

ISE



Industrial and
Systems Engineering

A Reduced-Space Algorithm for Minimizing ℓ_1 -Regularized Convex Functions

TIANYI CHEN

Dept. of Applied Mathematics and Statistics, Johns Hopkins University

FRANK E. CURTIS

Dept. Industrial and Systems Engineering, Lehigh University

DANIEL P. ROBINSON

Dept. of Applied Mathematics and Statistics, Johns Hopkins University

COR@L Technical Report 16T-002



LEHIGH
UNIVERSITY.

COR@L
COMPUTATIONAL OPTIMIZATION
RESEARCH AT LEHIGH 

A Reduced-Space Algorithm for Minimizing ℓ_1 -Regularized Convex Functions

TIANYI CHEN^{*1}, FRANK E. CURTIS^{†2}, AND DANIEL P. ROBINSON^{‡3}

¹Dept. of Applied Mathematics and Statistics, Johns Hopkins University

²Dept. Industrial and Systems Engineering, Lehigh University

³Dept. of Applied Mathematics and Statistics, Johns Hopkins University

February 19, 2016

Abstract

We present a new method for minimizing the sum of a differentiable convex function and an ℓ_1 -norm regularizer. The main features of the new method include: (i) an evolving set of indices corresponding to variables that are predicted to be nonzero at a solution (i.e., the support); (ii) a reduced-space subproblem defined in terms of the predicted support; (iii) conditions that determine how accurately each subproblem must be solved, which allow for Newton, Newton-CG, and coordinate-descent techniques to be employed; (iv) a computationally practical condition that determines when the predicted support should be updated; and (v) a reduced proximal gradient step that ensures sufficient decrease in the objective function when it is decided that variables should be added to the predicted support. We prove a convergence guarantee for our method and demonstrate its efficiency on a large set of model prediction problems.

1 Introduction

In this paper, we propose, analyze, and provide the results of numerical experiments for a new method for solving ℓ_1 -norm regularized convex optimization problems of the form

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \ F(x), \text{ where } F(x) := f(x) + \lambda \|x\|_1, \quad (1)$$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice continuously differentiable convex function, and $\lambda > 0$ is a weighting parameter. A necessary and sufficient optimality condition for (1) is

$$0 \in \partial F(x) = \nabla f(x) + \lambda \partial \|x\|_1 \quad (2)$$

with ∂F and $\partial \|\cdot\|_1$ denoting the subdifferentials of F and $\|\cdot\|_1$, respectively. Our method for solving (1) generates a sequence of iterates such that any limit point of the sequence satisfies (2). It is applicable when only first-order derivative information is computed, but is most effective when one can at least approximate second-order derivative matrices, e.g., using limited-memory quasi-Newton techniques.

Problems of the form (1) routinely arise in statistics, signal processing, and machine learning applications, and are usually associated with data fitting or maximum likelihood estimation. A popular setting is binary

^{*}Email: tchen59@jhu.edu

[†]Email: frank.e.curtis@gmail.com

[‡]Email: daniel.p.robinson@gmail.com

classification using logistic regression (where f is a logistic cost function), although instances of such problems also arise when performing multi-class logistic regression and profit regression. Instances of (1) also surface when using LASSO or elastic-net formulations to perform data analysis and discovery, such as in unsupervised subspace clustering on data drawn from a union of subspaces.

1.1 Literature review and our key contributions

Popular first-order optimization methods for solving (1) include ISTA, FISTA, and SpARSA [2, 13]. Second-order methods have also been proposed, which can roughly be split into the classes of proximal-Newton methods [4, 8, 10, 11, 14] and orthant-based methods [1, 3, 9]. Proximal-Newton methods solve problem (1) by minimizing a sequence of subproblems formed as the sum of a quadratic approximation to f and the nonsmooth ℓ_1 -norm regularizer. For example, the state-of-the-art software LIBLINEAR, which implements newGLMNET [14], uses a coordinate descent algorithm to approximately minimize each piecewise quadratic subproblem. Orthant-based methods, on the other hand, minimize smooth quadratic approximations to (1) over a sequence of orthants in \mathbb{R}^n until a solution is found. Of particular interest is the recently proposed orthant-based method OBA [9] in which every iteration consists of a corrective cycle of orthant predictions and subspace minimization steps. OBA was shown to be slower than LIBLINEAR when the Hessian matrices were diagonally dominant, but faster otherwise, at least on the collection of test problems considered in [9].

Since LIBLINEAR and OBA are the most relevant to the algorithm described in this paper, let us discuss their respective advantages and disadvantages in more detail. The key advantage of LIBLINEAR is its use of a coordinate descent (CD) algorithm to approximately minimize the piecewise quadratic subproblem. The use of CD means that one should expect excellent performance on problems whose Hessian matrices are strongly diagonally dominant. This expectation was confirmed, as mentioned above, by the OBA paper [9]. For some problems encountered in model prediction, e.g., when using logistic regression to perform classification, the Hessians are often strongly diagonally dominant, at least after certain data scaling techniques are used. However, not all prediction problems have such nice diagonal dominance properties, and in some instances the user would prefer to avoid discovering a proper scaling for their data. In these latter cases, the OBA method is typically superior.

Another potential advantage of the OBA method is its use of an active-set strategy that uses subproblems that are smaller in dimension than the ambient space. For many ℓ_1 -norm regularized prediction problems, the number of nonzero components in a solution is a small percentage of the ambient dimension, and thus OBA spends most of its time solving small dimensional problems. This is an advantage, at least when the zero and nonzero structure of the solution is quickly identified.

We have the perspective that both LIBLINEAR and OBA are valuable state-of-the-art algorithms that complement each other. Our *fast reduced space algorithm* (FaRSA) is designed to capitalize on the advantages of both while avoiding their disadvantages. The following bulleted points summarize our key contributions.

- (i) We present a new active-set line search method that utilizes reduced-space subproblems, approximate solutions of which can be computed efficiently.
- (ii) Unlike the active-set OBA method, our method does not require the computation of an ISTA step during each iteration to ensure convergence. We achieve convergence by combining a new projected backtracking line search procedure, an approximate subspace minimization scheme, and a mechanism for determining when the support of the solution estimate should be updated.
- (iii) Our framework is flexible. In particular, we introduce a new set of conditions that signal how accurately each subproblem should be solved and allow for various subproblem solvers to be used. In so doing, our method easily accommodates a Newton-CG subproblem solver as in OBA and a CD solver as in LIBLINEAR. Interestingly, this allows for multiple subproblem solvers to be used in parallel, thus allowing for numerical performance that can be as good as either LIBLINEAR and OBA regardless of whether the problem Hessians are strongly diagonally dominant.
- (iv) As demonstrated in the numerical experiments described in this paper, the practical performance of our method is state-of-the-art.

We end this review by remarking that our proposed algorithm has similarities with the iterative method that one would obtain using the following procedure: (i) at a given iterate x_k , construct a quadratic model of f and recast the minimization of this model plus the regularization term $\lambda\|x\|_1$ into a bound-constrained quadratic optimization problem (similarly to the procedure in SpaRSA), (ii) approximately solve this subproblem using the techniques developed in [5, 6, 7] (see also [12]), and (iii) translate the resulting solution back into the space of x variables to produce a trial step from x_k , call it d_k . Indeed, our initial developments of this work was based on these ideas. However, the algorithm proposed in this paper involves some deviations and enhancements from this starting point.

1.2 Notation

Let $\mathcal{I} \subseteq \{1, 2, \dots, n\}$ denote an index set of variables. For any $v \in \mathbb{R}^n$, we let $[v]_{\mathcal{I}}$ denote the subvector of v consisting of elements of v with indices in \mathcal{I} . Similarly, for any symmetric matrix $M \in \mathbb{R}^{n \times n}$, we let $[M]_{\mathcal{I}, \mathcal{I}}$ denote the submatrix of M consisting of the rows and columns of M that correspond to the index set \mathcal{I} . For any vector v , we let $\text{sgn}(v)$ denote the vector of the same length as v whose i th component is 0 when $[v]_i = 0$, is 1 when $[v]_i > 0$, and is -1 when $[v]_i < 0$. For any vector v , we let $\|v\|_1$ and $\|v\|$ denote its ℓ_1 -norm and ℓ_2 -norm, respectively.

2 Algorithm FaRSA

Crucial to our algorithm is the manner in which we handle the zero and nonzero components of a solution estimate. In order to describe the details of our approach, we first define the index sets

$$\mathcal{I}^0(x) := \{i : [x]_i = 0\}, \quad \mathcal{I}^+(x) := \{i : [x]_i > 0\}, \quad \text{and} \quad \mathcal{I}^-(x) := \{i : [x]_i < 0\}.$$

We call $\mathcal{I}^0(x)$ the set of *zero variables*, $\mathcal{I}^+(x)$ the set of *positive variables*, $\mathcal{I}^-(x)$ the set of *negative variables*, and the union of $\mathcal{I}^-(x)$ and $\mathcal{I}^+(x)$ the set of *nonzero variables* at x . We use these sets to define measures of optimality corresponding to the zero and nonzero variables at x . Respectively, these measures are as follows:

$$\begin{aligned} [\beta(x)]_i &:= \begin{cases} [\nabla f(x)]_i + \lambda & \text{if } i \in \mathcal{I}^0(x) \text{ and } [\nabla f(x)]_i + \lambda < 0, \\ [\nabla f(x)]_i - \lambda & \text{if } i \in \mathcal{I}^0(x) \text{ and } [\nabla f(x)]_i - \lambda > 0, \\ 0 & \text{otherwise;} \end{cases} \\ [\phi(x)]_i &:= \begin{cases} 0 & \text{if } i \in \mathcal{I}^0(x), \\ \min\{[\nabla f(x)]_i + \lambda, \max\{[x]_i, [\nabla f(x)]_i - \lambda\}\} & \text{if } i \in \mathcal{I}^+(x) \text{ and } [\nabla f(x)]_i + \lambda > 0, \\ \max\{[\nabla f(x)]_i - \lambda, \min\{[x]_i, [\nabla f(x)]_i + \lambda\}\} & \text{if } i \in \mathcal{I}^-(x) \text{ and } [\nabla f(x)]_i - \lambda < 0, \\ [\nabla f(x)]_i + \lambda \cdot \text{sgn}(x)_i & \text{otherwise.} \end{cases} \end{aligned}$$

The following result shows that the functions β and ϕ together correspond to a valid optimality measure for problem (1).

Lemma 2.1. *Let \mathcal{S} be an infinite set of positive integers such that $\{x_k\}_{k \in \mathcal{S}} \rightarrow x_*$. Then, x_* is a solution to (1) if and only if $\{\beta(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$ and $\{\phi(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$. Consequently, x_* is a solution to (1) if and only if $\|\beta(x_*)\| = \|\phi(x_*)\| = 0$.*

Proof. Suppose $\{\beta(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$ and $\{\phi(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$. Then, first, consider any i such that $[x_*]_i > 0$, which means that $[x_k]_i > 0$ for all sufficiently large $k \in \mathcal{S}$. We now consider two subcases. If $[\nabla f(x_k)]_i + \lambda \leq 0$ for infinitely many $k \in \mathcal{S}$, then it follows from the definition of $\phi(x_k)$ and $\{\phi(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$ that $[\nabla f(x_*)]_i + \lambda = 0$. On the other hand, if $[\nabla f(x_k)]_i + \lambda > 0$ for infinitely many $k \in \mathcal{S}$, then it follows from the definition of $\phi(x_k)$, $\{\phi(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$, and $[x_*]_i > 0$ that $[\nabla f(x_*)]_i + \lambda = 0$. By combining both cases, we have established

that $[\nabla f(x_*)]_i + \lambda = 0$, so that the i th component satisfies the optimality conditions (2). A similar argument may be used for the case when one considers i such that $[x_*]_i < 0$ to show that $[\nabla f(x_*)]_i - \lambda = 0$.

