

ASTRAL: AN ACTIVE SET ℓ_∞ -TRUST-REGION ALGORITHM FOR BOX CONSTRAINED OPTIMIZATION

LIANG XU¹ AND JAMES V. BURKE²

ABSTRACT. An algorithm for solving large-scale nonlinear optimization problems with simple bounds is described. The algorithm is an ℓ_∞ -norm trust-region method that uses both active set identification techniques as well as limited memory BFGS updating for the Hessian approximation. The trust-region subproblems are solved using primal-dual interior point techniques that exploit the structure of the limited memory Hessian approximation. A restart strategy ensures that the algorithm identifies the optimal constraints in finite number iterations under a standard nondegeneracy hypothesis. Local and global convergence properties are established, and the results of numerical tests are given.

1. INTRODUCTION

In [3], Burke and Weigmann proposed a strategy for implementing the limited memory BFGS matrix secant updating strategy in the context of a trust-region algorithm [19]. The algorithm was shown [15] to be very effective in improving the computational efficiency of numerical methods for wavefront reconstruction in optical imaging where it is common to have 2^{19} or more nonlinear unknowns. The trust-region approach in [3], builds on the line-search implementation of limited memory BFGS (L-BFGS) developed by Nocedal [16] and Byrd, Nocedal, and Schnabel [6]. In [5], Byrd, Lu, Nocedal, and Zhu extend the limited memory technology to problems having simple bound constraints again using a linesearch to adjust the step size. In this paper we extend the trust-region approach in [3] to problems with simple bound constraints. The algorithms are intended for large-scale nonlinear programs with bound constraints and are not recommended for problems of small to medium size, or for large-scale problems whose a structure can be exploited in algorithm design.

Trust-region algorithms typically have a greater cost per iteration due to the need to solve the trust-region subproblem. In this regard, the proposed method (ASTRAL) has a significantly greater linear algebra cost per iteration than the linesearch algorithm of Byrd-Lu-Nocedal-Zhu since a quadratic program (QP) is solved at each iteration. The underlying philosophy

Date: July 9, 2007.

1. Research supported in part by NSF grant DMS-0505712.

2. Research supported in part by NIH grant P41 EB-001975 and NSF grant DMS-0505712.

of trust-region methods is that by making more extensive use of accumulated local information a better step choice is possible thereby reducing the overall number of function and gradient calls required to achieve a desired level of accuracy. These methods work best in a context where function and gradient calls are costly in terms of CPU time. The primary goal in the design of ASTRAL is to mimic as closely as possible the linear algebra of the unconstrained limited memory BFGS (L-BFGS) update. In this way we hope to recover the numerical performance and efficiency of the L-BFGS algorithm in the presence of box constraints. We achieve this design goal by solving the underlying QP subproblems using interior point technology.

ASTRAL has the following basic features: (a) an ℓ_∞ -norm trust-region is used to conform with the geometry of the constraint region, (b) aggressive active set identification techniques are employed to reduce subproblem dimension when possible, (c) L-BFGS Hessian approximations are used to extract local curvature information while further reducing dimensionality, and (d) the QP subproblems are solved using an interior point methodology that exploits the underlying algebraic structure of the L-BFGS update. The end result is an ℓ_∞ -norm trust-region algorithm for large-scale box constrained optimization problems that compares favorably with the performance of the Byrd-Lu-Nocedal-Zhu algorithm L-BFGS-B.

Several algorithms have been developed to minimize a smooth nonlinear function subject to box constraints [5, 8, 14, 13, 17]. Although ASTRAL uses similar techniques, there are essential differences. These differences are designed to reduce the overall number of function evaluations. The methods [5, 8, 14, 13, 17] first compute an *active* face of the box. The methods in [5, 8, 14, 17] then compute a search direction or a trust-region step using this face, while [13] used a convergence criteria to determine if unconstrained minimization over the active face is warranted. The methods described in [5, 8, 14] compute a generalized Cauchy point (GCP) to identify the active face. The GCP is obtained by minimizing a quadratic approximation to the objective along the projected gradient path. Both [8] and [14] describe a trust-region algorithm based on the 2-norm with the GCP also satisfying the trust-region constraint. ASTRAL uses a single gradient projection to estimate the active set and does not require the computation of a GCP. In [13], a non-monotone backtracking line search is applied to the objective function to identify the active face. In [8, 10], the 2-norm trust-region step is computed by applying truncated conjugate gradient (CG) method to the quadratic approximation of the objective on the active face subject to the trust-region constraint. In [14], truncated CG is applied to the quadratic approximation of the objective on the subspace parallel to the active face subject to the 2-norm trust-region constraint, and recovers feasibility, if necessary, by searching along the projected path. The methods described in [8, 10] and [14] are not designed to solve the trust-region subproblem. The L-BFGS-B method in [5] computes a search direction by truncating the

unconstrained L-BFGS update relative to the subspace parallel to the active face. The method in [13] suggests the use of an unconstrained method applied to the objective function on the active face with an adaptive termination procedure. On the other hand, ASTRAL uses a simple method based on the current point to determine the active face and then uses an interior point method to compute an approximate solution to the ∞ -norm trust-region subproblem relative to this face.

The ASTRAL algorithm is given in the next section. In sections 3 and 4 we establish some of the key global and local convergence properties, respectively. The algorithm and its convergence theory follow a pattern established by Powell in [19] and has been used for similar algorithms from the literature [4, 7, 9, 18]. In section 5, we present a standard primal-dual interior point based method for solving the QP subproblems [22], and show why the structure of the L-BFGS update is well matched to the linear algebra required to implement the interior point algorithm. We conclude in Section 6 by presenting the results of our numerical experiments.

Most of the notation and terminology that we employ are standard (e.g. see [20, 18]). However, it is helpful to list a few of these before proceeding. Let S be a subset of \mathbb{R}^n . The closure of S is denoted $\text{cl}S$. The line segment connecting the two points $x, y \in S$ is given by

$$[x, y] = \{(1 - \lambda)x + \lambda y \mid 0 \leq \lambda \leq 1\} .$$

We say that S is *convex* if it contains all of the line segments connecting any two points in S . The set S is said to be *affine* if it is representable as the translate of a subspace. The *convex* and *affine hull* of S , denoted $\text{aff}(S)$ and $\text{co}(S)$, respectively, are the smallest convex and affine sets, respectively, containing S . If S is a cone, then the *lineality* of S , denoted $\text{lin}(S)$, is the largest subspace contained in S . The set S is said to be a (closed) *convex polyhedron* if can be represented as the convex hull of finitely many points. The *relative interior* of S , denoted $\text{ri}(S)$, is its interior relative to its affine hull. The *polar* of S is the set

$$S^\circ = \{z \mid \langle z, s \rangle \leq 1 \ \forall s \in S\} .$$

If S is convex, then the *tangent* and *normal* cone to S at a point $x \in S$ are given by

$$T_S(x) = \text{cl} \bigcup_{\lambda \geq 0} \lambda(C - x) \quad \text{and} \quad N_S(x) = T_S(x)^\circ ,$$

respectively. The closure operation in the definition of the tangent cone is superfluous if S is a closed convex polyhedron.

2. ASTRAL

Consider the box constrained problem

$$\mathcal{P} \quad \underset{x \in \Omega}{\text{minimize}} \ f(x) ,$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable and the set $\Omega \in \mathbb{R}^n$ is given by

$$\Omega = \{x \in \mathbb{R}^n \mid l \leq x \leq u\}$$

with the extended real-valued vectors l and u satisfying $l_j \in [-\infty, +\infty)$, $u_j \in (-\infty, +\infty]$, and $l_j \leq u_j$ for $j = 1, 2, \dots, n$. Constraint regions of this type are often referred to as generalized boxes. The proposed algorithm makes use of active set techniques based on gradient projections and binding constraints. Recall that the projection of a point $z \in \mathbb{R}^n$ onto the the set Ω is given by

$$P_\Omega(z) = \operatorname{argmin}_{x \in \Omega} \|z - x\|_2 .$$

It is easily verified that

$$(P_\Omega(z))_i = \begin{cases} z_i & \text{if } l_i \leq z_i \leq u_i; \\ l_i & \text{if } z_i < l_i; \\ u_i & \text{if } z_i > u_i. \end{cases}$$

Define the *binding constraints* at a point $x \in \Omega$ by

$$\mathcal{B}(x) = \left\{ i \mid \begin{array}{l} x_i = l_i \text{ and } (\nabla f(x))_i > 0, \text{ or} \\ x_i = u_i \text{ and } (\nabla f(x))_i < 0, \text{ or} \\ l_i = u_i \end{array} \right\},$$

and set $\mathcal{B}^c(x) = \{1, 2, \dots, n\} \setminus \mathcal{B}(x)$, $\nu(x) = |\mathcal{B}^c(x)|$, and $\Phi(x)$ equal to the $n \times \nu(x)$ matrix whose columns are those of the identity matrix corresponding to the non-binding constraints $\mathcal{B}^c(x)$. Note that $\Phi(x)\Phi(x)^T$ is the orthogonal projection onto the range of $\Phi(x)$. The projection leaves the components in $\mathcal{B}^c(x)$ unchanged and set the remaining components to zero.

Also define the *gradient projection* mapping at a point $x \in \Omega$, $p : \Omega \rightarrow \mathbb{R}^n$, by

$$p(x) = P_\Omega(x - \nabla f(x)) - x.$$

It is well known that $\bar{x} \in \Omega$ is a first-order stationary point for \mathcal{P} if and only if $p(\bar{x}) = 0$, or equivalently, $-\nabla f(x) \in N_\Omega(x)$. Moreover, if ∇f is Lipschitz on Ω , it follows from the non-expansiveness of P_Ω (in the 2-norm) that p is also Lipschitz on Ω .

