# A globally and quadratically convergent affine scaling method for linear $\ell_1$ problems

## Thomas F. Coleman*

*Computer Science Department and Center for Applied Mathematics, Cornell University, Ithaca, NY, USA*

## Yuying Li**

*Computer Science Department, Cornell University, Ithaca, NY, USA*

Recently, various interior point algorithms related to the Karmarkar algorithm have been developed for linear programming. In this paper, we first show how this "interior point" philosophy can be adapted to the linear $\ell_1$ problem (in which there are no feasibility constraints) to yield a globally and linearly convergent algorithm. We then show that the linear algorithm can be modified to provide a *globally* and ultimately *quadratically* convergent algorithm. This modified algorithm appears to be significantly more efficient in practise than a more straightforward interior point approach via a linear programming formulation: we present numerical results to support this claim.

## 1. Introduction

Let $A$ be an $n \times m$ matrix with $m > n$ and consider the overdetermined system

$$A^\mathrm{T} x \approx b,$$

where row $i$ of $A^T$ is denoted by $a_i^T$, for $i = 1, \ldots, m$. The linear $\ell_1$ problem is to find a vector $x$ which is a solution to

$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^m |a_i^T x - b_i|. \tag{1}$$

Note that the objective function is piecewise linear; it is not differentiable at any point $x$ such that $a_i^T x = b_i$, for some index $i$. Moreover, it is well known that problem (1) is equivalent to the following linear program (LP):

$$\min_{u,v,x} \sum_{i=1}^m (u_i + v_i)$$

$$\text{subject to} \quad a_i^T x - u_i + v_i = b_i, \quad i = 1, \ldots, m, \tag{2}$$

$$u_i \geqslant 0, \ v_i \geqslant 0, \quad i = 1, \ldots, m.$$

Thus, the $\ell_1$ problem can be solved by a general LP-solver. However, usually it is more efficient to use techniques especially tailored to problem (1), e.g., [2] and [3]. Most such methods are related to the simplex method in that they involve notions of an active set and pivoting or exchanging columns. These are finite algorithms.

Following Karmarkar's work [11] several alternative approaches to linear programming have recently been developed (e.g., [1, 10, 16, 19, 21]). Amongst the various types of approaches, affine scaling methods (e.g., [1, 21]) appear to represent a relatively practical approach. It is difficult to give a crisp accurate definition of an affine scaling method; however, the salient feature is the use of an *affine transformation*, defined locally (i.e., at the current point). The advent of affine scaling methods for linear programming raises the question: can such methods be tailored to the linear $\ell_1$ problem? Indeed, an interior point strategy applied to the dual of (2) has been suggested (see [12] and [17]). Such methods can be viewed as generating an *infinite* sequence of approximate solutions $\{x_k\}$.[1] However, despite numerous connections to Newton-like processes (e.g., [16]), affine scaling methods, such as Meketon's algorithm [12], do not possess second-order convergence, even in the limit.

In this paper we propose a globally convergent affine scaling algorithm, for the $\ell_1$ problem, which is ultimately quadratically convergent (in the nondegenerate case).

## 1.1. Notations and definitions

In this paper, the vector $r$ is defined to be the residual vector $b - A^T x$. We use the sign function $\text{sgn}(r)$, where $r$ is a vector, in the following sense: if $\sigma = \text{sgn}(r)$,

$$\sigma_i = \begin{cases} 1 & \text{if } r_i \geqslant 0, \\ -1 & \text{otherwise.} \end{cases}$$

---

[1] One can also take a finite view of such algorithms — assuming integer data and exact arithmetic. This view leads to a complexity analysis; e.g., is the number of steps bounded by a polynomial in the size of the problem? This is not our concern in this paper.

We use $e_i$ to denote the $i$th elementary vector. We let $\mathscr{A}(x)$ denote the indices of zero residuals at any point $x$, i.e.,

$$\mathscr{A}(x) = \{i \mid a_i^T x - b_i = 0\},$$

and $\mathscr{A}^c(x)$ denotes the complementary set. (We will suppress the argument when it is clear from context.) We use the subscripts 1 and 0 to denote the compressed vector whose components correspond to the nonactivities and activities respectively, e.g., $g_1 = g_{\mathscr{A}^c}$, $g_0 = g_{\mathscr{A}}$. Similarly, $A_0$ and $A_1$ denote the columns corresponding to zero residuals, $A_0^T x - b_0 = 0$, and nonzero residuals. In some cases this denotation will refer to the current point; otherwise, it will refer to the limit point. The usage will be clear from context.

In our presentation, the multiplication between two vectors is a componentwise operation. The operator $|\cdot|$ around a set, e.g., $|\mathscr{A}|$, denotes the cardinality of that set. Otherwise it denotes the componentwise absolute values of a number, vector or a matrix. The operator $\max(x, y)$ with two vectors as arguments defines a vector whose components are the maximum of the corresponding argument vectors. If the argument is a single vector, the result is the maximum entry; therefore $\max(x)$ of a vector $x$ denotes the maximum component of $x$. The operator $\text{null}(A)$ denotes the matrix whose columns form a basis for the null space of $A$, i.e., if $B = \text{null}(A)$, $AB = 0$ and $\text{rank}(B) = n - \text{rank}(A)$. The left arrow $x \leftarrow y$ denotes setting $y$ to $x$.

The dual of (2) is

$$\max_{\lambda} \quad b^T \lambda$$

$$\text{subject to} \quad A\lambda = 0, \tag{3}$$

$$-1 \le \lambda_i \le 1, \quad i = 1, \ldots, m.$$

**Definition 1.** We say an $\ell_1$ problem is *primal nondegenerate* if at any point $x$ the vectors $\{a_i : i \in \mathscr{A}(x)\}$ are linearly independent.

**Definition 2.** We say an $\ell_1$ problem is *dual nondegenerate* if, for any $\lambda$ satisfying $A\lambda = 0$, $|\{\lambda_i : |\lambda_i| = 1\}| \le m - n$.

The optimal solution to (1) can be characterized in various ways. For example, $x$ is optimal if and only if there exists $\lambda \in \mathbb{R}^{|\mathscr{A}|}$ such that

$$\sum_{i \in \mathscr{A}^c} \text{sgn}(a_i^T x - b_i) a_i = \sum_{i \in \mathscr{A}} \lambda_i a_i, \quad \text{where } -1 \le \lambda_i \le 1 \ \forall i \in \mathscr{A}. \tag{4}$$

Thus, $\lambda$ is clearly a feasible solution of (3). Note that, when $\mathscr{A}$ is empty, $\sum_{i \in \mathscr{A}^c} \text{sgn}(a_i^T x - b_i) a_i = 0$. It is easy to verify that the optimality condition (4) is

equivalent to the following: $x$ is optimal for (1) if and only if there exists $\lambda \in \mathbb{R}^m$ such that

$$(a_i^T x - b_i) * (\text{sgn}(a_i^T x - b_i) - \lambda_i) = 0, \quad i = 1, \ldots, m, \tag{5}$$

$$A\lambda = 0, \tag{6}$$

$$-1 \leq \lambda_i \leq 1 \quad \forall i \in \mathcal{A}. \tag{7}$$

This characterization is especially interesting because the nonlinear equation (5) suggests the possibility of a Newton method. We explore this possibility in Section 3.

**Remark.** Primal and dual nondegeneracy imply that the linear $\ell_1$ problem has a unique minimizer, $(x^*, \lambda^*)$, e.g., [15]. Moreover,

$$|\mathcal{A}(x^*)| = \text{rank}(A_0(x^*)) = n,$$

$$|\mathcal{A}^c(x^*)| = |\{\lambda_i^* : |\lambda_i^*| = 1\}| = m - n.$$

### 1.2. An equivalent formulation

Following Seneta and Steiger [18], we consider an alternative (but equivalent) formulation of the $\ell_1$ problem. Let $Z$ denote a matrix whose rows form a basis for the null space of $A$, i.e., $Z^T = \text{null}(A)$. Hence $Z$ has dimensions $(m - \text{rank}(A)) \times m$, $\text{rank}(Z) = m - \text{rank}(A)$, and $AZ^T = 0$. Then, defining $r = b - A^T x$, the linear $\ell_1$ problem is equivalent to the following constrained $\ell_1$ problem with $m$ variables $r$:

$$\min_{r \in \mathbb{R}^m} \quad \{\psi(r) = \|r\|_1\}$$

$$\text{subject to} \quad Zr = Zb. \tag{8}$$

This formulation of the $\ell_1$ problem is useful for the derivation of our algorithms; however, implementation of our methods, as discussed in Section 4, does not necessarily involve the computation of the matrix $Z$.

Once again optimality conditions can take different forms. For example, $r$ is a solution to (8) if and only if there exists $\mu \in \mathbb{R}^{|\mathcal{A}|}$ and $w \in \mathbb{R}^{m - |\mathcal{A}|}$ such that

$$\sum_{\mathcal{A}^c} \text{sgn}(r_i) e_i = \sum_{\mathcal{A}} e_i \mu_i + Z^T w,$$

$$Z(r - b) = 0, \tag{9}$$

$$-1 \leq \mu_i \leq 1 \quad \forall i \in \mathcal{A}.$$

Clearly the first system of equations in (9) imply that $|(Z^T w)_i| \leq 1 \ \forall i \in \mathcal{A}$. Therefore, if we define $\lambda = Z^T w$, an equivalent expression for optimality conditions (9) is

$$r_i(\text{sgn}(r_i) - \lambda_i) = 0, \quad i = 1, \ldots, m, \tag{10}$$

$$Z(r - b) = 0, \tag{11}$$

$$-1 \leq \lambda_i \leq 1, \quad i \in \mathcal{A}. \tag{12}$$

Note that if $x$ satisfies $r = b - A^T x$, then $(x, \lambda)$ satisfies (5)–(7).

This formulation is important to our design of a quadratically convergent algorithm: equations (10) and (11) lead to a Newton process, which we introduce in Section 3.

## 2. An affine scaling method

In this section we present a new affine scaling method for the linear $\ell_1$ problem. Our approach maintains feasibility with respect to the equality constraints in (8); the scaling involved in the computation of a descent direction is determined by the distance from the current point to the lines of nondifferentiability (i.e., $r_i = 0$). It is this notion that replaces the more standard "distance from the boundary" definition used in most interior methods (e.g., [1, 9, and 21]). However, since the differentiable region[2] in (8) is not connected, and since we do not know how to immediately identify a connected differentiable region adjacent to the optimal point, our algorithm must allow for the ability to cross lines of nondifferentiability.

Let $D$ be a positive diagonal matrix.

Assuming we are in the differentiable region, we can define a descent direction by the following "trust region problem":

$$\min_{d \in \mathbb{R}^m} \quad g^T d$$

$$\text{subject to} \quad Zd = 0, \tag{13}$$

$$\|D^{-1}d\|_2 \leq \delta,$$

where $\delta$ is a positive number reflecting the trust region size; $g = g(r)$ is the gradient of $\psi(r) = \sum |r_i|$; the matrix $D$ defines the shape of the ellipsoid created by $\|D^{-1}d\|_2 \leq \delta$. We choose $D = \text{diag}\{|r_i|^{1/2}\}$ and with this choice the ellipsoid is short in directions corresponding to components of $r_i$ close to zero, and long in directions corresponding to relatively large $|r_i|$. The solution to (13) is of the form $d_* = \alpha d$ where

$$d = -D^2(g - Z^T(ZD^2Z^T)^{-1}ZD^2g). \tag{14}$$

Here $\alpha$ is a function of $\delta$, for small $\delta > 0$.

Rather than choose $\delta$ a priori, we prefer to compute $d$ by (14) and then define $\alpha$ through minimizing a piecewise linear function $\psi(r + \alpha d)$ along the ray $d$ (allowing for the ability to cross lines of nondifferentiability). Hence we must determine all nonnegative breakpoints,

$$\mathcal{J} = \{a_i: \alpha_i > 0, \alpha_i = -r_i/d_i\}, \tag{15}$$

and using this information determine the minimizer of $\psi(r + \alpha d)$ with respect to $\alpha$:

$$\psi(r + \alpha_* d) = \min_{\alpha > 0} \psi(r + \alpha d) = \min_{\alpha \in \mathcal{J}} \psi(r + \alpha d). \tag{16}$$

---

[2] The differentiable region consists of all point $r \in \mathbb{R}^m$ such that $\prod_{i=1}^m r_i \neq 0$.

