10.4]. For the general case, i.e., when  $\mathcal{I}$  is equal to the union of a subset of  
 778  $\{\mathcal{G}_i\}_{i=1}^{n_g}$ , the result follows by using the same proof as for [2, Lemma 11.9].

779 **Proof of Lemma 2.4.** Denote  $g_i := \nabla_{\mathcal{G}_i} f(\bar{x})$ ,  $x_i = [\bar{x}]_{\mathcal{G}_i}$ , and  $s_i = [s(\bar{x}, \bar{\alpha})]_{\mathcal{G}_i}$ . Since  
 780  $f + r$  is differentiable with respect to the variables in  $\mathcal{G}_i$  at  $\bar{x}$  since  $[\bar{x}]_{\mathcal{G}_i} \neq 0$ , we have

$$781 \quad \|\nabla_{\mathcal{G}_i} (f + r)(\bar{x})\|_2^2 = \|g_i + \lambda_i x_i / \|x_i\|_2\|_2^2 = \|g_i\|_2^2 + 2\lambda_i \frac{g_i^T x_i}{\|x_i\|_2} + \lambda_i^2,$$

782 which means that it is sufficient to prove that

$$\|g_i\|_2^2 + 2\lambda_i \frac{g_i^T x_i}{\|x_i\|_2} + \lambda_i^2 \geq \|s_i\|_2^2.$$

783 Since  $x_i + s_i \neq 0$  by assumption, we know that  $s_i$  (see (2.3)) satisfies

$$784 \quad s_i = \left(1 - \frac{\bar{\alpha}\lambda_i}{\|x_i - \bar{\alpha}g_i\|_2}\right) (x_i - \bar{\alpha}g_i) - x_i$$

$$785 \quad = x_i - \bar{\alpha}g_i - \frac{\bar{\alpha}\lambda_i(x_i - \bar{\alpha}g_i)}{\|x_i - \bar{\alpha}g_i\|_2} - x_i = -\bar{\alpha} \left(g_i + \frac{\bar{\alpha}\lambda_i(x_i - \bar{\alpha}g_i)}{\|x_i - \bar{\alpha}g_i\|_2}\right)$$

786 so that

$$\|s_i\|_2^2 = \bar{\alpha}^2 \left(\|g_i\|_2^2 + 2\bar{\alpha}\lambda_i \frac{g_i^T (x_i - \bar{\alpha}g_i)}{\|x_i - \bar{\alpha}g_i\|_2} + \bar{\alpha}^2 \lambda_i^2\right).$$

Thus, it is sufficient to prove that

$$\|g_i\|_2^2 + 2\lambda_i \frac{g_i^T x_i}{\|x_i\|_2} + \lambda_i^2 \geq \bar{\alpha}^2 \left(\|g_i\|_2^2 + 2\bar{\alpha}\lambda_i \frac{g_i^T (x_i - \bar{\alpha}g_i)}{\|x_i - \bar{\alpha}g_i\|_2} + \bar{\alpha}^2 \lambda_i^2\right).$$

787 We consider two cases, and note that  $x_i \neq 0$  by assumption and that  $x_i - \bar{\alpha}g_i \neq 0$  as  
 788 a consequence of (2.3) and the assumption that  $x_i + s_i \neq 0$ .

789 *Case 1:*  $\bar{\alpha} = 1$ . In this case, the desired inequality simplifies to

$$790 \quad (\text{A.4}) \quad \frac{g_i^T x_i}{\|x_i\|_2} \geq \frac{g_i^T (x_i - g_i)}{\|x_i - g_i\|_2}.$$

791 We now consider the following two subcases.

792 *Case 1a:*  $g_i^T x_i \geq 0$ . The desired inequality clearly holds if  $g_i^T (x_i - g_i) \leq 0$ . Thus,  
 793 for the remainder of this subcase, we assume that  $g_i^T (x_i - g_i) > 0$ , which equivalently  
 794 means that  $g_i^T x_i > \|g_i\|_2^2$ , which implies that  $-2x_i^T g_i + \|g_i\|_2^2 < 0$ . It follows from this  
 795 inequality and the fact that  $(g_i^T x_i)^2 \leq \|g_i\|_2^2 \|x_i\|_2^2$  (by Cauchy-Schwarz) that

$$796 \quad (g_i^T x_i)^2 (-2x_i^T g_i + \|g_i\|_2^2) \geq (-2x_i^T g_i + \|g_i\|_2^2) \|g_i\|_2^2 \|x_i\|_2^2 = (\|g_i\|_2^4 - 2g_i^T x_i \|g_i\|_2^2) \|x_i\|_2^2.$$

798 We can now add the term  $(g_i^T x_i)^2 \|x_i\|_2^2$  to both sides to obtain

$$799 \quad (g_i^T x_i)^2 (\|x_i\|_2^2 - 2x_i^T g_i + \|g_i\|_2^2) \geq ((g_i^T x_i)^2 + \|g_i\|_2^4 - 2g_i^T x_i \|g_i\|_2^2) \|x_i\|_2^2,$$

which can be written equivalently as

$$(g_i^T x_i)^2 \|x_i - g_i\|_2^2 \geq (g_i^T x_i - \|g_i\|_2^2)^2 \|x_i\|_2^2 = (g_i^T (x_i - g_i))^2 \|x_i\|_2^2.$$