The ASTRAL Algorithm

Step 0: (Initialization) Let $\hat{\delta}_0 > 0$, $x^0 \in \Omega$, and $B_0 \in \mathbb{R}^{n \times n}$ be symmetric and positive definite. Set

$$g^0 = \nabla f(x^0) \quad \text{and} \quad p^0 = p(x^0),$$

and choose the following constants:

$$\begin{aligned} 0 < \beta_0 < \beta_1 < \beta_2 < 1, & \quad (\text{step acceptance parameters}) \\ 0 < \sigma_1 < \sigma_2 < 1 \leq \sigma_3, & \quad (\text{trust-region scaling parameters}) \\ 0 < \gamma < 1, \quad 0 \leq \tau_0 & \quad (\text{re-start parameters}) \end{aligned}$$

Set $\delta_0 = \min\{\max_i\{u_i - l_i\}, \hat{\delta}_0\}$ and $k = 0$.

Step 1: (Gradient Projection Re-Start) If $\|p^k\|_\infty = 0$, **STOP**; if $\|p^k\|_\infty > \tau_k$, go to Step 2; otherwise, set

$$\tau_{k+1} = \gamma\tau_k, \quad \hat{x} = P_\Omega(x^k - \nabla f(x^k)), \quad \hat{f} = f(\hat{x}), \quad \text{and } \hat{g} = \nabla f(\hat{x}).$$

If $\mathcal{B}(\hat{x}) \neq \mathcal{B}(x^k)$, set

$$x^{k+1} = \hat{x}, \quad f^{k+1} = \hat{f}, \quad g^{k+1} = \hat{g}, \quad \delta_{k+1} = \delta_k,$$

choose $B^{k+1} \in \mathbb{R}^{n \times n}$ to be symmetric and positive definite, compute

$$p^{k+1} = P_\Omega(x^{k+1} - g^{k+1}) - x^{k+1},$$

update k to $k + 1$, and go to Step 2.

Step 2 (Solve the Trust-Region Subproblem) Set

$$\begin{aligned} (l^k)_j &= \max\{(l - x^k)_j, -\delta_k\}, \quad (u^k)_j = \min\{(u - x^k)_j, \delta_k\} \\ &\text{for } j = 1, 2, \dots, n, \\ \Phi_k &= \Phi(x^k), \quad \tilde{B}_k = \Phi_k^T B_k \Phi_k, \quad \tilde{g}^k = \Phi_k^T g^k, \\ \tilde{l}^k &= \Phi_k^T l^k, \quad \text{and } \tilde{u}^k = \Phi_k^T u^k. \end{aligned}$$

Let \tilde{s} be the solution to the $\nu(x^k)$ -dimensional trust-region subproblem

$$(TR)_k \quad \begin{aligned} &\text{minimize} && \frac{1}{2} s^T \tilde{B}_k s + \tilde{g}^k s \\ &\text{subject to} && \tilde{l}^k \leq s \leq \tilde{u}^k \end{aligned}.$$

Set

$$q_k = \frac{1}{2} \tilde{s}^T \tilde{B}_k \tilde{s} + \tilde{g}^{kT} \tilde{s}, \quad \bar{s}^k = \Phi_k \tilde{s}, \quad \text{and } r_k = (f(x^k + \bar{s}^k) - f(x^k))/q_k.$$

Step 3: (Update the iterates) If $r_k \geq \beta_0$, set $s^k = \bar{s}^k$; otherwise, set $s^k = 0$. Update the iterate as follows:

$$x^{k+1} = x^k + s^k, \quad g^{k+1} = \nabla f(x^{k+1}), \quad p^{k+1} = p(x^{k+1}), \quad \tau_{k+1} = \tau_k,$$

and choose $B^{k+1} \in \mathbb{R}^{n \times n}$ to be symmetric and positive definite.

Step 4: (Update the trust-region radius) Set

$$\delta_{k+1} = \begin{cases} \sigma_1 \delta_k, & \text{if } r_k < \beta_0; \\ \sigma_2 \delta_k, & \text{if } \beta_0 \leq r_k < \beta_1; \\ \delta_k, & \text{if } \beta_1 \leq r_k < \beta_2; \\ \min(\sigma_3 \delta_k, \delta_0), & \text{if } r_k \geq \beta_2, \end{cases}.$$

Step 5: Update k to $k + 1$ and return to Step 1.

The re-start procedure in Step 1 is included to aid in the identification of the active constraints at the solution.

3. GLOBAL CONVERGENCE

First consider the $\tau_0 = 0$ case. In this case the algorithm only terminates if $\|p^k\|_\infty = 0$ and a gradient projection re-start is never executed in Step 1. Global convergence is established by adapting techniques of Powell [19] and others from the literature (see for example [4, 7, 9, 18]). For the remainder of this section $\{x^k\}$ denotes a sequence generated by ASTRAL with $\tau_0 = 0$. Our goal is to show that either $f(x^k) \downarrow -\infty$ or $p^k \rightarrow 0$. For this we make the following blanket assumptions.

- A1. There is a constant $M \geq 1$ such that $\|B_k\| \leq M$ for all k .
- A2. ∇f is Lipschitz continuous on Ω with Lipschitz constant L .

lemma 1. For all $k = 1, 2, \dots$,

$$-q_k \geq \frac{1}{2} \|p^k\|_\infty \min \left\{ \frac{\|p^k\|_\infty}{M}, \delta_k \right\}.$$

Proof. If $p^k = 0$, the inequality holds trivially since $s = 0$ is always feasible in $(TR)_k$, so assume that $p^k \neq 0$. We first note that $\|p^k\|_\infty = \|\Phi_k^T p^k\|_\infty$. Indeed, if $(x^k)_i = l_i$ and $(\nabla f(x^k))_i \geq 0$, then $(P_\Omega(x^k - \nabla f(x^k)))_i = l_i = (x^k)_i$, and so $(p^k)_i = 0$. Similarly, $(p^k)_i = 0$ when $(x^k)_i = u_i$ and $(\nabla f(x^k))_i \leq 0$. Therefore, $(p^k)_i = 0$ for $i \in \mathcal{B}(x^k)$, and so $(I - \Phi_k \Phi_k^T) p^k = 0$. Hence $p^k = \Phi_k \Phi_k^T p^k$, which implies $\|p^k\|_\infty = \|\Phi_k^T p^k\|_\infty$.

Set $d = \mu_k p^k$ where $\mu_k = \min\{1, \delta_k / \|p^k\|_\infty\}$. Then $x^k + d$ is in Ω , and $\|d\|_\infty \leq \delta_k$. Hence, for $0 \leq t \leq 1$, $s = t \Phi_k^T d$ is feasible for $(TR)_k$.

Set $\phi(t) = t \tilde{g}^k \Phi_k^T d + \frac{t^2}{2} d^T \Phi_k \tilde{B}_k \Phi_k^T d$. Since $\Phi_k \Phi_k^T d = d$, we have $\phi(t) = t g^k \Phi_k^T d + \frac{t^2}{2} d^T B_k d$. The definition of q_k implies $q_k \leq \phi(t)$ for all $0 \leq t \leq 1$. Since $d \neq 0$ and B_k is positive definite, we have

$$\operatorname{argmin}_{0 \leq t \leq 1} \phi(t) = \min \left\{ -\frac{g^k \Phi_k^T d}{d^T B_k d}, 1 \right\} = \bar{t},$$

which implies that

$$q_k \leq \phi(\bar{t}) \leq \begin{cases} -(g^k \Phi_k^T d)^2 / (2 d^T B_k d), & \text{when } -g^k \Phi_k^T d \leq d^T B_k d, \\ g^k \Phi_k^T d / 2, & \text{otherwise.} \end{cases}$$

Note that

$$\begin{aligned} -g^k \Phi_k^T d &= -\mu_k \langle g^k, p^k \rangle \\ &= \mu_k \langle -g^k, P_\Omega(x^k - g^k) - x^k \rangle \\ &= \mu_k [\langle (x^k - g^k) - P_\Omega(x^k - g^k), P_\Omega(x^k - g^k) - x^k \rangle \\ &\quad + \langle P_\Omega(x^k - g^k) - x^k, P_\Omega(x^k - g^k) - x^k \rangle] \\ &\geq \mu_k \|p^k\|_2^2 \geq \mu_k \|p^k\|_\infty^2 \end{aligned}$$

Consequently,

$$\begin{aligned}
 -q_k &\geq \min\{(g^k)^T d)^2/(2d^T B_k d), -g^k)^T d/2\} \\
 &\geq \frac{1}{2} \min\left\{\frac{\mu_k^2 \|p^k\|_2^4}{\mu_k^2 \|B_k\| \cdot \|p^k\|_2^2}, \mu_k \|p^k\|_\infty^2\right\} \\
 &\geq \frac{1}{2} \min\left\{\frac{\|p^k\|_\infty^2}{M}, \mu_k \|p^k\|_\infty^2\right\} \\
 &\geq \frac{1}{2} \|p^k\|_\infty \min\left\{\frac{\|p^k\|_\infty}{M}, \delta_k\right\},
 \end{aligned}$$

since B_k is positive definite and $M \geq 1$. \square

lemma 2. *If $\|p^k\|_\infty \geq \epsilon$ for all k , then δ_k does not converge to 0.*

Proof. If there is a subsequence $J \subset \mathbb{N}$ such that $\|p^k\|_\infty \leq \delta_k M$ for all $k \in J$, then $\delta_k \geq \|p^k\|_\infty / M \geq \epsilon / M$, and so δ_k does not converge to 0. Hence we can assume that $\|p^k\|_\infty \geq \delta_k M$ for all $k \geq k_0$ for some $k_0 \in \mathbb{N}$. In addition, we assume to the contrary that $\delta_k \rightarrow 0$.

By Taylor's theorem,

$$f(x^k + s^k) = f(x^k) + g^k)^T s^k + \int_0^1 [\nabla(f(x^k + ts^k) - \nabla f(x^k))]^T s^k dt,$$

and so

$$\begin{aligned}
 |f(x^k + s^k) - f(x^k) - q_k| &= \left| -\frac{1}{2} s^k)^T B_k s^k + \int_0^1 [\nabla(f(x^k + ts^k) - \nabla f(x^k))]^T s^k dt \right| \\
 &\leq c_1 \|s^k\|_2^2.
 \end{aligned}$$

where $c_1 = \frac{1}{2}(M + L)$.