It remains to consider i such that $[x_*]_i = 0$. We have four subcases to consider. First, if infinitely many $k \in \mathcal{S}$ satisfy $[x_k]_i = 0$ and $[\nabla f(x_k)]_i + \lambda < 0$, then it follows from the definition of $\beta(x_k)$ and $\{\beta(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$ that $[\nabla f(x_*)]_i + \lambda = 0$; a similar argument shows that if infinitely many $k \in \mathcal{S}$ satisfy $[x_k]_i = 0$ and $[\nabla f(x_k)]_i - \lambda > 0$, then $[\nabla f(x_*)]_i - \lambda = 0$. Second, if infinitely many $k \in \mathcal{S}$ satisfy $[x_k]_i = 0$ and $|\nabla f(x_k)_i| < \lambda$, then, trivially, $|\nabla f(x_*)_i| \leq \lambda$. Third, if infinitely many $k \in \mathcal{S}$ satisfy $[x_k]_i > 0$ and $[\nabla f(x_k)]_i + \lambda \leq 0$, then it follows from the definition of $\phi(x_k)$ and $\{\phi(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$ that $[\nabla f(x_*)]_i + \lambda = 0$; a similar argument shows that if infinitely many $k \in \mathcal{S}$ satisfy $[x_k]_i < 0$ and $[\nabla f(x_k)]_i - \lambda \geq 0$, then $[\nabla f(x_*)]_i - \lambda = 0$. Fourth, if infinitely many $k \in \mathcal{S}$ satisfy $[x_k]_i > 0$ and $[\nabla f(x_k)]_i + \lambda > 0$, then it follows from the definition of $\phi(x_k)$ and $\{\phi(x_k)\}_{k \in \mathcal{S}} \rightarrow 0$ that $|\nabla f(x_*)_i| \leq \lambda$; a similar argument shows that if infinitely many $k \in \mathcal{S}$ satisfy $[x_k]_i < 0$ and $[\nabla f(x_k)]_i - \lambda < 0$, then $|\nabla f(x_*)_i| \leq \lambda$. By combining these subcases, we conclude that $|\nabla f(x_*)_i| \leq \lambda$, so the i th component satisfies the optimality condition (2).

To prove the reverse implication, now suppose that x_* is a solution to problem (1). If $[x_*]_i > 0$, then $[\beta(x_k)]_i = 0$ for all sufficiently large $k \in \mathcal{S}$ and $\{\phi(x_k)_i\}_{k \in \mathcal{S}} \rightarrow 0$ since $[\nabla f(x_*)]_i + \lambda = 0$. If $[x_*]_i < 0$, then $[\beta(x_k)]_i = 0$ for all sufficiently large $k \in \mathcal{S}$ and $\{\phi(x_k)_i\}_{k \in \mathcal{S}} \rightarrow 0$ since $[\nabla f(x_*)]_i - \lambda = 0$. Finally, if $[x_*]_i = 0$, then $|\nabla f(x_*)_i| \leq \lambda$, which with the definitions of $\beta(x_k)$ and $\phi(x_k)$ implies that $\{\beta(x_k)_i\}_{k \in \mathcal{S}} \rightarrow 0$ and $\{\phi(x_k)_i\}_{k \in \mathcal{S}} \rightarrow 0$. \square

We now state our proposed method, FaRSA, as Algorithm 1. When considering a reduced-space subproblem defined by a chosen index set \mathcal{I}_k (see lines 7 and 14), the algorithm makes use of a quadratic model of the objective of the form (see line 10)

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

FaRSA also makes use of two line search subroutines, stated as Algorithms 2 and 3, the former of which employs the following projection operator dependent on x_k :

$$[\text{Proj}(y; x_k)]_i := \begin{cases} \max\{0, [y]_i\} & \text{if } \mathcal{I}^+(x_k), \\ \min\{0, [y]_i\} & \text{if } \mathcal{I}^-(x_k), \\ 0 & \text{if } \mathcal{I}^0(x_k). \end{cases}$$

FaRSA computes a sequence of iterates $\{x_k\}$. During each iteration, the sets $\mathcal{I}^0(x_k)$, $\mathcal{I}^+(x_k)$, and $\mathcal{I}^-(x_k)$ are identified, which are used to define $\beta(x_k)$ and $\phi(x_k)$. We can see in line 4 of Algorithm 1 that when both $\|\beta(x_k)\|$ and $\|\phi(x_k)\|$ are less than a prescribed tolerance $\epsilon > 0$, it returns x_k as an approximate solution to (1); this is justified by Lemma 2.1. Otherwise, it proceeds in one of two ways depending on the relative sizes of $\|\beta(x_k)\|$ and $\|\phi(x_k)\|$. We describe these two cases next.

- (i) The relationship $\|\beta(x_k)\| \leq \gamma \|\phi(x_k)\|$ indicates that significant progress toward optimality can still be achieved by reducing F over the current set of nonzero variables at x_k ; lines 7–12 are designed for this purpose. In line 7, a subset \mathcal{I}_k of variables are chosen such that the norm of $\phi(x_k)$ over that subset of variables is at least proportional to the norm of $\phi(x_k)$ over the full set of variables. This allows control over the size of the subproblem, which may be as small as one-dimensional. Note that for $i \in \mathcal{I}_k$, it must hold that $[\phi(x_k)]_i \neq 0$, which in turn means that $i \notin \mathcal{I}^0(x_k)$, i.e., the i th variable is nonzero. This means that the reduced space subproblem to minimize $m_k(d)$ over $d \in \mathbb{R}^{|\mathcal{I}_k|}$ is aimed at minimizing F over the variables in \mathcal{I}_k . Our analysis does not require an exact minimizer of m_k . Rather, we allow for the computation of any direction \bar{d}_k that satisfies the conditions in line 10, namely $g_k^T \bar{d}_k \leq g_k^T d_k^R$ and $m_k(\bar{d}_k) \leq m_k(0)$, where the reference direction d_k^R is computed in line 9 by minimizing m_k along the steepest descent direction. The first condition imposes how much descent is required by the search direction \bar{d}_k , while the second condition ensures that the model is reduced at least as much as a zero step. It will be shown (see Lemma 3.7) that the second condition ensures that \bar{d}_k is bounded by a multiple of $\|g_k\|$. Such conditions are satisfied by a Newton step, by any

Algorithm 1 FaRSA for solving problem (1).

```

1: Input:  $x_0$ 
2: Constants:  $\{\eta_\phi, \eta_\beta, \xi\} \subset (0, 1]$ ,  $\eta \in (0, 1/2]$ , and  $\{\gamma, \epsilon\} \subset (0, \infty)$ 
3: for  $k = 0, 1, 2, \dots$  do
4:   if  $\max\{\|\beta(x_k)\|, \|\phi(x_k)\|\} \leq \epsilon$  then
5:     Return the (approximate) solution  $x_k$  of problem (1).
6:   if  $\|\beta(x_k)\| \leq \gamma\|\phi(x_k)\|$  then [ $k \in \mathcal{S}_\phi$ ]
7:     Choose any  $\mathcal{I}_k \subseteq \{i : [\phi(x_k)]_i \neq 0\}$  such that  $\|[\phi(x_k)]_{\mathcal{I}_k}\| \geq \eta_\phi\|\phi(x_k)\|$ .
8:     Set  $H_k \leftarrow [\nabla^2 F(x_k)]_{\mathcal{I}_k \mathcal{I}_k}$  and  $g_k \leftarrow [\nabla F(x_k)]_{\mathcal{I}_k}$ .
9:     Compute the reference direction
           
$$d_k^R \leftarrow -\alpha_k g_k, \text{ where } \alpha_k \leftarrow \|g_k\|^2 / (g_k^T H_k g_k).$$

10:    Compute any direction  $\bar{d}_k$  that satisfies the inequalities
           
$$g_k^T \bar{d}_k \leq g_k^T d_k^R \text{ and } m_k(\bar{d}_k) \leq m_k(0).$$

11:    Set  $[d_k]_{\mathcal{I}_k} \leftarrow \bar{d}_k$  and  $[d_k]_i \leftarrow 0$  for  $i \notin \mathcal{I}_k$ .
12:    Use Algorithm 2 to compute  $x_{k+1} \leftarrow \text{LINESEARCH-}\phi(x_k, d_k, \mathcal{I}_k, \eta, \xi)$ .
13:  else [ $k \in \mathcal{S}_\beta$ ]
14:    Choose any  $\mathcal{I}_k \subseteq \{i : [\beta(x_k)]_i \neq 0\}$  such that  $\|[\beta(x_k)]_{\mathcal{I}_k}\| \geq \eta_\beta\|\beta(x_k)\|$ .
15:    Set  $[d_k]_{\mathcal{I}_k} \leftarrow -[\beta(x_k)]_{\mathcal{I}_k}$  and  $[d_k]_i \leftarrow 0$  for  $i \notin \mathcal{I}_k$ .
16:    Use Algorithm 3 to compute  $x_{k+1} \leftarrow \text{LINESEARCH-}\beta(x_k, d_k, \eta, \xi)$ .

```

Newton-CG iterate, and asymptotically by CD iterates. Once \bar{d}_k is obtained, the search direction d_k in the full space is obtained by filling its elements that correspond to the index set \mathcal{I}_k with the elements from \bar{d}_k , and setting the complementary set of variables to zero (see line 11). With the search direction d_k computed, we call Algorithm 2 in line 12, which performs a (non-standard) backtracking projected line search. This line search procedure makes use of the projection operator $\text{Proj}(\cdot; x_k)$. This operator projects vectors onto the orthant inhabited by x_k , a feature shared by OBA. The while-loop that starts in line 3 of Algorithm 2 checks whether the trial point y_j decreases the objective function F relative to its value at x_k when $\text{sgn}(y_j) \neq \text{sgn}(x_k)$. If the line search terminates in this while-loop, then this implies that at least one component of x_k that was nonzero has become zero for $x_{k+1} = y_j$. Since the dimension of the reduced space will therefore be reduced during the next iteration (provided line 6 of Algorithm 1 tests true), the procedure only requires $F(x_{k+1}) \leq F(x_k)$ instead of a more traditional sufficient decrease condition, e.g., one based on the Armijo condition. If line 7 of Algorithm 2 is reached, then the current trial iterate y_j satisfies $\text{sgn}(y_j) = \text{sgn}(x_k)$, i.e., the trial iterate has entered the same orthant as that inhabited by x_k . Once this has occurred, the method could then perform a standard backtracking Armijo line search as stipulated in the loop starting at line 12. For the purpose of guaranteeing convergence, however, the method first checks whether the largest step along d_k that stays in the same orthant as x_k (see lines 8 and 9) satisfies the Armijo sufficient decrease condition (see line 10). (This aspect makes our procedure different from a standard backtracking scheme.) If Algorithm 2 terminates in line 5 or 11, then at least one nonzero variable at x_k will have become zero at x_{k+1} , which we indicate by saying $k \in \mathcal{S}_\phi^{\text{ADD}} \subseteq \mathcal{S}_\phi$. Otherwise, if Algorithm 2 terminates in line 14, then x_{k+1} and x_k are housed in the same orthant and sufficient decrease in F was achieved (i.e., the Armijo condition in line 13 was satisfied). Since sufficient decrease has been achieved in this case, we say that $k \in \mathcal{S}_\phi^{\text{SD}} \subseteq \mathcal{S}_\phi$.

