

COMPUTING A TRUST REGION STEP FOR A PENALTY FUNCTION*

THOMAS F. COLEMAN†‡ AND CHRISTIAN HEMPEL†

Abstract. The problem of minimizing a quadratic function subject to an ellipsoidal constraint when the matrix involved is the Hessian of a quadratic penalty function (i.e., a function of the form $p(x) = f(x) + (1/2\mu)c(x)^T c(x)$) is considered. Most applications of penalty functions require $p(x)$ to be minimized for values of μ decreasing to zero. In general, as μ tends to zero the nature of finite precision arithmetic causes a considerable loss of information about the null space of the constraint gradients when $\nabla^2 p(x)$ is formed. This loss of information renders ordinary trust region Newton's methods unstable and degrades the accuracy of the solution to the trust region problem. The algorithm of Moré and Sorensen [*SIAM J. Sci. Statist. Comput.*, 4 (1983), pp. 553-572] is modified so as to be more stable and less sensitive to the nature of finite precision arithmetic in this situation. Numerical experiments clearly demonstrate the stability of the proposed algorithm.

Key words. equality constrained optimization, trust region, penalty function, extended system

AMS(MOS) subject classifications. 49D30, 49D37, 65K05, 65K10

1. Introduction. The generic trust region problem is

$$(1.1) \quad \begin{aligned} &\text{minimize} && \psi(s) = \frac{1}{2}s^T Hs + g^T s, \\ &\text{subject to} && \|s\| \leq \Delta \end{aligned}$$

where Δ is a positive number, $g \in R^n$, and $H \in R^{n \times n}$ is a symmetric matrix. We will concern ourselves with the case where H is the Hessian of an L_2 penalty function. That is, H is the Hessian of a function of the form

$$(1.2) \quad p(x) = f(x) + \frac{1}{2\mu} c(x)^T c(x)$$

where $f: R^n \rightarrow R$, $c: R^n \rightarrow R^l$ for $l \leq n$, and $\mu > 0$. For such functions the quadratic $\psi(s)$ is meant to serve as a local model at a current point x_c . Evaluating the first and second derivatives of $p(x)$ at x_c gives

$$(1.3) \quad g = \nabla p(x_c) = \nabla f(x_c) + \frac{1}{\mu} \nabla c(x_c) c(x_c)$$

and

$$(1.4) \quad H = \nabla^2 p(x_c) = \nabla^2 f(x_c) + \frac{1}{\mu} \sum_{i=1}^l \nabla^2 c_i(x_c) c_i(x_c) + \frac{1}{\mu} \nabla c(x_c) \nabla c(x_c)^T.$$

If we define $A = \nabla c(x_c)$ and

$$(1.5) \quad B = \nabla^2 f(x_c) + \frac{1}{\mu} \sum_{i=1}^l \nabla^2 c_i(x_c) c_i(x_c),$$

then

$$(1.6) \quad H = B + \frac{1}{\mu} AA^T.$$

* Received by the editors August 5, 1987; accepted for publication (in revised form) November 18, 1988. This research was supported by the Applied Mathematical Sciences Research Program (KC-04-02) of the Office of Energy Research of the U.S. Department of Energy under grant DE-FG02-86ER25013.A000.

† Center for Applied Mathematics, Cornell University, Ithaca, New York 14853.

‡ Department of Computer Science, Cornell University, Ithaca, New York 14853.

Matrices of the form (1.6) become ill conditioned as μ tends to zero. Also, when μ is small many of the digits of B will, in general, be lost when H is formed. This loss of information and ill-conditioning can have disastrous effects on methods that rely on having H . There is, as we shall see, a way to solve (1.1) that sidesteps these difficulties.

Our method is an adaptation of the method of Moré and Sorensen [1983] and retains their method's convergence properties. That is, our algorithm will compute a nearly optimal solution in finitely many iterations. More precisely, for $\sigma_1 \in (0, 1)$ the approximate solution s satisfies

$$(1.7) \quad \psi(s) - \psi(s^*) \leq \sigma_1(2 - \sigma_1)|\psi(s^*)| \quad \text{and} \quad \|s\| \leq (1 + \sigma_1)\Delta$$

where s^* is a solution to (1.1).

The paper is organized as follows. The next section briefly reviews the theoretical background needed to develop a suitable algorithm. Readers wishing a more thorough treatment should see Moré and Sorensen [1983]. The algorithm is constructed in § 3. Last, in § 4 we give a detailed report of numerical tests comparing TRSQPF with GQTPAR on trust region problems.

2. Structure of the problem. In this section we will show that the trust region problem usually reduces to a zero finding problem of a rational function with second-order poles. This perspective will enable us to construct an effective algorithm with good convergence properties. The next two lemmas provide necessary and sufficient conditions for a point $s \in R^n$ to be a solution to (1.1).

LEMMA 2.1. *If s is a solution to (1.1), then s is a solution to an equation of the form*

$$(2.2) \quad (H + \lambda I)s = -g$$

with $H + \lambda I$ positive semidefinite, $\lambda \geq 0$, and $\lambda(\Delta - \|s\|) = 0$.

LEMMA 2.3. *Let $s \in R^n$ satisfy (2.2) with $H + \lambda I$ positive semidefinite.*

- (i) *If $\lambda = 0$ and $\|s\| < \Delta$, then s solves (1.1);*
- (ii) *s solves $\psi(s) = \min \{\psi(w) : \|w\| = \|s\|\}$;*
- (iii) *If $\lambda \geq 0$ and $\|s\| = \Delta$, then s solves (1.1).*

For proofs of these lemmas see Sorensen [1982]. Moré and Sorensen [1983] also show that (1.1) has no solutions s with $\|s\| = \Delta$ if and only if H is positive definite and $\|H^{-1}g\| < \Delta$. The important consequence of Lemma 2.3 is that when (1.1) has a solution on the boundary of $\{w : \|w\| \leq \Delta\}$, then that solution can be found by solving

$$(2.4) \quad \|s_\lambda\|^2 - \Delta^2 = 0, \quad \text{where } s_\lambda = -(H + \lambda I)^{-1}g.$$

Thus, except for a special case, (1.1) has been reduced to a zero-finding problem in one variable, λ . The next crucial observation, due to Reinsch [1967], [1971] and Hebden [1973], is that (2.4) is a rational function in λ with second-order poles on a subset of the negatives of the eigenvalues of the matrix H . This follows easily from the Real-Schur decomposition:

$$(2.5) \quad H = Q\Lambda Q^T \quad \text{where } \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) \text{ and } Q^T Q = I.$$

Moreover, we assume that the eigenvalues are ordered so that $\lambda_1 \leq \dots \leq \lambda_n$. We now have

$$(2.6) \quad \|s_\lambda\|^2 = \|Q(\Lambda + \lambda I)^{-1}Q^T g\|^2 = \sum_{i=1}^n \frac{\gamma_i^2}{(\lambda_i + \lambda)^2}$$

where γ_i is the i th component of $Q^T g$. Let S_1 denote the eigenspace corresponding to the smallest eigenvalue λ_1 of H . If g has a component in S_1 , then (2.4) has a solution

in $(-\lambda_1, \infty)$, and this is called the "general case." If, however, g is perpendicular to S_1 , then (2.4) may not have a solution in $(-\lambda_1, \infty)$, and this leads to numerical difficulties. This situation is called the "hard case." As Moré and Sorensen point out, the characteristic difficulty of the hard case is that $\|s_\lambda\| < \Delta$ whenever $H + \lambda I$ is positive definite. In the hard case, a solution to (1.1) can be obtained by solving

$$(2.7) \quad (H - \lambda_1 I)s = -g$$

for s with $\|s\| < \Delta$ and by determining an eigenvector $z \in S_1$. Then

$$(2.8) \quad (H - \lambda_1 I)(s + \tau z) = -g$$

and $\|s + \tau z\| = \Delta$ for some τ . Lemma 2.3 now shows that $s + \tau z$ solves (1.1).

It is now easier to see how the ill-conditioning in matrices of the form (1.6) causes problems. Specifically, the solution to (2.2) becomes unstable as μ tends to zero. Gould [1986] offers a numerically stable way to solve such systems. Essentially the idea is to avoid using the true Hessian and instead compute the trust region steps with an extended matrix whose condition number is normally much better. If $H + \lambda I = B + (1/\mu)AA^T + \lambda I$, as in (1.6), we then set

$$(2.9) \quad X = \begin{pmatrix} B + \lambda I & A \\ A^T & -\mu I \end{pmatrix}$$

and if

$$(2.10) \quad X \begin{pmatrix} s_\lambda \\ r \end{pmatrix} = \begin{pmatrix} -\nabla f \\ -c \end{pmatrix}$$

for some $s_\lambda \in R^n$ and $r \in R^t$, then $(H + \lambda I)s = -\nabla f(x) - (1/\mu)Ac(x) = -g$. This way of solving (2.2) will play an instrumental role in the algorithm for solving (1.1). It is important to note that X does become ill conditioned as $\mu \rightarrow 0$ if A is nearly singular. While there is a way to cope with that situation we shall reserve the discussion for another paper. For now we assume that A has full column rank in a neighborhood of the minimum of $p(x)$. Every iteration of our algorithm will need to know if $H + \lambda I$ is positive definite. Therefore, the following definition and lemma are in order.