For $k \geq k_0$,

$$-q_k \geq \frac{1}{2} \|p^k\|_\infty \min\left\{\frac{\|p^k\|_\infty}{M}, \delta_k\right\} \geq \frac{1}{2} \delta_k \|p^k\|_\infty \geq \frac{1}{2} \delta_k \epsilon,$$

and so

$$|r_k - 1| = \frac{|f(x^k + s^k) - f(x^k) - q_k|}{-q_k} \leq \frac{2c_1 \|s^k\|_2^2}{\epsilon \delta_k} \leq \frac{2c_1}{\epsilon} \delta_k$$

But then, for all k sufficiently large, $|r_k - 1| < 1 - \beta_2$, or $r_k > \beta_2$, so that $\delta_{k+1} = \min(\sigma_3 \delta_k, \delta_0)$. This contradicts our working assumption that $\delta_k \rightarrow 0$. \square

lemma 3. *If $\{f(x^k)\}$ is bounded below and $\|p^k\|_\infty > \epsilon$ for all k , then $\delta_k \rightarrow 0$ and $\{x^k\}$ converges.*

Proof. Set $U = \{k \mid \frac{f(x^{k+1}) - f(x^k)}{q_k} \geq \beta_1\}$ and $U_1 = \{k \mid k \in U, \delta_k \geq \frac{\|p^k\|_\infty}{M}\}$. For $k \in U_1$,

$$f(x^k) - f(x^{k+1}) \geq \beta_1 (-q_k) \geq \frac{\beta_1}{2M} \|p^k\|_\infty^2 \geq c_2 \epsilon^2$$

where $c_2 = \frac{\beta_1}{2M}$. Since $\{f(x^k)\}$ is non-increasing and bounded below, $f(x^k)$ is convergent, and so

$$\sum_{k \in U_1} c_2 \epsilon^2 \leq \sum_{k \in U_1} (f(x^k) - f(x^{k+1})) \leq \sum_k (f(x^k) - f(x^{k+1})) < \infty.$$

Therefore U_1 must be finite. Hence, there is a $k_0 \in \mathbb{N}$ such that $\delta_k \leq \|p^k\|_\infty / M$ for all $k \in \tilde{U} = \{k \geq k_0 \mid k \in U\}$. Consequently, for all $k \in \tilde{U}$,

$$f(x^k) - f(x^{k+1}) \geq \beta_1(-q_k) \geq \frac{\beta_1 \epsilon}{2} \delta_k.$$

Write $c_3 = \beta_1 \epsilon / 2$. Then

$$\sum_{k \in \tilde{U}} \delta_k \leq \frac{1}{c_3} \sum_{k \in \tilde{U}} (f(x^k) - f(x^{k+1})) < \infty.$$

For $k \notin U$, $\delta_{k+1} \leq \sigma_2 \delta_k$. Writing $\tilde{U} = \{k_1, k_2, \dots\}$ where $k_1 < k_2 < \dots$, we obtain for each $k_j \in \tilde{U}$ that

$$\sum_{k_j < k < k_{j+1}} \delta_k \leq \sum_{i=0}^{k_{j+1}-k_j} \sigma_2^i \delta_{k_j} \leq \frac{1}{1-\sigma_2} \delta_{k_j}.$$

Therefore

$$\sum_{k_0 \leq k} \|x^{k+1} - x^k\|_\infty \leq \sum_{k_0 \leq k} \delta_k \leq \frac{2}{1-\sigma_2} \sum_{k_j \in U} \delta_{k_j} < \infty,$$

which implies that the sequence $\{x^k\}$ is Cauchy and $\delta_k \rightarrow 0$. \square

We now prove the main global convergence result.

Theorem 4. *Assume A1 and A2 hold and let $\{x^k\}$ be a sequence generated by *ASTRAL* with $\tau_0 = 0$. Then either the algorithm terminates finitely at a first-order stationary point of \mathcal{P} , or $f(x^k) \downarrow -\infty$, or $\lim_{k \rightarrow \infty} p^k = 0$.*

Proof. We assume that $\{x^k\}$ is not finitely terminating and that $\{f(x^k)\}$ is bounded below, and show that $\lim_{k \rightarrow \infty} p^k = 0$. First note that these two assumptions and a combination of the previous two lemmas imply that $\liminf_{k \rightarrow \infty} \|p^k\|_\infty = 0$.

Let $k_0 \in \mathbb{N}$ be given and set $\epsilon_0 = \|p^{k_0}\|_\infty / 2$. Since the algorithm is not finitely terminating, $\epsilon_0 > 0$. If $L_1 > 1$ is an ∞ -norm Lipschitz constant for $p(x)$ on Ω , then $\|p(x)\| \geq \epsilon_0$ whenever $x \in \mathbb{B}(x^{k_0}, \rho_0) \cap \Omega$ where $\rho_0 = \epsilon_0 / L_1$ and $\mathbb{B}(x^{k_0}, \rho_0) = \{x \mid \|x - x^{k_0}\| \leq \rho_0\}$. Hence, if $\{x^k \mid k \geq k_0\} \subset \mathbb{B}(x^{k_0}, \rho_0) \cap \Omega$, then $\|p^k\|_\infty > \epsilon_0$ for $k > k_0$ contradicting the fact that $\liminf_{k \rightarrow \infty} p^k = 0$.

Let $l \geq k_0$ be such that x^{l+1} is the first iterate after x^{k_0} lying outside of $\mathbb{B}(x^{k_0}, r) \cap \Omega$ (in particular, $x^l \neq x^{l+1}$). Then

$$\begin{aligned} f(x^{k_0}) - f(x^{l+1}) &= \sum_{k=k_0}^l (f(x^k) - f(x^{k+1})) \\ &\geq \sum_{k=k_0, x^k \neq x^{k+1}}^l -\beta_0 q_k \\ &\geq \frac{\beta_0}{2} \sum_{k=k_0, x^k \neq x^{k+1}}^l \|p^k\|_\infty \min\left\{\frac{\|p^k\|_\infty}{M}, \delta_k\right\}. \end{aligned}$$

Case 1. If $\delta_k \leq \epsilon_0 \leq \|p^k\|_\infty$ for all $k_0 \leq k \leq l+1$ with $x^k \neq x^{k+1}$, then

$$\begin{aligned} f(x^{k_0}) - f(x^{l+1}) &\geq \frac{\beta_0}{2M} \epsilon_0 \sum_{k=k_0, x^k \neq x^{k+1}}^l \delta_k \\ &\geq \frac{\beta_0}{2M} \epsilon_0 \rho_0 \\ &\geq \frac{\beta_0}{2L_1 M} \epsilon_0^2. \end{aligned}$$

Case 2. If $\epsilon_0 < \delta_k \leq \|p^k\|_\infty$ for some $k_0 \leq k \leq l+1$ with $x^k \neq x^{k+1}$, then

$$f(x^{k_0}) - f(x^{l+1}) \geq \frac{\beta_0}{2M} \epsilon_0^2.$$

Case 3. If $\delta_k \geq \|p^k\|_\infty$ some $k_0 \leq k \leq l+1$ with $x^k \neq x^{k+1}$, then

$$f(x^{k_0}) - f(x^{l+1}) \geq \frac{\beta_0}{2M} \epsilon_0^2.$$

Therefore,

$$f(x^{k_0}) - f(x^{l+1}) \geq \frac{\beta_0}{2ML_1} \epsilon_0^2 = \frac{\beta_0}{4ML_1} \|p^{k_0}\|_\infty^2.$$

Now as $k_0 \rightarrow \infty$, we must have $\|p^{k_0}\|_\infty \rightarrow 0$ since the decreasing sequence $\{f(x^k)\}$ is bounded below. Therefore, $p^k \rightarrow 0$. \square

The following corollary is an immediate consequence of the convergence theorem and the continuity of $p(x)$ on Ω .

Corollary 5. *Under the assumptions of the previous theorem, any cluster point of the sequence $\{x^k\}$ is a first-order stationary point for \mathcal{P} .*

4. LOCAL CONVERGENCE

In this section we make use of three additional assumptions at a local solution \bar{x} to \mathcal{P} .

A3 (Non-Degeneracy) $-\nabla f(\bar{x}) \in \text{ri}(N_\Omega(\bar{x}))$.

A4 (Second-Order Sufficiency) $\Phi(\bar{x})^T \nabla^2 f(\bar{x}) \Phi(\bar{x})$ is positive definite.

A5 (Uniform Positive Definiteness) There exists $\rho > 0$ such that

$$\rho \|\Phi(\bar{x})^T s\|_2^2 \leq s^T \Phi(\bar{x}) \Phi(\bar{x})^T B_k \Phi(\bar{x}) \Phi(\bar{x})^T s \quad \text{for all } s \in \mathbb{R}^n.$$

The non-degeneracy condition in A3 is closely related to the strict complementary slackness condition of constrained optimization [2, Lemma 3.2] and is easily shown to be equivalent to the following componentwise conditions:

$$(1) \quad \begin{aligned} (\nabla f(\bar{x}))_i &> 0, & \text{whenever } x_i = l_i < u_i, \\ (\nabla f(\bar{x}))_i &< 0, & \text{whenever } x_i = u_i > l_i, \text{ and} \\ (\nabla f(\bar{x}))_i &= 0, & \text{whenever } l_i < x_i < u_i. \end{aligned}$$

The second-order sufficiency condition in A4 simply states that the Hessian is positive definite in the subspace spanned by the non-binding variables, i.e. $\text{Ran}(\Phi(\bar{x}))$. The uniform positive definite hypothesis A5 is used to assure the convergence of the iterates under assumptions A3 and A4.