- (ii) When $\|\beta(x_k)\| > \gamma\|\phi(x_k)\|$, progress toward optimality is best achieved by freeing at least one variable that is currently set to zero; lines 14–16 are designed for this purpose. Since $\|\beta(x_k)\|$ is relatively large, in line 14 of Algorithm 1 a subset \mathcal{I}_k of variables is chosen such that the norm of $\beta(x_k)$ over that subset

Algorithm 2 A line search procedure for computing x_{k+1} when $k \in \mathcal{S}_\phi$.

```

1: procedure  $x_{k+1} = \text{LINESEARCH}_\phi(x_k, d_k, \mathcal{I}_k, \eta, \xi)$ 
2:   Set  $j \leftarrow 0$  and  $y_0 \leftarrow \text{Proj}(x_k + d_k; x_k)$ .
3:   while  $\text{sgn}(y_j) \neq \text{sgn}(x_k)$  do
4:     if  $F(y_j) \leq F(x_k)$  then
5:       return  $x_{k+1} \leftarrow y_j$ .  $[k \in \mathcal{S}_\phi^{\text{ADD}}]$ 
6:     Set  $j \leftarrow j + 1$  and then  $y_j \leftarrow \text{Proj}(x_k + \xi^j d_k; x_k)$ .
7:   if  $j \neq 0$  then
8:     Set  $\alpha_B \leftarrow \text{argsup} \{ \alpha > 0 : \text{sgn}(x_k + \alpha d_k) = \text{sgn}(x_k) \}$ .
9:     Set  $y_j \leftarrow x_k + \alpha_B d_k$ .
10:    if  $F(y_j) \leq F(x_k) + \eta \alpha_B [\nabla F(x_k)]_{\mathcal{I}_k}^T [d_k]_{\mathcal{I}_k}$  then
11:      return  $x_{k+1} \leftarrow y_j$ .  $[k \in \mathcal{S}_\phi^{\text{ADD}}]$ 
12:  loop
13:    if  $F(y_j) \leq F(x_k) + \eta \xi^j [\nabla F(x_k)]_{\mathcal{I}_k}^T [d_k]_{\mathcal{I}_k}$  then
14:      return  $x_{k+1} \leftarrow y_j$ .  $[k \in \mathcal{S}_\phi^{\text{SD}}]$ 
15:    Set  $j \leftarrow j + 1$  and then  $y_j \leftarrow x_k + \xi^j d_k$ .

```

Algorithm 3 A line search procedure for computing x_{k+1} when $k \in \mathcal{S}_\beta$.

```

1: procedure  $x_{k+1} = \text{LINESEARCH}_\beta(x_k, d_k, \eta, \xi)$ 
2:   Set  $j \leftarrow 0$  and  $y_0 \leftarrow x_k + d_k$ .
3:   while  $F(y_j) > F(x_k) - \eta \xi^j \|d_k\|^2$  do
4:     Set  $j \leftarrow j + 1$  and then  $y_j \leftarrow x_k + \xi^j d_k$ .
5:   return  $x_{k+1} \leftarrow y_j$ .

```

of variables is at least proportional to the norm of $\beta(x_k)$ over the full set of variables. Similar to the previous case, this allows control over the size of the subproblem, which in the extreme case may be one-dimensional. If $i \in \mathcal{I}_k$, then $[\beta(x_k)]_i \neq 0$, which in turn means that $i \in \mathcal{I}^0(x_k)$, i.e., the i th variable has the value zero. The components of $\beta(x_k)$ that correspond to \mathcal{I}_k are then used to define the search direction d_k in line 15. With the search direction d_k computed, Algorithm 3 is called in line 16, which performs a standard backtracking Armijo line search to obtain x_{k+1} . If a unit step length is taken, i.e., if $x_{k+1} = x_k + d_k$, then x_{k+1} can be interpreted as the iterate that would be obtained by taking a *reduced ISTA step* in the space of variables indexed by \mathcal{I}_k . (For additional details, see Lemma A.1 in the appendix.)

3 Convergence Analysis

Our analysis uses the following assumption that is assumed to hold throughout this section.

Assumption 3.1. *The function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, twice continuously differentiable, and bounded below on the level set $\mathcal{L} := \{x \in \mathbb{R}^n : F(x) \leq F(x_0)\}$. The gradient function $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz continuous on \mathcal{L} with Lipschitz constant L . The Hessian function $\nabla^2 f : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$ is uniformly positive definite and bounded on \mathcal{L} , i.e., there exist positive constants θ_{\min} and θ_{\max} such that*

$$\theta_{\min} \|v\|^2 \leq v^T H(x) v \leq \theta_{\max} \|v\|^2 \quad \text{for all } \{x, v\} \subset \mathbb{R}^n.$$

Our analysis uses the index sets (already shown in Algorithms 1–2)

$$\begin{aligned}\mathcal{S}_\phi &:= \{k : \text{lines 7–12 in Algorithm 1 are performed during iteration } k\}, \\ \mathcal{S}_\phi^{\text{ADD}} &:= \{k \in \mathcal{S}_\phi : \text{sgn}(x_{k+1}) \neq \text{sgn}(x_k)\}, \\ \mathcal{S}_\phi^{\text{SD}} &:= \{k \in \mathcal{S}_\phi : \text{sgn}(x_{k+1}) = \text{sgn}(x_k)\}, \text{ and} \\ \mathcal{S}_\beta &:= \{k : \text{lines 14–16 in Algorithm 1 are performed during iteration } k\}.\end{aligned}$$

We start with a lemma that establishes an important identity for iterations in \mathcal{S}_β .

Lemma 3.2. *If $k \in \mathcal{S}_\beta$, then (\mathcal{I}_k, d_k) in lines 14 and 15 of Algorithm 1 yield*

$$[d_k]_{\mathcal{I}_k} = -[\nabla f(x_k) + \lambda \cdot \text{sgn}(x_k + \xi^j d_k)]_{\mathcal{I}_k} \text{ for any integer } j. \quad (3)$$

Consequently, the right-hand side of (3) has the same value for any integer j .

Proof. We prove that (3) holds for an arbitrary element of \mathcal{I}_k . To this end, let j be any integer and $i \in \mathcal{I}_k \subseteq \{\ell : [\beta(x_k)]_\ell \neq 0\}$, where \mathcal{I}_k is defined in line 14. It follows from the definition of \mathcal{I}_k , the definition of d_k in line 15, and $i \in \mathcal{I}_k$ that

$$[d_k]_i = \begin{cases} -([\nabla f(x_k)]_i + \lambda) & \text{if } [\nabla f(x_k)]_i + \lambda < 0, \\ -([\nabla f(x_k)]_i - \lambda) & \text{if } [\nabla f(x_k)]_i - \lambda > 0, \end{cases} \quad (4)$$

so that $[d_k]_i \neq 0$. Also, since $[x_k]_i = 0$ for $i \in \mathcal{I}_k$, we know that $[x_k + \xi^j d_k]_i \neq 0$. Thus, we need only consider the following two cases.

Case 1: Suppose $[x_k + \xi^j d_k]_i > 0$. In this case, the right-hand-side of (3) is equal to $-([\nabla f(x_k)]_i + \lambda)$. As for the left-hand-side, since $[x_k]_i = 0$ and $[x_k + \xi^j d_k]_i > 0$, we have $0 < [x_k + \xi^j d_k]_i = \xi^j [d_k]_i$, which combined with $\xi^j > 0$ means that $[d_k]_i > 0$. This fact and (4) gives $[d_k]_i = -([\nabla f(x_k)]_i + \lambda)$, so (3) holds.

Case 2: Suppose $[x_k + \xi^j d_k]_i < 0$. In this case, the right-hand-side of (3) is equal to $-([\nabla f(x_k)]_i - \lambda)$. As for the left-hand-side, since $[x_k]_i = 0$ and $[x_k + \xi^j d_k]_i < 0$ we have $0 > [x_k + \xi^j d_k]_i = \xi^j [d_k]_i$, which when combined with $\xi^j > 0$ means that $[d_k]_i < 0$. This fact and (4) gives $[d_k]_i = -([\nabla f(x_k)]_i - \lambda)$, so (3) holds. \square

We can now establish a bound for a decrease in the objective when $k \in \mathcal{S}_\beta$.

Lemma 3.3. *If $k \in \mathcal{S}_\beta$, then d_k in line 15 of Algorithm 1 yields*

$$F(x_k + \xi^j d_k) \leq F(x_k) - \frac{\xi^j}{2} \|d_k\|^2 \text{ for any integer } j \text{ with } 0 \leq \xi^j \leq \frac{1}{L}.$$

Proof. Let j be any integer with $0 \leq \xi^j \leq \frac{1}{L}$ and let $y_j := x_k + \xi^j d_k$. By Lipschitz continuity of the gradient function ∇f , we have

$$\begin{aligned}f(y_j) &\leq f(x_k) + \xi^j \nabla f(x_k)^T d_k + \frac{L}{2} \|\xi^j d_k\|^2 \\ &\leq f(x_k) + \xi^j \nabla f(x_k)^T d_k + \frac{\xi^j}{2} \|d_k\|^2.\end{aligned} \quad (5)$$

It then follows from (5), convexity of both f and $\lambda \|\cdot\|_1$, the fact that $\text{sgn}(y_j) \in \partial \|y_j\|_1$, the definition of d_k

(in particular that $[d_k]_i = 0$ for $i \notin \mathcal{I}_k$), and Lemma 3.2 that the following holds for all $z \in \mathbb{R}^n$:

$$\begin{aligned}
& F(y_j) \\
&= f(y_j) + \lambda \|y_j\|_1 \\
&\leq f(x_k) + \xi^j \nabla f(x_k)^T d_k + \frac{\xi^j}{2} \|d_k\|^2 + \lambda \|y_j\|_1 \\
&\leq f(z) + \nabla f(x_k)^T (x_k - z) + \xi^j \nabla f(x_k)^T d_k + \frac{\xi^j}{2} \|d_k\|^2 + \lambda \|z\|_1 + \lambda \cdot \text{sgn}(y_j)^T (y_j - z) \\
&\leq F(z) + [\nabla f(x_k) + \lambda \cdot \text{sgn}(y_j)]^T (x_k - z) + \xi^j [\nabla f(x_k) + \lambda \cdot \text{sgn}(y_j)]^T d_k + \frac{\xi^j}{2} \|d_k\|^2 \\
&= F(z) + [\nabla f(x_k) + \lambda \cdot \text{sgn}(y_j)]^T (x_k - z) - \xi^j \|d_k\|^2 + \frac{\xi^j}{2} \|d_k\|^2 \\
&= F(z) + [\nabla f(x_k) + \lambda \cdot \text{sgn}(y_j)]^T (x_k - z) - \frac{\xi^j}{2} \|d_k\|^2.
\end{aligned} \tag{6}$$

The desired result follows by considering $z = x_k$ in (6). \square

We now show that Algorithm 3 called in line 16 of Algorithm 1 is well defined, and that it returns x_{k+1} yielding sufficient decrease in the objective function.