We refer to $\alpha_*$ as the optimal breakpoint. The point $\alpha_*$ is computed by considering each break-point in $\mathcal{J}$ in turn, adjusting the gradient $\nabla \psi(r)$ to reflect a step just beyond the breakpoint and then determining if $d$ continues to be a descent direction for $\psi(r)$. For example, if $\alpha_\diamond$ is the smallest positive breakpoint, then a step just beyond this point yields the following gradient:

$$g^+ = g - 2g_\diamond e_\diamond, \quad \text{where } \alpha_\diamond = \min\{\alpha_i : \alpha_i \in \mathcal{J}\}. \tag{17}$$

If $(g^+)^T d < 0$ the next breakpoint is considered, etc. Since we want to stay in the differentiable region of $\psi(r)$, we cannot step all the way to the minimizer. Therefore we take as our steplength

$$\alpha = \alpha_\sharp + \tau(\alpha_* - \alpha_\sharp), \quad \text{where } \alpha_\sharp = \max_{\mathcal{J} \cup \{0\}} \{\alpha_i : 0 \le \alpha_i < \alpha_*\} \tag{18}$$

and $0 < \tau < 1$. Here, if $\alpha_*$ is not the smallest positive breakpoint, $\alpha_\sharp$ is the largest breakpoint smaller than the optimal breakpoint $\alpha_*$; otherwise, $\alpha_\sharp = 0$.

In summary, we have the following line search procedure for $\psi(r)$:

**Piecewise Linear Line Search Procedure.** Given $0 < \tau < 1$ and $d$.
    *Step 1.* Compute $\mathcal{J} = \{\alpha_i : \alpha_i = -r_i/d_i, \ r_i d_i < 0\}$.
    *Step 2.* Find the optimal breakpoint $\alpha_*$ such that

$$\psi(r + \alpha_* d) = \min_{\alpha > 0} \psi(r + \alpha d) = \min_{\alpha \in \mathcal{J}} \psi(r + \alpha d).$$

    *Step 3.* Let $\alpha_\sharp = \max_{\mathcal{J} \cup \{0\}} \{\alpha_i : 0 \le \alpha_i < \alpha_*\}$. Set the stepsize

$$\alpha = \alpha_\sharp + \tau(\alpha_* - \alpha_\sharp).$$

We now present the linearly convergent algorithm. Let $r^0$ be an initial differentiable point satisfying $Zr^0 = Zb$; $k \leftarrow 0$; $\tau \in (0, 1)$.

**Algorithm 1.**
    *Step 1.* Define $D^k = \text{diag}\{|r_i^k|^{1/2}\}$ and $g^k = \text{sgn}(r^k)$.
    *Step 2.* Compute

$$d^k = -(D^k)^2\{g^k - (Z^k)^T[Z^k(D^k)^2(Z^k)^T]^{-1} Z^k(D^k)^2 g^k\}.$$

    *Step 3.* Apply the Piecewise Linear Line Search Procedure with the constant $\tau$ to determine $\alpha^k$. Then

$$r^{k+1} \leftarrow r^k + \alpha^k d^k, \quad k \leftarrow k + 1.$$

**Remarks.** (1) Note that since the linesearch maintains nonzero residual components, we have $D^k > 0$; therefore, $d^k$ is well-defined and the linesearch guarantees $\psi(r_{k+1}) < \psi(r_k)$.

(2) We have also considered using $D^k = \text{diag}\{|r_i^k|\}$ as our scaling matrix. Certainly this choice is more consistent with other methods, e.g., [1, 9, 21]. However, we have not had success with $D^k = \text{diag}\{|r_i^k|\}$ in two aspects. First, as we will show the choice $D^k = \text{diag}\{|r_i^k|^{1/2}\}$ leads to a hybrid algorithm that smoothly connects the linearly convergent process to a second-order Newton method; we have been unable to devise a similar hybrid scheme using $D^k = \text{diag}\{|r_i^k|\}$. Second, it is not clear that the choice $D^k = \text{diag}\{|r_i^k|\}$ leads to a globally convergent process — certainly our convergence proof does not apply in this case. The difficulty appears to be that the choice $D^k = \text{diag}\{|r_i^k|\}$ produces directions too nearly tangential to the near activities thus inhibiting the algorithm's ability to cross lines of nondifferentiability. This is not an issue in an interior point method since, in that setting, constraints are never crossed due to the maintenance of feasibility.

## 3. A local Newton process

A Newton process for the $\ell_1$ problem can be defined by considering conditions (10) and (11). Most importantly, consider (10) which we rewrite as

$$\text{diag}(r)(\text{sgn}(r) - Z^T w) = 0. \tag{19}$$

In general this system is not differentiable due to the discontinuities caused when $r_i = 0$, for some $i$. However, system (19) can certainly be differentiated at any point with no zero residual. Moreover, in a neighbourhood of a solution $r^*$, $\text{sgn}(r_i^*)$ will remain constant for any $i \in \mathscr{A}^c(r^*)$. Therefore, in a neighbourhood of $r^*$, discontinuities will be due only to the active equations, $i \in \mathscr{A}(r^*)$. As we formally establish in Section 6, these discontinuities do not impede the local quadratic convergence behaviour of a Newton process. Therefore, define $g = \text{sgn}(r)$, $D_r = \text{diag}(r)$, and $D_\lambda = \text{diag}(g - Z^T w)$; differentiate (19) and (11) to define a Newton correction for (19),

$$\begin{bmatrix} D_\lambda & -D_r Z^T \\ Z & 0 \end{bmatrix} \begin{bmatrix} \Delta r \\ \Delta w \end{bmatrix} = \begin{bmatrix} -D_r(g - Z^T w) \\ 0 \end{bmatrix}. \tag{20}$$

Simple algebraic manipulation yields

$$\Delta r = -A^T (A D_r^{-1} D_\lambda A^T)^{-1} A g. \tag{21}$$

Note that the step (14) used in the linearly convergent algorithm is equivalent to

$$d = -A^T (A D^{-2} A^T)^{-1} A g. \tag{22}$$

It is this similarity in form that yields a smooth transition from the linearly convergent algorithm in the previous section to the Newton method. Before presenting our new algorithm, it is worth noting that the matrix $A D_r^{-1} D_\lambda A^T$ is positive definite when $(r, w)$ is close enough to the solution $(r^*, w^*)$ and $\prod r_i \neq 0$.

**Lemma 1.** *Assume $(r^*, w^*)$ is a solution and the $\ell_1$ problem is both primal and dual non-degenerate. Then there exists a neighbourhood of $(r^*, w^*)$ such that when $(r, w)$ is within this neighbourhood, and no component of $r$ is equal to zero, the matrix $AD_r^{-1}D_\lambda A^T$ is positive definite.*

**Proof.** First we note that provided the diagonal of $D_r$ is zero-free,

$$AD_r^{-1}D_\lambda A^T = A_0 D_{r0}^{-1} D_{\lambda 0} A_0^T + A_1 D_{r1}^{-1} D_{\lambda 1} A_1^T,$$

where the zero subscript refers to the active set at $r^*$, i.e., $r_i^* = 0$, and the unity subscript refers to the inactive set at $r^*$. Hence, by definition, $D_{r1}^*$ is zero-free and since $D_{\lambda 1}^* = 0$, by complementary slackness (19), it follows that

$$A_1 (D_{r1}^*)^{-1} D_{\lambda 1}^* A_1^T = 0.$$

Primal and dual nondegeneracy yield $|\mathcal{A}(r^*)| = n$, $\mathcal{A}^c(r^*) = m - n$. Therefore

$$|\lambda_i^*| < 1 \quad \forall i \in \mathcal{A}(r^*).$$

Hence, for $i \in \mathcal{A}(r^*)$, there exists a neighbourhood around $(r^*, w^*)$ such that $((\mathrm{sgn}(r_i) - \lambda_i)/r_i)$ is positive and strictly bounded away from zero for all $i \in \mathcal{A}^*$ and $r_i \neq 0$. Since $A_0$ is full rank (by primal nondegeneracy) the result follows.    $\square$

Lemma 1 implies the surprising result that Newton directions (21), for the nonlinear system (19), are descent directions for $\psi(r)$ in a neighbourhood of the solution.

It is also the case that as the iterates converge, for each $i \in \mathcal{A}(r^*)$,

$$(D_\lambda D_r^{-1})_i \to +\infty.$$

This unbounded behaviour may appear to be cause for concern; however, in Sections 5 and 6 we prove that theoretically this is not a problem. Moreover, our numerical experience suggests that reliable numerical performance can be attained provided sufficient care is taken with respect to implementation details.

## 4. A hybrid method

The Newton and linear directions, introduced in the previous sections, have similar forms but differ in their definitions of diagonal matrix. Consider now a third choice:

$$D^2 = |D_r D_\theta^{-1}| \quad \text{where } D_r = \mathrm{diag}(r), \ D_\theta = \theta \, \mathrm{diag}(g) + (1 - \theta) D_\lambda, \qquad (23)$$

and $0 \leq \theta \leq 1$. Clearly if $\theta = 1$ then $D^2 = |D_r|$ and $d$ is the linear step (22) as in Algorithm 1. On the other hand, as $\theta \to 0$,

$$A^T D^{-2} A^T \to AD_r^{-1} D_\lambda A^T$$

(which, by Lemma 1, is positive definite in a neighbourhood of the solution), and so the Newton direction (21) is approached asymptotically (i.e., as $\theta \to 0$).

Our remaining task is to define $\theta \in [0, 1]$ so that $\theta \to 0$ if and only if $(r, \lambda) \to (r^*, \lambda^*)$. Our idea is to let $\theta$ encapsulate the optimality conditions. One possible choice is

$$\theta = \frac{\max\{\max(|r_i(g_i - \lambda_i)|/\psi(r^0)), \max\{\max\{|\lambda| - e, 0\}\}}{\gamma + \max\{\max(|r_i(g_i - \lambda_i)|/\psi(r^0)), \max\{\max\{|\lambda| - e, 0\}\}\}}, \tag{24}$$

where $0 < \gamma < 1$. Clearly $\theta$ is bounded above by one; assuming $Z(r - b) = 0$ then $\theta = 0$ if and only if $(r, \lambda) = (r^*, \lambda^*)$. The following result indicates that, unless $(r, \lambda)$ is optimal, $D_\theta$ has a zero-free diagonal. Consequently, when $\theta \neq 0$ and no component of $r$ is zero, $D^2$ is a *positive* diagonal matrix and so $-A^T(AD^{-2}A^T)^{-1}Ag$ is a descent direction for $\psi(r)$.

**Lemma 2.** *Suppose $0 < \gamma < 1$. Assume $\theta$ is defined by (24). Then $D_\theta$ satisfies*[3]

$$(1 - \gamma)\theta I \leq |D_\theta| \leq (2 - (1 - \gamma)\theta)I.$$

**Proof.** By definition of the matrix $D_\theta$, we can write

$$D_\theta = \text{diag}(g) - (1 - \theta)\,\text{diag}(Z^T w).$$

Consider the $i$th diagonal element of $D_\theta$. If $|\lambda_i = Z_i^T w| \leq 1$, it is clear that

$$(1 - \gamma)\theta \leq \theta \leq |g_i - (1 - \theta)\lambda_i| \leq 2 - \theta \leq 2 - (1 - \gamma)\theta,$$

since $0 < \gamma < 1$. If $|\lambda_i| > 1$ and $\theta$ is defined by (24), we have

$$\theta = \frac{\max\{\max(|r_i(g_i - \lambda_i)|/\psi(r^0)), \max\{\max\{|\lambda| - e, 0\}\}}{\gamma + \max\{\max(|r_i(g_i - \lambda_i)|/\psi(r^0)), \max\{\max\{|\lambda| - e, 0\}\}\}}.$$

Hence

$$(1 - \theta)\max\{\max(|r_i(g_i - \lambda_i)|/\psi(r^0))), \max\{\max\{|\lambda| - e, 0\}\} = \gamma\theta.$$

Therefore

$$(1 - \theta)(|\lambda_i| - 1) \leq \gamma\theta.$$

Thus

$$|\lambda_i| \leq 1 + \frac{\gamma\theta}{1 - \theta} = \frac{1 - (1 - \gamma)\theta}{1 - \theta},$$

which yields

$$(1 - \theta)|\lambda_i| \leq 1 - (1 - \gamma)\theta.$$

Therefore

$$(1 - \gamma)\theta \leq |g_i - (1 - \theta)\lambda_i| \leq 2 - (1 - \gamma)\theta,$$

and the result is established.  □

[3] Recall: If $M = (m_{ij})$ is a matrix, $|M|$ is the matrix obtained by replacing $m_{ij}$ with $|m_{ij}|$ for all $i, j$.