DEFINITION 2.11. The inertia of an $n \times n$ symmetric matrix H is the triple $\text{In}(H) = (a, b, c)$ where a, b, c are, respectively, the number of positive, negative, and zero eigenvalues of H , respectively.

LEMMA 2.12. $\text{In}(H + \lambda I) = (a, b, c)$ if and only if $\text{In}(X) = (a, b + t, c)$.

Proof. See Gould [1986] for the proof. \square

Note that Lemma 2.12 implies that the extended matrix X has at least t negative eigenvalues and so cannot be positive definite. Thus we simply cannot use the Cholesky factorization to solve (2.10). Fortunately, the symmetric indefinite factorization of Bunch and Kaufman [1977] is ideally suited to this situation. The Bunch-Kaufman factorization factors the extended system as $X = US^{-1}SDSS^{-1}U^T$ where U is a product of unit upper triangular and permutation matrices. That is,

$$U = P_{n+t}U_{n+t}^{-1} \cdots P_1U_1^{-1} \quad \text{and} \quad U_k^{-1} = \begin{pmatrix} I_{k-1} & -m_k & 0 \\ 0 & I_l & 0 \\ 0 & 0 & I_{n+t-k} \end{pmatrix}$$

(with $l = 1$ or 2) and D is a block diagonal matrix with blocks of order one or two, and S is an optional diagonal scaling matrix chosen so that the entries of US^{-1} are bounded. Moreover, it turns out that the inertia of X can be determined easily from

D. This property along with Lemma 2.12 will be very useful to our algorithm. For more details on the Bunch-Kaufman factorization, see the Linpack Users' Guide (Dongarra et al. [1979] and Bunch and Kaufman [1977]).

3. The algorithm. As in § 2 the development here will be much the same as that of Moré and Sorensen [1983] but will include those modifications necessary to make the algorithm perform well on penalty functions. The algorithmic development begins with a search for a good way to solve

$$(3.1) \quad \Phi(\lambda) = \|s_\lambda\|^2 - \Delta^2 = 0 \quad \text{where } s_\lambda = -(H + \lambda I)^{-1}g \text{ and } \lambda > -\lambda_1(H).$$

It is well known that $\Phi(\lambda)$ is convex and decreasing on $(-\lambda_1, \infty)$. When Newton's method is applied to such functions, λ_+ (the updated value of λ) is always less than λ^* , and if $\lambda > \lambda^*$ then λ_+ may be much less than λ^* . Dennis and Schnabel [1983] have attempted to overcome the problem of excessively short or long steps by modeling (3.1) with a single term rational function and have derived the following update:

$$(3.2) \quad \lambda_+ = \lambda + \left(\frac{\|s_\lambda\|^2}{\Delta} \right) \left[\frac{\|s_\lambda\| - \Delta}{s_\lambda^T (H + \lambda I)^{-1} s_\lambda} \right].$$

Moré and Sorensen [1983] obtain the same update by applying Newton's method to the function

$$(3.3) \quad \phi(\lambda) = \frac{1}{\Delta} - \frac{1}{\|s_\lambda\|} = 0.$$

It follows from (2.6) that $\phi(\lambda)$ is convex and twice continuously differentiable on the interval $(-\lambda_1, \infty)$. In practice we have observed that (3.1) works well when $-\lambda_1 < \lambda < \lambda^*$ but it still generally overcorrects when $\lambda > \lambda^*$. This is not surprising since a one-term approximation to (3.1) is valid only near $-\lambda_1$. The following example typifies the difficulties encountered when trying to compute trust region steps for the L_2 penalty function. Let $\mu = 10^{-2}$, $\Delta = 1$, $c = 1$, and

$$(3.4) \quad B = \begin{pmatrix} -.5 & 1.5 \\ 1.5 & -.5 \end{pmatrix}, \quad A = \begin{pmatrix} .5 \\ 1 \end{pmatrix}, \quad \text{and hence } H = \begin{pmatrix} 24.5 & 51.5 \\ 51.5 & 99.5 \end{pmatrix}.$$

Moreover, if we let $\nabla f = [-3, 2]^T$, then $g = [47, 102]^T$. Then eigenvalues of H are $\lambda_1 = -1.7064$ and $\lambda_2 = 125.7064$ and the value of λ^* for the trust region problem defined by these quantities is 9.5377. It follows that for $\lambda > 21$ equation (3.2) will return a value of $\lambda_+ < -\lambda_1 = 1.7064$. When that happens, λ must be updated in some other way. While the interval $(1.706, 21)$ may appear to provide a relatively large radius of convergence for Newton's method it should also be noted that a commonly used starting value

$$(3.5) \quad \lambda = \frac{\|g\|}{\Delta} = 112.3076$$

is much larger than 21. Therefore, to construct an efficient algorithm for solving trust region problems in this setting we must derive a better starting value for λ and devise a more robust way to update λ .

If a method is to produce a good update for λ when $\lambda > \lambda^*$ it must in some way pay attention to the curvature of ϕ . With this goal in mind we will try to construct a meaningful quadratic model of ϕ and then let λ_+ be the left root of the model. For the same reason Newton's method failed when $\lambda > \lambda^*$, the usual quadratic model (i.e., the first three terms of a Taylor series) also often produces a poor update in this setting.

That is, $\phi''(\lambda)$ tends to zero as λ goes to infinity (actually, $\phi''(\lambda) = O(1/\lambda)$ as $\lambda \rightarrow \infty$) so that any quadratic model built on local information at λ will not have enough curvature and will therefore greatly underestimate λ^* . A more successful approach is to construct a quadratic model with information acquired at points left and right of λ^* . As we shall see, when $\lambda > \lambda^*$ the algorithm has on hand a lower bound λ_S for $-\lambda_1(H)$ and an approximate unit eigenvector \hat{z} corresponding to $\lambda_1(H)$. Moreover, $\phi(-\lambda_1) = 1/\Delta$ and $\phi'(-\lambda_1^+) = -1/|v_1^T g| \phi'(-\lambda_1^+)$ denotes the derivative from the right, where v_1 is the exact eigenvector corresponding to λ_1 . Thus it seems reasonable to use λ_S and \hat{z} in scheme for updating λ . We are now poised to state the basic algorithm.

ALGORITHM 3.6. Let $\lambda \geq 0$ with $H + \lambda I$, as in (1.6), positive definite, and $\Delta > 0$ be given.

- (1) Form X and compute its Bunch-Kaufman factorization.
- (2) Solve $X \begin{pmatrix} s \\ r \end{pmatrix} = \begin{pmatrix} -v^T f \\ -c \end{pmatrix}$, as in (2.9).
- (3) Solve $X \begin{pmatrix} q \\ w \end{pmatrix} = \begin{pmatrix} s \\ 0 \end{pmatrix}$.
- (4) **If** $\|s\| > \Delta$ **then**
 - Update λ via (3.2)
 - Elseif** $\phi(\lambda) > -1/\Delta$ **then**
 - Determine the quadratic, $\chi(\lambda)$, which interpolates the points $(\lambda_S, 1/\Delta)$, $(\lambda, \phi(\lambda))$ and satisfies $\chi'(\lambda) = \phi'(\lambda)$.
 - Let λ_{1+} be the left root of $\chi(\lambda)$.
 - Let λ_{2+} be the left root of the best local quadratic approximation to ϕ at λ (i.e., the first three terms of a Taylor series).
 - If** $\lambda_{2+} > \lambda_S$ **then**
 - Let $\lambda_+ = \min\{\lambda_{1+}, \lambda_{2+}\}$.
 - Else**
 - $\lambda_+ = \lambda_{1+}$
 - Else**
 - Determine the quadratic, $\chi(\lambda)$, which interpolates the points $(\lambda_S, 1/\Delta)$, $(\lambda, \phi(\lambda))$ and satisfies $\chi'(\lambda_S) = -1/|\hat{z}^T g|$.
 - Let λ_+ be the left root of $\chi(\lambda)$.

Note that where possible we have used the best local model of ϕ . This has been done to ensure quadratic convergence near λ^* . In practice we have seen that Algorithm 3.6 is usually robust enough to converge from almost any reasonable starting value. However, to ensure convergence Algorithm 3.6 must, especially in the hard case, be safeguarded. Also, trust region methods usually require only an approximate solution. Thus, it remains to develop good convergence criteria. We shall postpone these matters until we have dealt effectively with the hard case.

In § 2 we have indicated that an eigenvector corresponding to λ_1 is required to solve the hard case. Moré and Sorensen [1983] show that a good approximate eigenvector can provide an acceptable inexact solution.

LEMMA 3.7. Let $0 < \sigma < 1$ be given and suppose that for $\lambda \geq 0$, $H + \lambda I$ is positive definite and $(H + \lambda I)s_\lambda = -g$. Let $z \in R^n$ satisfy

$$(3.8) \quad \|s_\lambda + z\| = \Delta \quad \text{and} \quad z^T(H + \lambda I)z \leq \sigma[s_\lambda^T(H + \lambda I)s_\lambda + \lambda\Delta^2];$$

then

$$(3.9) \quad -\psi(s_\lambda + z) \geq \frac{1}{2}(1 - \sigma)[s_\lambda^T(H + \lambda I)s_\lambda + \lambda\Delta^2] \geq (1 - \sigma)|\psi(s^*)|$$

where $\psi(s^*)$ is the optimal value of (1.1).