Non-degeneracy has a number of important consequences that are best visualized using the *face* geometry of Ω . Recall that a subset F of a convex set $C \in \mathbb{R}^n$ is said to be a *face* of C if whenever $x, y \in C$ are such that $\text{ri}([x, y]) \cap F \neq \emptyset$, then $x, y \in F$. It is well-known that faces of a convex set are themselves convex and the relative interiors of all of the faces of a closed convex set form a partition of the set [20]. Therefore, every point x in C can be associated with a unique face of C , denoted $F_C(x)$, where $x \in \text{ri}(F_C(x))$. It is easily shown [2, Theorem 2.3] that for every face F of C we have

$$N_C(x) = N_C(\bar{x}) \quad \text{and} \quad T_C(x) = T_C(\bar{x}) \quad \text{for all } x, \bar{x} \in \text{ri}(F).$$

Hence, we define $N_C(F) = N_C(\bar{x})$ and $T_C(F) = T_C(\bar{x})$ for any $\bar{x} \in \text{ri}(F)$. If it is further assumed that C is polyhedral, then [2, Theorem 2.6]

$$(2) \quad x + \text{lin}(T_C(F)) = \text{aff}(F) \quad \forall x \in \text{ri}(F),$$

where recall that the lineality $\text{lin}(T_C(F))$ is the largest subspace contained in $T_C(F)$. From this it follows that [2, Theorem 2.8]

$$(3) \quad \text{int}(F + N_C(F)) = \text{ri}(F) + \text{ri}(N_C(F)) \neq \emptyset.$$

The geometric notion of a face is closely connected to active indices. For example, it is easily shown [2, Equation 2.7] that for $x \in \Omega$

$$F_\Omega(x) = \text{cl}\{y \in C \mid \mathcal{A}_l(y) = \mathcal{A}_l(x) \text{ and } \mathcal{A}_u(y) = \mathcal{A}_u(x)\},$$

where

$$\mathcal{A}_l(x) = \{i \mid x_i = l_i\} \quad \text{and} \quad \mathcal{A}_u(x) = \{i \mid x_i = u_i\}.$$

Hence, in particular, $\mathcal{A}^c(y) = \mathcal{A}^c(x)$ for all $y \in \text{ri}(F_\Omega(x))$, where

$$\mathcal{A}^c(x) = \{i \mid l_i < x_i < u_i\},$$

and so

$$\text{aff}(F_\Omega(y)) = \text{aff}(F_\Omega(x)) = x + \text{Ran}(\Phi(x)) \quad \text{for all } y \in \text{ri}(F_\Omega(x)).$$

Another important object in our convergence analysis is the *projected gradient*:

$$\nabla_\Omega f(x) = P_{T_\Omega(x)}(-\nabla f(x)).$$

The first-order optimality condition $-\nabla f(x) \in N_\Omega(x)$ is equivalent to the condition $\nabla_\Omega f(x) = 0$. However, the projected gradient is problematic since it is not a continuous function of x . Nonetheless, it is still quite useful. For example, in [2, Corollary 3.6] it is shown that if $\{x^k\}$ is any sequence converging to a non-degenerate solution \bar{x} to \mathcal{P} , then the active constraints at \bar{x} are identified in a finite number of iterations if and only if $\nabla_\Omega f(x^k) \rightarrow 0$. Insight into why this is true is given in the next 2 lemmas which expand on the local face geometry developed in [2].

lemma 6. *Let C be a nonempty polyhedral convex subset of \mathbb{R}^n , and let $\bar{x} \in C$ and $\bar{w} \in \text{ri}(N_C(\bar{x}))$. Set $F = F_C(\bar{x})$, $T = T_C(F)$, $N = N_C(F)$, $A = \text{aff}(F)$ and $S = \text{lin}(T_C(F))$. Then there is an $\epsilon > 0$ such that for all $v \in (\bar{x} + \bar{w}) + \epsilon\mathbb{B}$*

$$(4) \quad P_C(v) = P_A(v) = \bar{x} + P_S(v - \bar{x}),$$

$$(5) \quad P_T(v - \bar{x}) = P_S(v - \bar{x}),$$

and

$$(6) \quad P_C(v) \in \text{ri}(F) \quad \text{and} \quad v - P_C(v) \in \text{ri}(N)$$

with

$$(7) \quad \|P_C(v) - \bar{x}\| \leq \epsilon \quad \text{and} \quad \|(v - P_C(v)) - \bar{w}\| \leq \epsilon.$$

Proof. We do not prove the second equivalence in (4) as it is an elementary property of affine sets. By (3), there is an $\epsilon > 0$ such that $(\bar{x} + \bar{w}) \in \text{int}(F + N_C(F))$. Reduce ϵ if necessary so that

$$(\bar{x} + \epsilon\mathbb{B}) \cap A \subset \text{ri}(F) \quad \text{and} \quad (\bar{w} + \epsilon\mathbb{B}) \cap \text{aff}(N) \subset \text{ri}(N).$$

Let $v \in (\bar{x} + \bar{w}) + \epsilon\mathbb{B}$, and write

$$\hat{x} = P_C(v) \quad \text{and} \quad \hat{w} = v - P_C(v).$$

Since both P_C and $I - P_C$ are contractions,

$$\|\hat{x} - \bar{x}\| = \|P_C(v) - P_C(\bar{x} + \bar{w})\| \leq \|v - (\bar{x} + \bar{w})\| \leq \epsilon$$

and

$$\|\hat{w} - \bar{w}\| = \|(v - P_C(v)) - (\bar{x} + \bar{w} - P_C(\bar{x} + \bar{w}))\| \leq \|v - (\bar{x} + \bar{w})\| \leq \epsilon.$$

Hence, since $P_C(x + w) = x$ for all $x \in F$ and $w \in N$, we must have $\hat{x} = P_C(v) \in \text{ri}(F)$ and $\hat{w} = v - P_C(v) \in \text{ri}(N)$ by the choice of ϵ . This establishes (6) and (7).

By [2, Lemma 2.7] and (2),

$$A - \bar{x} = S = \text{span}(N)^\perp,$$

and so

$$\hat{x} - \bar{x} \in S \quad \text{and} \quad \hat{w} \in S^\perp.$$

Therefore,

$$\bar{x} - P_A(v) = \bar{x} - (\bar{x} + P_S(v - \bar{x})) = -P_S(\hat{x} - \bar{x} + \hat{w}) = \bar{x} - \hat{x} = \bar{x} - P_C(v),$$

and so $P_C(v) = P_A(v)$ establishing the first equation in (4).

It remains only to establish (5). Since $S \subset T$, (5) follows once it is shown that $P_T(v - \bar{x}) \in S$. Using the notation above, $v - \bar{x} = (\hat{x} - \bar{x}) + \hat{w}$ with $\hat{x} - \bar{x} \in T$ and $\hat{w} \in N$, and by (4), $\langle \hat{x} - \bar{x}, \hat{w} \rangle = 0$. Hence, $(\hat{x} - \bar{x}) + \hat{w}$ is the Moreau decomposition of $v - \bar{x}$, giving $P_T(x - \bar{x} + w) = \hat{x} - \bar{x}$. This proves (5) since we have also shown that $P_S(v - \bar{x}) = \hat{x} - \bar{x}$. \square

lemma 7. *If \bar{x} is a non-degenerate stationary point for \mathcal{P} , then there exists $\epsilon > 0$ such that*

- (a) $\mathcal{A}_l(x) = \mathcal{A}_l(\bar{x})$ and $\mathcal{A}_u(x) = \mathcal{A}_u(\bar{x})$ for all $x \in (\bar{x} + \epsilon\mathbb{B}) \cap F_\Omega(\bar{x})$,
- (b) $\mathcal{A}^c(x) = \mathcal{B}^c(x) = \mathcal{A}^c(\bar{x})$ for all $x \in (\bar{x} + \epsilon\mathbb{B}) \cap F_\Omega(\bar{x})$,
- (c) $\Phi(x) = \Phi(\bar{x})$ for all $x \in (\bar{x} + \epsilon\mathbb{B}) \cap F_\Omega(\bar{x})$,
- (d) $\text{Ran}(\Phi(x)) = \text{aff}(F_\Omega(\bar{x})) - \bar{x}$ for all $x \in (\bar{x} + \epsilon\mathbb{B}) \cap F_\Omega(\bar{x})$,
- (e) $P_\Omega(x - \nabla f(\bar{x})) \in F_\Omega(\bar{x})$ for all $x \in \bar{x} + \epsilon\mathbb{B}$,
- (f) $\nabla_\Omega f(x) = -\Phi(x)\Phi(x)^T \nabla f(x)$ for all $x \in (\bar{x} + \epsilon\mathbb{B}) \cap F_\Omega(\bar{x})$, and
- (g) $P_\Omega(x - \nabla f(x)) = x + \nabla_\Omega f(x)$ for all $x \in (\bar{x} + \epsilon\mathbb{B}) \cap F_\Omega(\bar{x})$.
- (h) $p(x) = \nabla_\Omega f(x)$ for all $x \in (\bar{x} + \epsilon\mathbb{B}) \cap F_\Omega(\bar{x})$.

Proof. Set $F = F_\Omega(\bar{x})$ and $\Phi = \Phi(\bar{x})$, and let $\epsilon > 0$ be chosen so that (4)-(7) in Lemma 6 hold with $C = \Omega$ and $\bar{w} = -\nabla f(\bar{x})$. Then (e)-(h) follow since $\Phi\Phi^T$ is the orthogonal projection onto $\text{aff}(F_\Omega(\bar{x})) - \bar{x}$, and (f) and (g) are equivalent by definition.