Therefore, $D_\theta$ is nonsingular when $\theta \neq 0$; consequently, the definition of $\theta$ leads to a descent direction. Moreover, $\theta$ induces a smooth transition from the linear step to the Newton step.

### 4.1. Computation of the hybrid step

We consider three different methods for computing the hybrid step.

#### 4.1.1. Full space implementation

The hybrid step can be viewed as an approximate Newton step. Assuming $(r, w)$ is our current guess, where $Zr = Zb$, the approximate Newton step, in analogy with (20), can be expressed as

$$\begin{bmatrix} \text{diag}(g)|D_\theta| & -D_r Z^T \\ Z & 0 \end{bmatrix} \begin{bmatrix} d \\ d_w \end{bmatrix} = \begin{bmatrix} -D_r(g - Z^T w) \\ 0 \end{bmatrix}. \tag{25}$$

Note that $d$ is the hybrid direction defined above. The dimension of this system is somewhat daunting, $(2m - n) \times (2m - n)$, and can be reduced by implicitly satisfying $Zd = 0$:

$$\text{Solve}[\text{diag}(g)|D_\theta|A^T, -D_r Z^T]\begin{bmatrix} d_x \\ d_w \end{bmatrix} = -D_r(g - Z^T w) \tag{26}$$

and then

$$d \leftarrow -A^T d_x. \tag{27}$$

After solving this full-space system and obtaining $(d_x, d_w)$, the dual variables, $\lambda$, can be updated $\lambda^+ \leftarrow \lambda + D_r Z^T d_w$; the primal variables, $r$, are updated using our line search procedure along direction $d$, $r^+ \leftarrow r + \alpha d$.

Note that as $\theta \to 0$, i.e., as we converge to the solution, the matrix in (26) approaches $(D_{\lambda*}A^T, -D_{r*}Z^T)$ which, under nondegeneracy assumptions, is full rank (and bounded). This is the advantage of using the full-space implementation: the limit matrix is well-behaved. The major disadvantage is that the system is still large, $m \times m$, and requires the computation of $Z$.

#### 4.1.2. Reduced space implementations

It is possible to reduce the dimension of the system to be solved (and still compute the same correction) by realizing that what is needed is either $d_x$ and $Z^T d_w$ or $w$ and $A^T d_x$ but not both $d_x$ and $d_w$. More generally, consider the following simple observations.

**Lemma 3.** *Let* $M = [M_1, M_2]$ *be full rank,* $U = \text{null}(M_2^T)$ [4], $V = \text{null}(M_1^T)$. *Then the matrix* $[U, V]$ *has full rank.* □

**Theorem 4.** *Assume* $M = [M_1, M_2]$ *has full rank,* $U = \text{null}(M_2^T)$, $V = \text{null}(M_1^T)$. *Then the system*

$$Mv = b$$

---

[4] Recall: the notation $B = \text{null}(C)$ indicates that $B$ is a matrix whose columns form a basis for the null space of $C$.

*is equivalent to*

$$\begin{bmatrix} U^T M_1 & 0 \\ 0 & V^T M_2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} U^T b \\ V^T b \end{bmatrix}.$$

**Proof.** From Lemma 3, $[U, V]$ has full rank. Hence $Mv = b$ is equivalent to

$$[U, V]^T Mv = [U, V]^T b. \qquad \square$$

**Corollary 5.** *To obtain* $[v_1, M_2 v_2]$ *such that*

$$[M_1, M_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = b,$$

*the following algorithm suffices*:
1. *Determine* $U = \mathrm{null}(M_2^T)$;
2. *Solve* $U^T M_1 v_1 = U^T b$;
3. $M_2 v_2 \leftarrow b - M_1 v_1$. $\quad \square$

**Range space implementation.** If we define $M_1 = \mathrm{diag}(g)|D_\theta|A^T$ and $M_2 = -D_r Z^T$ then we can solve (26) using Corollary 5:

$$\text{Solve } D^{-1}A^T d_x \stackrel{\text{l.s.}}{=} Dg \tag{28}$$

and

$$d \leftarrow -A^T d_x, \quad \lambda^+ \leftarrow g + D^{-2}d. \tag{29}$$

The main computational work is the solution of the linear least squares problem of order $m \times n$. Note that $Z$ is not required; this approach is particularly attractive when $n$ is small or $A$ is sparse. The primal variables are updated via a linesearch on $\psi(r)$: $r^+ \leftarrow r + \alpha d$.

**Null space implementation.** A null space implementation is obtained from Corollary 5 by letting $M_1 = -D_r Z^T$ and $M_2 = \mathrm{diag}(g)|D_\theta|A^T$. This implementation first requires the computation of $Z$ and then

$$\text{Solve } DZ^T w^+ \stackrel{\text{l.s.}}{=} Dg \tag{30}$$

and

$$A^T d_x \stackrel{\text{def}}{=} d \leftarrow -D^2(g - Z^T w^+). \tag{31}$$

In this case an $m \times (m - n)$ least squares system is solved at each iteration; hence, this approach is attractive when $n$ is large. Note that $Z$ is computed only once; moreover, it is often possible to compute a sparse $Z$ given that $A$ is sparse (e.g., [6, 7]).

We conclude this section with a presentation of the simple hybrid algorithm. Bear in mind that there are a number of numerical concerns, such as unequal row scaling and stopping criteria, that must be taken care of before this algorithm can be reliably used. These issues are briefly discussed in Section 7.

When applying the linesearch procedure, instead of using a constant $\tau$, $\tau^k = \max\{\tau, 1 - \theta^k\}$ is used in order to obtain final quadratic convergence. Thus, we have

$$\alpha^k = \alpha_{\#}^k + \tau^k(\alpha_*^k - \alpha_{\#}^k), \quad \tau^k = \max\{\tau, 1 - \theta^k\}, \tag{32}$$

where $\alpha_*^k$ and $\alpha_{\#}^*$ are as defined in (16) and (18). Notice that the subscripts $\#$ and $*$ depend on $k$ in general.

Let $r^0$ be an initial differentiable point satisfying $Z(r^0 - b) = 0$; $k \leftarrow 0$; Compute an initial point $w^0$.

**Algorithm 2.**
  *Step 1.* Compute $\theta^k$ from (24). Define $D^k$ from (23) and $g^k \leftarrow \text{sgn}(r^k)$.
  *Step 2.* Compute $d^k$ and $w^{k+1}$ using one of the three methods above.
  *Step 3.* Do a line search on the piecewise linear function $\psi(\alpha)$ (as described in Section 2 and using (32)) to determine $\alpha^k$,

$$r^{k+1} \leftarrow r^k + \alpha^k d^k, \quad k \leftarrow k + 1.$$

## 5. Global convergence

In this section, we establish global convergence of the linear and hybrid methods, Algorithms 1 and 2. We make the following global assumption: *The $n \times m$ matrix $A$ has full row rank $n$.* Let $P^k$ be the orthogonal projector onto $\text{null}(ZD^k)$, i.e.,

$$P^k = I - D^k Z^T (Z(D^k)^2 Z^T)^{-1} Z D^k.$$

Both algorithms use a diagonal matrix $D^k \stackrel{\text{def}}{=} \text{diag}(s_i^k)$ to yield the search direction $d^k$:

$$d^k = -D^k P^k D^k g^k$$

$$= -(D^k)^2 (g^k - Z^T w^{k+1}), \tag{33}$$

where $w^{k+1}$ is the least squares solution to

$$D^k Z^T w^{k+1} \stackrel{\text{l.s.}}{=} D^k g^k,$$

and

$$s_i^k = \begin{cases} |r_i^k|^{1/2} & \text{for Algorithm 1,} \\ |r_i^k|^{1/2} |\theta^k g_i^k + (1 - \theta^k)(g_i^k - Z_i^T w^k)|^{-1/2} & \text{for Algorithm 2.} \end{cases}$$

The first major step in the convergence proof is to show that $\|P^k D^k g^k\| \to 0$. This is established in Lemma 10 after several preliminary results.

**Lemma 6.** *Assume $\{d^k\}$ is defined by Algorithm 1 or Algorithm 2. Then*

$$\lim_{k \to \infty} \alpha_\diamond^k \|P^k D^k g^k\|_2^2 = 0, \tag{34}$$

*where $\alpha_\diamond^k = \min\{\alpha_i^k : \alpha_i^k > 0\}$ corresponds to the first positive breakpoint in the direction $d^k$ (17). Moreover,*

$$\lim_{k \to \infty} \sum_{g_i^k \neq g_i^{k+1}} (\alpha_i^k - \alpha_\diamond^k) g_i^k d_i^k = 0.$$

**Proof.** Since $\|r^k\|_1$ is monotonically decreasing and bounded below, $\|r^k\|_1$ converges; therefore,

$$\lim_{k \to \infty} (\|r^k\|_1 - \|r^{k+1}\|_1) = 0. \tag{35}$$

From (33),

$$g^{k^T} d^k = -\|P^k D^k g^k\|_2^2;$$

therefore, since $r^{k+1} = r^k + \alpha^k d^k$,

$$\|r^k\|_1 - \|r^{k+1}\|_1$$
$$= g^{k^T} r^k - g^{k+1^T} r^k - \alpha^k g^{k+1^T} d^k$$
$$= (g^k - g^{k+1})^T r^k + \alpha^k (g^k - g^{k+1})^T d^k + \alpha^k \|P^k D^k g^k\|_2^2.$$

However, if $g_i^k \neq g_i^{k+1}$ then $g_i^k = -g_i^{k+1}$; hence,

$$\|r^k\|_1 - \|r^{k+1}\|_1$$
$$= \sum_{g_i^k \neq g_i^{k+1}} 2(g_i^k r_i^k + \alpha^k g_i^k d_i^k) - (\alpha^k - \alpha_\diamond^k) g^{k^T} d^k + \alpha_\diamond^k \|P^k D^k g^k\|_2^2. \tag{36}$$

For $g_i^k \neq g_i^{k+1}$, recall $\alpha_i^k$ is the breakpoint for equation $i$ in direction $d_i^k$: $\alpha_i^k = -r_i^k / d_i^k$. Therefore, for $g_i^k \neq g_i^{k+1}$,

$$r_i^k + \alpha^k d_i^k = (\alpha^k - \alpha_i^k) d_i^k. \tag{37}$$

Consequently, substituting (37) into (36),

$$\|r^k\|_1 - \|r^{k+1}\|_1$$

$$= \sum_{g_i^k \neq g_i^{k+1}} 2(\alpha^k - \alpha_i^k) g_i^k d_i^k - (\alpha^k - \alpha_\diamond^k) g^{k^T} d^k + \alpha_\diamond^k \|P^k D^k g^k\|_2^2$$

$$= -(\alpha^k - \alpha_\diamond^k) \left( g^{k^T} d^k - 2 \sum_{g_i^k \neq g_i^{k+1}} g_i^k d_i^k \right)$$

$$\quad - 2 \sum_{g_i^k \neq g_i^{k+1}} (\alpha_i^k - \alpha_\diamond^k) g_i^k d_i^k + \alpha_\diamond^k \|P^k D^k g^k\|_2^2.$$

Next we establish that each term in the expression above is positive. The linesearch stops before optimality is reached in the direction $d^k$; hence, $g^{k+1^T} d^k < 0$. But

$g^{k+1^T}d^k = g^{k^T}d^k - 2\sum_{g_i^k \neq g_i^{k+1}} g_i^k d_i^k$ and so it follows that

$$\left( g^{k^T}d^k - 2 \sum_{g_i^k \neq g_i^{k+1}} g_i^k d_i^k \right) < 0,$$

and therefore the first term is nonnegative. But, for $g_i^k \neq g_i^{k+1}$, $\alpha_\diamond^k \leq \alpha_i^k$ and so

$$\sum_{g_i^k \neq g_i^{k+1}} (\alpha_i^k - \alpha_\diamond^k) g_i^k d_i^k \leq 0,$$

which yields a nonnegative second term. Thus

$$0 \leq \alpha_\diamond^k \| P^k D^k g^k \|_2^2 \leq \| r^k \|_1 - \| r^{k+1} \|_1,$$

and therefore, applying (35),

$$\lim_{k \to \infty} \alpha_\diamond^k \| P^k D^k g^k \|_2^2 = 0.$$

Similarly,

$$\lim_{k \to \infty} \sum_{g_i^k \neq g_i^{k+1}} (\alpha_i^k - \alpha_\diamond^k) g_i^k d_i^k = 0. \qquad \square$$

Recall: We use the subscripts 1 and 0 to denote the components corresponding to the non-activities and activities respectively. So, for example, $Z_1$ is the submatrix of $Z$ consisting of columns of $Z$ corresponding to the nonzero components of $r$.