Lemma 3.4. *If $k \in \mathcal{S}_\beta$, then x_{k+1} satisfies*

$$F(x_{k+1}) \leq F(x_k) - \kappa_\beta \max\{\|\beta(x_k)\|^2, \gamma^2 \|\phi(x_k)\|^2\}, \tag{7}$$

where $\kappa_\beta := \eta \eta_\beta \min\{1, \xi/L\}$.

Proof. Let j be any integer with $0 \leq \xi^j \leq \frac{1}{L}$ and let $y_j := x_k + \xi^j d_k$. It follows from Lemma 3.3 and the fact that $\eta \in (0, 1/2]$ in Algorithm 1 that

$$F(y_j) \leq F(x_k) - \frac{\xi^j}{2} \|d_k\|^2 \leq F(x_k) - \eta \xi^j \|d_k\|^2,$$

It follows from this inequality that Algorithm 3 will return the vector $x_{k+1} = x_k + \xi^{j^*} d_k$ with $\xi^{j^*} \geq \min\{1, \xi/L\}$ when called in line 16 of Algorithm 1. Using this bound, line 3 of Algorithm 3, and lines 15 and 14 of Algorithm 1, we have

$$\begin{aligned}
F(x_{k+1}) &\leq F(x_k) - \eta \xi^{j^*} \|d_k\|^2 \leq F(x_k) - \eta \min\{1, \xi/L\} \|d_k\|^2 \\
&= F(x_k) - \eta \min\{1, \xi/L\} \|\beta(x_k)_{\mathcal{I}_k}\|^2 \leq F(x_k) - \eta \eta_\beta \min\{1, \xi/L\} \|\beta(x_k)\|^2.
\end{aligned}$$

The inequality (7) follows from the definition of κ_β , the previous inequality, and the fact that the inequality in line 6 of Algorithm 1 must not hold since line 16 is assumed to be reached. \square

We now show that the index set \mathcal{S}_β must be finite.

Lemma 3.5. *The index set \mathcal{S}_β must be finite, i.e., $|\mathcal{S}_\beta| < \infty$.*

Proof. To derive a contradiction, suppose that $|\mathcal{S}_\beta| = \infty$, which also means that Algorithm 1 does not terminate finitely. Since Algorithm 1 does not terminate finitely, we know from line 4 of Algorithm 1 that $\max\{\|\beta(x_k)\|, \|\phi(x_k)\|\} > \epsilon$ for all $k \geq 0$. Combining this inequality with Lemma 3.4 and the fact that $F(x_{k+1}) \leq F(x_k)$ for all $k \notin \mathcal{S}_\beta$ (as a result of Algorithm 2 called in line 12 of Algorithm 1), we may

conclude for any nonnegative integer ℓ and $\kappa_\beta > 0$ defined in Lemma 3.4 that

$$\begin{aligned}
F(x_0) - F(x_{\ell+1}) &= \sum_{k=0}^{\ell} [F(x_k) - F(x_{k+1})] \\
&\geq \sum_{k \in \mathcal{S}_\beta, k \leq \ell} [F(x_k) - F(x_{k+1})] \\
&\geq \sum_{k \in \mathcal{S}_\beta, k \leq \ell} \kappa_\beta \max\{\|\beta(x_k)\|^2, \gamma^2 \|\phi(x_k)\|^2\} \\
&\geq \sum_{k \in \mathcal{S}_\beta, k \leq \ell} \kappa_\beta \min\{1, \gamma^2\} \epsilon^2.
\end{aligned}$$

Rearranging the previous inequality shows that

$$\begin{aligned}
\lim_{l \rightarrow \infty} F(x_{l+1}) &\leq \lim_{\ell \rightarrow \infty} \left[F(x_0) - \sum_{k \in \mathcal{S}_\beta, k \leq \ell} \kappa_\beta \min\{1, \gamma^2\} \epsilon^2 \right] \\
&= F(x_0) - \sum_{k \in \mathcal{S}_\beta} \kappa_\beta \min\{1, \gamma^2\} \epsilon^2 = -\infty,
\end{aligned}$$

which contradicts Assumption 3.1. Thus, we conclude that $|\mathcal{S}_\beta| < \infty$. \square

To prove that Algorithm 1 terminates finitely with an approximate solution to problem (1), all that remains is to prove that the set \mathcal{S}_ϕ is finite. To establish that $\mathcal{S}_\phi \equiv \mathcal{S}_\phi^{\text{ADD}} \cup \mathcal{S}_\phi^{\text{SD}}$ is finite, we proceed by showing individually that both $\mathcal{S}_\phi^{\text{ADD}}$ and $\mathcal{S}_\phi^{\text{SD}}$ are finite. We begin with the set $\mathcal{S}_\phi^{\text{ADD}}$.

Lemma 3.6. *The set $\mathcal{S}_\phi^{\text{ADD}}$ is finite, i.e., $|\mathcal{S}_\phi^{\text{ADD}}| < \infty$.*

Proof. To derive a contradiction, suppose that $|\mathcal{S}_\phi^{\text{ADD}}| = \infty$, which in particular means that Algorithm 1 does not terminate finitely. Since Lemma 3.5 shows that \mathcal{S}_β is finite, we may also conclude that there exists an iteration k_1 such that $k \in \mathcal{S}_\phi = \mathcal{S}_\phi^{\text{ADD}} \cup \mathcal{S}_\phi^{\text{SD}}$ for all $k \geq k_1$.

We proceed by making two observations. First, if the i th component of x_k becomes zero for some iteration $k \geq k_1$, it will remain zero for the remainder of the iterations. This can be seen by using lines 11 and 7 of Algorithm 1 and the definition of $\phi(x_k)$ to deduce that if $[d_k]_i \neq 0$, then $i \in \mathcal{I}_k \subseteq \{\ell : [\phi(x_k)]_\ell \neq 0\} \subseteq \mathcal{I}^+(x_k) \cup \mathcal{I}^-(x_k)$ for all $k \geq k_1$; equivalently, if $i \in \mathcal{I}^0(x_k)$, then $[d_k]_i = 0$. The second observation is that at least one nonzero component of x_k becomes zero at x_{k+1} for each $k \in \mathcal{S}_\phi^{\text{ADD}}$. This can be seen by construction of Algorithm 2 when it is called in line 12 of Algorithm 1. Together, these observations contradict $|\mathcal{S}_\phi^{\text{ADD}}| = \infty$, since at most n variables may become zero. Thus, we must conclude that $|\mathcal{S}_\phi^{\text{ADD}}| < \infty$. \square

To establish that $\mathcal{S}_\phi^{\text{SD}}$ is finite, we require the following two lemmas. The first lemma gives a bound on the size of \bar{d}_k that holds whenever $k \in \mathcal{S}_\phi$.

Lemma 3.7. *If $k \in \mathcal{S}_\phi$, then $\|\bar{d}_k\| \leq (2/\theta_{\min}) \|g_k\|$ where $\theta_{\min} > 0$ is defined in Assumption 3.1.*

Proof. Let $k \in \mathcal{S}_\phi$ so that \bar{d}_k is computed in line 10 of Algorithm 1, and let d_k^N be the Newton step satisfying $H_k d_k^N = -g_k$ with H_k and g_k defined in line 8 of Algorithm 1. It follows that

$$\|d_k^N\| \leq \|H_k^{-1}\| \|g_k\|. \quad (8)$$

Let us also define the quadratic function $\bar{m}_k(d) := m_k(d_k^N + d)$ and the associated level set $\mathcal{L}_k := \{d : \bar{m}_k(d) \leq 0\}$. We then see that

$$(\bar{d}_k - d_k^N) \in \mathcal{L}_k \quad (9)$$

since $\bar{m}_k(\bar{d}_k - d_k^N) = m_k(\bar{d}_k) \leq m_k(0) = 0$, where we have used the condition $m_k(\bar{d}_k) \leq m_k(0)$ that is required to hold in line 10 of Algorithm 1.

We are now interested in finding a point in \mathcal{L}_k with largest norm. To characterize such a point, we consider the optimization problem

$$\underset{d \in \mathbb{R}^n}{\text{maximize}} \quad \frac{1}{2} \|d\|^2 \quad \text{subject to } d \in \mathcal{L}_k. \quad (10)$$

It is not difficult to prove that a global maximizer of problem (10) is $d_* := \alpha_* v$ with $\alpha_*^2 := (-g_k^T d_k^N)/\theta$, where (v, θ) with $\|v\| = 1$ is an eigenpair corresponding to the left-most eigenvalue $\theta \geq \theta_{\min}$ of H_k . Thus, it follows that $\|d\|^2 \leq \|d_*\|^2$ for all $d \in \mathcal{L}_k$. Combining this with (9), the definition of d_* , and (8) shows that

$$\|\bar{d}_k - d_k^N\|^2 \leq \|d_*\|^2 = \alpha_*^2 \|v\|^2 = \frac{-g_k^T d_k^N}{\theta} \leq \frac{\|g_k\| \|d_k^N\|}{\theta} \leq \frac{\|H_k^{-1}\| \|g_k\|^2}{\theta} = \left(\frac{\|g_k\|}{\theta} \right)^2.$$

By combining the previous inequality with the triangle inequality and (8), we obtain

$$\|\bar{d}_k\| \leq \|\bar{d}_k - d_k^N\| + \|d_k^N\| \leq \frac{\|g_k\|}{\theta} + \frac{\|g_k\|}{\theta} = \frac{2\|g_k\|}{\theta} \leq \frac{2\|g_k\|}{\theta_{\min}},$$

which complete the proof. \square

The next result establishes a bound on the decrease in F when $k \in \mathcal{S}_\phi^{\text{SD}}$.

Lemma 3.8. *If $k \in \mathcal{S}_\phi^{\text{SD}}$, then x_{k+1} satisfies*

$$F(x_{k+1}) \leq F(x_k) - \kappa_\phi \max\{\gamma^{-2} \|\beta(x_k)\|^2, \|\phi(x_k)\|^2\}, \quad (11)$$

where $\kappa_\phi := \eta_\phi^2 \min \left\{ \frac{\eta}{\theta_{\max}}, \frac{\eta \xi (1-\eta) \theta_{\min}^2}{2\theta_{\max}^3} \right\} > 0$.