Proof. See Moré and Sorensen [1983] for the proof. \square

It follows from Lemma 3.7 that $|\psi(s_\lambda + z) - \psi(s^*)| \leq \sigma |\psi(s^*)|$ and so $s_\lambda + z$ is a nearly optimal solution to (1.1). As Moré and Sorensen suggest, anytime we have $\|s_\lambda\| < \Delta$ we try to satisfy (3.8) with $z = \tau \hat{z}$ where $\|s_\lambda + \tau \hat{z}\| = \Delta$ and \hat{z} is an approximate eigenvector of unit norm corresponding to $-\lambda_1$. Furthermore, our complete algorithm (to be given later) will require that \hat{z} satisfy $\hat{z}^T(H + \lambda I)\hat{z} \rightarrow 0$ as $\lambda \rightarrow -\lambda_1^+$. Such a vector \hat{z} with $\|\hat{z}\| = 1$ can be obtained with the following algorithm.

ALGORITHM 3.10. Given $H + \lambda I$ positive definite, the extended system X , and its $US^{-1}SDSS^{-1}U^T$ factorization, the following algorithm computes a vector \hat{z} such that $\|\hat{z}\| = 1$ and $\hat{z}^T(H + \lambda I)\hat{z}$ is as small as possible.

- (1) If $U = P_n U_n^{-1} \cdots P_1 U_1^{-1}$, let $P = P_1 P_2 \cdots P_n$ and let $\tilde{U} = PU$. Note that \tilde{U} is upper triangular.
- (2) Compute the QR factorization of $\tilde{U}DS$. We hasten to point out that since D is tridiagonal, the product $\tilde{U}DS$ is upper Hessenberg and so its QR factorization can be computed in $O(n+t)$ flops.
- (3) Solve $Rw = e$ where e is a vector of the form $(\pm 1, \cdots, \pm 1)^T$ using the following procedure
 $\theta_i = 0$ for $i = 1, \cdots, n+t$
For $k = n+t, \cdots, 1$

$$w_k^+ = (1 + \theta_k) / r_{kk}$$

$$w_k^- = (-1 + \theta_k) / r_{kk}$$

$$v^+ = |w_k^+| + \sum_{i=1}^{k-1} |\theta_i + r_{ik} w_k^+|$$

$$v^- = |w_k^-| + \sum_{i=1}^{k-1} |\theta_i + r_{ik} w_k^-|$$

If $v^+ \geq v^-$

then $w_k = w_k^+$

else $w_k = w_k^-$

$\theta_i = \theta_i + r_{ik} w_k$ for $i = 1, \cdots, k-1$

- (4) Solve $S^{-1}\tilde{U}^T x = w$.
- (5) Solve $Py = x$.
- (6) If $y = (z_1, z_2)^T$ with $z_1 \in R^n$ and $z_2 \in R'$ we then set $\hat{z} = z_1 / \|z_1\|$.

As we said before the only property of \hat{z} required by the complete algorithm is that $\hat{z}^T(H + \lambda I)\hat{z}$ approach zero as λ approaches $-\lambda_1$. The following lemma confirms this.

LEMMA 3.11. *If $\lambda \rightarrow -\lambda_1^+$, then for the \hat{z} from Algorithm 3.10, $\hat{z}^T(H + \lambda I)\hat{z} \rightarrow 0$.*

Proof. See Appendix A for the proof. \square

Last, given \hat{z} from Algorithm 3.10 the proof of Lemma 3.7 shows that the correct choice of τ is

$$(3.12) \quad \tau = \frac{\Delta^2 - \|s_\lambda\|^2}{s_\lambda^T \hat{z} + \operatorname{sgn}(s_\lambda^T \hat{z}) [(s_\lambda^T \hat{z})^2 + (\Delta^2 - \|s_\lambda\|^2)]^{1/2}}$$

Algorithm 3.10 will not, in general, be used at every iteration where $\|s_\lambda\| < \Delta$, but rather to generate starting values for inverse iteration. If at some iteration we have $\|s_\lambda\| < \Delta$ and we have an approximate eigenvector \hat{z} of unit norm from a previous iteration (where $\lambda_p \cong -\lambda_1$), we then solve $(H + \lambda I)y = \hat{z}$ and if

$$\hat{y}^T(H + \lambda I)\hat{y} < \varepsilon \hat{z}^T(H + \lambda_p I)\hat{z} \quad \text{where } \hat{y} = y / \|y\|$$

and $\varepsilon \in (0, 1)$ is some fixed constant, we then take \hat{y} to be our new approximate eigenvector. Otherwise we resort to Algorithm 3.10. In practice Algorithm 3.10 is seldom used more than twice within each instance of (1.1).

As it stands it is possible, although unlikely, for Algorithm 3.6 to return a value of λ_+ that is less than $-\lambda_1$. To prevent this and to guarantee that a solution will be found in finitely many iterations, the algorithm must be safeguarded. Moré and Sorensen's safeguarding scheme is perfectly adequate for this setting and so we use it. An interval containing λ^* , $[\lambda_L, \lambda_U]$, and a lower bound λ_S , for $-\lambda_1$ are maintained at all times. At the beginning of each iteration the value of λ_+ inherited from the previous iteration is checked against these parameters by the following procedure.

Safeguard λ :

- (1) $\lambda = \max(\lambda, \lambda_L)$,
- (2) $\lambda = \min(\lambda, \lambda_U)$,
- (3) If $\lambda \leq \lambda_S$ then $\lambda = \max\{.001\lambda_U, (\lambda_L, \lambda_U)^{1/2}\}$.

As Moré and Sorensen point out, the third step is probably the most important. Without it, we could not prove that the algorithm will yield an approximate solution in finitely many iterations. The procedure for updating λ_L , λ_U , and λ_S is straightforward (for details see Moré and Sorensen [1983]).

Update λ_L , λ_U , and λ_S :

- (1) **If $H + \lambda I$ is nonsingular then**
 Compute \hat{z} as described above
If $H + \lambda I$ is positive definite then
 $\lambda_S = \max\{\lambda_S, \lambda - \hat{z}^T(H + \lambda I)\hat{z}\}$
Else
 $\lambda_S = \max\{\lambda_S, \lambda, \lambda - \hat{z}^T(H + \lambda I)\hat{z}\}$
- (2) **If $\lambda \in (-\lambda_1, \infty)$ and $\phi(\lambda) < 0$ then**
 $\lambda_U = \min(\lambda_U, \lambda)$
Else
 $\lambda_L = \max(\lambda_L, \lambda)$
- (3) **Let $\lambda_L = \max(\lambda_L, \lambda_S)$.**

Initial values for the safeguarding parameters are:

$$\lambda_S = \max\{-h_{ii}, \text{for } i = 1, \dots, n, -\|B\|_1 \text{ (if } t < n)\},$$

$$\lambda_L = \max\left\{0, \lambda_S, \frac{\|g\|}{\Delta} - \|B\|_1 - \frac{1}{\mu} \|A\|_1 \|A\|_\infty\right\},$$

$$\lambda_U = \min\left\{\frac{\|g\|}{\Delta} + \|B\|_1 + \frac{1}{\mu} \|A\|_1 \|A\|_\infty, \frac{\|g\|}{\Delta} + \|B\|_1 \text{ (if } t < n)\right\}$$

where h_{ii} is the i th diagonal element of $H = B + (1/\mu)AA^T$.

Since each iteration of the algorithm will require $O(n+t)^3$, it is imperative that we try to minimize the number of iterations needed to find a solution. As we saw earlier in the example, a poor initial value for λ will certainly hinder progress toward a solution. It has been our experience that it is generally worthwhile expending some effort to compute a better starting value for λ . In other words, we were unable to find a starting value for λ that was both inexpensive and guaranteed to be close to the solution. The initial value described below proved itself computationally cost effective

when compared to less expensive initial values. To begin the derivation we observe that (2.5) implies

$$(3.13) \quad \|s_\lambda\|^2 \leq \sum_{i=1}^{n-t} \frac{\gamma_i^2}{(\lambda_1 + \lambda)^2} + \sum_{i=n-t+1}^n \frac{\gamma_i^2}{\lambda_{n-t+1}^2}$$

and if reasonable estimates for $\sum_{i=1}^{n-t} \gamma_i^2$, $\sum_{i=n-t+1}^n \gamma_i^2$, λ_1 , and λ_{n-t+1} can be found, then we may be able to use (3.13) to compute a suitable initial value for λ . The estimates should, if possible, improve as μ tends to zero since the accuracy of the initial values becomes increasingly critical as μ becomes small. To secure such estimates we let $A = QR$ where Q is orthogonal and R is upper-triangular. Note that since Q is orthogonal the last $n-t$ columns $[q_{t+1}, \dots, q_n]$ form an orthonormal basis for the null space of A^T (computing a QR factorization of A is even more useful when A is nearly rank deficient in a neighborhood of the solution but this is the topic of another paper). We now make the following approximations:

$$\sum_{i=1}^{n-t} \gamma_i^2 \approx \sum_{j=t+1}^n (q_j^T g)^2 \equiv \delta$$

and

$$\sum_{i=n-t+1}^n \gamma_i^2 \approx \|g\|^2 - \delta \equiv \zeta.$$

Unfortunately, we could not find estimates for λ_1 and λ_{n-t+1} that improve as μ tends to zero. However, it is easy to see that if $\lambda_1(H) < 0$, then $\lambda_1(B) < \lambda_1(H)$. To compute $\lambda_1(B)$ we reduce B to tridiagonal form with Householder transformations. The characteristic polynomial is then easily formed, and its smallest root can be readily found via bisection. For details see Golub and Van Loan [1983]. If A has full rank it can be shown that as μ tends to zero, $\lambda_{n-t+1}(H)$ converges to $(1/\mu)\lambda_{n-t+1}(AA^T)$. Furthermore,

$$(3.14) \quad \lambda_{n-t+1}(AA^T) \leq \frac{1}{\mu} \min \{\bar{r}_{ii}^2, \bar{a}_{ii}\} \equiv \rho$$

where \bar{r}_{ii} is the element of smallest absolute value along the diagonal of R , and \bar{a}_{ii} is the smallest diagonal entry of AA^T . And so we use (3.14) as our approximation for λ_{n-t+1} . Substituting δ , ζ , $\lambda_1(B)$, and ρ for their respective quantities in (3.13) and solving for λ gives

$$(3.15) \quad \lambda \approx \left[\frac{\delta^2}{\Delta^2 - \zeta^2 / \rho^2} \right]^{1/2} + |\lambda_1(B)|.$$

It is certainly possible for the denominator in (3.15) to be negative and in that case we recommend using $\|B\|_1$ or $\|g\|/\Delta$ for a starting value. Formula (3.15) is most likely to fail when Δ is small or when ρ is small, but in these situations λ^* (the value of λ corresponding to s^*) is generally larger. In other words, λ^* is usually inversely proportional to Δ and ρ . The final ingredient of the iteration is the convergence criteria. Here again we invoke, without modification, Moré and Sorensen's work. For the convergence parameter $\sigma_1 \in (0, 1)$, and a $\lambda \geq 0$ such that $H + \lambda I$ is positive definite, the candidate step s_λ , computed by Algorithm 3.6, is checked for near optimality with the following scheme.

Convergence Criteria:

- (1) If $|\Delta - \|s_\lambda\|| \leq \sigma_1 \Delta$ or $\|s_\lambda\| \leq \Delta$, $\lambda = 0$ then
Terminate with $s = s_\lambda$ as the approximate solution.
- (2) If $\|s_\lambda\| < \Delta$ then
Use Algorithm 3.10 to compute \hat{z} and τ
If s_λ , \hat{z} , and τ satisfy (3.8) then
Terminate with $s = s_\lambda + \tau \hat{z}$ as the approximate solution.

More and Sorensen [1983] show that these criteria guarantee that if the algorithm terminates, then the approximate solution s satisfies

$$\psi(s) - \psi(s^*) \leq \sigma_1(2 - \sigma_1)|\psi(s^*)| \quad \text{and} \quad \|s\| \leq (1 + \sigma_1)\Delta,$$

and thus s is a nearly optimal solution of (1.1). It should be noted that (3.8) provides termination criteria that is applicable to both the general and the hard cases. In our tests we observe that when a very inexact solution (i.e., $\sigma_1 = .02$) is sought, termination usually occurs on (3.8) even in the general case. We have now discussed all the ingredients of the iterative scheme for solving problem (1.1). The following algorithm summarizes these ingredients and defines a typical iteration.

ALGORITHM 3.16.

- (1) Safeguard λ .
- (2) Form X and compute its $US^{-1}SDSS^{-1}U^T$ factorization. If $H + \lambda I$ is not positive definite then go to (5).
- (3) Solve $(H + \lambda I)s_\lambda = -g$ as in Algorithm 3.6.
- (4) If $\|s_\lambda\| < \Delta$ use Algorithm 3.10 to compute τ and \hat{z} .
- (5) Update λ_L , λ_U , and λ_S .
- (6) Check the convergence criteria.
- (7) If $H + \lambda I$ is positive definite and $g \neq 0$ then update λ via Algorithm 3.6, otherwise set $\lambda = \lambda_S$.

Last, we must now show that the above algorithm will produce a nearly optimal solution in a finite number of steps.

THEOREM 3.17. *For $\phi(\lambda)$ as in (3.3) and any $\varepsilon > 0$ Algorithm 3.16 will, in finitely many iteration, return a value of λ in $[-\lambda_1, \infty)$ such that one of the following hold:*

- (i) $\lambda \in (\lambda^* - \varepsilon, \lambda^* + \varepsilon)$, where $\lambda^* \in (-\lambda_1, \infty)$ and $\phi(\lambda^*) = 0$, in the event of the general case.
- (ii) $\lambda \in [-\lambda_1, -\lambda_1 + \varepsilon)$ in the event of the hard case.
- (iii) $\lambda \in [0, \varepsilon)$ in the event of the positive definite case.

Proof. (i) General Case. The proof is by contradiction. Suppose that the length of the interval of uncertainty is bounded away from zero. Then the bisection step of the safeguarding scheme can only be invoked finitely many times. If during some iteration a value of λ is found such that $\phi(\lambda) > 0$ and $\lambda > -\lambda_1$, then the convexity of ϕ implies that Newton's method will proceed monotonically to λ^* and a nearly optimal solution will be found in finitely many more iterations.

Now suppose that the length of the interval of uncertainty does not go to zero and the algorithm never finds a value of λ for which $\phi(\lambda) > 0$ and $\lambda > -\lambda_1$. This implies that after a finite number of iterations we have a sequence $\{\lambda_k\}_N^\infty$ that satisfies $\lambda_k > \lambda_{k+1} > \lambda^* > -\lambda_1$ and $\phi(\lambda_k) < 0$. To show that the sequence $\{\lambda_k\}_N^\infty$ converges, we will compare our method to one that we know converges: the secant method. In particular, we will show that $\lambda_k - \lambda_{k+1}$ is larger for our method than for a suitable and obviously convergent secant method. This result, under all of the assumptions we have already made, will prove the convergence.

Suppose that the exact value of $-\lambda_1$ is known. Then since $\phi(-\lambda_1) = 1/\Delta$ we could update λ_k by the secant method applied to the points $(-\lambda_1, 1/\Delta)$ and $(\lambda_k, \phi(\lambda_k))$ and thereby obtain

$$\lambda_a^+ = \frac{\lambda_k + \lambda_1}{1 - \Delta\phi(\lambda_k)} - \lambda_1.$$

Since ϕ is convex on $(-\lambda_1, \infty)$, it is easy to see that updating λ_k by this formula will generate a sequence converging to λ^* . Now consider the update

$$\lambda_b^+ = \frac{\lambda_k + \lambda_S}{1 - \Delta\phi(\lambda_k)} + \lambda_S,$$

which was obtained by applying the secant formula to the points $(\lambda_S, 1/\Delta)$, $(\lambda_k, \phi(\lambda_k))$. Since $\lambda_S \leq -\lambda_1$ it easily follows that $\lambda_b^+ < \lambda_a^+$. Let χ be any convex quadratic that interpolates the points $(\lambda_S, 1/\Delta)$, $(\lambda_k, \phi(\lambda_k))$. Since $1/\Delta > 0$ and $\phi(\lambda_k) < 0$, it easily follows that the leftmost root r_l satisfies $\lambda_S < r_l < \lambda_b^+ < \lambda_a^+$. If, on the other hand, χ is concave, then r_l clearly satisfies $r_l < \lambda_S < \lambda_b^+ < \lambda_a^+$. In either case we have $r_l < \lambda_b^+ < \lambda_a^+$, and taking $\lambda_{k+1} = r_l$ it immediately follows that $\lambda_k - \lambda_{k+1} > \lambda_k - \lambda_a^+$, as required. This proves convergence for the general case.

(ii) Hard Case. In this case we have $\phi(-\lambda_1) \leq 0$ and again we assume that the interval of uncertainty remains bounded away from zero and so after finitely many iterations, we have a sequence $\{\lambda_k\}_N^\infty$ that satisfies $\lambda_k > \lambda_{k+1} > -\lambda_1$ and $\phi(\lambda_k) < 0$. As in the general case we will show that our method converges by comparing it to an obviously convergent variant of the secant method. Suppose that the exact values of $-\lambda_1$ and $\phi(-\lambda_1 + \delta)$, where $\delta \in (0, \varepsilon)$, are known. Updating λ_k by applying the secant method to the points $(-\lambda_1, 1/\Delta)$ and $(\lambda_k, \phi(-\lambda_1 + \delta))$ yields

$$\lambda_a^+ = \frac{\lambda_k + \lambda_1}{1 - \Delta\phi(-\lambda_1 + \delta)} - \lambda_1.$$

Even though the points $(-\lambda_1, 1/\Delta)$ and $(\lambda_k, \phi(-\lambda_1 + \delta))$ do not lie on the function ϕ , it is easy to see that the above formula will generate a sequence that converges monotonically to $-\lambda_1$. Consider the following update:

$$(3.18) \quad \lambda_b^+ = \frac{\lambda_k + \lambda_S}{1 - \Delta\phi(\lambda_k)} - \lambda_S.$$

Since $\lambda_S \leq -\lambda_1$ and $\phi(\lambda_k) < \phi(-\lambda_1 + \delta)$, for $\lambda_k > -\lambda_1 + \delta$, we clearly have $\lambda_b^+ < \lambda_a^+$. Thus for any sequence $\{\lambda_k\}_N^\infty$ obtained from updating λ_k by (3.18) it follows that for some finite $K \geq N$ we have $\lambda_K \leq -\lambda_1 + \delta$. But since δ was arbitrary, the sequence $\{\lambda_k\}_N^\infty$ generated by (3.18) converges to $-\lambda_1$.