801 After taking the square root of both sides, we obtain (A.4).

802 *Case 1b:*  $g_i^T x_i < 0$ . Using  $g_i^T x_i < 0$  and  $(g_i^T x_i)^2 \leq \|g_i\|_2^2 \|x_i\|_2^2$  (by Cauchy-Schwarz),  
803 we have

$$804 (g_i^T x_i)^2 (-2x_i^T g_i + \|g_i\|_2^2) \leq (-2x_i^T g_i + \|g_i\|_2^2) \|g_i\|_2^2 \|x_i\|_2^2 = (\|g_i\|_2^4 - 2g_i^T x_i \|g_i\|_2^2) \|x_i\|_2^2.$$

806 We can now add the term  $(g_i^T x_i)^2 \|x_i\|_2^2$  to both sides to obtain

$$807 (g_i^T x_i)^2 (\|x_i\|_2^2 - 2x_i^T g_i + \|g_i\|_2^2) \leq ((g_i^T x_i)^2 + \|g_i\|_2^4 - 2g_i^T x_i \|g_i\|_2^2) \|x_i\|_2^2,$$

which can be written equivalently as

$$(g_i^T x_i)^2 \|x_i - g_i\|_2^2 \leq (g_i^T x_i - \|g_i\|_2^2)^2 \|x_i\|_2^2 = (g_i^T (x_i - g_i))^2 \|x_i\|_2^2.$$

After taking the square root of both sides and rearranging, we obtain

$$\frac{|g_i^T x_i|}{\|x_i\|_2} \leq \frac{|g_i^T (x_i - g_i)|}{\|x_i - g_i\|_2}.$$

809 Combining this result with  $0 > g_i^T x_i \geq g_i^T (x_i - g_i)$  gives (A.4), as claimed.

810 *Case 2:*  $\bar{\alpha} \in (0, 1)$ . The proof follows from Case 1 and [2, Theorem 10.9], which in  
811 our notation from (2.2) proves that  $\|s(\bar{x}, \bar{\alpha})\|_2 \leq \|s(\bar{x}, 1)\|_2$  when  $\bar{\alpha} \in (0, 1)$ .

812

## REFERENCES

- 813 [1] F. BACH, R. JENATTON, J. MAIRAL, AND G. OBOZINSKI, *Optimization with sparsity-inducing*  
814 *penalties*, Foundations and Trends® in Machine Learning, 4 (2012), pp. 1–106.
- 815 [2] A. BECK, *First-order methods in optimization*, vol. 25, SIAM, 2017.
- 816 [3] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse*  
817 *problems*, SIAM Journal on Imaging Sciences, 2 (2009), pp. 183–202.
- 818 [4] D. P. BERTSEKAS, *Convex optimization theory*, Athena Scientific, Belmont, Ma., 2009.
- 819 [5] T. CHEN, F. E. CURTIS, AND D. P. ROBINSON, *A reduced-space algorithm for minimizing  $\ell_1$ -*  
820 *regularized convex functions*, SIAM Journal on Optimization, 27 (2017), pp. 1583–1610.
- 821 [6] T. CHEN, F. E. CURTIS, AND D. P. ROBINSON, *FaRSA for  $\ell_1$ -regularized convex optimization:*  
822 *local convergence and numerical experience*, 33 (2018), pp. 396–415.
- 823 [7] P. COMBETTES AND J.-C. PESQUET, *Proximal splitting methods in signal processing*, in Fixed-  
824 *point Algorithms for Inverse Problems in Science and Eng.*, Springer, 2011, pp. 185–212.
- 825 [8] A. R. CONN, N. I. M. GOULD, AND PH. L. TOINT, *Trust-Region Methods*, Society for Industrial  
826 *and Applied Mathematics (SIAM)*, Philadelphia, PA, 2000.
- 827 [9] F. E. CURTIS AND D. P. ROBINSON, *Exploiting negative curvature in deterministic and stochas-*  
828 *tic optimization*, Mathematical Programming, 176 (2019), pp. 69–94.
- 829 [10] I. DAUBECHIES, M. DEFRISE, AND C. MOL, *An iterative thresholding algorithm for linear inverse*  
830 *problems with a sparsity constraint*, Comm. Pure Appl. Math., 58 (2004), pp. 1413–1457.
- 831 [11] R. S. DEMBO, S. C. EISENSTAT, AND T. STEIHAUG, *Inexact newton methods*, SIAM Journal on  
832 *Numerical analysis*, 19 (1982), pp. 400–408.
- 833 [12] D. DONOHO, *Denoising by soft-thresholding*, Trans. Inform. Theory, 41 (1995), pp. 613–627.
- 834 [13] R.-E. FAN, K.-W. CHANG, C.-J. HSIEH, X.-R. WANG, AND C.-J. LIN, *Liblinear: A library for*  
835 *large linear classification*, J. Mach. Learn. Res., 9 (2008), pp. 1871–1874.
- 836 [14] R.-E. FAN, K.-W. CHANG, C.-J. HSIEH, X.-R. WANG, AND C.-J. LIN, *LIBLINEAR: A library*  
837 *for large linear classification*, JMLR, 9 (2008), pp. 1871–1874.
- 838 [15] M. A. T. FIGUEIREDO, R. D. NOWAK, AND S. J. WRIGHT, *Gradient projection for sparse*  
839 *reconstruction: Application to compressed sensing and other inverse problems*, IEEE J.  
840 *Selected Topics Signal Process.*, 1 (2007), pp. 586–597.