**Lemma 7.** *Assume primal nondegeneracy; let $r$ be an arbitrary point. Then $Z_1$ is of full row rank. Moreover, if we assume that $(r, \lambda = Z^T w)$ satisfy complementary slackness (9), then $\lambda$ is a dual basic solution. In addition, if we assume dual nondegeneracy, then $|\mathscr{A}(r)| = n$.*

**Proof.** Assume that there exist $y_1$ such that

$$Z_1^T y_1 = 0. \tag{38}$$

From $AZ^T = 0$, we have

$$A_0 Z_0^T + A_1 Z_1^T = 0.$$

Hence

$$A_0 Z_0^T y_1 + A_1 Z_1^T y_1 = 0.$$

From (38), we have $A_0 Z_0^T y_1 = 0$; from the primal nondegeneracy assumption, $A_0$ is of full column rank. Therefore $Z_0^T y_1 = 0$. Thus

$$Z^T y_1 = 0.$$

But $Z^T$ has full column rank; therefore, $y_1 = 0$ and we have established that $Z_1$ is of full row rank. Now, let us further assume complementary slackness (9); hence,

$$g_1 = Z_1^T w$$

where $\lambda = Z^T w$. But $AZ^T = 0$ which yields $A\lambda = 0$ and so we obtain dual feasibility (i.e., $\lambda$ is a dual basic point). Finally, dual nondegeneracy and full row rank of $Z_1$ result in $|\mathscr{A}| = n$. $\quad \square$

**Lemma 8.** *Under primal and dual nondegeneracy assumptions, $J$ is nonsingular at any point $(r, \lambda = Z^\mathrm{T} w)$, where*

$$J = \begin{bmatrix} \mathrm{diag}(g)|D_\theta| & -D_r Z^\mathrm{T} \\ Z & 0 \end{bmatrix}.$$

**Proof.** For any $y = [y_0^\mathrm{T}, y_1^\mathrm{T}]^\mathrm{T}$ such that $Jy = 0$, we have

$$\mathrm{diag}(g)|D_\theta|y_0 - D_r Z^\mathrm{T} y_1 = 0, \tag{39}$$

$$Zy_0 = 0. \tag{40}$$

Hence $y_0 = A^\mathrm{T} z$, for some $z$. Thus

$$\mathrm{diag}(g)|D_\theta|A^\mathrm{T} z - D_r Z^\mathrm{T} y_1 = 0. \tag{41}$$

First assume $\theta \neq 0$. Then, by Lemma 2, $D_\theta$ is nonsingular. Multiplying the above by the full rank matrix $Z \, \mathrm{diag}(g)|D_\theta^{-1}|$ yields

$$Z \, \mathrm{diag}(g)|D_\theta|^{-1} D_r Z^\mathrm{T} y_1 = 0,$$

which implies

$$Z_1 \, \mathrm{diag}(g_1)|D_{\theta 1}|^{-1} D_{r1} Z_1^\mathrm{T} y_1 = 0.$$

From the primal nondegeneracy assumption and Lemma 7, it is clear that $y_1 = 0$. From (41), we have $z = 0$. Thus, $y_0 = 0$. Hence $J$ is nonsingular.

Second, assume $\theta = 0$. Complementary slackness and (41) yield

$$|D_{\lambda 0}|A_0^\mathrm{T} z = 0 \quad \text{and} \quad D_{r1} Z_1^\mathrm{T} y_1 = 0.$$

From primal and dual nondegeneracy and Lemma 7, $|D_{\lambda 0}|A_0^\mathrm{T}$ and $D_{r1} Z_1^\mathrm{T}$ are nonsingular, hence

$$y_1 = 0, \quad z = 0.$$

Therefore $y_0 = 0$ and so $J$ is nonsingular. $\quad\square$

The next result establishes that, in Algorithm 2, the multipliers are bounded in size.

**Lemma 9.** *Assume that an $\ell_1$ problem is both primal and dual nondegenerate and $\{\lambda^k = Z^\mathrm{T} w^k\}$ is obtained by Algorithm 2. Then $\{\lambda^k\}$ and $\{w^k\}$ are bounded: there exists $M > 0$ such that*

$$\|\lambda^k\| \leq M \quad \text{and} \quad \|w^k\| \leq M.$$

**Proof.** Using (26), $(d_x^k, w^{k+1})$ is the solution to the following linear system:

$$[\mathrm{diag}(g^k)|D_\theta^k|A^\mathrm{T}, -D_r^k Z]\begin{bmatrix} d_x^k \\ w^{k+1} \end{bmatrix} = -D_r^k g^k. \tag{42}$$

Following Lemma 2, $\{D_\theta^k\}$ is always bounded. Thus the coefficient matrix is always bounded. Moreover, following Lemma 2, the coefficient matrix in (42) is nonsingular at any limit point (this is easily proved along the lines of the first half of the proof of Lemma 8). It follows immediately that $\{\lambda^{k+1}\}$ and $\{w^{k+1}\}$ are bounded. $\quad\square$

We can now state the first major result.

**Lemma 10.** *Assume $\{d^k\}$ is defined by Algorithm 1 or 2; assume primal and dual nondegeneracy. Then*

$$\lim_{k\to\infty} \|P^k D^k g^k\|_2 = 0, \qquad \lim_{k\to\infty} d^k = 0 \quad and \quad \lim_{k\to\infty} D_r^k(g^k - \lambda^{k+1}) = 0.$$

**Proof.** Using Lemma 6, we know that

$$\lim_{k\to\infty} \alpha_\diamond^k \|P^k D^k g^k\|_2^2 = 0. \tag{43}$$

(Recall: $\alpha_\diamond^k$ is the step to the first breakpoint from $r^k$ in the direction $d^k$.) From Lemma 9, there exists $M > 0$, such that

$$\|(g^k - Z^T w^{k+1})\| \le M. \tag{44}$$

Now assume there exists a subsequence, which we denote with the use of an under-score, satisfying,

$$\{-\|\underline{P}^k\underline{D}^k\underline{g}^k\|_2^2\} \to c_1 < 0. \tag{45}$$

We now prove that the corresponding subsequence of first breakpoints, $\{\underline{\alpha}_\diamond^k\}$ is bounded away from zero, i.e., $\underline{\alpha}_\diamond^k > c_4$, for some $c_4 > 0$, and this will lead to the obvious contradiction.

From the definition of $\theta^k$ and (45), there exists $c_2 \in (0, 1)$ such that $\underline{\theta}^k > c_2$. To see this assume the contrary: a subsequence $\{\hat{\theta}^k\}$ of $\{\underline{\theta}^k\}$ converges to zero. Then

$$\hat{D}_r^k(\hat{g}^k - Z^T \hat{w}^k) \to 0,$$

which implies, by (25), that $\hat{d}^k \to 0$. But $\hat{d}^k = \hat{D}^k \hat{P}^k \hat{D}^k \hat{g}^k$ and so (45) is contradicted.

Recall, from Lemma 2, that for any $i$,

$$|\underline{\theta}^k \underline{g}_i^k + (1 - \underline{\theta}^k)(\underline{g}_i^k - \underline{\lambda}_i^k)| \ge (1 - \gamma)\underline{\theta}^k \ge (1 - \gamma)c^2 > 0. \tag{46}$$

However,

$$\underline{\alpha}_\diamond^k = -\frac{\underline{r}_\diamond^k}{\underline{d}_\diamond^k} = \frac{|\underline{\theta}^k \underline{g}_\diamond^k + (1 - \underline{\theta}^k)(\underline{g}_\diamond^k - \underline{\lambda}_\diamond^k)|}{(|\underline{g}_\diamond^k - Z_\diamond^T \underline{w}^{k+1}|)} \quad \text{(using (31))}$$

$$\ge \frac{c_3}{M} \quad \text{(using (44) and (46))}$$

$$\overset{\text{def}}{=} c_4 > 0.$$

Therefore, by Lemma 6, $\|P^k D^k g^k\| \to 0$.

Using (33), we have

$$\lim_{k \to \infty} D^k(g^k - \lambda^{k+1}) = 0 \quad \text{and} \quad \lim_{k \to \infty} d^k = 0.$$

Hence $\lim_{k \to \infty} D_r^k(g^k - \lambda^{k+1}) = 0$ follows from the boundedness of $D_\theta^k$ (Lemma 2). $\square$

The next result says that if the primal variables converge then so do the duals. The point of convergence, $(\bar{r}, \bar{\lambda} = Z^T \bar{w})$, satisfies all optimality conditions except possibly for $|\bar{\lambda}_0| < e$.

**Lemma 11.** *Suppose $\{r^k\}$ and $\{w^k\}$ are obtained by Algorithm 1 or 2 and assume $\{r^k\} \to \bar{r}$, a limit point. Further, assume primal and dual nondegeneracy. Then $\{w^k\}$ converges; i.e., there exists a point $\bar{w}$ such that*

$$\lim_{k \to \infty} w^k = \bar{w}, \tag{47}$$

*where $w^{k+1}$ is defined from*

$$D^k Z^T w^{k+1} \overset{\text{l.s.}}{=} D^k g^k. \tag{48}$$

*Moreover,*

$$|\mathscr{A}(\bar{r})| = n, \quad A\bar{\lambda} = 0, \quad |\bar{\lambda}_1| = e, \quad \text{where } \bar{\lambda} = Z^T \bar{w}. \tag{49}$$

**Proof.** Using Lemma 10, we know that

$$D_r^k(g^k - Z^T w^{k+1}) \to 0.$$

Without loss of generality, assume

$$D_{\bar{r}} = \begin{bmatrix} D_{\bar{r}1} & 0 \\ 0 & D_{\bar{r}0} \end{bmatrix}, \qquad Z^T = \begin{bmatrix} Z_1^T \\ Z_0^T \end{bmatrix},$$

where $D_{\bar{r}1}$ is a diagonal matrix with zero-free diagonal, and $D_{\bar{r}0}$ is the zero matrix. But $D^k = |D_r^k[D_\theta^k]^{-1}|$ and so

$$\lim_{k \to \infty} [D_{\theta_1}^k]^{-1}(g_1^k - Z_1^T w^{k+1}) = 0.$$

But, Lemma 2 established the boundedness of $|D_\theta^k|$; hence, there exists $\bar{M} > 0$ such that $|[D_{\theta_1}^k]^{-1}| > \bar{M}$. Therefore, for any convergent subsequence of $w^k$, we have

$$\bar{g}_1 = Z_1^T \bar{w}, \tag{50}$$

where $\bar{w}$ is the limit point of the subsequence. (Note: $\bar{g} = \text{sgn}(\bar{r})$.) By primal nondegeneracy and Lemma 7, $Z_1^T$ is of full column rank; therefore, the solution to

(50) is unique. Since the limit point of any convergent subsequence of $\{w^k\}$ is a solution of (50), and since by Lemma 9 $\{w^k\}$ is bounded, $\{w^k\}$ converges to $\bar{w}$.