(iii) Positive-Definite Case. The proof for this case is identical to the proof of the hard case except that $-\lambda_1$ should be replaced with zero. \square

4. Computational results. In this section we report the results of some tests in which we compare the performance of TRSQPF (our algorithm) with that of GQTPAR (Moré and Sorensen's algorithm) on individual trust region problems where g and H are of the form (1.3) and (1.6), respectively. The tests covered a wide range of values of n , t , and μ and for each triple (n, t, μ) Tables 1-12 below summarize the outcome of 10 different trust region problems. Individual problems were created by specifying an $n \times n$ diagonal matrix for D_B , an $n \times t$ diagonal matrix for D_A , an n vector for g , and a t vector for c . By appropriately choosing the entries in each of these quantities we can construct any of the four main types of trust region problems: general case, hard case, positive definite case, and the saddle-point case.

TABLE 1

General case										
			GQTPAR				TRSQPF			
n	t	μ	# crash	min	max	mean	min	max	mean	adjst. mean
20	5	10^{-2}	0	5	9	6.4	1	5	2.9	3.8
		10^{-5}	0	5	10	7.4	2	4	2.7	3.6
		10^{-9}	0	4	11	8.8	2	4	2.8	3.7
		10^{-12}	0	4	13	9.0	2	4	2.7	3.6
		10^{-16}	8	5	6	5.5	2	5	3.4	4.3
	10	10^{-2}	0	3	9	5.7	1	7	3.9	4.5
		10^{-5}	0	4	10	7.2	1	4	2.5	3.1
		10^{-9}	0	5	14	9.2	2	4	2.9	3.5
		10^{-12}	0	3	11	6.7	2	5	3.0	3.6
		10^{-16}	9	8	8	8.0	2	4	3.1	3.7
	15	10^{-2}	0	3	8	5.3	1	5	3.4	4.0
		10^{-5}	0	6	11	7.4	1	7	3.8	4.4
		10^{-9}	0	5	13	9.4	2	4	2.6	3.2
		10^{-12}	0	2	9	6.1	1	4	2.7	3.3
		10^{-16}	9	7	7	7.0	1	4	3.0	3.6

TABLE 2

General case										
			GQTPAR				TRSQPF			
n	t	μ	# crash	min	max	mean	min	max	mean	adjst. mean
40	10	10^{-2}	0	3	9	6.7	1	5	3.3	4.5
		10^{-5}	0	5	10	7.3	1	5	3.2	4.4
		10^{-9}	0	5	10	7.9	2	4	3.2	4.4
		10^{-12}	0	3	12	7.6	2	4	3.3	4.5
		10^{-16}	7	2	9	6.0	2	4	3.2	4.4
	20	10^{-2}	0	3	7	5.6	1	6	3.3	4.2
		10^{-5}	0	5	9	7.2	2	6	3.3	4.2
		10^{-9}	0	5	12	8.9	2	4	2.8	3.7
		10^{-12}	0	3	12	6.9	2	4	3.6	4.5
		10^{-16}	7	7	10	8.7	2	5	2.6	3.5
	30	10^{-2}	0	4	9	5.2	2	5	3.2	3.6
		10^{-5}	0	4	10	7.5	1	6	4.2	4.6
		10^{-9}	0	6	13	9.4	2	4	2.8	3.2
		10^{-12}	0	3	12	8.5	2	4	3.2	3.6
		10^{-6}	10	-	-	-	2	5	2.8	3.2

TABLE 3

General case										
<i>n</i>	<i>t</i>	μ	GQTPAR				TRSQPF			adjst. mean
			# crash	min	max	mean	min	max	mean	
80	20	10^{-2}	0	4	9	6.4	1	6	3.9	5.5
		10^{-5}	0	5	10	7.0	1	5	3.7	5.3
		10^{-9}	0	4	12	7.6	2	4	3.3	4.9
		10^{-12}	0	3	10	7.6	2	5	3.6	5.2
		10^{-16}	5	8	13	6.3	2	5	3.8	5.4
		10^{-2}	0	3	9	5.5	1	5	3.5	4.7
	40	10^{-5}	0	4	11	7.5	4	7	5.3	6.5
		10^{-9}	0	5	13	9.0	3	4	3.3	4.5
		10^{-12}	0	2	12	7.0	2	5	3.5	4.7
		10^{-16}	6	5	8	6.8	2	5	3.2	4.4
		10^{-2}	0	3	6	4.6	1	6	3.6	4.5
		10^{-5}	0	3	10	8.1	3	6	5.2	6.1
	60	10^{-9}	0	5	11	8.7	2	6	3.2	4.1
		10^{-12}	0	3	11	7.5	2	6	3.2	4.1
		10^{-16}	7	2	7	5.3	2	4	2.5	3.4

TABLE 4

Hard case										
<i>n</i>	<i>t</i>	μ	GQTPAR				TRSQPF			adjst. mean
			# crash	min	max	mean	min	max	mean	
20	5	10^{-2}	0	4	6	5.4	2	4	2.8	3.8
		10^{-5}	0	4	11	7.8	3	5	3.3	4.2
		10^{-9}	0	5	11	8.9	2	6	3.5	4.4
		10^{-12}	0	3	12	8.1	2	5	3.5	4.4
		10^{-16}	9	2	2	2.0	2	5	3.6	4.5
		10^{-2}	0	3	9	5.9	2	5	3.9	4.5
	10	10^{-5}	0	4	11	7.5	2	4	2.9	3.5
		10^{-9}	0	5	13	9.3	2	5	3.3	3.9
		10^{-12}	0	3	12	7.8	3	6	4.2	4.8
		10^{-16}	8	6	7	6.5	3	5	3.8	4.4
		10^{-2}	0	3	8	5.7	1	4	2.3	2.9
		10^{-5}	0	4	10	7.7	1	6	3.8	4.4
	15	10^{-9}	0	6	10	8.5	2	7	4.6	5.2
		10^{-12}	0	4	12	8.9	2	6	4.4	5.0
		10^{-16}	9	7	7	7.0	2	7	4.1	4.7

TABLE 5

Hard case											
<i>n</i>	<i>t</i>	μ	GQTPAR				TRSQPF			adjst. mean	
			# crash	min	max	mean	min	max	mean		
40	10	10^{-2}	0	3	9	5.5	1	5	3.6	4.8	
		10^{-5}	0	5	10	7.3	1	4	3.0	4.2	
		10^{-9}	0	4	12	8.1	2	4	3.5	4.7	
		10^{-12}	0	2	12	7.3	1	6	3.3	4.5	
	20	10^{-16}	6	7	11	8.8	2	5	3.7	4.9	
		10^{-2}	0	1	8	5.0	1	7	3.3	4.2	
		10^{-5}	0	4	9	6.5	2	6	3.9	4.8	
		10^{-9}	0	5	14	8.8	2	4	3.0	3.9	
		10^{-12}	0	2	13	8.0	2	5	3.5	4.4	
		10^{-16}	7	9	11	9.7	3	5	3.7	4.6	
		30	10^{-2}	0	3	8	5.4	1	6	3.6	4.0
			10^{-5}	0	4	10	7.5	2	7	5.0	5.4
10^{-9}	0		6	12	9.4	3	6	3.9	4.3		
10^{-12}	0		3	11	7.9	3	7	4.0	4.4		
	10^{-16}	8	5	11	8.0	2	6	4.5	4.9		

TABLE 6

Hard case										
<i>n</i>	<i>t</i>	μ	GQTPAR				TRSQPF			adjst. mean
			# crash	min	max	mean	min	max	mean	
80	20	10^{-2}	0	3	7	5.1	1	6	3.5	5.1
		10^{-5}	0	4	11	7.5	2	5	3.3	4.9
		10^{-9}	0	5	12	8.5	2	4	3.5	5.1
		10^{-12}	0	3	11	7.3	2	4	3.4	5.0
		10^{-16}	5	2	9	8.8	2	4	3.0	4.6
	40	10^{-2}	0	3	7	5.5	1	4	3.7	4.9
		10^{-5}	0	4	11	7.2	2	7	4.4	5.6
		10^{-9}	0	5	13	9.4	2	4	3.3	4.5
		10^{-12}	0	3	11	7.3	2	4	3.4	4.6
		10^{-16}	8	4	8	6.0	2	4	3.3	4.5
	60	10^{-2}	0	3	7	4.4	1	6	2.9	3.8
		10^{-5}	0	4	10	7.6	2	7	4.8	5.7
		10^{-9}	0	5	11	8.2	2	5	3.7	4.6
		10^{-12}	0	3	11	8.3	2	4	3.3	4.2
		10^{-16}	3	6	9	7.8	2	5	3.5	4.4