Proof. Let $k \in \mathcal{S}_\phi^{\text{SD}}$. We consider two cases. First, suppose that $j = 0$ when line 7 in Algorithm 2 is reached. In this case, it follows by construction of Algorithm 2 that $\text{sgn}(y_0) = \text{sgn}(x_k + d_k) = \text{sgn}(x_k)$, i.e., the full step d_k and the vector x_k are contained in the same orthant. Consequently, the loop that starts in line 12 is simply a backtracking Armijo line search. Thus, if

$$\xi^j \in \left(0, \frac{2(\eta-1)[\nabla F(x_k)]_{\mathcal{I}_k}^T [d_k]_{\mathcal{I}_k}}{\theta_{\max} \|[d_k]_{\mathcal{I}_k}\|^2} \right] \equiv \left(0, \frac{2(\eta-1)g_k^T \bar{d}_k}{\theta_{\max} \|\bar{d}_k\|^2} \right], \quad (12)$$

then, by well known properties of twice continuously differentiable functions with Lipschitz continuous gradients, we have that

$$\begin{aligned} F(x_k + \xi^j d_k) &\leq F(x_k) + \xi^j [\nabla F(x_k)]_{\mathcal{I}_k}^T [d_k]_{\mathcal{I}_k} + \frac{1}{2} \xi^{2j} \theta_{\max} \|[d_k]_{\mathcal{I}_k}\|^2 \\ &\leq F(x_k) + \xi^j [\nabla F(x_k)]_{\mathcal{I}_k}^T [d_k]_{\mathcal{I}_k} + \xi^j (\eta-1) [\nabla F(x_k)]_{\mathcal{I}_k}^T [d_k]_{\mathcal{I}_k} \\ &= F(x_k) + \eta \xi^j [\nabla F(x_k)]_{\mathcal{I}_k}^T [d_k]_{\mathcal{I}_k}, \end{aligned}$$

i.e., the inequality in line 13 will hold whenever (12) holds. On the other hand, suppose that $j > 0$ when line 7 in Algorithm 2 is reached. Then, since $k \in \mathcal{S}_\phi^{\text{SD}}$, we may conclude that

$$F(x_k + \alpha_B d_k) > F(x_k) + \eta \alpha_B [\nabla F(x_k)]_{\mathcal{I}_k}^T [d_k]_{\mathcal{I}_k} = F(x_k) + \eta \alpha_B g_k^T \bar{d}_k \quad (13)$$

in line 10, because otherwise we would have $k \in \mathcal{S}_\phi^{\text{ADD}} = \mathcal{S}_\phi \setminus \mathcal{S}_\phi^{\text{SD}}$. Since no points of non-differentiability of $\|\cdot\|_1$ exist on the line segment connecting x_k to $x_k + \alpha_B d_k$ (which follows by the definition of α_B in line 8 of Algorithm 2), we can conclude for the same reason that we acquired (12) that (13) implies

$$\alpha_B > \frac{2(1-\eta)|g_k^T \bar{d}_k|}{\theta_{\max} \|\bar{d}_k\|^2}.$$

Combining these two cases, we have that the line search procedure in Algorithm 2 will terminate with $x_{k+1} = x_k + \xi^j d_k$ where

$$\xi^j \geq \min \left\{ 1, \frac{2\xi(1-\eta)|g_k^T \bar{d}_k|}{\theta_{\max} \|\bar{d}_k\|^2} \right\} \quad \text{and} \quad F(x_{k+1}) \leq F(x_k) + \eta \xi^j g_k^T \bar{d}_k. \quad (14)$$

Let us now consider two cases. First, suppose that $\xi^j = 1$ is returned from the line search, i.e., $j = 0$. Then, it follows from (14), lines 10 and 9 of Algorithm 1, the Cauchy-Schwarz inequality, and Assumption 3.1 that

$$\begin{aligned} F(x_k) - F(x_{k+1}) &\geq -\eta \xi^j g_k^T \bar{d}_k = \eta |g_k^T \bar{d}_k| \geq \eta |g_k^T d_k^R| \\ &= \eta \alpha_k \|g_k\|^2 = \eta \frac{\|g_k\|^4}{g_k^T H_k g_k} \geq \frac{\eta}{\theta_{\max}} \|g_k\|^2. \end{aligned} \quad (15)$$

Now suppose that $\xi^j < 1$. Then, it follows from (14), the inequality $|g_k^T \bar{d}_k| \geq \|g_k\|^2 / \theta_{\max}$ established while deriving (15), and Lemma 3.7 that

$$\begin{aligned} F(x_k) - F(x_{k+1}) &\geq -\eta \xi^j g_k^T \bar{d}_k = \eta \xi^j |g_k^T \bar{d}_k| \geq \frac{2\eta \xi(1-\eta) |g_k^T \bar{d}_k|^2}{\theta_{\max} \|\bar{d}_k\|^2} \\ &\geq \frac{2\eta \xi(1-\eta) \theta_{\min}^2 \|g_k\|^4}{4\theta_{\max}^3 \|g_k\|^2} = \left(\frac{\eta \xi(1-\eta) \theta_{\min}^2}{2\theta_{\max}^3} \right) \|g_k\|^2. \end{aligned} \quad (16)$$

Combining (15) and (16) for the two cases establishes that

$$\begin{aligned} F(x_k) - F(x_{k+1}) &\geq \min \left\{ \frac{\eta}{\theta_{\max}}, \frac{\eta \xi(1-\eta) \theta_{\min}^2}{2\theta_{\max}^3} \right\} \|g_k\|^2 \\ &= \min \left\{ \frac{\eta}{\theta_{\max}}, \frac{\eta \xi(1-\eta) \theta_{\min}^2}{2\theta_{\max}^3} \right\} \|\nabla F(x_k)\|_{\mathcal{I}_k}^2 \\ &\geq \min \left\{ \frac{\eta}{\theta_{\max}}, \frac{\eta \xi(1-\eta) \theta_{\min}^2}{2\theta_{\max}^3} \right\} \|\phi(x_k)\|_{\mathcal{I}_k}^2 \\ &\geq \eta_{\phi}^2 \min \left\{ \frac{\eta}{\theta_{\max}}, \frac{\eta \xi(1-\eta) \theta_{\min}^2}{2\theta_{\max}^3} \right\} \|\phi(x_k)\|^2 \quad \text{for } k \in \mathcal{S}_{\phi}^{\text{SD}}, \end{aligned}$$

where we have also used the condition in lines 7 of Algorithm 1 and the definition of $\phi(x_k)$. The inequality (11) follows from the previous inequality and the fact that $\|\beta(x_k)\| \leq \gamma \|\phi(x_k)\|$ for all $k \in \mathcal{S}_{\phi} \subseteq \mathcal{S}_{\phi}^{\text{SD}}$ as can be seen by line 6 of Algorithm 1. \square

We may now establish finiteness of the index set $\mathcal{S}_{\phi}^{\text{SD}}$.

Lemma 3.9. *The index set $\mathcal{S}_{\phi}^{\text{SD}}$ is finite, i.e., $|\mathcal{S}_{\phi}^{\text{SD}}| < \infty$.*

Proof. To derive a contradiction, suppose that $|\mathcal{S}_{\phi}^{\text{SD}}| = \infty$, which means that Algorithm 1 does not terminate finitely. Thus, it follows from line 4 of Algorithm 1 that $\max\{\|\beta(x_k)\|, \|\phi(x_k)\|\} > \epsilon$ for all $k \geq 0$. Also, it follows from Lemmas 3.5 and 3.6 that there exists an iteration number k_1 such that $k \in \mathcal{S}_{\phi}^{\text{SD}}$ for all $k \geq k_1$. Thus, with Lemma 3.8, we have for all $\ell \geq k_1$ that

$$\begin{aligned} F(x_{k_1}) - F(x_{\ell+1}) &= \sum_{k=k_1}^{\ell} [F(x_k) - F(x_{k+1})] \\ &= \sum_{k \in \mathcal{S}_{\phi}^{\text{SD}}, k_1 \leq k \leq \ell} [F(x_k) - F(x_{k+1})] \\ &\geq \sum_{k \in \mathcal{S}_{\phi}^{\text{SD}}, k_1 \leq k \leq \ell} \kappa_{\phi} \max\{\gamma^{-2} \|\beta(x_k)\|^2, \|\phi(x_k)\|^2\} \\ &\geq \sum_{k \in \mathcal{S}_{\phi}^{\text{SD}}, k_1 \leq k \leq \ell} \kappa_{\phi} \min\{\gamma^{-2}, 1\} \epsilon^2. \end{aligned}$$

Rearranging the previous inequality shows that

$$\begin{aligned}\lim_{\ell \rightarrow \infty} F(x_{\ell+1}) &\leq \lim_{\ell \rightarrow \infty} \left[F(x_{k_1}) - \sum_{k \in \mathcal{S}_\phi^{\text{SD}}, k_1 \leq k \leq \ell} \kappa_\phi \min\{\gamma^{-2}, 1\} \epsilon^2 \right] \\ &= F(x_{k_1}) - \sum_{k \in \mathcal{S}_\phi^{\text{SD}}, k_1 \leq k} \kappa_\phi \min\{\gamma^{-2}, 1\} \epsilon^2 = -\infty,\end{aligned}$$

which contradicts Assumption 3.1. Thus, we conclude that $|\mathcal{S}_\phi^{\text{SD}}| < \infty$. \square

We now prove our first main convergence result.

Theorem 3.10. *Algorithm 1 terminates finitely.*

Proof. Since each iteration number k generated in the algorithm is an element of $\mathcal{S}_\beta \cup \mathcal{S}_\phi^{\text{ADD}} \cup \mathcal{S}_\phi^{\text{SD}}$, the result follows by Lemmas 3.5, 3.6, and 3.9. \square

Our final convergence result states what happens when the finite termination criterion is removed from Algorithm 1.

Theorem 3.11. *Let x_* be the unique solution to problem (1). If ϵ in the finite termination condition in line 4 of Algorithm 1 is replaced by zero, then either:*

- (i) *there exists an iteration k such that $x_k = x_*$; or*
- (ii) *infinitely many iterations $\{x_k\}$ are computed and they satisfy*

$$\lim_{k \rightarrow \infty} x_k = x_*, \quad \lim_{k \rightarrow \infty} \varphi(x_k) = 0, \quad \text{and} \quad \lim_{k \rightarrow \infty} \beta(x_k) = 0.$$

Proof. If case (i) occurs, then there is nothing left to prove. Thus, for the remainder of the proof, we assume that case (i) does not occur. Since case (i) does not occur, we know that Algorithm 1 performs an infinite sequence of iterations. Let us then define the set $\mathcal{S} := \mathcal{S}_\beta \cup \mathcal{S}_\phi^{\text{SD}}$, which must be infinite (since any consecutive subsequence of iterations in $\mathcal{S}_\phi^{\text{ADD}}$ must be finite by the finiteness of n). It follows from (7) for $k \in \mathcal{S}_\beta$, (11) for $k \in \mathcal{S}_\phi^{\text{SD}}$, and Assumption 3.1 (specifically, the assumption that f is bounded below over \mathcal{L}) that

$$\lim_{k \in \mathcal{S}} \max\{\|\beta(x_k)\|, \|\varphi(x_k)\|\} = 0.$$

Combining this with Assumption 3.1 and Lemma 2.1 gives

$$\lim_{k \in \mathcal{S}} x_k = x_*. \tag{17}$$

Now, we claim that the previous limit holds over all iterations. To prove this by contradiction, suppose that there exists an infinite $\mathcal{K} \subseteq \mathcal{S}_\phi^{\text{ADD}}$ and a scalar $\varepsilon > 0$ with

$$\|x_k - x_*\| \geq \varepsilon \quad \text{for all } k \in \mathcal{K}. \tag{18}$$

From Assumption 3.1, we conclude that there exists $\delta > 0$ such that

$$\text{if } F(x) \leq F(x_*) + \delta, \text{ then } \|x - x_*\| < \varepsilon. \tag{19}$$

Moreover, from (17) and Assumption 3.1, there exists a smallest $k_S \in \mathcal{S}$ such that

$$F(x_{k_S}) \leq F(x_*) + \delta. \tag{20}$$

There then exists a smallest $k_K \in \mathcal{K}$ such that $k_K > k_S$. Since, by construction, $\{F(x_k)\}_{k \geq 0}$ is monotonically decreasing, we may conclude with (20) that

$$F(x_{k_K}) \leq F(x_{k_S}) \leq F(x_*) + \delta. \tag{21}$$

Combining (21) and (19), we deduce that $\|x_{k_K} - x_*\| < \varepsilon$, which contradicts (18) since $k_K \in \mathcal{K}$. This completes the proof. \square

4 Numerical Results

In this section, we present results when employing an implementation of FaRSA to solve a collection of ℓ_1 -norm regularized logistic regression problems. Such problems routinely arise in the context of model prediction, making the design of advanced optimization algorithms that efficiently and reliably solve them paramount in big data applications. We first describe the datasets considered in our experiments, then describe some details of our implementation (henceforth simply referred to as FaRSA), and then present the results of our experiments.