Non-degeneracy and the continuity of ∇f imply that we may reduce $\epsilon > 0$ if necessary so that

$$x - \nabla f(x) \in \text{int}(F + N_\Omega(F)) \quad \text{whenever } x \in (\bar{x} + \epsilon\mathbb{B}) .$$

(a)-(d) By [2, Equation 2.7], a point $x \in \Omega$ satisfies $\mathcal{A}_l(x) = \mathcal{A}_l(\bar{x})$ and $\mathcal{A}_u(x) = \mathcal{A}_u(\bar{x})$ if and only if $x \in \text{ri}(F)$. Thus, in particular, $\mathcal{A}^c(x) = \mathcal{A}^c(\bar{x})$ for all $x \in \text{ri}(F)$. Next note that

$$\mathcal{B}^c(x) = \mathcal{A}^c(x) \cup \mathcal{B}_1^c(x) \cup \mathcal{B}_2^c(x),$$

where

$$\mathcal{B}_1^c(x) = \{i \mid x_i = l_i < u_i \text{ and } (\nabla f(x))_i \leq 0\}$$

and

$$\mathcal{B}_2^c(x) = \{i \mid x_i = u_i > l_i \text{ and } (\nabla f(x))_i \geq 0\}.$$

But, by (1) and the continuity of ∇f , there is an $\hat{\epsilon} > 0$ such that the sets $\mathcal{B}_1^c(x)$ and $\mathcal{B}_2^c(x)$ are both empty for $x \in (\bar{x} + \hat{\epsilon}\mathbb{B}) \cap F$, hence, if $\hat{\epsilon} < \epsilon$, we simply reset ϵ to $\hat{\epsilon}$. In particular, (c) and (d) follow. \square

The following lemma parallels a similar result in [4]. It shows that any non-degenerate stationary point of \mathcal{P} at which the second-order sufficiency condition is satisfied is an attractor for the ASTRAL algorithm.

lemma 8. *Suppose \bar{x} is a non-degenerate stationary point for \mathcal{P} that satisfies the second-order sufficient condition (A4). Then given $\rho > 0$ there exists a $\bar{\epsilon} > 0$ such that for any symmetric matrix $B \in \mathbb{R}^{n \times n}$, $x \in (\bar{x} + \bar{\epsilon}\mathbb{B}_2) \cap F_\Omega(\bar{x})$, and $\hat{x} \in F_\Omega(\bar{x})$ satisfying*

- (a) $\rho \|z\|_2^2 \leq z^T \Phi(x)^T B \Phi(x) z$ for all $z \in \mathbb{R}^{\nu(x)}$, where $\nu(x) = n - |\mathcal{B}(x)|$,
- (b) $\nabla f(x)^T s + \frac{1}{2} s^T B s \leq 0$, and
- (c) $f(\hat{x}) \leq f(x)$,

where $s = \hat{x} - x$, it must be the case that $\hat{x} \in (\bar{x} + \bar{\epsilon} \mathbb{B}_2) \cap F_\Omega(\bar{x})$.

Proof. Set $F = F_\Omega(\bar{x})$, $\Phi = \Phi(\bar{x})$, and $\Pi = \Phi \Phi^T$. Recall that Π is the orthogonal projector onto $\text{Ran}(\Phi)$. With no loss in generality, we may reduce ρ if necessary so that the hypotheses, Lemma 7, and (A4-A5) imply the existence of an $\epsilon > 0$ such that $(\bar{x} + \epsilon \mathbb{B}_2) \cap F \subset \text{ri}(F)$ and

$$\begin{aligned} \Phi(x) &= \Phi, \quad \nabla_\Omega f(x) = -\Pi \nabla f(x), \\ \rho \|s\|_2^2 &\leq \min\{s^T \Pi B \Pi s, s^T \Pi \nabla^2 f(x) \Pi s\}, \text{ and} \\ f(x) &\geq f(\bar{x}) + \nabla f(\bar{x})^T (x - \bar{x}) + \frac{\rho}{2} \|x - \bar{x}\|_2^2 \end{aligned}$$

for all $x \in (\bar{x} + \epsilon \mathbb{B}_2) \cap F$. The final inequality above can be simplified to

$$f(x) \geq f(\bar{x}) + \frac{\rho}{2} \|x - \bar{x}\|_2^2$$

since $\nabla_\Omega f(x) = -\Pi \nabla f(x)$ and so

$$\nabla f(\bar{x})^T (x - \bar{x}) = \nabla f(\bar{x})^T \Pi (x - \bar{x}) = -\nabla_\Omega f(\bar{x})^T (x - \bar{x}) = 0.$$

Moreover, since $\nabla_\Omega f(x) = -\Pi \nabla f(x)$ is continuous on $x \in (\bar{x} + \epsilon \mathbb{B}_2) \cap F$ and $\nabla_\Omega f(\bar{x}) = 0$, we may choose $\epsilon > 0$ so that $\|\nabla_\Omega f(x)\|_2 \leq \rho\epsilon/4$ whenever $x \in (\bar{x} + \frac{1}{2}\epsilon \mathbb{B}_2) \cap F$.

Set $\bar{\epsilon} = \frac{1}{2}\epsilon$ and suppose $x \in (\bar{x} + \bar{\epsilon} \mathbb{B}) \cap F$ and $\hat{x} \in F$ satisfy (a)-(c). Then (a) implies that $\hat{x} - x = s = \Pi s$ satisfies

$$\frac{1}{2} s^T B s = \frac{1}{2} s^T \Pi B \Pi s \geq \frac{\rho}{2} \|\Phi^T s\|^2 = \frac{\rho}{2} \|s\|_2^2.$$

This inequality and (b) give

$$\frac{\rho}{2} \|s\|_2^2 \leq \frac{1}{2} s^T B s \leq -s^T \nabla f(x) = -s^T \Pi \nabla f(x) \leq \|s\|_2 \|\nabla_\Omega f(x)\|_2,$$

and so $\frac{\rho}{2} \|s\|_2 \leq \|\nabla_\Omega f(x)\|_2$, which in turn implies that $\|s\|_2 \leq \frac{\epsilon}{2}$. Therefore $\|\hat{x} - \bar{x}\|_2 \leq \epsilon$. Finally, (c) implies that

$$\begin{aligned} f(x) &\geq f(\hat{x}) \\ &\geq f(\bar{x}) + \frac{\rho}{2} \|\hat{x} - \bar{x}\|_2^2 \\ &= f(x) + \nabla f(x)^T (\bar{x} - x) + \frac{1}{2} (\bar{x} - x)^T \nabla^2 f(z) (\bar{x} - x) \\ &\geq f(x) - \nabla_\Omega f(x)^T (\bar{x} - x) + \frac{\rho}{2} \|\hat{x} - \bar{x}\|_2^2, \end{aligned}$$

for some z between x and \bar{x} . Since $\hat{x} \in (\bar{x} + \epsilon \mathbb{B}_2) \cap F$ and $x \in (\bar{x} + \frac{\epsilon}{2} \mathbb{B}_2) \cap F$, this inequality gives

$$\|\hat{x} - \bar{x}\|_2 \leq \frac{1}{2} \epsilon = \bar{\epsilon}.$$

□

The main local convergence result follows.

Theorem 9. *Suppose A1-A2 hold and let $\{x^k\}$ be a sequence generated by *ASTRAL* with $\tau_0 > 0$. Let \bar{x} is a cluster point of the subsequence $J = \{k \mid \|p^{k-1}\| \leq \tau_{k-1}\}$. By Corollary 5, \bar{x} is a first-order stationary point of \mathcal{P} , that is $-\nabla f(\bar{x}) \in N_\Omega(\bar{x})$. If, in addition, \bar{x} and the sequence of matrices $\{B_k\}$ satisfy A3-A5, then a gradient projection re-start ($x^{k+1} = \hat{x}$ in Step 1 of *ASTRAL*) occurs at most finitely many times, $x^k \rightarrow \bar{x}$, and $\nabla_\Omega f(x^k) \rightarrow 0$.*

Proof. Let $\bar{\epsilon} > 0$ be chosen so that the conclusions of Lemmas 7 and 8 hold at \bar{x} with ρ as in (A5) uniformly over $\{B_k\}$. Set $F = F_\Omega(\bar{x})$, $\Phi = \Phi(\bar{x})$, and $\Pi = \Phi\Phi^T$. Since \bar{x} is a cluster point for $\{x^k\}_J$, by Part (e) of Lemma 7, there is a $k_0 \in J$ such that $x^{k_0} \in (\bar{x} + \bar{\epsilon}\mathbb{B}_2) \cap F$, and, by construction,

$$\nabla f(x^{k_0})^T s^{k_0} + \frac{1}{2}(s^{k_0})^T \nabla^2 f(x^{k_0}) s^{k_0} \leq 0 \quad \text{and} \quad f(x^{k_0+1}) \leq f(x^{k_0})$$

(note that it is possible that $s^{k_0} = 0$ and $x^{k_0+1} = x^{k_0}$). Hence, by Lemma 8, $x^{k_0+1} \in (\bar{x} + \bar{\epsilon}\mathbb{B}_2) \cap F$, and $\mathcal{B}^c(x^{k_0+1}) = \mathcal{B}^c(\bar{x}) = \mathcal{A}^c(\bar{x})$. By induction, we have $x^k \in (\bar{x} + \bar{\epsilon}\mathbb{B}_2) \cap F$ for $k \geq k_0$. Thus, the gradient projection re-start is never activated in Step 1 of the algorithm for $k \geq k_0$. By Theorem 4, $p^k \rightarrow 0$, hence, by Part (g) of Lemma 7, $\nabla_\Omega f(x^k) = p^k \rightarrow 0$ since $\{x^k \mid k \geq k_0\} \subset (\bar{x} + \bar{\epsilon}\mathbb{B}_2) \cap F$. Finally, since the same argument holds for all $0 < \epsilon \leq \bar{\epsilon}$ (for increasing values k_0), we find that $x^k \rightarrow \bar{x}$. \square

The fact that $\nabla_\Omega f(x^k) \rightarrow 0$ is significant since it is shown in [2, Theorem 3.4] that this occurs if and only if the sets $\mathcal{A}_l(\bar{x})$ and $\mathcal{A}_u(\bar{x})$ are identified after a finite number of iterations. Thus, if we take $B_k = \nabla^2 f(x^k)$, the algorithm should locally behave like Newton's method and achieve quadratic convergence. This is the content of our final convergence result whose proof parallels that of a similar result in [18].