TABLE 7

Positive-definite case										
<i>n</i>	<i>t</i>	μ	GQTPAR				TRSQPF			adjst. mean
			# crash	min	max	mean	min	max	mean	
20	5	10^{-2}	0	2	5	4.0	1	3	2.0	2.9
		10^{-5}	0	4	5	4.3	1	2	1.9	2.8
		10^{-9}	0	3	5	4.2	1	4	2.2	3.1
		10^{-12}	0	2	4	3.4	2	3	2.5	3.4
		10^{-16}	10	-	-	-	2	4	2.4	3.3
	10	10^{-2}	0	2	7	4.3	2	5	3.5	4.1
		10^{-5}	0	2	5	4.2	2	4	2.4	3.0
		10^{-9}	0	4	5	4.3	2	5	2.5	3.1
		10^{-12}	0	2	4	3.0	1	3	2.9	3.5
		10^{-16}	8	7	17	12.0	2	3	2.2	2.8
	15	10^{-2}	0	2	6	4.0	2	5	2.8	3.4
		10^{-5}	0	2	5	3.4	1	4	2.6	3.2
		10^{-9}	0	2	4	3.1	2	2	2.0	2.6
		10^{-12}	0	2	5	2.7	2	3	2.1	2.7
		10^{-16}	8	6	6	6.0	2	3	2.2	2.8

TABLE 8

Positive-definite case										
<i>n</i>	<i>t</i>	μ	GQTPAR				TRSQPF			adjst. mean
			# crash	min	max	mean	min	max	mean	
40	10	10^{-2}	0	4	5	4.9	1	5	3.5	4.7
		10^{-5}	0	4	5	4.4	2	3	2.1	3.3
		10^{-9}	0	4	5	4.8	2	3	2.4	3.6
		10^{-12}	0	2	5	3.6	2	4	2.3	3.5
		10^{-16}	4	2	11	6.5	1	2	1.9	3.1
	20	10^{-2}	0	3	6	4.7	1	5	3.1	4.0
		10^{-5}	0	4	5	4.8	2	3	2.4	3.3
		10^{-9}	0	4	5	4.6	1	3	2.0	2.9
		10^{-12}	0	2	5	3.8	2	6	2.6	3.5
		10^{-16}	7	2	7	4.0	1	3	2.2	3.1
	30	10^{-2}	0	2	6	3.9	1	3	2.1	2.5
		10^{-5}	0	3	5	3.9	2	4	3.3	3.7
		10^{-9}	0	2	5	3.5	2	5	2.4	2.8
		10^{-12}	0	2	4	3.4	2	3	2.4	2.8
		10^{-16}	10	-	-	-	2	3	2.4	2.8

TABLE 9

Positive-definite case												
			GQTPAR				TRSQPF					
n	t	μ	# crash	min	max	mean	min	max	mean	adjst. mean		
80	20	10^{-2}	0	3	6	5.0	1	5	2.8	4.4		
		10^{-5}	0	5	5	5.0	2	3	2.1	3.7		
		10^{-9}	0	4	5	4.9	2	2	2.0	3.6		
		10^{-12}	0	2	5	3.5	1	3	2.0	3.6		
	40	10^{-16}	10^{-2}	5	3	9	7.8	2	3	2.2	3.8	
			10^{-5}	0	2	6	4.9	1	5	3.7	4.9	
			10^{-9}	0	4	5	4.7	1	5	3.2	4.4	
			10^{-12}	0	4	5	4.7	2	3	2.1	3.3	
		60	10^{-16}	10^{-2}	0	2	5	3.4	1	3	2.1	3.3
				10^{-5}	9	5	5	5.0	2	3	2.4	3.6
				10^{-9}	0	2	8	4.2	1	5	3.1	4.0
				10^{-12}	0	2	6	4.9	2	5	3.7	4.6
	10^{-16}	10^{-2}	0	2	5	4.2	2	4	2.3	3.2		
		10^{-5}	0	2	5	4.2	2	4	2.3	3.2		
		10^{-9}	0	2	5	3.5	1	3	2.1	3.0		
		10^{-12}	6	8	11	9.5	2	3	2.1	3.0		

TABLE 10

Saddle-point case											
			GQTPAR				TRSQPF				
n	t	μ	# crash	min	max	mean	min	max	mean	adjst. mean	
20	5	10^{-2}	0	4	12	7.4	1	4	1.3	2.1	
		10^{-5}	0	8	17	10.2	1	5	1.7	2	
		10^{-9}	0	9	18	12.4	1	3	1.4	2.3	
		10^{-12}	0	8	19	13.9	1	2	1.1	2.0	
		10^{-16}	10	-	-	-	1	4	1.9	2.8	
		10	10^{-2}	0	8	14	10.6	1	6	2.9	3.5
	10^{-5}		0	6	13	9.7	1	11	2.5	3.1	
	10^{-9}		0	10	17	13.3	1	7	2.1	2.7	
	10^{-12}		1	11	16	13.1	1	7	2.7	3.3	
	15	10^{-16}	10^{-2}	9	12	12	12.0	1	6	2.8	3.4
			10^{-5}	0	6	12	9.7	1	8	5.3	5.9
		10^{-16}	10^{-2}	0	2	15	11.8	3	10	5.4	6.0
			10^{-5}	0	2	17	11.0	1	6	4.1	4.7
			10^{-9}	0	2	17	11.0	1	6	4.1	4.7
			10^{-12}	0	12	18	15.3	1	6	3.7	4.3
	10^{-16}	10	-	-	-	1	6	3.4	4.0		

TABLE 11

Saddle-point case										
<i>n</i>	<i>t</i>	μ	GQTPAR				TRSQPF			adjst. mean
			# crash	min	max	mean	min	max	mean	
40	10	10^{-2}	0	6	16	9.5	1	2	1.3	2.5
		10^{-5}	0	6	15	11.0	1	2	1.1	2.3
		10^{-9}	0	8	15	12.2	1	5	1.4	2.6
		10^{-12}	0	9	17	11.4	1	4	1.3	2.5
	20	10^{-16}	9	13	13	13.0	1	5	2.3	3.5
		10^{-2}	0	7	10	7.8	1	5	1.9	2.8
		10^{-5}	0	7	19	11.1	1	4	2.0	2.9
		10^{-9}	0	7	18	13.1	1	6	2.1	3.0
	30	10^{-12}	0	10	16	13.4	1	5	2.5	3.4
		10^{-16}	9	14	14	14.0	1	9	2.8	3.7
		10^{-2}	0	7	15	11.5	1	6	4.5	4.9
		10^{-5}	0	7	14	11.4	1	6	4.1	4.5
	10^{-9}	0	10	15	12.8	1	7	3.6	4.0	
	10^{-12}	0	8	16	13.3	1	6	3.7	4.1	
	10^{-16}	9	16	16	16.0	1	6	3.0	3.4	

TABLE 12

Saddle-point case										
<i>n</i>	<i>t</i>	μ	GQTPAR				TRSQPF			adjst. mean
			# crash	min	max	mean	min	max	mean	
80	20	10^{-2}	0	7	14	10.1	1	3	1.4	3.0
		10^{-5}	0	8	13	10.9	1	3	1.5	3.1
		10^{-9}	0	8	14	11.6	1	4	1.6	3.2
		10^{-12}	0	10	16	13.5	1	2	1.1	2.7
	40	10^{-16}	8	10	12	11.0	1	1	1.0	2.6
		10^{-2}	0	9	16	11.2	1	4	2.1	3.3
		10^{-5}	0	7	15	11.9	1	4	1.8	3.0
		10^{-9}	0	11	17	13.3	1	3	1.5	2.7
	60	10^{-12}	0	9	18	13.3	1	4	1.6	2.8
		10^{-16}	8	12	16	13.3	1	5	1.9	3.1
		10^{-2}	0	6	13	9.1	1	5	3.2	4.1
		10^{-5}	0	6	15	10.2	1	5	2.9	3.8
	10^{-9}	0	11	18	13.7	1	6	3.4	4.3	
	10^{-12}	0	11	18	14.9	1	7	3.0	3.9	
	10^{-16}	10	-	-	-	1	5	2.4	3.3	

An instance of the general case is formed by choosing the entries of D_B and D_A as uniformly distributed random numbers in $(-1, 1) \setminus 0$ and $(-1, -2) \cup (.2, 1)$, respectively. The entries of D_A were so chosen to avoid the case where A is nearly rank deficient. The entries of g were chosen as uniformly distributed random numbers in $(-1, 1)$. However, the entries of the constraint vector c were restricted to the interval $(-\mu_{\text{prev}}, \mu_{\text{prev}})$ where μ_{prev} is the next largest value of μ used in the testing. For example, if the current value of μ is 10^{-12} , then μ_{prev} is 10^{-9} . We feel this scaling of c provides a more realistic test because when TRSQPF is used in the context equality constrained minimization, the magnitudes of the constraints are, at any given iteration, usually of the order of the previous value of μ .

To form an instance of the hard case we set to zero the component of g corresponding to the smallest element of $\{D_{B_{t+1,t+1}}, \dots, D_{B_{nn}}\}$. If none of those entries of D_B are negative, then $D_{B_{nn}}$ is chosen from the interval $(-1, 0)$, and g_n is set to zero. We need only consider the last $n-t$ elements of D_B because for the values of μ that we use $\{10^{-2}, 10^{-5}, 10^{-9}, 10^{-12}, 10^{-16}\}$ and our choice for the entries of D_A ensure that the first t eigenvalues $D_B + (1/\mu)D_AD_A^T$ are positive. An instance of the positive-definite case was formed by choosing $[D_{B_{t+1,t+1}}, \dots, D_{B_{nn}}]$ from the interval $(0, 1)$. The saddle-point case is somewhat more involved and will be described a little later.