Finally, to show (49) note that

$$[A_0, A_1]\bar{\lambda} = 0, \quad \text{where } \bar{\lambda}^T \stackrel{\text{def}}{=} (\bar{w}^T Z_0, \bar{w}^T Z_1)$$

since $AZ^T = 0$. But, by (50), the magnitude of each component of $\bar{\lambda}_1$ is unity; therefore, by Lemma 7, $A_0$ is order $n$. $\square$

Theorem 14 below is the next major result: it says that if the primal variables converge, they converge to the optimal point. Before proving this theorem two technical results are provided, the second of which we state without proof (the proof is relatively straightforward and unenlightening).

**Lemma 12.** *Assume $\{d^k\}$ is obtained by either Algorithm 1 or Algorithm 2. Suppose there exist vectors $u$ and $v$ such that*

$$A_1 u = -A_0 v.$$

*Then, if $g_1^k = u$,*

$$g_1^{k^T} d_1^k = -v^T d_0^k. \tag{51}$$

**Proof.** From the definition of the direction $d^k$ (27), we know there exists $d_x^k$ such that

$$d^k = -A^T d_x^k.$$

Thus

$$\begin{aligned}
g_1^{k^T} d_1^k &= g_1^T A_1^T d_x^k \quad \text{(since } g_1^k = g_1) \\
&= -v^T A_0^T d_x^k \quad \text{(from (51))} \\
&= -v^T d_0^k. \quad \square
\end{aligned}$$

**Lemma 13.** *Suppose $\theta \in (0, 1)$ and $-1 < \beta < 1 < \eta$. Then,*

$$\frac{1 + (1 - \theta)\beta}{1 + \beta} > \frac{1 + (1 - \theta)\eta}{1 + \eta}. \quad \square$$

**Theorem 14.** *Assume $\{r^k\}$ is obtained from Algorithm 1 or Algorithm 2. Assume primal and dual nondegeneracy. If the sequence $\{r^k\}$ converges to a point $r^*$, then $r^*$ is optimal.*

**Proof.** From Lemma 11, $|\mathcal{A}(r^*)| = n$, $\{w^k\} \to w^*$ and if $\lambda^*$ is defined $\lambda^* = Z^T w^*$ then $A\lambda^* = 0$ with $|\lambda_1^*| = e$. Therefore, to establish optimality we need only show that $\|\lambda_0^*\|_\infty \le 1$.

Assume the contrary, i.e., for some $i$ with $r_i^* = 0$, $\lambda_i^* > 1$. (Note that this implies: $\lim \inf_{k \to \infty} \theta^k > 0$). Then, there exists $k_1$ such that for $k > k_1$,

$$|\lambda_i^k| > 1 \quad \text{for all } i \text{ such that } |\lambda_i^*| > 1,$$

and

$$\lambda_j^k r_j^k > 0 \quad \text{for all } j \text{ such that } |\lambda_j^*| = 1.$$

If there exists $|\lambda_i^*| > 1$ and $\lambda_i^{k+1} r_i^k > 0$ where $k > k_1$, from

$$r_i^{k+1} = r_i^k - \alpha^k s_i^{k^2} (g_i^k - \lambda_i^{k+1})$$

it follows that $|r_i^{k+1}| > |r_i^k| > 0$, for all $k > k_1$. (To see this note that $\text{sgn}(r_i^k) = -\text{sgn}(\alpha^k s_i^{k^2} (g_i^k - \lambda_i^{k+1}))$). Hence, we see immediately that $\{r_i^k\} \nrightarrow 0$ which contradicts that $\{r_i^k\} \to r_i^* = 0$.

Now assume there does not exist $k_2 > k_1$ such that $r_i^{k_2} \lambda_i^{k_2+1} > 0$, for any $|\lambda_i^*| > 1$. In Algorithm 1, the breakpoint corresponding to equation $l$ is given by

$$\alpha_l^k = \frac{1}{g_l^k (g_l^k - \lambda_l^{k+1})},$$

whereas for Algorithm 2,

$$\alpha_l^k = \frac{r_l^k |g_l^k - (1 - \theta^k) \lambda_l^k|}{|r_l^k|(g_l^k - \lambda_l^{k+1})} = \frac{|g_l^k - (1 - \theta^k) \lambda_l^k|}{g_l^k (g_l^k - \lambda_l^{k+1})}.$$

This latter expression can be simplified if we consider two separate cases. For $|\lambda_j^*| < 1$, we have

$$\alpha_j^k = \frac{g_j^k (g_j^k - (1 - \theta^k) \lambda_j^k)}{g_j^k (g_j^k - \lambda_j^{k+1})} = \frac{1 - (1 - \theta^k) g_j^k \lambda_j^k}{1 - g_j^k \lambda_j^{k+1}}.$$

Hence, for $l \in \mathscr{A}^c$, $\alpha_l^k \to \infty$. For $|\lambda_i^*| > 1$, since $g_i^k = -\text{sign}(\lambda_i^{k+1}) = -\text{sign}(\lambda_i^*)$, we have

$$\alpha_i^k = \frac{|g_i^k - (1 - (1 - \theta^k) \lambda_i^k|}{(1 - \lambda_i^{k+1} g_i^k)} = \frac{1 + (1 - \theta^k) |\lambda_i^k|}{1 + |\lambda_i^{k+1}|}.$$

Using Lemma 13, we know that

$$\frac{1 + (1 - \theta^*) |\lambda_i^*|}{1 + |\lambda_i^*|} < \frac{1 + (1 - \theta^*) \lambda_j^*}{1 + \lambda_j^*}$$

and

$$\frac{1 + (1 - \theta^*) |\lambda_i^*|}{1 + |\lambda_i^*|} < \frac{1 - (1 - \theta^*) \lambda_j^*}{1 - \lambda_j^*},$$

since $\theta^* \neq 0$.

Hence, when $k$ is sufficiently large, we obtain, for both algorithms,

$$\alpha_i^k < \alpha_j^k \quad \text{for any } j \neq i, |\lambda_j^*| < 1. \tag{52}$$

Now let $g^+$ denote the new gradient after passing all the breakpoints with $|\lambda_i^*| > 1$. Then

$$g^{+^T}d^k = -\sum_{|\lambda_i^*|>1, i \in \mathscr{A}} g_i^k d_i^k + \sum_{|\lambda_j^*|<, j \in \mathscr{A}} g_j^{k^T} d_j^k + g_1^{k^T} d_1^k$$

$$= -\sum_{|\lambda_i^*|>1, i \in \mathscr{A}} g_i^k d_i^k + \sum_{|\lambda_j^*|<1, j \in \mathscr{A}} g_j^{k^T} d_j^k - \lambda_0^{*^T} d_0^k \quad \text{(using Lemma 12)}$$

$$= \sum_{|\lambda_i^*|>1, i \in \mathscr{A}} (-\lambda_i^* - g_i^k) d_i^k + \sum_{|\lambda_j^*|<1, j \in \mathscr{A}} (g_j^k - \lambda_j^*) d_j^k.$$

But

$$-\lambda_i^* d_i^k < g_i^k d_i^k < 0 \quad \text{for } |\lambda_i^*| > 1, r_i^* = 0,$$

and

$$-|\lambda_j^* d_j^k| > g_j^k d_j^k \quad \text{for } |\lambda_j^*| < 1, r_j^* = 0.$$

Hence

$$g^{+^T} d^k < 0.$$

Therefore, the linesearch will yield $r_i^{k+1} \lambda_i^{k+2} > 0$, a contradiction. $\quad\square$

**Lemma 15.** *Assume* (1) *is both primal nondegenerate and dual nondegenerate and* $\{(r^k, \lambda^k)\}$ *is obtained by either Algorithm 1 or Algorithm 2. Then* $\{r^k\}$ *and* $\{w^k : Z^T w^k = \lambda^k\}$ *have a finite number of limit of points. Moreover, for any convergent subsequence* $\{(r^k, \lambda^k)\}$ *with* $\{\alpha^k \|d^k\|\}$ *having a positive limit, we have*

$$\lim_{k \to \infty} D_r^k (g^k - \lambda^k) = 0, \quad \text{where } \lambda^k = Z^T w^k.$$

**Proof.** It is clear that $\{r^k\}$ is bounded. From Lemma 9, $\{\lambda^k\}$ is bounded. Moreover, from Lemma 10,

$$\lim_{k \to \infty} D_r^k (g^k - \lambda^{k+1}) = 0. \tag{53}$$

The nondegeneracy assumptions imply that the set of points $\mathscr{V}$ satisfying complementary slackness, i.e., $\mathscr{V} = \{(r, w) : D_r(g - \lambda) = 0, \lambda = Z^T w\}$, is a finite set. Hence, any limit point of $\{(r^k, Z^T w^{k+1})\}$ belongs to $\mathscr{V}$. Thus, $\{r^k\}$ and $\{\lambda^{k+1}\}$ have finite number of limit points. This implies the sequence $\{(r^k, \lambda^k)\}$ has a finite number of limit points.

Now assume a subsequence $\{(r^k, \lambda^k)\}$ converging to $(\bar{r}, \bar{\lambda})$ with $\{\alpha^k \|d^k\|\}$ having a positive limit. (To simplify the notation, we use subscripts $k$ to denote a subsequence as well but mention it is a subsequence explicitly.) Following primal nondegeneracy assumption, the corresponding subsequence has $\alpha^k \|d_0^k\| \geq c_0 > 0$ for some $c_0 > 0$ (otherwise $\alpha^k d_0^k = -A_0^T(\alpha^k d_x^k)$ converging to zero leads to $\alpha^k d_x^k$ converging to zero and thus $\alpha^k d^k = -\alpha^k A^T d_x^k$ converging to zero).

Following Lemma 6, $\lim_{k\to\infty} \alpha_\diamond^k \|P^k D^k g^k\|_2^2 = 0$. The nondegeneracy assumptions imply that $\lim_{k\to\infty} \alpha_\diamond^k d_0^k = 0$ and subsequently $\lim_{k\to\infty} \alpha_\diamond^k d_x^k = 0$. Hence, $\lim_{k\to\infty} \alpha_\diamond^k d_\diamond^k = 0$. Immediately, we have $\lim_{k\to\infty} r_\diamond^k = 0$. Using (15) and (33),

$$\alpha_i^k = -\frac{r_i^k}{d_i^k} = -\frac{\delta_i^k}{g_i^k - \lambda_i^{k+1}}, \quad 1 \leq i \leq m. \tag{54}$$

Moreover, under the nondegeneracy assumptions, the corresponding sequence $\{|g_\diamond^k - \lambda_\diamond^{k+1}|\}$ must be bounded away from zero for $k$ sufficiently large. Hence, in light of Lemma 2, $\{\alpha_\diamond^k\}$ is bounded. Since $\lim_{k\to\infty} d^k = 0$ (Lemma 10), it follows that

$$\lim_{k\to\infty} \sum_{g_i^k \neq g_i^{k+1}} \alpha_\diamond^k g_i^k d_i^k = 0.$$

Using Lemma 6,

$$\lim_{k\to\infty} \sum_{g_i^k \neq g_i^{k+1}} (\alpha_i^k - \alpha_\diamond^k) g_i^k d_i^k = 0.$$

Hence

$$\lim_{k\to\infty} \sum_{g_i^k \neq g_i^{k+1}} \alpha_i^k g_i^k d_i^k = 0.$$

For $\alpha_\sharp^k > 0$, $g_\sharp^k \neq g_\sharp^{k+1}$. (Recall: $\alpha_\sharp^k$ is the largest breakpoint smaller than the optimal breakpoint $\alpha_*^k$.) Thus

$$\lim_{k\to\infty} r_\sharp^k = \lim_{k\to\infty} \alpha_\sharp^k d_\sharp^k = 0.$$

Therefore, we have $r_1^{k+1} r_1^k > 0$ for $k$ sufficiently large. Applying the same arguments of the boundedness of $\{\alpha_\diamond^k\}$, we conclude $\{\alpha_\sharp^k\}$ is also bounded. Again from $\lim_{k\to\infty} d^k = 0$,

$$\lim_{k\to\infty} \alpha_\sharp^k d_1^k = 0.$$

Using (18) or (32), we have

$$r_1^{k+1} = r_1^k + \rho^k \alpha_\sharp^k d_1^k + \tau^k \alpha_*^k d_1^k.$$

Hence

$$\min(|r_1^{k+1}|) = \min(|r_1^k + \rho^k \alpha_\sharp^k d_1^k + \tau^k \alpha_*^k d_1^k|)$$

$$\geq \min(\tfrac{1}{2}(g_1^k r_1^k + \tau^k \alpha_*^k g_1^k d_1^k)) \quad \text{(for $k$ sufficiently large)}.$$