- 841 [16] J. FRIEDMAN, T. HASTIE, AND R. TIBSHIRANI, *Regularization paths for generalized linear models*  
842 *via coordinate descent*, Journal of Statistical Software, 33 (2010), p. 1.
- 843 [17] G. N. GRAPIGLIA AND Y. NESTEROV, *Accelerated regularized newton methods for minimizing*  
844 *composite convex functions*, SIAM Journal on Optimization, 29 (2019), pp. 77–99.
- 845 [18] N. KESKAR, J. NOCEDAL, F. OZTOPRAK, AND A. WÄCHTER, *A second-order method for convex*  
846  *$\ell_1$ -regularized optimization with active-set prediction*, Optimization Methods and Software,  
847 31 (2016), pp. 605–621.
- 848 [19] J. D. LEE, Y. SUN, AND M. A. SAUNDERS, *Proximal newton-type methods for minimizing*  
849 *composite functions*, SIAM Journal on Optimization, 24 (2014), pp. 1420–1443.
- 850 [20] Q. LIN, Z. LU, AND L. XIAO, *An accelerated randomized proximal coordinate gradient method*  
851 *and its application to regularized empirical risk minimization*, SIAM Journal on Optimiza-  
852 tion, 25 (2015), pp. 2244–2273.
- 853 [21] J. LIU, S. JI, AND J. YE, *SLEP: Sparse Learning with Efficient Projections*, Arizona State  
854 University, 2009, <http://www.public.asu.edu/~jye02/Software/SLEP>.
- 855 [22] J. LIU AND S. J. WRIGHT, *Asynchronous stochastic coordinate descent: Parallelism and con-*  
856 *vergence properties*, SIAM Journal on Optimization, 25 (2015), pp. 351–376.
- 857 [23] S. MA, X. SONG, AND J. HUANG, *Supervised group Lasso with applications to microarray data*  
858 *analysis*, BMC bioinformatics, 8 (2007), p. 60.
- 859 [24] J. L. MORALES, *A numerical study of limited memory BFGS methods*, Applied Mathematics  
860 Letters, 15 (2002), pp. 481–487.
- 861 [25] B. S. MORDUKHOVICH AND N. M. NAM, *An easy path to convex analysis and applications*,  
862 vol. 6, Morgan & Claypool Publishers, 2013.
- 863 [26] Y. NESTEROV, *A method of solving a convex programming problem with convergence rate*  
864  *$\mathcal{O}(1/k^2)$* , Soviet Mathematics Doklady, 27(2) (1983), pp. 372–376.
- 865 [27] Y. NESTEROV, *Gradient methods for minimizing composite functions*, Mathematical Program-  
866 ming, 140 (2013), pp. 125–161.
- 867 [28] J. NUTINI, M. SCHMIDT, AND W. HARE, *Active-set complexity of proximal gradient: How long*  
868 *does it take to find the sparsity pattern?*, Optimization Letters, 13 (2019), pp. 645–655.
- 869 [29] P. RICHTÁRIK AND M. TAKÁČ, *Parallel coordinate descent methods for big data optimization*,  
870 Mathematical Programming, 156 (2016), pp. 433–484.
- 871 [30] R. TAPPENDEN, P. RICHTÁRIK, AND J. GONDZIO, *Inexact coordinate descent: complexity and*  
872 *preconditioning*, J. of Optimization Theory and Applications, 170 (2016), pp. 144–176.
- 873 [31] S. J. WRIGHT, R. D. NOWAK, AND M. A. FIGUEIREDO, *Sparse reconstruction by separable*  
874 *approximation*, IEEE Transactions on Signal Processing, 57 (2009), pp. 2479–2493.
- 875 [32] Y. YANG AND H. ZOU, *A fast unified algorithm for solving group-lasso penalize learning prob-*  
876 *lems*, Statistics and Computing, 25 (2015), pp. 1129–1141.
- 877 [33] G.-X. YUAN, C.-H. HO, AND C.-J. LIN, *An improved GLMNET for  $\ell_1$ -regularized logistic*  
878 *regression*, Journal of Machine Learning Research, 13 (2012), pp. 1999–2030.
- 879 [34] M. YUAN AND Y. LIN, *Model selection and estimation in regression with grouped variables*, J.  
880 of the Royal Statistical Society: Series B (Statistical Methodology), 68 (2006), pp. 49–67.
- 881 [35] Y. ZENG AND P. BREHENY, *Overlapping group logistic regression with applications to genetic*  
882 *pathway selection*, Cancer Informatics, 15 (2016), pp. CIN–S40043.