4.1 Datasets

We tested FaRSA on ℓ_1 -norm regularized logistic regression problems using 31 datasets (see Table 1), 19 of which are available only after standard scaling practices have been applied. For the remaining 12 datasets, we considered both unscaled and scaled versions, where, for each, the scaling technique employed is described in the last column of Table 1. A checkmark in the “Unscaled” column indicates that we were able to obtain an unscaled version of that dataset.

Most of the datasets in Table 1 can be obtained from the LIBSVM repository.¹ From this repository, we excluded all regression and multiple-class (greater than two) instances, except for mnist since it is such a commonly used dataset. Since mnist is for digit classification, we transformed it for binary classification by assigning the digits 0–4 to the label -1 , and the digits 5–9 to the label 1 . The remaining datasets were binary classification examples from which we removed HIGGS, kdd2010(algebra), kdd2010(bridge to algebra), epsilon, url, and webspam since insufficient computer memory was available. (All experiments were conducted on a 64-bit machine with an Intel I7 4.0GHz CPU and 16GB of main memory.) Finally, for the adult data (a1a–a9a) and webpage data (w1a–w8a) we only used the largest instances, namely problems a9a and w8a. This left us with our final subset of datasets from LIBSVM.

In addition, we also tested FaRSA on three other datasets: synthetic, gene-ad, and pathway-ad. The synthetic set is a randomly generated non-diagonally dominant dataset created by the authors of OBA. The sets gene-ad and pathway-ad are datasets related to Alzheimer’s Disease. They were obtained by preprocessing the sets GSE4226² and GSE4227³ using the method presented in [15], and merging the results into the single dataset: gene-ad. The gene data (gene-ad) was converted to pathway data (pathway-ad) using the ideas described in [15]. The union of these three datasets and those from the LIBSVM repository comprised our complete test set.

For the unscaled datasets (see column 4 in Table 1), we adopted standard scaling techniques. For problems scaled into $[-1, 1]$ a simple linear scaling transformation was used. For problem mnist, which was scaled into $[0, 1]$, we used a common converting method in image processing. Specifically, we defined

$$I(i, j) = \frac{P(i, j)}{2^b}, \quad (22)$$

where $P(i, j)$ is the given unscaled integer pixel value satisfying

$$P(i, j) \in \{0, 1, 2, \dots, 2^b - 1\},$$

b is the intensity resolution ($b = 8$ for the mnist dataset), and (i, j) range over the size of the image. The scaled pixel values are then given by the values $I(i, j) \in [0, 1]$.

4.2 Implementation details

We developed a preliminary MATLAB implementation of FaRSA that we are happy to provide upon request. In this section, we describe the algorithm-specific choices made to obtain the results that we present.

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

²<http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE4226>

³<http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE4227>

Table 1: Data sets.

Dataset	# of Samples	# of Features	Unscaled	Scaling Used
fourclass	862	2	✓	into [-1,1]
svmguide1	3089	4	✓	into [-1,1]
cod-rna	59535	8		
breast-cancer	683	10		
australian	690	14		
SUSY	5000000	18	✓	into [-1,1]
splice	1000	60		
heart	270	13		
german.numer	1000	24	✓	into [-1,1]
diabetes	768	8	✓	into [-1,1]
liver-disorders	345	6	✓	into [-1,1]
w8a	49749	300		
madelon	2000	500	✓	into [-1,1]
a9a	32561	123		
mnist	30001	784	✓	into [0,1]
skin-nonskin	245057	3	✓	into [-1,1]
sonar	208	60		
ijcnn1	49990	22		
svmguide3	1243	22		
synthetic	5000	5000		
gisette	6000	5000		
pathway-ad	278	71		
real-sim	72309	20958		
covtype.binary	581012	8		
mushrooms	8124	112		
rcv1.binary	20242	47236		
leukemia	34	7129		
duke-breast-cancer	38	7129	✓	into [-1,1]
gene-ad	71	17375	✓	into [-1,1]
colon-cancer	62	2000	✓	into [-1,1]
news20	19996	1355191		

First, the weighting parameter in (1) was defined as

$$\lambda = \frac{1}{\# \text{ of Samples}}.$$

For determining the iteration type, we chose $\gamma = 1$ in line 6 of Algorithm 1 so that no preference was given to iterations being in either \mathcal{S}_ϕ or \mathcal{S}_β .

For any $k \in \mathcal{S}_\phi$, we made the simple choice of $\mathcal{I}_k = \{i : [\phi(x_k)]_i \neq 0\}$. This made the inequality in line 7 satisfied for any $\eta_\phi \in (0, 1]$, making the choice of this parameter irrelevant. (In a more sophisticated implementation, one might consider other choices of \mathcal{I}_k , say to adaptively control $|\mathcal{I}_k|$, to improve efficiency.) With this choice for \mathcal{I}_k made, Algorithm 1 allows for great flexibility in obtaining a search direction that satisfies the conditions in line 10 (see (iii) in Section 1.1 for additional comments). For our tests, we applied the linear-CG method to the system $H_k d = -g_k$ defined by the terms constructed in line 8, except that we added a diagonal matrix with entires 10^{-8} to H_k (an approach also adopted by OBA and LIBLINEAR). As discussed in Section 2, the conditions that are required to be satisfied by the trial step will hold if CG is terminated during any iteration. To help limit the number of backtracking steps required by the subsequent backtracking line search, we terminated CG as soon as one of three conditions was satisfied. To describe

these conditions, we let d_j denote the j th CG iteration, $r_j = \|H_k d_j + g_k\|$ denote the j th CG residual, and v_j denote the number of components in $x_k + d_j$ that fall into a different orthant than x_k . With these definitions, we terminated CG as soon as one of the following was satisfied:

$$\begin{aligned} r_j &\leq \max\{10^{-1}r_0, 10^{-12}\}, \\ v_j &\geq \max\{10^3, 10^{-1}|\mathcal{I}_k|\}, \text{ or} \\ \|d_j\| &\geq \delta_{k,\phi} := \max\{10^{-3}, \min\{10^3, 10\|x_{k_\phi(k)+1} - x_{k_\phi(k)}\|\}\}, \end{aligned}$$

where $k_\phi(k) := \max\{\bar{k} : \bar{k} \in \mathcal{S}_\phi \text{ and } \bar{k} < k\}$. This first condition is a standard requirement of asking the residual to be reduced by a fraction of the initial residual. We used the second condition to trigger termination when a CG iterate predicted that “too many” of the variables at $x_k + d_j$ are in the “wrong” orthant. Finally, the third condition ensured that the size of the trial step was moderate, thus functioning as an implicit trust-region constraint; this condition was motivated by the well-known fact that CG iterations $\{d_j\}$ are monotonically increasing in norm.

When $k \in \mathcal{S}_\beta$, we again made the simplest choice of $\mathcal{I}_k = \{i : [\beta(x_k)]_i \neq 0\}$, making the choice of $\eta_\beta \in (0, 1]$ irrelevant in our tests (though adaptive choices of \mathcal{I}_k might be worthwhile in a more sophisticated implementation). Since there is no natural scaling for the direction $\beta(x_k)$ because it is based on first derivative information only, it is important from a practical perspective to adaptively scale the direction. Therefore, in line 15, we used the alternative safeguarded direction defined by

$$[d_k]_{\mathcal{I}_k} = -\delta_{k,\beta} \frac{[\beta(x_k)]_{\mathcal{I}_k}}{\|[\beta(x_k)]_{\mathcal{I}_k}\|}, \quad (23)$$

where

$$\delta_{k,\beta} := \max\{10^{-5}, \min\{1, \|x_{k_\beta(k)+1} - x_{k_\beta(k)}\|\}\}$$

with $k_\beta(k) := \max\{\bar{k} : \bar{k} \in \mathcal{S}_\beta \text{ and } \bar{k} < k\}$. Since this is a safeguarded scaling of the d_k defined in line 15, it is fully covered by the theory that we developed in Section 3.

During each iteration, the values $\eta = 10^{-2}$ and $\xi = 0.5$ were used during the line search regardless of whether it was the line search performed by Algorithm 2 when called by Algorithm 1 (line 12) or if it was the line search performed by Algorithm 3 when called by Algorithm 1 (line 16). The starting point x_0 was chosen as the zero vector for all problems, and the termination tolerance, maximum allowed number of iterations, and maximum allowed time limit values were chosen to be $\epsilon = 10^{-6}$, 1000, and 10 minutes, respectively.

4.3 Test results

The output from FaRSA for the problems corresponding to the scaled and unscaled datasets in our experiments are summarized in Tables 2 and 3, respectively. These tables focus on the computational time in seconds and percentage of zeros (sparsity) in the computed solutions. For comparison purposes, we also provide the output from the OBA solver whose MATLAB implementation was graciously provided by the authors. For a fair comparison, we used the same stopping tolerance value of $\epsilon = 10^{-6}$ for OBA and made no modifications to their code. The numbers reported for each problem (named according to the corresponding dataset) are the averages from running each problem instance 10 times. We do not provide the final objective values since they were the same for FaRSA and OBA on all problems that were successfully solved by both algorithms. We use red numbering to indicate that an average CPU time was relatively lower for an algorithm, or if the average percentage of zeros in the solution was relatively larger for an algorithm.

We can observe from Table 2 that FaRSA performed better than OBA on 26 of the 31 (83.87%) scaled test problems. OBA is faster than FaRSA only on problems cod-rna, SUSY, synthetic, gene-ad, and colon-cancer. However, FaRSA is between 3 and 8 times faster than OBA on problems a9a, mnist, pathway-ad, covtype.binary, and news20, and between 1 and 3 times faster than OBA on the remaining 21 problems. In terms of sparsity, the two algorithms are comparable. Although not presented in the table, we find it interesting to note that FaRSA required an average of 34.12 iterations to solve the problems, with, on

Table 2: CPU time and sparsity for FaRSA and OBA on scaled problem variants.