Using (54) and the definition of $\alpha_*^k$ (16), it is easy to verify

$$\alpha_*^k g_1^k d_1^k \geq \alpha_1^k g_1^k d_1^k.$$

Hence, for $k$ sufficiently large,

$$\min(|r_1^{k+1}|) \geq \min(\tfrac{1}{2}|r_1^k + \tau^k \alpha_1^k d_1^k|) = \min(\tfrac{1}{2}\rho^k |r_1^k|).$$

We prove by contradiction that the corresponding $D_r^k(g^k - \lambda^k)$ converges to zero. Assume otherwise. Then there exists $c > 0$ and a subsequence with $\|D_r^k(g^k - \lambda^k)\| > c$

for $k$ sufficiently large. Let $\rho^k = 1 - \tau^k$ where either $\tau^k = \tau$ (for Algorithm 1) or $\tau^k = \max(\tau, 1 - \theta^k)$ as in (32). Hence there exists $c_1 > 0$ with $\rho^k > c_1$. Therefore

$$\min(|r_1^{k+1}|) \geq c_2 > 0. \tag{55}$$

Since $\|r_0^{k+1} - r_0^k\| = \alpha^k \|d_0^k\| \geq c_0 > 0$, we conclude that $\max(|r_0^{k+1}|)$ is bounded away from zero. From (53) and nondegeneracy assumptions, zero is a limit point of the corresponding subsequence $\min(|r_1^{k+1}|)$. But this contradicts (55). In conclusion, the corresponding $D_r^k(g^k - \lambda^k)$ converges to zero. $\square$

**Theorem 16.** *If an $\ell_1$ problem is both primal nondegenerate and dual nondegenerate, then the sequence $\{(r^k, \lambda^k)\}$, generated by either Algorithm 1 or 2, is convergent.*

**Proof.** From Lemma 15, $\{r^k\}$ and $\{w^k\}$ have finite number of limits points. Let $\mathscr{S}_r$, $\mathscr{S}_w$ and $\mathscr{S}$ denote the sets of limit points of $\{r^k\}$, $\{w^k\}$ and $\{(r^k, w^k)\}$ respectively. We prove, by contradiction, that $\{(r^k, \lambda^k)\}$ can have exact one limit point.

Suppose $\{(r^k, w^k)\}$ has another limit point, e.g., $(\hat{r}, \hat{w})$. From the nondegeneracy assumptions, $\bar{w} \neq \hat{w}$ and $\bar{r} \neq \hat{r}$.

Define $\delta = \min\{\delta_1, \delta_2\} > 0$ where

$$\delta_1 = \tfrac{1}{2} \min\{\|\bar{w} - \check{w}\|: \check{w} \in \mathscr{S}_w, \check{w} \neq \bar{w}\}, \qquad \delta_2 = \tfrac{1}{2} \min\{\|\bar{r} - \check{r}\|: \check{r} \in \mathscr{S}_r, \check{r} \neq \bar{r}\}.$$

Let

$$C_\delta = \{(r, w): \|\bar{w} - w\| < \tfrac{1}{2}\delta \text{ and } \|\bar{r} - r\| < \tfrac{1}{2}\delta\}.$$

Since $(\bar{r}, \bar{w})$ is a limit point, $(r^k, w^k)$ belongs to $C_\delta$ infinitely often. In order for the sequence $\{(r^k, w^k)\}$ to have another limit point $(\check{r}, \check{w})$, $\{(r^k, w^k)\}$ must leave $C_\delta$ infinitely often. Let $\mathscr{K}$ denote the subsequence corresponding to all such points which leave $C_\delta$, i.e.,

$$\mathscr{K} = \{(r^k, w^k): (r^k, w^k) \in C_\delta \text{ but either } \|r^{k+1} - \bar{r}\| \geq \tfrac{1}{2}\delta \text{ or } \|w^{k+1} - \bar{w}\| \geq \tfrac{1}{2}\delta\}.$$

Assume there exists a subsequence of $\mathscr{K}$ with $\alpha^k d^k$ converging to a positive limit. From Lemma 15, the corresponding $\{D_r^k(g^k - \lambda^k)\}$ converges to zero. Using (20), the corresponding $\{d_w^k\}$ also converges to zero. Thus, there exists $k_1 > 0$ such that, when $k > k_1$, $\|d_w^k\| \leq \tfrac{1}{2}\delta$. Hence

$$\|\bar{w} - w^{k+1}\| \leq \|\bar{w} - w^k\| + \|d_w^k\| \leq \delta.$$

Hence $w^{k+1} \in \{w: \|w - \bar{w}\| \leq \delta\}$. If there exists a infinite subsequence with $\|r^{k+1} - \bar{r}\| \geq \tfrac{1}{2}\delta$, there exists a limit point $\hat{r}$ of $r^{k+1}$ with $\hat{r} \neq \bar{r}$. But the corresponding $\{w^{k+1}\}$ has a limit point $\hat{w} \in \{w: \|w - \bar{w}\| \leq \delta\}$ (Recall that $\{w^k\}$ and $\{r^k\}$ are bounded). This means that $\hat{w} = \bar{w}$ by definition of $\delta$. But this is impossible because $(\hat{r}, \hat{w}) \in \mathscr{S}$ but $\hat{r} \neq \bar{r}$ and $\hat{w} = \bar{w}$. Thus, for $k$ sufficiently large $\|w^{k+1} - \bar{w}\| \geq \tfrac{1}{2}\delta$. However, this indicates the existence of a limit point $\check{w}$ with $\|\check{w} - \bar{w}\| \leq \delta$ and $\|\check{w} - \bar{w}\| \geq \tfrac{1}{2}\delta$. Again, this is impossible from the definition of $\delta$.

Consequently, for $(r^k, w^k) \in \mathcal{H}$, $\alpha^k d^k$ converges to zero. This implies that for $k$ sufficiently large, $\|r^{k+1} + r^k\| < \frac{1}{2}\delta$ and $\|w^{k+1} - \bar{w}\| \geq \frac{1}{2}\delta$. But this indicates the existence of a limit point $(\hat{r}, \hat{w})$ with $\hat{r} = \bar{r}$ but $\hat{w} \neq \bar{w}$. This is again a contradiction.

We conclude that $\{(r^k, \lambda^k)\}$ can have only one limit point. $\square$

It is now clear that under nondegeneracy assumptions, Algorithms 1 and 2 generate points that converge to the optimum point. This follows from Lemma 11, and Theorems 14, 16.

## 6. Quadratic convergence

In this section we establish that Algorithm 2 produces a sequence $\{(r^k, w^k)\}$ that converges to $(r^*, w^*)$ at a *quadratic* rate. The primary difficulty in establishing this is the discontinuous nature of the first derivatives of $\psi$. This problem is circumvented by considering a finite set $\mathcal{F}$ of systems of nonlinear equations, where each system in $\mathcal{F}$ has the following form:

$$\hat{F}_\nu(y) \stackrel{\text{def}}{=} D_r(g(\nu) - Z^\mathrm{T} w) = 0, \qquad Zr = 0, \tag{56}$$

where

$$y = \begin{bmatrix} r \\ w \end{bmatrix}, \qquad g(\nu)_i = \begin{cases} \operatorname{sign}(r_i^*) & \text{if } r_i^* \neq 0, \\ \nu_i & \text{if } r_i^* = 0, \end{cases}$$

and $\nu_i$ can be either 1 or $-1$. Note that $y^* = (r^*, w^*)$ is a solution to each system; each system is continuously differentiable in a neighbourhood of $y^*$.

The Newton step at $y^k$, for each of the above systems $F_\nu$, is defined by

$$J_\nu^k d_N^k = -F_\nu(y^k)$$

where

$$J_\nu^k = \begin{bmatrix} \operatorname{diag}(g(\nu) - \lambda^k) & -D_r^k Z^\mathrm{T} \\ Z & 0 \end{bmatrix}, \quad \lambda_k = Z^\mathrm{T} w^k,$$

and

$$F_\nu(y^k) = \begin{bmatrix} \hat{F}_\nu(y^k) \\ 0 \end{bmatrix}. \tag{57}$$

Note that the hybrid step $d^k$ satisfies

$$B_{\nu^k} d^k = -F_{\nu^k}(y^k), \tag{58}$$

where

$$B_{\nu^k} = \begin{bmatrix} \operatorname{diag}(g^k)|D_\theta^k| & -D_r^k Z^\mathrm{T} \\ Z & 0 \end{bmatrix}, \tag{59}$$

and $v^k = g^k$. It is clear that $F_{\nu^k} \in \mathcal{F}$ (i.e., $d^k$ is a Newton-like step for some $F_{\nu^k} \in \mathcal{F}$). Of course $F_{\nu^k} \neq F_{\nu^{k+1}}$, in general, and therefore quadratic convergence is not automatic; however, a slight modification of Theorem 3.4 in [8] yields a viable approach.

**Theorem 17.** *Let $\mathcal{F} = \{F_{\nu} : \mathbb{R}^m \to \mathbb{R}^m\}$ be a finite set of functions satisfying:*

1. *each $F_{\nu}$ is continuously differentiable in an open convex set $D$;*

2. *there is a $y^*$ in $D$ such that $F_{\nu}(y^*) = 0$ and $\nabla F_{\nu}(y^*)$ is nonsingular for each $F_{\nu} \in \mathcal{F}$;*

3. *$\|F_{\nu}(y) - F_{\nu}(y^*)\| \leq c\|y - y^*\|$ for all $y \in D$, $F_{\nu} \in \mathcal{F}$ and some $c \geq 0$.*

*Let $\{B^k\}$ in $\mathcal{L}(\mathbb{R}^m)$ be a sequence of nonsingular matrices. Suppose that for some $y^0$ in $D$ the sequence*

$$y^{k+1} = y^k - (B^k)^{-1} F_{\nu^k}(y^k), \quad k = 0, 1, \ldots,$$

*remains in $D$, $y^k \neq y^*$ for $k > 0$, and converges to $y^*$. Moreover, assume*

$$\|B^k - \nabla F_{\nu^k}(y^*)\| = O(\|y^k - y^*\|). \tag{60}$$

*Then $\{y^k\}$ converges quadratically to $y^*$.* $\quad\square$

We now show that Algorithm 2 can be described in a manner consistent with Theorem 17 and therefore quadratic convergence is achieved. Specifically, (60) must be established: the next four results establish several preliminary bounds.

First we show that the change in the dual variables, $\lambda$, is bounded by the distance to the solution. Recall the full space implementation (27) of Algorithm 2: the hybrid step $d^k$ solves

$$[\operatorname{diag}(g^k)|D_\theta^k|A^T, -D_r^k Z^T] \begin{bmatrix} d_x^k \\ d_w^k \end{bmatrix} = -D_r^k(g^k - Z^T w^k), \tag{61}$$

and $d^k = A^T d_x^k$. We denote the $m \times m$ matrix $C^k$ by

$$C^k = [\operatorname{diag}(g^k)|D_\theta^k|A^T, -D_r^k Z^T].$$

**Lemma 18.** *Assume that $\{(r^k, w^k)\}$ is any subsequence, convergent to $y^* = (r^*, w^*)$, obtained by Algorithm 2. Then,*

$$\|\lambda^{k+1} - \lambda^k\| = \|d_w^k\| = O(\|y^k - y^*\|), \tag{62}$$

*where $\lambda = Z^T w$.*

**Proof.** Since

$$\begin{bmatrix} d_r^k \\ d_w^k \end{bmatrix} = -C^{k^{-1}} D_r^k(g^k - Z^T w^k),$$

we immediately obtain

$$\|w^{k+1} - w^k\| \leq L\|C^{k^{-1}}\| \|y^k - y^*\|$$

and since $\{\|(C^k)^{-1}\|\}$ is bounded, the result ensues. $\quad\square$

Next we show that $\theta$ is bounded by the distance to the solution.