Corollary 10. *Suppose $B_k = \nabla^2 f(x^k)$ whenever $\Phi(x^k)^T \nabla^2 f(x^k) \Phi(x^k)$ is positive definite and the assumptions of Theorem 9 hold with the further assumption that $\nabla^2 f(\bar{x})$ is locally Lipschitz at \bar{x} with local Lipschitz constant L_2 . Then the *ASTRAL* iterates $\{x^k\}$ converge quadratically to \bar{x} .*

Proof. Set $F = F_\Omega(\bar{x})$, $\Phi = \Phi(\bar{x})$, $\Pi = \Phi\Phi^T$, and $\nu = \nu(\bar{x}) = |\mathcal{B}^c(\bar{x})|$. We begin by showing that the trust-region constraint is eventually inactive.

By Theorem 9, there exists a $k_0 > 0$, such that for $k > k_0$, we have $x^k \in F$, and the gradient projection re-start is never activated in Step 1 of the algorithm. Therefore $\Pi s^k = s^k$ for $k > k_0$. Moreover, there exist $\epsilon > 0$, $\rho > 0$, such that

$$\rho \|z\|_2^2 \leq z^T \Phi^T \nabla^2 f(x) \Phi z \quad \text{for all } z \in \mathbb{R}^\nu \text{ and } x \in (\bar{x} + \epsilon\mathbb{B}) \cap F$$

and, by Lemma 7,

$$p(x) = \nabla_\Omega f(x) = -\Pi \nabla f(x) \quad \text{for all } x \in (\bar{x} + \epsilon\mathbb{B}_2) \cap F.$$

We may assume that k_0 is sufficiently large to ensure that $x^k \in (\bar{x} + \epsilon \mathbb{B}_2) \cap F$ for all $k > k_0$, and so

$$\frac{\rho}{2} \|s^k\|_2^2 \leq \frac{1}{2} s^{kT} \nabla^2 f(x^k) s^k \leq \|\nabla_{\Omega} f(x)\|_2 \|s^k\| \quad \text{for all } k \geq k_0.$$

Hence $\frac{\rho}{2} \|s\|_2 \leq \|\nabla_{\Omega} f(x^k)\|_2$, which implies

$$\frac{\rho}{2\sqrt{n}} \|s\|_{\infty} \leq \|\nabla_{\Omega} f(x^k)\|_{\infty} = \|p^k\|_{\infty}.$$

By Lemma 1, we now have

$$\begin{aligned} -(\nabla f(x^k)^T s^k + \frac{1}{2} s^{kT} \nabla^2 f(x^k) s^k) &\geq \frac{1}{2} \|p^k\|_{\infty} \min\left\{\frac{\|p^k\|_{\infty}}{M}, \delta_k\right\} \\ &\geq \frac{\rho}{4\sqrt{n}} \|s^k\|_{\infty} \min\left(\frac{\rho \|s^k\|_{\infty}}{2M\sqrt{n}}, \|s^k\|_{\infty}\right) \\ &\geq c \|s^k\|_{\infty}^2. \end{aligned}$$

where $c = \frac{\rho}{4\sqrt{n}} \min(\frac{\rho}{2M\sqrt{n}}, 1)$. Consequently

$$\begin{aligned} |r_k - 1| &= \left| \frac{f(x^k + s^k) - (f(x^k) + \nabla f(x^k)^T s^k + \frac{1}{2} s^{kT} \nabla^2 f(x^k) s^k)}{-(\nabla f(x^k)^T s^k + \frac{1}{2} s^{kT} \nabla^2 f(x^k) s^k)} \right| \\ &\leq \frac{1}{2} \left| \frac{\int_0^1 (s^k)^T \nabla^2 f(x + ts^k) s^k - (s^k)^T \nabla^2 f(x^k) s^k dt}{c \|s^k\|_{\infty}^2} \right| \\ &\leq \frac{L_2 \|s^k\|_2^3}{4c \|s^k\|_{\infty}^2}. \end{aligned}$$

Hence $|r_k - 1| \rightarrow 0$ since $s^k \rightarrow 0$, and so the trust region bound becomes inactive after finite number iterations. Therefore, by Theorem 9, there is a k_0 such that the trust-region radius is inactive, $s^k \in \text{Ran}(\Phi)$, and $x^k \in \text{ri}(F)$ for all $k \geq k_0$. Hence, for all $k \geq k_0$, $x^{k+1} = x^{k_0} + \Phi z^{k+1}$, where $z^{k_0} = 0$ and z^{k+1} solves the problem

$$\min_{z \in \mathbb{R}^{\nu}} (\Phi^T \nabla f(x^k))^T (z - z^k) + \frac{1}{2} (z - z^k)^T (\Phi^T \nabla^2 f(x^k) \Phi) (z - z^k).$$

That is, the iterates reduce to the Newton iterates applied to the function $\phi(z) = f(x^{k_0} + \Phi z)$ initiated at $z^{k_0} = 0$. Since the iterates converge to a stationary point \bar{z} (necessarily unique) at which $\bar{x} = x^{k_0} + \Phi \bar{z}$, $\nabla \phi(\bar{z}) = \Phi^T \nabla f(\bar{x}) = 0$, and $\nabla^2 \phi(\bar{z}) = \Phi^T \nabla^2 f(\bar{x}) \Phi$ is positive definite, it must be the case that $z^k \rightarrow \bar{z}$ quadratically. Consequently, $x^k \rightarrow \bar{x}$ quadratically as well. \square

5. L-BFGS AND AN INTERIOR POINT METHOD FOR THE QP
SUBPROBLEMS $(TR)_k$

The key to the success of `ASTRAL` comes from the observation that limited memory BFGS updating nicely dovetails with the interior point methodology for solving the QP subproblems $(TR)_k$. In this section we illustrate this observation in some detail.

Our subproblem template has the form

$$(8) \quad \begin{aligned} \mathcal{Q} \quad & \text{minimize} && \frac{1}{2}w^T\Phi^TB\Phi w + k^T w \\ & \text{subject to} && 0 \leq w \leq h. \end{aligned}$$

where $k = \Phi^TB\Phi\Phi^Tl + \Phi^Tg$. So the first step in solving the subproblems $(TR)_k$ is to recast them into this form. The Hessian approximation $B \approx \nabla^2 f(x)$ is based on L-BFGS technology [6, 16, 21]. The user chooses an integer m and obtains an approximate Hessian by updating an initial Hessian approximation using the difference of m previous successive iterates and gradients. More specifically, suppose $\{x^k\}$ is the sequence of iterates produced by `ASTRAL` with associated gradient sequence $g^k = \nabla f(x^k)$. Define

$$s^k = x^{k+1} - x^k \text{ and } y^k = g^{k+1} - g^k$$

for all k . Select iterations k_1, k_2, \dots, k_m (usually the m most recent iterates $k_1 = k - m + 1, \dots, k_m = k$) and set

$$S = [s^{k_1}, \dots, s^{k_m}], \quad Y = [y^{k_1}, \dots, y^{k_m}], \quad \text{and } S^TY = L + D + R$$

where L is strictly lower triangular, D is diagonal, and R is strictly upper triangular. Under the assumption that $s^{k_i T} y^{k_i} > 0$ for $i = 1, \dots, m$ (D is positive definite), the L-BFGS Hessian approximation at x^k is given by

$$(9) \quad B = \lambda I - \Psi\Gamma^{-1}\Psi^T$$

for some $\lambda > 0$, where

$$\Psi = [Y, \lambda S] \quad \text{and} \quad \Gamma = \begin{bmatrix} -D & L^T \\ L & \lambda S^T S \end{bmatrix}.$$

The positive definiteness of D guarantees the positive definiteness of B and the invertibility of Γ [5]. In our implementation, we use

$$\lambda = \frac{\|y^{k_m}\|^2}{y^{k_m T} s^{k_m}}$$

as suggested by Shanno and Phua in [21].

In the remainder of this section we focus on the interplay between the L-BFGS update and the linear algebra associated with the Newton step in a primal-dual version of the interior point method. It is this interplay that provides the greatest efficiencies in our implementation. The remaining details of the interior point algorithm we employ are given in [22].

The relaxed Karush-Kuhn-Tucker (KKT) conditions for \mathcal{Q} can be written as

$$(10) \quad \begin{bmatrix} \Phi^T B \Phi & I \\ -I & 0 \end{bmatrix} \begin{bmatrix} w \\ v \end{bmatrix} + \begin{bmatrix} k \\ h \end{bmatrix} = \begin{bmatrix} z \\ u \end{bmatrix},$$

$$(11) \quad ZW e = t e,$$

$$(12) \quad UV e = t e,$$

$$(13) \quad u, v, w, z \geq 0.$$

where $e = [1, 1, \dots, 1]^T$ is the vector of all ones of the appropriate dimension and U, V, W, Z are diagonal matrices of u, v, w, z respectively. The KKT conditions are obtained by taking $t = 0$.

The set of componentwise strictly positive solutions to (10)-(13) for $t > 0$ is called the *central path*, and a primal-dual interior point algorithm for \mathcal{Q} is an algorithm that attempts to approximately follow this central path as $t \downarrow 0$ to a solution to \mathcal{Q} [23]. All points u, v, w, z on the central path are necessarily componentwise strictly positive due to the equations (11) and (12). At each iteration one computes a Newton step for the system of equations (10)-(13) and computes a steplength in that direction which, among other things, preserves the strict positivity of all of the iterates. After some reformulation, the equations for the Newton step can be written as

$$(14) \quad \begin{aligned} \Delta w &= -w + (\Phi^T B \Phi + U^{-1}V + W^{-1}Z)^{-1}(z - v - k - tU^{-1}e \\ &\quad + U^{-1}Vh + tW^{-1}e), \\ \Delta u &= -w + h - u - \Delta w, \\ \Delta v &= tU^{-1}e - v - U^{-1}V\Delta s, \\ \Delta z &= -W^{-1}Z\Delta w + tW^{-1}e - z. \end{aligned}$$

Here we use the fact that the matrix B is positive definite which implies that the matrix $\Phi^T B \Phi + U^{-1}V + W^{-1}Z$ is also positive definite. In this formulation, the major numerical difficulty is the computation of Δw . That is, we need to develop an efficient method for solving equations involving the matrix $\Phi^T B \Phi + U^{-1}V + W^{-1}Z$. We do this by combining the Sherman-Morrison-Woodbury formula with some of the linear algebra techniques developed in [6].