The next step in constructing a trust region problem is to scramble the structure of D_B and D_A . We set $B = QD_BQ^T$ and $A = QD_AZ$ where Q and Z are, respectively, $n \times n$ and $t \times t$ orthogonal matrices. Both Q and Z are of the form $Q_1Q_2Q_3$ or $Z_1Z_2Z_3$ where

$$(4.1) \quad Q_i = I - 2 \frac{v_i v_i^T}{v_i^T v_i} \quad \text{and} \quad Z_i = I - 2 \frac{w_i w_i^T}{w_i^T w_i}, \quad i = 1, 2, 3.$$

The entries of v_i and w_i are random numbers uniformly distributed in $(-1, 1)$. We are now better able to describe the saddle-point case. Quite simply, an instance of the saddle-point case can be created by setting $g = -(1/\mu)Ac$. Last, the values of Δ were random numbers chosen from the interval $(0, 10)$.

For the purposes of the testing we did not use the initial values for the safeguarding parameters described in § 3. Instead we used the following cruder bounds:

$$\begin{aligned} \lambda_S &= \max \{-h_{ii}, i = 1, \dots, n, \}, \\ \lambda_B &= \max \left\{ 0, \lambda_S, \frac{\|g\|}{\Delta} - \|B\|_1 - \frac{1}{\mu} \|A\|_1 \|A\|_\infty \right\}, \\ \lambda_U &= \frac{\|g\|}{\Delta} + \|B\|_1 + \frac{1}{\mu} \|A\|_1 \|A\|_\infty. \end{aligned}$$

These initial values were used to make for a more stringent test of the rest of the algorithm.

The charts also deserve some explanation. For both GQTPAR and TRSQPF we report the minimum, maximum, and average number of iterations used to solve 10 distinct trust region problems for each triple (n, t, μ) . For GQTPAR we also note the number of problems (out of 10) on which it ran for 20 iterations. Moré and Sorensen interpret this as a failure to solve the problem and so we report it under the #crash column. The subsequent entries in the min, max, and mean columns were computed from the results of those problems on which GQTPAR halted in less than 20 iterations. Furthermore, to account for the work done by TRSQPF to compute the initial value for α , we add an appropriate amount to the average number of iterations. For specific values of n and t , the number of floating point operations (flops) required to compute

the initial α is

$$(4.2) \quad \frac{2}{3}n^3 + t^2\left(n - \frac{t}{3}\right)$$

while the work count per iteration of TRSQPF is given by the formula

$$(4.3) \quad \frac{(n+t)^3}{6} + 7(n+t)^2.$$

The adjustment is made by evaluating the above formulas and dividing (4.3) into (4.2). For example, for $n = 40$ and $t = 20$, the ratio of (4.2) to (4.3) is .9 and this is the amount we add to the average number of iterations. Thus, the average number flops per problem can be obtained by multiplying the adjusted mean number of iterations times the number of operations per iteration (i.e., equation (4.3)). In all problems $\sigma_1 = 0.01$ and $\sigma_2 = 0.0$.

To conclude this section we would like to offer explanations for some of the results presented in the Tables 1-12. When comparing average numbers of iterations we can see that TRSQPF compares quite favorably with GQTPAR in the general and hard cases. On the basis of some informal experimentation we believe that both the new starting value for λ (i.e., equation (3.21)) and the more complicated way of updating λ (i.e., Algorithm 3.6) are responsible. In the hard case TRSQPF is aided further by the condition estimator put forth in Algorithm 3.10. By contrast, the average numbers of iterations in the positive-definite case are nearly the same except when $\mu = 10^{-16}$. In the positive-definite case both algorithms proceed until a candidate value of λ is found that is greater than $-\lambda_1$ and less than zero. When that happens, the safeguarding scheme resets λ to zero, and the problem is solved. Since neither algorithm is especially well tailored to the positive-definite case, neither is especially adept at it, and so the results are unsurprisingly similar. The saddle-point case is where TRSQPF registers the sharpest improvement over GQTPAR. The saddle-point case is, essentially, a search for the direction of most negative curvature. It stands to reason then that the new initial value for λ and the condition estimator in Algorithm 3.10 are the primary factors contributing to TRSQPF's good performance in the saddle-point case.

In all of the testing presented here, we set $\sigma_1 = 0.01$ and $\sigma_2 = 0.0$ that essentially require both algorithms to return a solution in which the first two significant digits of the trust region step are correct. We feel this represents a realistic test because within the context of a global algorithm to minimize penalty functions and ultimately solve equality constrained minimization problems, the accuracy requested of a trust region solver is typically kept fixed. Despite this, however, some informal tests were conducted in which a much more accurate solution was requested of both TRSQPF and GQTPAR. The results indicated that for a given trust region problem of the type considered in this paper, the maximum number of correct significant digits in the trust region step GQTPAR could return was for the most part determined by the value of μ and the machine precision. For example, if μ is 10^{-9} and the machine precision is on the order of 10^{-17} , then in general at most the first $17-9=8$ digits of the trust region step would be correct. TRSQPF, however, would in general return a solution correct to nearly full precision.

5. Concluding remarks. Our objective in this work has been to develop a technique for the solution of the model trust region problem (1.1) under the assumption that the objective function is the L_2 penalty function. In particular, we have tailored the Moré and Sorensen [1983] algorithm with this specific structure in mind. Our numerical

experiments indicate that it can be beneficial to use this modified algorithm. This is particularly true when increased accuracy is demanded, the penalty parameter is quite small, or the current quadratic exhibits a saddle point. It should be noted that in this context algorithm GQTPAR requires additional preprocessing since the product AA^T must be explicitly determined. On the other hand, TRSQPF requires computation of the Bunch-Parlett factorization of the extended system that is more expensive than the Cholesky factorization of the Hessian matrix.

With respect to the method presented in this paper there are several avenues of research we intend to pursue. Among the most obvious possible things we could do is to extend the existing method to the sparse case. At first glance it might seem that the technology used to exploit sparsity in the Cholesky factorization could be adapted to the Bunch-Kaufman factorization. The latter, however, utilizes a complicated pivoting scheme that is required to guarantee stability. It is not clear that we can exploit sparsity and maintain stability in the Bunch-Kaufman factorization. Another possible endeavor would be to develop an implementation of the Bunch-Kaufman factorization on a parallel architecture. Here again the pivoting scheme would no doubt hinder our efforts. Throughout this paper we have assumed that the matrix A is of full column rank. The case where A is rank deficient (or nearly rank deficient) will be the focus of some future work.

The most important remaining project, and the one we plan to dispatch first, is to effectively use the algorithm proposed in this paper in a global method for solving equality constrained optimization problems. Existing methods suffer from often having to decrease μ slowly and consequently can be inefficient. Our goal is to overcome such problems.

Appendix A. We shall now restate and prove Lemma 3.11.

LEMMA 3.11. *If $\lambda \rightarrow -\lambda_1^+$, then for the \hat{z} from Algorithm 3.10, $\hat{z}^T(H + \lambda I)\hat{z} \rightarrow 0$.*

Proof. The proof is basically a suitable generalization of the proof given by Moré and Sorensen [1983]. Unfortunately, our use of the extended system and the Bunch-Kaufman factorization complicates matters. Therefore, to facilitate understanding the proof we begin by outlining the proof strategy.

- (I) Show that as $\lambda \rightarrow -\lambda_1^+$ and $H + \lambda I$ becomes singular, then the extended system X (equation (2.9)) also becomes singular.
- (II) For the vector $\hat{y} = y/\|y\|$ (where y is computed by Algorithm 3.10), prove that $\hat{y}^T X \hat{y} \rightarrow 0$ as $\lambda \rightarrow -\lambda_1^+$.
- (III) If $y^T = (z_1^T, z_2^T)$ where $z_1 \in R^n$ and $z_2 \in R^l$ and $\hat{z} = z_1/\|z_1\|$, show that $\hat{z}^T(H + \lambda I)\hat{z} \rightarrow 0$ as $\lambda \rightarrow -\lambda_1^+$.

Let v_1 be an eigenvector of H corresponding to λ_1 . The $(H + \lambda I)v_1 = (\lambda_1 + \lambda)v_1$ and so $\|(H + \lambda I)v_1\| \rightarrow 0$ as $\lambda \rightarrow -\lambda_1^+$. We now observe

$$(A1) \quad \begin{pmatrix} B + \lambda I & A \\ A^T & -\mu I \end{pmatrix} \begin{pmatrix} v_1 \\ (1/\mu)A^T v_1 \end{pmatrix} = \begin{pmatrix} (\lambda_1 + \lambda)v_1 \\ 0 \end{pmatrix}.$$

This shows that some singular value of X tends to zero as $\lambda \rightarrow -\lambda_1$. Since $\|X\|_2$ is bounded in a neighborhood of $-\lambda_1$ we may conclude that the determinant of X goes to zero, and hence X becomes singular, as $\lambda \rightarrow -\lambda_1$.