**Lemma 19.** *Suppose $\{\theta^k\}$ is defined as in (24). Then,*

$$\theta^k \leq L_1 \| y^k - y^* \|. \tag{63}$$

**Proof.** By definition of $\theta^k$,

$$\theta^k \leq L_0 \sqrt{\sum_{i=1}^m r_i^{k^2}(g_i^k - \lambda_i^k)^2}.$$

Thus, from (57) and (56), we have

$$\theta^k \leq L_0 \| F_{\nu^k}(y^k) \|_2$$

$$= L_0 \| F_{\nu^k}(y^k) - F_{\nu^k}(y^*) \|_2$$

$$\leq L_1 \| y^k - y^* \|_2,$$

for positive constants $L_0, L_1$. The last inequality comes from continuity and finiteness of $\mathscr{F}$. $\square$

**Lemma 20.** *Assume that an $\ell_1$ problem is primal and dual nondegenerate and that the sequence $\{(r_k, w_k)\}$ is generated by Algorithm 2. Then, for each $i$ such that $r_i^* = 0$,*

$$\alpha_i^k - 1 = \mathrm{O}(\| y^k - y^* \|). \tag{64}$$

*Moreover, the sequence of optimal stepsizes, $\{\alpha_*^k\}$, converges to unity.*

**Proof.** First we establish (64). For any $i$ corresponding to an activity at the solution,

$$\alpha_i^k - 1 = \frac{-r_i^k}{d_i^k} - 1$$

$$= \left( \frac{r_i^k(g_i^k - (1 - \theta^k)\lambda_i^k)}{r_i^k(g_i^k - \lambda_i^{k+1})} - 1 \right) \quad \text{(by (31))}$$

$$= \frac{\lambda_i^{k+1} - \lambda_i^k + \theta^k \lambda_i^k}{g_i^k - \lambda_i^{k+1}}.$$

But since $i$ corresponds to an active constraint, then, by dual nondegeneracy, $g_i^k - \lambda_i^{k+1}$ is bounded from zero. Therefore, since $\theta^k \to 0$ and Lemma 18 holds, we have $\alpha_i^k - 1 = \mathrm{O}(\| y^k - y^* \|)$ and $\lim_{k \to \infty} \alpha_i^k = 1$.

Next we establish that the optimal stepsize, $\alpha_*^k$, converges to unity. First we note that the stepsizes to the breakpoints of the inactive constraints go to infinity. That is, if $i$ corresponds to an inactive constraint at the solution,

$$\lim_{k \to \infty} \frac{-r_i^k}{d_i^k} = \infty. \tag{65}$$

This is true because $d_i^k \to 0$ but $r_i^k$ is bounded from zero when $i$ corresponds to an inactivity at the solution. (To see that $d_i^k \to 0$, notice that by convergence $\theta^k \to 0$ which implies $D_r^k(g^k - Z^T w^k) \to 0$; the result follows from Lemma 8 and (25).) Next we prove that the stepsizes to the breakpoints of the activities all go to unity. To establish this we first prove that[5]

$$|g_1^{k^T} d_1^k| < |g_0^{k^T} d_0^k|.$$

To see this recall, from dual feasibility and dual nondegeneracy, that

$$g_1^* A_1 = -A_0 \lambda_0^*, \quad \text{where } \max\{|\lambda_0^*|\} < 1.$$

Moreover, for $k$ sufficiently large, it is clear that $g_1^k = g_1^*$. Hence, from Lemma 12, we have

$$g_1^{k^T} d_1^k = -\lambda_0^{*^T} d_0^k.$$

Since $|\lambda_0^*| < 1$, we have

$$-g_0^{k^T} d_0^k \le \sum_{i \in \mathcal{A}} |d_i^k| \quad \text{for sufficiently large } k.$$

Therefore,

$$\frac{|g_1^{k^T} d_1^k|}{|g_0^T d_0^k|} \le \max\{|\lambda_0^*|\} < 1.$$

Thus, when $k$ is sufficiently large,

$$\frac{|g_1^{k^T} d_1^k|}{|g_0^T d_0^k|} < 1,$$

which implies

$$|g_1^{k^T} d_1^k| \le |g_0^{k^T} d_0^k| = -g_0^{k^T} d_0^k \quad \text{for sufficiently large } k. \tag{66}$$

Now let $g^{k^+}$ denote the gradient after crossing all the breakpoints corresponding to the activities. Then, by (66),

$$g^{k^{+T}} d^k = -g_0^{k^T} d_0^k + g_1^{k^T} d_1^k > 0 \quad \text{for sufficiently large } k.$$

This, together with (65), implies that the optimal steplength equals one of the steplengths corresponding to an activity. Thus, the optimal steplength, $\alpha_k^*$, converges to unity.   $\square$

**Lemma 21.** *Assume $\{(r^k, w^k)\}$ is generated by Algorithm 2 where $\{\alpha^k\}$ is the stepsize; assume primal and dual nondegeneracy. Then*

$$|\alpha^k - 1| \le L_2 \|y^k - y^*\|. \tag{67}$$

---

[5] Notation: In this section subscript 0 refers to components corresponding to activities at the solution, i.e., $r_i^* = 0$. Subscript 1 refers to components corresponding to inactivities at the solution, i.e., $r_i^* \ne 0$.

**Proof.** From the computation of the steplength (32), convergence, and Lemma 19, it is clear that

$$\alpha^k = \alpha^k_* - \theta^k \alpha^k_* + \mathrm{O}(\|y^k - y^*\|) \quad \text{for } k \text{ sufficiently large.}$$

Furthermore, from the proof of Lemma 20, $* \in \mathcal{A}$.

From (62),

$$|\alpha^k_i - 1| = \mathrm{O}(\|y^k - y^*\|), \tag{68}$$

therefore,

$$|\alpha^k - 1| = |\alpha^k_i - 1 - \theta^k \alpha^k_i| + \mathrm{O}(\|y^k - y^*\|)$$

$$\leq |\alpha^k_i - 1| + \theta^k \alpha^k_i + \mathrm{O}(\|y^k - y^*\|)$$

$$\leq L_2 \|y^k - y^*\| \quad \text{(using (68) and Lemma 19).} \qquad \square$$

Denote $B^k = B_{\nu^k} S^k$ where

$$S^k = \mathrm{diag}\left( \begin{bmatrix} 1/\alpha^k\, e_m \\ e_{m-n} \end{bmatrix} \right),$$

where $B_{\nu^k}$ is defined as in (59). But $y^{k+1} = y^k + S^{k-1} d^k$ where $d^k$ is defined by (58); therefore, from Lemma 21, we have

$$\|S^k - I\| = \mathrm{O}(\|y^k - y^*\|).$$

The hybrid step defined by Algorithm 2 satisfies

$$B^k(y^{k+1} - y^k) = -F_{\nu^k}(y^k).$$

**Lemma 22.** *Assume $\{y^k\}$ is obtained through Algorithm 2; assume primal and dual nondegeneracy. Then*

$$\|B^k - \nabla F_{\nu^k}(y^*)\| \leq L \|y^k - y^*\|$$

*for some $F_{\nu^k} \in \mathcal{F}$.*

**Proof.** From continuity and Lemma 21, it is clear that

$$B^k - J^*_{\nu^k} = (B_{\nu^k} - J^k_{\nu^k}) S^k + (J^k_{\nu^k} - J^*_{\nu^k}) S^k + (S^k - I) J^*_{\nu^k}$$

$$= \mathrm{O}(\|B_{\nu^k} - J^k_{\nu^k}\|) + \mathrm{O}(\|y^k - y^*\|) + \mathrm{O}(\|y^k - y^*\|).$$

But,

$$B_{\nu^k} - J^k_{\nu^k} = \begin{bmatrix} \mathrm{diag}(g^k)|D^k_\theta| - D^k_\lambda & 0 \\ 0 & 0 \end{bmatrix}.$$

Consider the $(1, 1)$-block of $B_{\nu^k} - J_{\nu^k}^k$. It is clear that, for $k$ sufficiently large, $g_0^k = \text{sign}(g_0^k - \lambda_0^k)$, $g_1^k = \lambda_1^*$. Hence

$$\|\text{diag}(g^k)|g^k - (1 - \theta^k)\lambda^k| - (g^k - \lambda^k)\|_2$$

$$= \|\text{diag}(g_0^k)|g_0^k - (1 - \theta^k)\lambda_0^k| - (g_0^k - \lambda_0^k)\|_2$$

$$+ \|\text{diag}(g_1^k)|g_1^k - (1 - \theta^k)\lambda_1^k| - (g_1^k - \lambda_1^k)\|_2$$

$$\leq \|\theta^k \lambda_0^k\|_2 + \|\lambda_1^* - \lambda_1^k\|_2 + \theta^k\|\lambda_1^k\|_2 + \|\lambda_1^* - \lambda_1^k\|_2.$$

Thus, by Lemma 19,

$$\|\text{diag}(g^k)|g^k - (1 - \theta^k)\lambda^k - (g^k - \lambda^k)\|_2 = O(\|y^k - y^*\|_2).$$

Therefore,

$$\|B^k - J_{\nu_k}^*\| = O(\|y^k - y^*\|). \qquad \square$$

The assumptions of Theorem 17 are now established; quadratic convergence of Algorithm 2 follows immediately.

**Theorem 23.** *Suppose the $\ell_1$ problem is primal and dual nondegenerate. Assume the sequence $\{(r^k, w^k)\}$ is obtained from the Algorithm 2. Then $y^k = \left[\begin{smallmatrix} r^k \\ w^k \end{smallmatrix}\right]$ converges quadratically to $y^* = \left[\begin{smallmatrix} r^* \\ w^* \end{smallmatrix}\right]$.* $\square$

## 7. Numerical testing

In this section we provide numerical results concerning Algorithm 2, the hybrid method (*New*). In particular, we compare our range-space implementation of this algorithm — see Section 4 — with our implementation of the interior point algorithm described in [12]. We denote the latter method by *Dual*. Our comparisons are based on the number of iterations: since the dominant work in both implementations is the solution of a linear least squares problem at each iteration, of size $m \times n$, the number of iterations accurately reflects the overall relative computational cost.

Our stopping criterion for both algorithms is based on the satisfaction of the optimality conditions:

$$\|D_r^k(g^k - \lambda^{k+1})\|_\infty + \|\max(|\lambda| - e, 0)\|_\infty < 10^{-13},$$

where machine precision on our system is approximately $10^{-16}$. The origin is a natural starting point for algorithm *Dual*; the starting point for *New* is computed as follows:

$$r^0 \leftarrow b - A^T x^0, \quad \text{where } A^T x^0 \overset{\text{l.s.}}{=} b,$$

$$\lambda^0 \leftarrow \frac{\tau}{\max(|r|)} * r^0.$$

The settings of the parameters for Algorithm 2, *New*, are:

$$\tau \leftarrow 0.975, \quad \gamma \leftarrow 0.99.$$

Algorithm *Dual* also uses $\tau = 0.975$ in the linesearch algorithm; there are no other parameters.

We have implemented the methods in PRO-MATLAB [13] using SUN 3/50 and 3/160 work-stations. The linear least squares subproblems are solved with the orthogonal QR-factorization using row interchanges for greater stability (row interchanges are advisable when a least squares problem involves widely varying row scalings [20]). No account was made of sparsity in our experiments.

Our experiments are not exhaustive; our purpose here is to determine the viability of our approach. Specifically, does the ultimate quadratic convergence appear to yield significantly fewer iterations in practise?

We have generated several kinds of test problems. First, since $\ell_1$ minimization is often used in a function approximation context [15], we have tried several such problems. We have also generated several random problems of varying dimensions. Finally, despite lack of theory for degenerate problems, we have experimented with degenerate and near-degenerate problems.

**Problem 1.** Approximate following $f(z)$, evaluated for $z = 0\,(1/m)\,1$ by a polynomial of degree $n - 1$:

$$\phi(z) = \sum_{j=1}^{n} \alpha_j z^{j-1}.$$

$n = 5$, $f_1(z) = \exp(z)$: see Table 1. $n = 6$, $f_2(z) = \sin(z)$: see Table 2.