Recall that the matrix B is positive definite and the matrix Γ is assumed invertible. Definition (9) gives

$$\Phi^T B \Phi + U^{-1}V + W^{-1}Z = (\lambda I + U^{-1}V + W^{-1}Z) - (\Phi^T \Psi) \Gamma^{-1} (\Phi^T \Psi)^T,$$

where $\lambda I + U^{-1}V + W^{-1}Z$ is positive definite and diagonal. Denote this diagonal matrix by Ξ , and write

$$\begin{aligned} S_\Phi &= \Phi^T S, \\ Y_\Phi &= \Phi^T Y, \\ \Psi_\Phi &= \Phi^T \Psi. \end{aligned}$$

Applying the Sherman-Morrison-Woodbury formula gives

$$(\Phi^T B \Phi + U^{-1}V + W^{-1}Z)^{-1} = \Xi^{-1} + \Xi^{-1} \Psi_\Phi (\Gamma - \Psi_\Phi^T \Xi^{-1} \Psi_\Phi)^{-1} (\Psi_\Phi)^T \Xi^{-1}$$

whenever the matrix $\Gamma - \Psi_\Phi^T \Xi^{-1} \Psi_\Phi$ is invertible.

We now adapt the approach developed in [6] to evaluate the Newton step (14).

Proposition 11. *If the matrix Γ is invertible, then $\Gamma - \Psi_\Phi^T \Xi^{-1} \Psi_\Phi$ is invertible if and only if $\Phi^T B \Phi + U^{-1}V + W^{-1}Z = \Xi - \Psi_\Phi \Gamma^{-1} \Psi_\Phi^T$ is invertible.*

Proof. Let us first assume that $\Xi - \Psi_\Phi \Gamma^{-1} \Psi_\Phi^T$ is invertible, and let $v \in \mathbb{R}^{2m}$ be such that $(\Gamma - \Psi_\Phi^T \Xi^{-1} \Psi_\Phi)v = 0$. Then $\Gamma v = \Psi_\Phi^T \Xi^{-1} \Psi_\Phi v$ so that $v = \Gamma^{-1} \Psi_\Phi^T \Xi^{-1} \Psi_\Phi v$, and consequently, $\Psi_\Phi v = \Psi_\Phi \Gamma^{-1} \Psi_\Phi^T \Xi^{-1} \Psi_\Phi v$. Set $w = \Xi^{-1} \Psi_\Phi v$. We have $\Xi w = \Psi_\Phi \Gamma^{-1} \Psi_\Phi^T w$, and so $(\Xi - \Psi_\Phi \Gamma^{-1} \Psi_\Phi^T)w = 0$. The non-singularity of $\Xi - \Psi_\Phi \Gamma^{-1} \Psi_\Phi^T$ implies that $w = 0$, and so $\Psi_\Phi v = 0$. Therefore $\Gamma v = \Psi_\Phi^T \Xi^{-1} \Psi_\Phi v = 0$, which gives $v = 0$. Hence $\Gamma - \Psi_\Phi^T \Xi^{-1} \Psi_\Phi$ must be non-singular.

The reverse implication is established in the same way. \square

Although the matrix $\Gamma - \Psi_\Phi^T \Xi^{-1} \Psi_\Phi$ may be indefinite, a triangular factorization can be obtained from the Cholesky factors of a related matrix. To see this, write $\Gamma - \Psi_\Phi^T \Xi^{-1} \Psi_\Phi$ in its explicit form,

$$\Gamma - \Psi_\Phi^T \Xi^{-1} \Psi_\Phi = \begin{bmatrix} -D - Y^T \Theta Y & L^T - \lambda Y^T \Theta S \\ L - \lambda S^T \Theta Y & \lambda S^T S - \lambda^2 S^T \Theta S \end{bmatrix}.$$

where $\Theta = \Phi \Xi^{-1} \Phi^T$.

Note $D + Y^T \Theta Y$ is positive definite and so there exists a lower triangular matrix M , such that

$$D + Y^T \Theta^{-1} Y = M M^T.$$

Set

$$\hat{L} = L - \lambda S^T \Theta Y, \quad \hat{D} = D + Y^T \Theta Y, \quad \text{and} \quad \hat{W} = \lambda S^T (I - \lambda \Theta) S.$$

Proposition 12. *The $m \times m$ matrix*

$$(15) \quad \hat{W} + \hat{L} \hat{D}^{-1} \hat{L}^T = \begin{bmatrix} S^T & \hat{L} \end{bmatrix} \begin{bmatrix} \lambda(I - \lambda \Theta) & 0 \\ 0 & \hat{D}^{-1} \end{bmatrix} \begin{bmatrix} S \\ \hat{L}^T \end{bmatrix}$$

is positive definite.

Proof. The matrix $I - \lambda\Theta$ positive definite since $\Theta = \Phi(\lambda I + U^{-1}V + W^{-1}Z)^{-1}\Phi^T$. Hence the block diagonal matrix appearing on the right hand side of (15) is positive definite since \hat{D} is positive definite. Thus $\hat{W} + \hat{L}\hat{D}^{-1}\hat{L}^T$ is positive semidefinite. We now show that $\hat{W} + \hat{L}\hat{D}^{-1}\hat{L}^T$ is nonsingular. For this, suppose that the vector ω satisfies $(\hat{W} + \hat{L}\hat{D}^{-1}\hat{L}^T)\omega = 0$. Then both $S\omega = 0$ and $\hat{L}^T\omega = 0$ since the block diagonal matrix in (15) is positive definite. But then $0 = L^T\omega - \lambda Y^T\Theta S\omega = L^T\omega$ and $0 = Y^T S\omega = L^T\omega + R^T\omega = R^T\omega$. Since R is an upper triangular matrix with diagonal entries greater than 0, we must have $\omega = 0$. \square

If we now let J be a lower triangular matrix such that

$$JJ^T = \hat{W} + \hat{L}\hat{D}^{-1}\hat{L}^T,$$

then

$$\Gamma - \Psi^T\Xi^{-1}\Psi = \begin{bmatrix} M & 0 \\ -\hat{L}M^{-T} & J \end{bmatrix} \begin{bmatrix} -M^T & M^{-1}\hat{L}^T \\ 0 & J^T \end{bmatrix}.$$

That is, we can obtain a triangular factorization for $\Gamma - \Psi^T\Xi^{-1}\Psi$ by computing two $m \times m$ Cholesky factorizations. ASTRAL uses this factorization to obtain the update Δw in (14), which in turn allows us to easily compute Δu , Δv , and Δz .

6. NUMERICAL EXPERIMENTS

This section compares f77 implementations of ASTRAL and L-BFGS-B [5, 17] on a Linux workstation. The implementation of L-BFGS-B is downloaded from <http://www.ece.northwestern.edu/nocedal/lbfgsb.html>. The memory size equals 5 in our numerical experiments. The experiments were conducted in two runs. In the first run, we execute both ASTRAL and L-BFGS-B to get *fbest*, which is the best function value one can get from these two algorithms using the following four stopping criteria:

$$\begin{aligned} \|\nabla_{\Omega} f(x^k)\|_{\infty} &\leq 10^{-5} \\ \frac{|f^k - f^{k-1}|}{\max(1, |f^k|, |f^{k-1}|)} &\leq \epsilon_{mach} * 10^7 \\ nf &\leq 1000 \\ \|s^k\|_2 &\leq 10^{-10}, \end{aligned}$$

where ϵ_{mach} is machine- ϵ and nf is the total number of function evaluations. In the second run, we execute these two algorithms and the stopping criteria becomes

$$\frac{|f(x^k) - fbest|}{\max(|fbest|, 1)} \leq 10^{-5} \text{ and } nf \leq 1000.$$

In ASTRAL, we take

$$\begin{aligned}\beta_0 &= 10^{-5}, \quad \beta_1 = 0.2, \quad \beta_2 = 0.8, \\ \sigma_1 &= 0.25, \quad \sigma_2 = 0.5, \quad \sigma_3 = 2, \\ \tau_0 &= 10^{-3}, \quad \gamma = 0.5.\end{aligned}$$

The trust-region subproblem is solved using interior-point framework described in [22] with stopping criteria

$$\|r; cv\|_2 \leq 10^{-8}$$

where r is the residual of the Newton equation and cv is the complementarity vector. The homotopy parameter for the central path pc is updated by

$$pc = \begin{cases} 1, & \text{if } \alpha < 0.2; \\ \min(\sum(cv), 1 - \alpha) & \text{if } 0.2 \leq \alpha < 0.8; \\ \min(\sum(cv)/\nu, 1 - \alpha) & \text{if } \alpha \geq 0.8 \end{cases}$$

where α is the stepsize taken on the previous step.

Test set consists of 23 problems from CUTER [1, 12]. All problems have dimension at least 1000 with 21 of them having dimension at least 10,000. Both algorithms converge to the same local minimizer. The CPU time in seconds as well as the number of function and gradient evaluations from both algorithms are posted at www.math.washington.edu/1xu/research.