The proof in this section relies heavily on the results of Bunch and Kaufman [1977]. They show that if $\tilde{M} = \tilde{U}S^{-1}$ then $|\tilde{M}_{ij}| \leq 1/(1 - .6404)$, and if $\tilde{D} = SDS$ then $|\tilde{D}_{ij}| \leq (2.57)^{n+l-1} \max_{1 \leq i, j \leq n+l} |X_{ij}|$. Again, since the entries of X are bounded in a neighborhood of $-\lambda_1$, it follows that the entries of $\tilde{U}DS$ are also bounded. The orthogonality of Q further implies that the entries of R remain bounded as $\lambda \rightarrow -\lambda_1^+$.

Let $\rho > 0$ satisfy

$$(A2) \quad \sum_{j=1}^{n+t} \sum_{i=1}^j |r_{ij}| \leq \rho.$$

It follows from the way we compute w_k^+ and w_k^- in Algorithm 3.10 that

$$(A3) \quad \max(|w_k^+|, |w_k^-|) \geq \frac{1}{|r_{kk}|}$$

and for our choice of w_k we now have

$$(A4) \quad \frac{1}{|r_{kk}|} \leq |w_k| + \sum_{i=1}^{k-1} |\theta_i + r_{ik} w_k|.$$

This together with equation (A2) implies

$$(A5) \quad \frac{1}{|r_{kk}|} \leq (1+\rho) \|w\| \quad \text{or} \quad \|w\| \geq \frac{1}{(1+\rho)|r_{kk}|},$$

which holds for all k and so

$$(A6) \quad \frac{1}{\|w\|} \leq \min(|r_{kk}|: k=1, \dots, n+t)(1+\rho) \quad \text{or} \quad \|w\| \geq \frac{1}{\min(|r_{kk}|)(1+\rho)}.$$

Now since the entries of $S^{-1}\tilde{U}^T$ are bounded there exists a $\delta > 0$ such that $\|S^{-1}\tilde{U}^T\|_2 \leq \delta$. This along with equation (A6) and the orthogonality of P now imply

$$\|y\| = \|x\| \geq \frac{1}{\delta} \|w\| \geq \frac{1}{\delta(1+\rho) \min(|r_{kk}|)}.$$

If $\hat{y} = y/\|y\|$ then

$$\hat{y}^T X \hat{y} = \frac{\hat{y}^T e}{\|y\|} \leq \hat{y}^T e [\delta(1+\rho) \min(|r_{kk}|)]$$

and by the Cauchy-Schwartz inequality

$$\hat{y}^T X \hat{y} \leq (n+t)^{1/2} \delta(1+\rho) \min(|r_{kk}|).$$

It remains to show that $\min(|r_{kk}|) \rightarrow 0$ as $\lambda \rightarrow -\lambda_1^+$. To do so we must investigate, in some detail, the pivoting strategy and the choice of entries for S used by Bunch and Kaufman [1977]. For $l = n+t, \dots, 1$, let $A_{ij}^{(l)}$ denote the ij th entry in the $l \times l$ unfactored matrix remaining at the l th iteration of the Bunch-Kaufman factorization. At this point S_{ll} is chosen according to the following rule:

$$S_{ll} = \begin{cases} 1 & \text{if } (|A_{ll}^{(l)}| \geq \lambda \lambda^{(l)}) \text{ or } (|A_{rr}^{(l)}| \geq \lambda \sigma^{(l)}) \text{ or } (\lambda^{(l)} \geq \sigma^{(k)}), \\ \min(\sqrt{\lambda \sigma^{(l)}}/|A_{ll}^{(l)}|, \sigma^{(l)}/\lambda^{(l)}) & \text{otherwise} \end{cases}$$

where $\lambda = .6404$ and $\lambda^{(l)}$ is the absolute value of the largest off-diagonal element in the last column of $A^{(l)}$. That is,

$$\lambda^{(l)} = \max_{1 \leq i < l} |A_{il}^{(l)}|$$

and if r is the least integer such that $|A_{rl}^{(l)}| = \lambda^{(l)}$ then $\sigma^{(l)}$ is the largest off-diagonal element in the r th column of $A^{(l)}$,

$$\sigma^{(l)} = \max_{\substack{1 \leq m \leq l \\ m \neq r}} |A_{mr}^{(l)}|.$$

It should be noted that $S_{ii} \neq 1$ only when a 2×2 pivot is used or when $\sigma^{(i)} > \lambda^{(i)}$, in which case a 1×1 pivot is used but no permutation is applied to $A^{(i)}$. We now observe that since \tilde{U} is a product Gauss transformations, it has determinant one. It follows that $|\det X| = |\det D|$ and $|\det R| = |\det DS|$ and hence if the determinant of X is small then so is the determinant of D . More specifically, either one of the 1×1 blocks of D must be small, or one of the 2×2 blocks must have a small determinant. For the case where $S_{ii} \neq 1$ and $\sigma^{(i)} > \lambda^{(i)}$ then a 1×1 pivot is used and $D_{ii} = A_{ii}^{(i)}$. Furthermore,

$$(DS)_{ii} \leq \sqrt{\lambda \sigma^{(i)} |A_{ii}^{(i)}|}$$

and since $\sigma^{(i)}$ and λ are bounded, this shows that if a 1×1 block in D has small magnitude then so must the corresponding entry in DS . The case where $S_{ii} \neq 1$ and a 2×2 block is used is more complicated. When this happens

$$\begin{pmatrix} D_{l-1,l-1} & D_{l-1,l} \\ D_{l,l-1} & D_{l,l} \end{pmatrix} = \begin{pmatrix} A_{rr}^{(i)} & \lambda^{(i)} \\ \lambda^{(i)} & A_{ll}^{(i)} \end{pmatrix}$$

and the conditions of the pivoting scheme imply

$$-\lambda^{(i)2} \left[\frac{\lambda |A_{rr}^{(i)}|}{\sigma^{(i)}} + 1 \right] \leq A_{ll}^{(i)} A_{rr}^{(i)} - \lambda^{(i)2} \leq -(1 - \lambda^2) \lambda^{(i)2} \leq 0$$

where the center term is the determinant of the 2×2 pivot above. The important implication of equation (A18) is that if the center term is small then so is $\lambda^{(i)}$. Moreover, if the last column of the 2×2 pivot is scaled by S_{ii} then the determinant of the new 2×2 block satisfies

$$S_{ii} |A_{ll}^{(i)} A_{rr}^{(i)} - \lambda^{(i)2}| \leq \lambda^{(i)} [\lambda |A_{rr}^{(i)}| + \sigma^{(i)}].$$

Taken together, these observations imply that $|\det DS| \rightarrow 0$ as $|\det D| \rightarrow 0$ and by virtue of equation (A15) we now know that $\min(r_{kk}) \rightarrow 0$ as $\lambda \rightarrow -\lambda_1^+$. This completes the proof of the second part.

To prove the last part we recall that $Xy = PQe$ and if we let $PQe = (u_1, u_2)^T$, then

$$(A7) \quad \begin{pmatrix} B + \lambda I & A \\ A^T & -\mu I \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$

from which it follows that

$$\hat{z}^T (H + \lambda I) \hat{z} = \left(\hat{z}^T u_1 + \frac{1}{\mu} \hat{z}^T A u_2 \right) / \|z_1\|.$$

It is easy to see that the numerator of the second term is bounded and hence all we need to show is that $\|z_1\|$ must grow as $\|y\|$ becomes large. Multiplying out equation (A7) gives

$$(B + \lambda I) z_1 + A z_2 = u_1$$

and since A has full column rank and $\|u_1\|$ is bounded, we deduce that $\|z_1\| \rightarrow \infty$ as $\|z_2\| \rightarrow \infty$. This completes the proof. \square

REFERENCES

- J. R. BUNCH AND L. KAUFMAN [1977], *Some stable methods for calculating inertia solving symmetric linear systems*, Math. Comp., 31, pp. 163-179.
 J. E. DENNIS, JR. AND R. B. SCHNABEL [1983], *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Prentice-Hall, Englewood Cliffs, NJ.

- J. J. DONGARRA, J. R. BUNCH, C. B. MOLER, AND G. W. STEWART [1979], *Linpack Users' Guide*, Society for Industrial and Applied Mathematics, Philadelphia, PA.
- G. H. GOLUB AND C. F. VAN LOAN [1983], *Matrix Computations*, The Johns Hopkins University Press, Baltimore, MD.
- N. I. M. GOULD [1986], *On the accurate determination of search directions for simple differentiable penalty functions*, I.M.A. J. Numer. Anal., 6, pp. 357-372.
- M. D. HEBDEN [1973], *An algorithm for minimization using exact second derivatives*, Atomic Energy Research Establishment report T.P. 515, Harwell, England.
- J. J. MORÉ AND D. C. SORENSEN [1983], *Computing a trust region step*, SIAM J. Sci. Statist. Comput., 4, pp. 553-572.
- C. H. REINSCH [1967], *Smoothing by spline functions*, Numer. Math., 10, pp. 177-183.
- [1971], *Smoothing by spline functions II*, Numer. Math., 16, pp. 451-454.
- D. C. SORENSEN [1982], *Newton's method with a model trust region modification*, SIAM J. Numer. Anal., 19, pp. 409-426.