Table 1

Number of steps

| $m$ | Dual | New |
|-----|------|-----|
| 100 | 18 | 8 |
| 200 | 19 | 11 |
| 400 | 21 | 10 |
| 600 | 21 | 12 |

Table 2

Number of steps

| $m$ | Dual | New |
|-----|------|-----|
| 100 | 15 | 9 |
| 200 | 16 | 8 |
| 400 | 16 | 9 |
| 600 | 17 | 9 |

**Problem 2.** Approximate following $f(z)$, evaluated for $z = 0\,(1/m)\,1$ by a polynomial of degree 10:

$$\phi(z) = \sum_{j=1}^{n} \alpha_j z^{j-1}.$$

Let $x = (\alpha_1, \ldots, \alpha_n)$.

$$f_3(z) = \exp(z) + \begin{cases} 1 & \text{if } 0.1 < z \leq 0.2, \\ 0 & \text{otherwise.} \end{cases}$$

See Tables 3 and 4.

Table 3

Number of steps

| $m$ | Dual | New |
| --- | --- | --- |
| 40 | 30 | 13 |
| 100 | 27 | 12 |
| 200 | 29 | 12 |
| 500 | 33 | 12 |

Table 4

Number of steps

| $m$ | $n$ | $f(z)$ | New |
| --- | --- | --- | --- |
| 1000 | 5 | $f_1(z)$ | 13 |
| 1000 | 6 | $f_2(z)$ | 10 |
| 800 | 10 | $f_3(z)$ | 12 |
| 1000 | 10 | $f_3(z)$ | 15 |

**Problem 3.** Random problem. See Tables 5–8.

Table 5

$m = 50$

Number of steps

| $n$ | Dual | New |
| --- | --- | --- |
| 10 | 25 | 9 |
| 20 | 25 | 12 |
| 30 | 25 | 10 |
| 40 | 23 | 9 |

Table 6

$m = 100$

Number of steps

| $n$ | Dual | New |
| --- | --- | --- |
| 10 | 25 | 10 |
| 20 | 26 | 13 |
| 40 | 26 | 13 |
| 50 | 26 | 11 |
| 70 | 25 | 11 |
| 90 | 23 | 10 |

Table 7

$m = 150$

Number of steps

| $n$ | Dual | New |
| --- | --- | --- |
| 10 | 26 | 10 |
| 20 | 27 | 13 |
| 50 | 26 | 17 |
| 70 | 27 | 11 |
| 90 | 27 | 13 |
| 135 | 24 | 12 |
| 140 | 23 | 9 |

Table 8

$m = 200$

Number of steps

| $n$ | Dual | New |
| --- | --- | --- |
| 10 | 26 | 13 |
| 20 | 27 | 17 |
| 50 | 27 | 16 |
| 70 | 27 | 19 |
| 100 | 26 | 15 |
| 140 | 26 | 14 |
| 160 | 25 | 14 |
| 190 | 24 | 9 |

**Problem 4.** Degenerate problems. In theory, the convergence and the convergence and the convergence rate of the algorithm is not guaranteed under the degenerate cases. Our experiments, however, indicate that the performance of the algorithm is not affected.

*Dual degenerate problems.* At the solution, the multipliers of *ndeg* activities are equal to either 1 or −1. $m = 100$. See Table 9.

Table 9

Number of steps

| $n$ | *ndeg* | *Dual* | *New* |
|-----|--------|--------|-------|
| 50  | 5      | 27     | 14    |
| 50  | 10     | 27     | 13    |
| 70  | 15     | 25     | 11    |
| 70  | 20     | 27     | 10    |

*Primal degenerate problems.* At the solution, there are *ndeg* activities which are redundant in forming a maximum linearly independent set of activities. $m = 100$. See Table 10.

Table 10

Number of steps

| $n$ | *ndeg* | *Dual* | *New* |
|-----|--------|--------|-------|
| 20  | 5      | 25     | 9     |
| 40  | 10     | 25     | 12    |
| 70  | 15     | 24     | 10    |
| 90  | 9      | 9      | 2     |

**Problem 5.** Primal near-degenerate problem. These problems are primal nearly degenerate in the sense that there are *ndeg* residuals that are active up to a tolerance of $\delta$. We believe that the degeneracy does affect the local quadratic convergence region. However, the algorithm still performs reasonably well. $m = 50$, *ndeg* = 2, $\delta = 10^{-9}$: see Table 11. $m = 100$: see Table 12.

**Remarks.** In general, the quadratically convergent Algorithm 2 (*New*) provides a clear advantage over the (linear) interior point method (*Dual*) on our test collection. The number of iterations required by *New* is always fewer, usually by a factor between two and three. *New* wins by the smallest margin on near-degenerate problems.

Table 11

| | Number of steps | |
|---|---|---|
| $n$ | Dual | New |
| 10 | 25 | 14 |
| 20 | 25 | 15 |
| 30 | 25 | 16 |
| 40 | 25 | 18 |

Table 12

| | Number of steps | |
|---|---|---|
| $n$ | Dual | New |
| 10 | 25 | 16 |
| 20 | 26 | 21 |
| 40 | 26 | 21 |

The number of iterations required by *New* is not quite as consistent, over different problems and different dimensions, as the number required by *Dual*. For example, on the random problems, *Dual* varies between 23 and 27 iterations whereas *New* ranges between 9 and 17. This greater variation is probably due to different radii of convergence for Newton's method. Nevertheless, the number of required iterations is almost never more than 15 and is considerably more consistent than number of iterations required by a finite pivoting algorithm.

## 8. Conclusions

We have presented a new global and quadratic algorithm for the linear $\ell_1$ problem. Computationally, the algorithm appears to be consistently superior to the interior point approach [12], typically requiring a factor of two to three fewer iterations. (The iterations are comparable in cost and so the overall running times compare in a similar way.) To support this claim consider Table 13 in which, on a representative example, we compare values of $\theta^k$ for the two algorithms.

Comparing the values of $\theta$, it is difficult to distinguish between the two algorithms for the first eight iterations; shortly after that second-order convergence kicks in for our new algorithm whereas *Dual* continues to exhibit linear convergence. This example is typical.

The quadratic convergence property stems from the nonlinear expression for complementary slackness and the resulting Newton-like process; the integration of the Newton process with a (linear) global affine scaling method comes in part through the use of a "square-root" scaling matrix to define the linear algorithm (Algorithm 1). The use of this scaling is compatible with the fact that "constraint lines", i.e., $r_i = 0$, must be crossed, in general.

The linear $\ell_1$ algorithm proposed here is probably most applicable to large-scale problems. There are several reasons for this. First, Algorithm 2 (*New*) seems fairly insensitive to problem size: it appears that about 15 iterations will be required, regardless of dimension. Therefore, this method will be increasingly attractive, relative to finite pivoting algorithms, as the dimension grows. Second, for sparse problems the technology required to exploit sparsity is only that required for solving

Table 13

Change of $\theta$ for random $(60, 12)$

| Iteration | New | Dual |
|---|---|---|
| 1 | 6.68708 e −02 | 5.66975 e −01 |
| 2 | 2.94760 e −01 | 2.70565 e −01 |
| 3 | 1.37327 e −01 | 1.44938 e −01 |
| 4 | 7.80604 e −02 | 5.23090 e −02 |
| 5 | 2.49524 e −02 | 6.83901 e −02 |
| 6 | 8.45633 e −03 | 7.05475 e −02 |
| 7 | 8.48194 e −03 | 3.17295 e −03 |
| 8 | 3.72855 e −04 | 6.34264 e −04 |
| 9 | 1.54264 e −04 | 1.57652 e −02 |
| 10 | 2.52256 e −05 | 1.12097 e −02 |
| 11 | 1.90326 e −09 | 2.49742 e −04 |
| 12 | 1.33204 e −12 | 3.60838 e −06 |
| 13 | | 1.46895 e −06 |
| 14 | | 1.55962 e −07 |
| 15 | | 4.92358 e −08 |
| 16 | | 1.06442 e −08 |
| 17 | | 2.72495 e −09 |
| 18 | | 8.21490 e −10 |
| 19 | | 1.74925 e −10 |
| 20 | | 5.06344 e −11 |
| 21 | | 1.27783 e −11 |
| 22 | | 3.32228 e −12 |

a sequence of sparse linear least squares problems: this is a heavily studied problem with techniques available and research ongoing. Finally, the new method is attractive for parallel computation in that the number of outer iterations (sequential steps) is modest; moreover, parallel techniques for solving linear least-squares problems have been developed (e.g., [5] and [4]) and research is continuing in this area.

## Acknowledgement

## References

[1] E. Barnes, "A variation on Karmarkar's algorithm for solving linear programming problems," *Mathematical Programming* 36 (1986) 174–182.

[2] I. Barrodale and F. Roberts, "An improved algorithm for discrete $\ell_1$ linear approximation," *SIAM Journal on Numerical Analysis* 10 (1972 839–848.

[3] R.H. Bartels, A.R. Conn and J.W. Sinclair, "Minimization techniques for piecewise differentiable functions: the $\ell_1$ solution to an overdetermined linear system," *SIAM Journal Numerical Analysis* 15 (1978) 224–240.

[4] C. Bischof, "QR factorization algorithms for coarse-grained distributed systems," Technical Report 88-939, Cornell University (Ithaca, NY, 1988).

[5] T.F. Coleman and P. Plassmann, "Solution of nonlinear least-squares problems on a multiprocessor," in: G. van Zee and J. van de Vorst, eds., *Parallel Computing 1988, Shell Conference Proceedings,* Lecture Notes in Computer Science No. 384 (Springer, Berlin, 1989).

[6] T.F. Coleman and A. Pothen, "The null space problem I: Complexity," *SIAM Journal on Algebraic and Discrete Methods* 7 (1987) 527–537.

[7] T.F. Coleman and A. Pothen, "The null space problem II: Algorithms," *SIAM Journal on Algebraic and Discrete Methods* 8 (1987) 544–563.

[8] J.E. Dennis, Jr. and J.J. Moré, "Quasi-Newton methods, motivation and theory," *SIAM Review* 19 (1977) 46–89.

[9] I. Dikin, "Iterative solution of problems of linear and quadratic programming," *Doklady Akademiia Nauk SSSR* 174 (1967) 747–748.

[10] P. Gill, W. Murray, M. Saunders, J. Tomlin and M. Wright, "On projected Newton barrier methods for linear programming and an equivalence to Karmarkar's projective method," *Mathematical Programming* 36 (1986) 183–209.

[11] N. Karmarkar, "A new polynomial-time algorithm for linear programming," *Combinatorica* 4 (1984) 373–395.

[12] M.S. Meketon, "Least absolute value regression," Technical Report, AT&T Bell Laboratory (Murray Hill, NJ, 1987).

[13] C.B. Moler, J. Little, S. Bangert and S. Kleiman, *ProMatlab User's Guide* (MathWorks, Sherborn, MA, 1987).

[14] J.M. Ortega and W.C. Rheinboldt, *Iterative Solution of Nonlinear Equations in Several Variables* (Academic Press, New York, 1970).

[15] M.R. Osborne, *Finite Algorithms in Optimization and Data Analysis* (Wiley, New York, 1985).

[16] J. Renegar, "A polynomial-time algorithm, based on Newton's method, for linear programming," *Mathematical Programming* 40 (1988) 59–93.

[17] S.A. Ruzinsky and E.T. Olsen, "$\ell_1$ and $\ell_\infty$ minimization via a variant of Karmarkar's algorithm," *IEEE Transactions on Acoustics Speech and Signal Processing* 37 (1989) 245–253.

[18] E. Seneta and W.L. Steiger, "A new lad curve-fitting algorithm: Slightly overdetermined equation systems in $\ell_1$," *Discrete Applied Mathematics* 7 (1984) 79–91.

[19] M. Todd, "Polynomial algorithms for linear programming," in: H. Eiselt and G. Pederzoli, eds., *Advances in Optimization and Control* (Springer, Berlin, 1988) pp. 49–66.

[20] C. Van Loan, "On the method of weighting for equality-constrained least squares problems," *SIAM Journal on Numerical Analysis* 22 (1985) 851–864.

[21] R.J. Vanderbei, M.S. Meketon and B.A. Freedman, "A modification of Karmarkar's linear programming algorithm," *Algorithmica* 1 (1986) 395–407.