As stated in the introduction, our goal is to develop an algorithm that is efficient in its use of function and gradient evaluations at the expense of the more intensive linear algebra required to solve the trust-region subproblems. For this reason we first compare the performance of the two algorithms with respect to the sum of all function and gradient evaluations prior to termination. We do not give a higher weight to gradient evaluations. Relative performance is illustrated using the performance profile technique described in Dolan and Moré[11], i.e., for each algorithm, we plot the proportion P of the problems for which method is within factor τ of the smallest sum of the number of function and gradient evaluations(on a log scale). In Figure 1, we compare the performance of ASTRAL and L-BFGS-B using the 23 CUTER problems. Since the curve w.r.t ALTRAL is on the top of L-BFGS-B, it implies that these 23 test problems favor ALSTRAL over L-BFGS-B with respect to the number of function and gradient evaluations. To see that the profile is not biased by the requested precision, we also include the graph for the performance profile with termination criteria

$$\frac{|f(x^k) - f_{best}|}{\max(|f_{best}|, 1)} \leq 10^{-3} \text{ and } nf \leq 1000.$$

Next we compare the performance of these two algorithms with respect to CPU time. Note that the solution time consists of two part:

1. The time is used to evaluate the functions and gradients.
2. The remaining time is dominated by the cost of the linear algebra.

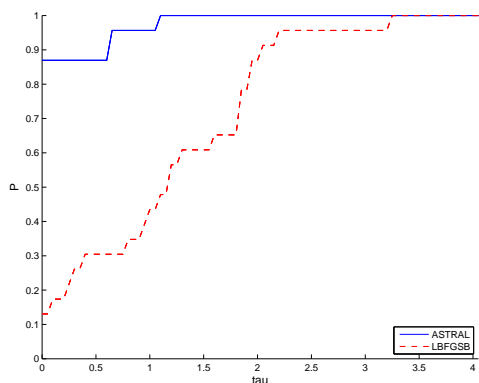


Figure 1a

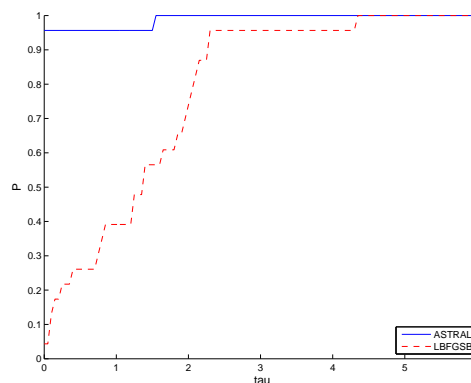


Figure 1b

Figure 1a. Performance profiles, sum of the function and gradient evaluations, relative accuracy 10^{-5} . Figure 1b. Performance profiles, sum of the function and gradient evaluations, relative accuracy 10^{-3} .

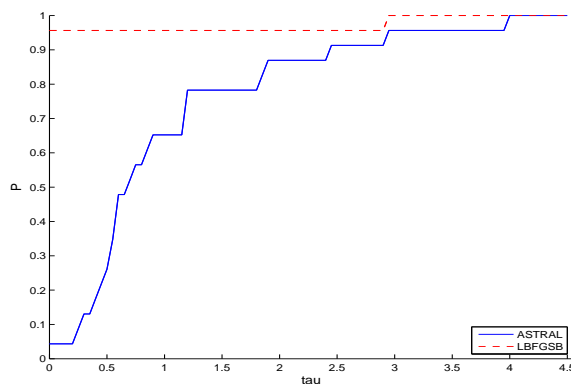


Figure 2. Performance profiles, CPU time.

For each test problem, we first evaluate the average cost per function and gradient evaluation as follows: randomly generate 50 feasible points, then evaluate the function and gradient 20 times for each of the points, record the total CPU time for all 50 points, and set the cost per evaluation to 1/1000 of the total CPU time. The average cost of each function evaluation in our test set is 0.002s, and the average cost of each gradient evaluation is 0.008s. We show three figures for comparing the algorithms based on CPU time. Figure 2 plots the performance profiles of the true CPU time used by both algorithms. Figures 3a and 3b present performance profiles with the CPU time devoted to function and gradient evaluations adjusted based on the computed average cost per test problem previously computed. In Figure 3a the CPU time for function and gradient evaluations is set to 10 times the computed average, and in Figure 3b, the CPU time for function and gradient evaluations is 20 times the computed average. Observe that when the cost in CPU time for function and gradient evaluations is relatively

small compared to linear algebra, L-BFGS-B outperforms ASTRAL in our test set. But when the cost of CPU time in function and gradient evaluations increases, ASTRAL shows comparable results.

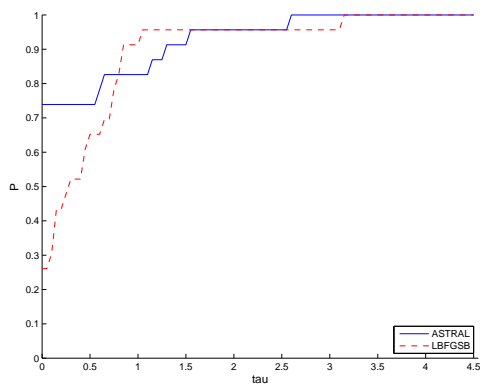


Figure 3a

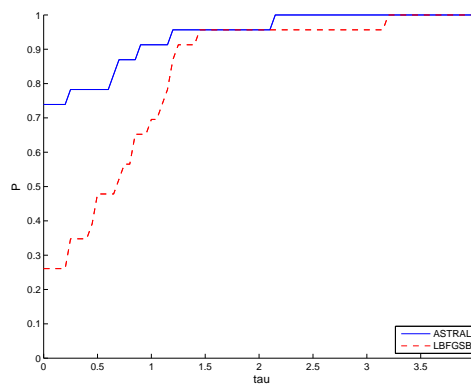


Figure 3b

Figure 3a. Performance profiles, function and gradient evaluation is 10 times the original cost. Figure 3b. Performance profiles, function and gradient evaluation is 20 times the original cost.

REFERENCES

- [1] I. Bongartz, A.R. Conn, N.I.M. Gould, and Ph.L. Toint. Cute: Constrained and unconstrained testing environment. *ACM Transactions on Mathematical Software*, 21(1):123–160, 1995.
- [2] James V. Burke and Jorge J. Moré. On the identification of active constraints. *SIAM J. NUMER. ANAL.*, 25(5), October 1988.
- [3] James V. Burke and A. Weigmann. Notes on limited memory bfgs updating in a trust-region framework. Technical report, University of Washington, Seattle, WA 98195, USA, 1997.
- [4] J.V. Burke, J.J. Moré, and G. Toraldo. Convergence properties of trust region methods for linear and convex constraints. *Math. Program.*, 47(3):305–336, 1990.
- [5] Richard H. Byrd, Peihuang Lu, Jorge Nocedal, and Ciyou Zhu. A limited memory algorithm for bound constrained optimization. *SIAM Journal on Scientific and Statistical Computing*, 16(5):1190–1208, 1995.
- [6] Richard H. Byrd, Jorge Nocedal, and Robert B. Schnabel. Representations of quasi-newton matrices and their use in limited memory methods. *Mathematical Programming*, 63(4):129–156, 1994.
- [7] A.R. Conn, N.I.M. Gould, and Ph. L. Toint. Global convergence of a class of trust region algorithms for optimization with simple bounds. *SIAM Journal on Numerical Analysis*, 25(2):433–460, 1988.
- [8] A.R. Conn, N.I.M. Gould, and Ph. L. Toint. Testing a class of methods for solving minimization problems with simple bounds on variables. *Mathematics of Computation*, 50(180):399–430, 1988.
- [9] A.R. Conn, N.I.M. Gould, and Ph. L. Toint. *Trust-Region Methods*. Society for Industrial and Applied Mathematics and Mathematical Programming Society, 2000.
- [10] A.R. Conn, N.I.M. Gould, and Ph.L. Toint. Testing a class of methods for solving minimization problems with simple bounds on the variables. *Mathematics of Computation*, 50, 1988.

- [11] E. D. Dolen and J.J. Moré. Benchmarking optimization software with performance profiles. Technical report, Mathematics and Computer Science, Argonne National Laboratory, Argonne, Illinois, USA, 2001.
- [12] N. I. M. Gould, D. Orban, and Ph.L. Toint. General cutter documentation. Technical report, CERFACS, 2002.
- [13] William W. Hager and Hongchao Zhang. A new active set algorithm for box constrained optimization. *SIAM Journal on Optimization*, 17(2):526–557, 2006.
- [14] Chih-Jen Lin and Jorge J. Moré. Newton’s method for large bound-constrained optimization problems. *SIAM Journal on Optimization*, 9(4):1100–1127, 1999.
- [15] D.R. Luke, J.V. Burke, and R.G. Lyon. Optical wavefront reconstruction: theory and numerical methods. *SIAM Rev.*, 44:196–224, 2002.
- [16] J. Nocedal. Updating quasi-newton matrices with limited storage. *Mathematics of Computation*, 35:773–782, 1980.
- [17] J. Nocedal, C. Zhu, R. Byrd, and P. Lu. Algorithm 778: L-bfgs-b, fortran routines for large scale bound constrained optimization. *ACM Transactions on Mathematical Software*, 23(4):550 – 560, 1997.
- [18] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 2006.
- [19] M.J.D. Powell. Convergence properties of a class of minimization algorithms. In O.L. Mangasarian, R.R. Meyer, and S.M. Robinson, editors, *Nonlinear Programming 2*, pages 1–27. Academic Press, 1975.
- [20] R.T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, NJ., 1970.
- [21] D.F. Shanno and K.H. Phua. Matrix conditioning and non-linear optimization. *Math. Programming*, (14):149–160, 1978.
- [22] Evangelia M. Simantiraki and David F. Shanno. An infeasible - interior - point method for linear complementarity problems. *SIAM J. OPTIM.*, 7(3):620 – 640, August 1997.
- [23] S. J. Wright. *Primal-Dual Interior-Point Methods*. SIAM, Philadelphia, USA, 1997.

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF WASHINGTON, SEATTLE, WA 98195, U.S.A.

E-mail address: lxu@math.washington.edu

1. DEPARTMENT OF MATHEMATICS, UNIVERSITY OF WASHINGTON, SEATTLE, WA 98195, U.S.A.

E-mail address: burke@math.washington.edu