Problems	Time (seconds)			% of zeros	
	FaRSA	OBA	OBA/FaRSA	FaRSA	OBA
fourclass	0.00326	0.00705	2.1626	0	0
svmguide1	0.0384	0.06457	1.6815	0	0
cod-rna	0.48762	0.18618	0.3818	0	0
breast-cancer	0.0089	0.03769	1.9674	0	0
australian	0.01443	0.0174	1.2058	0	0
SUSY	241.2437	205.1242	0.8502	0	0
splice	0.0101	0.01982	1.9624	5	5
heart	0.00706	0.01357	1.9221	7.7	7.7
german.numer	0.01159	0.02111	1.8214	8.3	8.3
diabetes	0.00581	0.00979	1.6850	12.5	12.5
liver-disorders	0.01254	max iter	Inf	16.7	—
w8a	0.97079	0.99154	1.0214	19.1	18.7
madelon	0.26604	0.37497	1.4094	19.8	19.8
a9a	0.78203	3.26994	4.1813	22.0	20.3
mnist	18.78034	54.432	3.0545	37.7	37.8
skin-nonskin	3.20594	ascent	Inf	41.7	—
sonar	0.02012	0.02938	1.4602	41.7	41.7
ijcnn1	0.06153	0.08178	1.3291	45.5	45.5
svmguide3	0.01856	0.03478	1.8739	45.5	45.5
synthetic	45.82688	20.42464	0.4457	57.4	49.5
gisette	13.30533	28.22136	2.1211	84.6	84.6
pathway-ad	0.16054	1.30585	8.1341	87.4	87.4
real-sim	2.3221	2.43764	1.0214	91.9	91.8
covtype.binary	1.49449	5.95536	3.9849	96.3	90.7
mushrooms	0.03089	0.05815	1.8825	97.3	97.3
rcv1.binary	0.39186	0.80563	1.7427	98.8	98.8
leukemia	0.09151	0.12086	1.3207	99.7	99.7
duke-breast-cancer	0.06227	0.10628	1.7068	99.7	99.7
gene-ad	0.21525	0.15943	0.7407	99.8	99.8
colon-cancer	0.04069	0.03905	0.9597	99.9	99.9
news20	6.09086	19.77945	3.2474	99.9	99.9

average, 29.93 of them being in \mathcal{S}_ϕ . This indicates that FaRSA quickly identifies the orthant that contains the optimal solution.

By turning our attention to Table 3, we see that the performance of both FaRSA and OBA deteriorates when the problems are unscaled. Moreover, OBA fails on problems skin-nonskin, gene-ad, and mnist because it generates iterates that increase the objective function; we denote these failures as “ascent” in the table. In theory, ascent is only possible for their method when their fixed estimate (10^8 in their code) of the Lipschitz constant for the gradient of f is not large enough. Although simple adaptive strategies could be used to avoid such issues, we made no such attempts because we did not want to make any edits to their code. Overall, FaRSA was able to solve 9 of the 12 unscaled problems, and OBA only performed better than FaRSA on a single test problem (colon-cancer).

The previous tables show that FaRSA efficiently and reliably obtains solutions that satisfy the stopping tolerance value of $\epsilon = 10^{-6}$. In practice, one sometimes only requests a low accuracy solution, often motivated by problems that may arise due to overfitting. To explore the performance of FaRSA for various stopping tolerance levels, we created the plots in Figures 1 and 2. Each plot shows the run time (y -axis) required to achieve the desired optimality accuracy (x -axis) for the stated problem. These figures show that the superior

Table 3: CPU time and sparsity for FaRSA and OBA on unscaled problem variants.

Problems	Time (seconds)			% of zeros	
	FaRSA	OBA	OBA/FaRSA	FaRSA	OBA
fourclass	0.00486	0.00775	1.5946	0	0
diabetes	0.01964	0.02159	1.0993	0	0
german.numer	0.03168	0.0564	1.7803	0	0
skin-nonskin	0.11378	ascent	Inf	0	—
madelon	6.55674	41.9519	6.3983	8	8
liver-disorders	0.00571	0.03277	5.7391	83.3	83.3
colon-cancer	0.08429	0.05364	0.6364	98.7	98.7
duke-breast-cancer	0.08936	0.12487	1.3973	99.7	99.7
gene-ad	4.62839	ascent	Inf	99.8	—
svmguide1	max iter	max iter	—	—	—
mnist	max time	ascent	—	—	—
SUSY	max iter	max iter	—	—	—

performance of FaRSA previously displayed for the stopping tolerance 10^{-6} also generally holds for larger stopping tolerances.

5 Conclusions

We presented a new reduced-space algorithm, FaRSA, for minimizing an ℓ_1 -norm regularized convex function. The method uses an adaptive condition to determine when the current reduced-space should be updated, which is itself based on measures of optimality in the current reduced space and its complement. Global convergence was established for our method, while numerical experiments on ℓ_1 -norm regularized logistic problems exhibited its practical performance. In particular, the experiments showed that FaRSA was generally superior to a recently proposed reduced-space orthant-based algorithm called OBA, regardless of the solution accuracy requested. Since OBA was shown in [9] to be better than the state-of-the-art solver used in LIBLINEAR when the second derivative matrices were not diagonally dominant, we expect that FaRSA will serve as a valuable data analysis tool. OBA and our preliminary implementation of FaRSA will often be outperformed by LIBLINEAR when the second derivative matrices are diagonally dominant. However, FaRSA was designed with great flexibility in how the subproblem solutions are obtained. Although our preliminary implementation invoked linear-CG as the subproblem solver, our framework also allows for coordinate-descent based algorithms to be used, such as those used in LIBLINEAR. We expect to provide such options as well as include features that control the subproblem size in a future release of our solver. We believe that once these enhancements have been made, FaRSA will be competitive with LIBLINEAR on all classes of problems, and superior when the second derivative matrices are not diagonally dominant.

Acknowledgments

We thank Nitish Keskar, Jorge Nocedal, Figen Öztoprak, and Andreas Wächter for providing the MATLAB code of their OBA algorithm, and for several discussions on their numerical experience with OBA. We also thank Qingsong Zhu for providing us the datasets gene-ad and pathway-ad used in Section 4.1.

References

- [1] Galen Andrew and Jianfeng Gao. Scalable training of l_1 -regularized log-linear models. In *Proceedings of the 24th international conference on Machine learning*, pages 33–40. ACM, 2007.

- [2] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [3] Richard H Byrd, Gillian M Chin, Jorge Nocedal, and Figen Oztoprak. A family of second-order methods for convex l1-regularized optimization. *Unpublished: Optimization Center: Northwestern University, Tech Report*, 2012.
- [4] Richard H Byrd, Jorge Nocedal, and Figen Oztoprak. An inexact successive quadratic approximation method for convex l-1 regularized optimization. *arXiv preprint arXiv:1309.3529*, 2013.
- [5] Zdenek Dostál. Box constrained quadratic programming with proportioning and projections. *SIAM Journal on Optimization*, 7(3):871–887, 1997.
- [6] Zdenek Dostál. A proportioning based algorithm with rate of convergence for bound constrained quadratic programming. *Numerical Algorithms*, 34(2):293–302, 2003.
- [7] Zdenek Dostál and Joachim Schoberl. Minimizing quadratic functions subject to bound constraints with the rate of convergence and finite termination. *Computational Optimization and Applications*, 30(1):23–43, 2005.
- [8] Cho-Jui Hsieh, Inderjit S Dhillon, Pradeep K Ravikumar, and Mátyás A Sustik. Sparse inverse covariance matrix estimation using quadratic approximation. In *Advances in Neural Information Processing Systems*, pages 2330–2338, 2011.
- [9] Nitish Shirish Keskar, Jorge Nocedal, Figen Oztoprak, and Andreas Waechter. A second-order method for convex ℓ_1 -regularized optimization with active set prediction. *arXiv preprint arXiv:1505.04315*, 2015.
- [10] Jason Lee, Yuekai Sun, and Michael Saunders. Proximal newton-type methods for convex optimization. In *Advances in Neural Information Processing Systems*, pages 836–844, 2012.
- [11] Katya Scheinberg and Xiaocheng Tang. Practical inexact proximal quasi-newton method with global complexity analysis. *arXiv preprint arXiv:1311.6547*, 2013.
- [12] Hassan Mohy ud Din and Daniel P. Robinson. A solver for nonconvex bound-constrained quadratic optimization. *SIAM Journal on Optimization*, 25(4):2385–2407, 2015.
- [13] Stephen J Wright, Robert D Nowak, and Mário AT Figueiredo. Sparse reconstruction by separable approximation. *Signal Processing, IEEE Transactions on*, 57(7):2479–2493, 2009.
- [14] Guo-Xun Yuan, Chia-Hua Ho, and Chih-Jen Lin. An improved glmnet for l1-regularized logistic regression. *The Journal of Machine Learning Research*, 13(1):1999–2030, 2012.
- [15] Qingsong Zhu, Evgeny Izumchenko, Alexander M Aliper, Evgeny Makarev, Keren Paz, Anton A Buzdin, Alex A Zhavoronkov, and David Sidransky. Pathway activation strength is a novel independent prognostic biomarker for cetuximab sensitivity in colorectal cancer patients. *Human Genome Variation*, 2, 2015.

A A Relationship between FaRSA and ISTA

The step in Line 16 of Algorithm 1 may be interpreted as a *reduced* ISTA [2] step. The next lemma makes this relationship precise.

Lemma A.1. For any k , let s_k be the full ISTA step defined by

$$s_k := \text{shrink}(x_k - \nabla f(x_k)) - x_k, \text{ where}$$

$$[\text{shrink}(x_k - \nabla f(x_k)) - x_k]_i := \begin{cases} -[\nabla f(x_k)]_i + \lambda & \text{if } [x_k - \nabla f(x_k)]_i < -\lambda, \\ -[x_k]_i & \text{if } [x_k - \nabla f(x_k)]_i \in [-\lambda, \lambda], \\ -[\nabla f(x_k)]_i - \lambda & \text{if } [x_k - \nabla f(x_k)]_i > \lambda. \end{cases}$$

Then, $s_k = -(\beta(x_k) + \phi(x_k))$.

Proof. Recall the definitions of the components of $\beta(x_k)$ and $\phi(x_k)$, which may be rewritten in a slightly more convenient form as follows:

$$[\beta(x_k)]_i := \begin{cases} [\nabla f(x_k)]_i + \lambda & \text{if } [x_k]_i = 0 \text{ and } [\nabla f(x_k)]_i + \lambda < 0, \\ [\nabla f(x_k)]_i - \lambda & \text{if } [x_k]_i = 0 \text{ and } [\nabla f(x_k)]_i - \lambda > 0, \\ 0 & \text{otherwise,} \end{cases}$$

$$[\phi(x_k)]_i := \begin{cases} 0 & \text{if } [x_k]_i = 0, \\ \min\{[\nabla f(x_k)]_i + \lambda, \max\{[x_k]_i, [\nabla f(x_k)]_i - \lambda\}\} & \text{if } [x_k]_i > 0 \text{ and } [\nabla f(x_k)]_i + \lambda > 0, \\ \max\{[\nabla f(x_k)]_i - \lambda, \min\{[x_k]_i, [\nabla f(x_k)]_i + \lambda\}\} & \text{if } [x_k]_i < 0 \text{ and } [\nabla f(x_k)]_i - \lambda < 0, \\ [\nabla f(x_k)]_i + \lambda \cdot \text{sgn}(x_k)_i & \text{otherwise.} \end{cases}$$

For any component i , we proceed by considering various cases and subcases.

Case 1: Suppose that

$$[x_k - \nabla f(x_k)]_i > \lambda, \text{ meaning that } [x_k]_i > [\nabla f(x_k)]_i + \lambda. \quad (24)$$

Subcase 1a: Suppose that $[x_k]_i > 0$ and $[\nabla f(x_k)]_i + \lambda > 0$, so $[\nabla f(x_k)]_i > -\lambda$. Then, $[\beta(x_k)]_i = 0$ and

$$[\phi(x_k)]_i = \min\{[\nabla f(x_k)]_i + \lambda, \max\{[x_k]_i, [\nabla f(x_k)]_i - \lambda\}\}. \quad (25)$$

By (24), it follows that $[x_k]_i > [\nabla f(x_k)]_i + \lambda > [\nabla f(x_k)]_i - \lambda$, which along with $[x_k]_i > 0$ means that the max in (25) evaluates as $[x_k]_i$. Then, again with (24), the min in (25) yields

$$[\phi(x_k)]_i = [\nabla f(x_k)]_i + \lambda = -[s_k]_i. \quad (26)$$

Subcase 1b: Suppose that $[x_k]_i > 0$ and $[\nabla f(x_k)]_i + \lambda \leq 0$, so $[\nabla f(x_k)]_i \leq -\lambda$. Then, $[\beta(x_k)]_i = 0$ and

$$[\phi(x_k)]_i = [\nabla f(x_k)]_i + \lambda = -[s_k]_i. \quad (27)$$

Subcase 1c: Suppose that $[x_k]_i = 0$ and $[\nabla f(x_k)]_i + \lambda \leq 0$, so $[\nabla f(x_k)]_i \leq -\lambda$. Then, $[\phi(x_k)]_i = 0$ and

$$[\beta(x_k)]_i = [\nabla f(x_k)]_i + \lambda = -[s_k]_i. \quad (28)$$

Subcase 1d: Suppose that $[x_k]_i < 0$ and $[\nabla f(x_k)]_i + \lambda < 0$, so $[\nabla f(x_k)]_i < -\lambda$. Then, $[\beta(x_k)]_i = 0$ and

$$[\phi(x_k)]_i = \max\{[\nabla f(x_k)]_i - \lambda, \min\{[x_k]_i, [\nabla f(x_k)]_i + \lambda\}\}. \quad (29)$$

By (24), it follows that $[x_k]_i > [\nabla f(x_k)]_i + \lambda$, which along with $[\nabla f(x_k)]_i + \lambda < 0$ means that the min in (29) evaluates as $[\nabla f(x_k)]_i + \lambda$. Then, since $[\nabla f(x_k)]_i + \lambda > [\nabla f(x_k)]_i - \lambda$, the max in (29) yields

$$[\phi(x_k)]_i = [\nabla f(x_k)]_i + \lambda = -[s_k]_i. \quad (30)$$

Since Subcases 1a–1d exhaust all possibilities under (24), we conclude from the results in (26), (27), (28), and (30) that for Case 1 we have $[s_k]_i = -[\beta(x_k) + \phi(x_k)]_i$.

Case 2: Suppose that

$$[x_k - \nabla f(x_k)]_i < -\lambda, \text{ meaning that } [x_k]_i < [\nabla f(x_k)]_i - \lambda. \quad (31)$$

We claim that the analysis for this case is symmetric to that in Case 1 above, from which we may conclude that for this case we again have $[s_k]_i = -[\beta(x_k) + \phi(x_k)]_i$.

Case 3: Suppose that

$$[x_k - \nabla f(x_k)]_i \in [-\lambda, \lambda], \text{ meaning that } [x_k]_i \in [\nabla f(x_k)]_i + [-\lambda, \lambda]. \quad (32)$$

Subcase 3a: Suppose that $[x_k]_i > 0$. Then, $[\beta(x_k)]_i = 0$ and, since $[x_k]_i > 0$ and (32) imply $[\nabla f(x_k)]_i > -\lambda$,

$$[\phi(x_k)]_i = \min\{[\nabla f(x_k)]_i + \lambda, \max\{[x_k]_i, [\nabla f(x_k)]_i - \lambda\}\}. \quad (33)$$

Since (32) also implies $[x_k]_i > [\nabla f(x_k)]_i - \lambda$, it follows along with $[x_k]_i > 0$ that the max in (33) evaluates as $[x_k]_i$. Then, since (32) implies $[x_k]_i < [\nabla f(x_k)]_i + \lambda$, the min in (33) yields

$$[\phi(x_k)]_i = [x_k]_i = -[s_k]_i. \quad (34)$$

Subcase 3b: Suppose that $[x_k]_i = 0$. Then, $[\phi(x_k)]_i = 0$ and, along under (32),

$$[\beta(x_k)]_i = -[s_k]_i = 0. \quad (35)$$

Subcase 3c: Suppose that $[x_k]_i < 0$. We claim that the analysis for this case is symmetric to that in Subcase 3.a, from which we may conclude that for this subcase we again have

$$[\phi(x_k)]_i = [x_k]_i = -[s_k]_i. \quad (36)$$

Since Subcases 1.a–1.d exhaust all possibilities under (32), we conclude from the results in (34), (35), and (36) that for Case 3 we have $[s_k]_i = -[\beta(x_k) + \phi(x_k)]_i$. The result follows as we have proved the desired result under all cases. \square

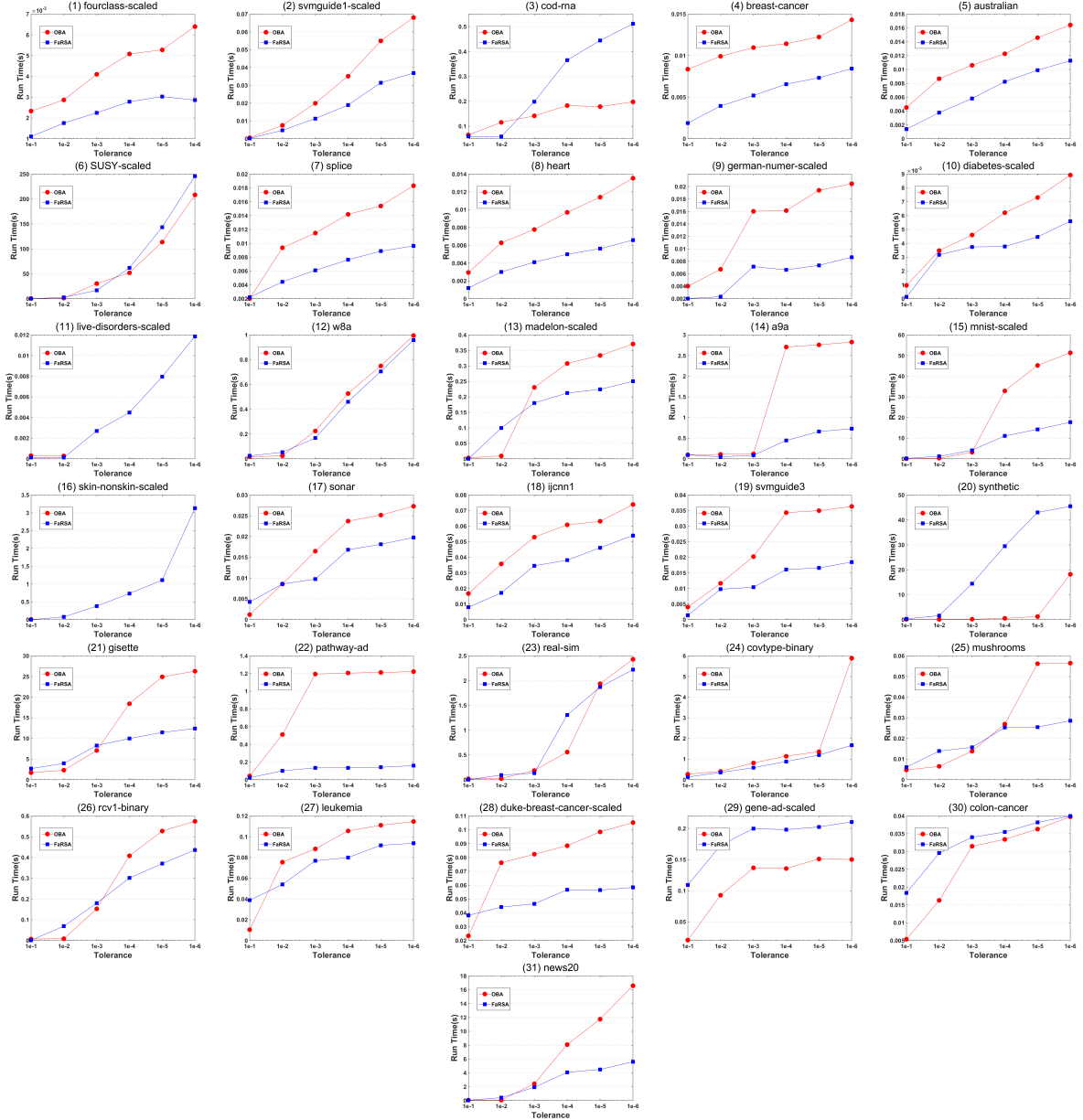


Figure 1: CPU time comparison between FaRSA and OBA for various stopping tolerances on the set of problems with scaled data.

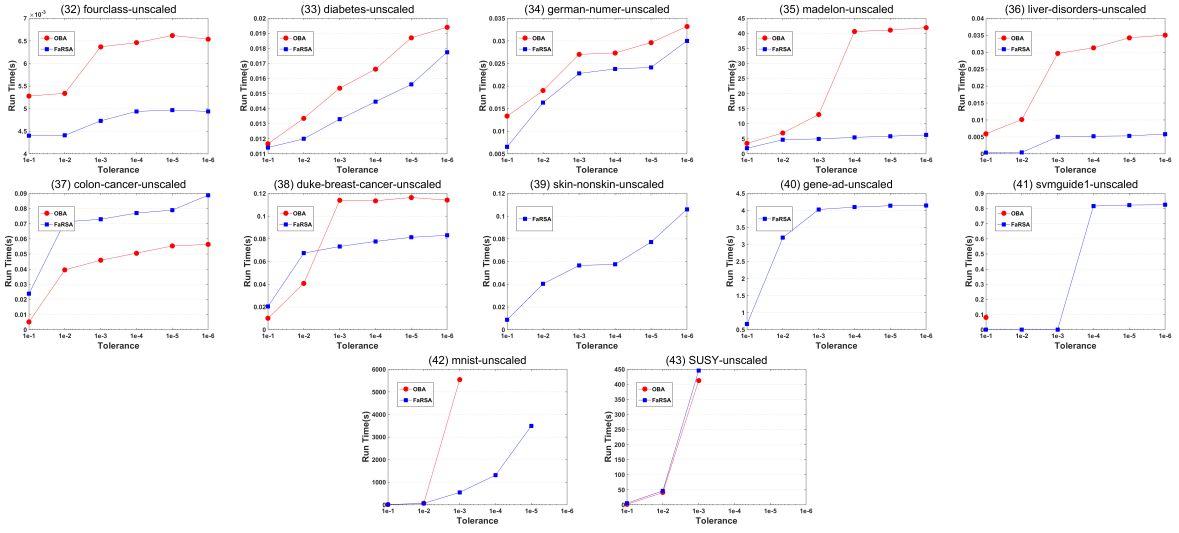


Figure 2: CPU time comparison between FaRSA and OBA for various stopping tolerances on the set of problems with unscaled data.