On the Convergence of the Newton/Log-Barrier Method^{*}

Stephen Wright[†]

Abstract

In the Newton/log-barrier method, Newton steps are taken for the log barrier function for a fixed value of the barrier parameter until a certain convergence criterion is satisfied. The barrier parameter is then decreased and the Newton process is repeated. A naive analysis indicates that Newton's method does not exhibit superlinear convergence to the minimizer of each instance of the log-barrier function until it reaches a very small neighborhood of the minimizer. By partitioning according to the subspace of active constraint gradients, however, we show that this neighborhood is actually quite large, thus explaining why reasonably fast local convergence can be attained in practice. Finally, we show that the overall convergence rate of the Newton/log-barrier algorithm is superlinear in the number of function/derivative evaluations, provided that the nonlinear program is formulated with a linear objective and the schedule for decreasing the barrier parameter is related in a certain way to the convergence criterion for each Newton process.

1 Introduction

We consider the nonlinear programming problem

min
$$f(x)$$
 subject to $c(x) \ge 0$, (1)

where $f : \mathbb{R}^n \to \mathbb{R}$ and $c : \mathbb{R}^n \to \mathbb{R}^m$ are smooth (twice Lipschitz continuously differentiable) functions. The logarithmic barrier function for (1) is

$$P(x;\mu) = f(x) - \mu \sum_{i=1}^{m} \ln c_i(x).$$
(2)

We denote by $x(\mu)$ a minimizer of $P(.;\mu)$ for $\mu > 0$ and assume that $x(\mu)$ exists for all sufficiently small μ . Methods based on (2) approximate $x(\mu)$ for a sequence of small, decreasing values of $\mu > 0$. Under certain conditions (see Fiacco and McCormick [3]), we have $\lim_{\mu \downarrow 0} x(\mu) = x^*$, where x^* is a local minimizer of (1).

The Newton/log-barrier method proceeds by fixing μ at a certain value and applying Newton's method to the unconstrained problem

$$\min_{x} P(x;\mu), \tag{3}$$

^{*}This research was supported by the Mathematics, Information, and Computational Sciences Division subprogram of the Office of Computational and Technology Research, U.S. Department of Energy, under Contract W-31-109-Eng-38.

[†]Mathematics and Computer Science Division, Argonne National Laboratory, 9700 South Cass Avenue, Argonne, Illinois 60439; wright@mcs.anl.gov



Figure 1: Typical sequence of Newton/log-barrier iterates

stopping the Newton iterations when some tolerance is satisfied. (Typically, the size of the Newton step or of the gradient $P_x(x;\mu)$ is required to fall below a certain threshold that depends on μ .) The barrier parameter μ is then decreased, and the process is repeated. A typical sequence of iterates is shown in Figure 1.

It is well known that the first Newton step for each value of μ —the one taken immediately after μ is decreased from its previous value—usually must be curtailed sharply to avoid leaving the feasible region. That is, a step length α considerably smaller than 1 is usually needed (see Conn, Gould, and Toint [1], Wright [4], and Wright and Jarre [6]). Often, however, steplengths of 1 can be taken safely after a few Newton iterations, yielding quadratic convergence toward $x(\mu)$. This phenomenon may not seem surprising, because it is also well known that, under typical nondegeneracy and second-order sufficiency conditions, the Hessian $P_{xx}(x;\mu)$ is positive definite at and near $x = x(\mu)$ for small values of μ . More investigation makes this simple explanation a little less satisfactory, however, because application of the standard theory seems to imply that the neighborhood within which quadratic convergence occurs becomes exceedingly small as μ approaches zero. In this paper, we show that this neighborhood is in fact not so small, suggesting that the number of Newton iterations required for each value of μ is not excessive.

We conclude the paper by using earlier results of Wright and Jarre [6] to show that, when the objective function f is linear and the convergence test for the Newton iterations at each value of μ has a certain form, the number of Newton iterates required for each μ value can be bounded above by a modest number. Moreover, μ can be decreased at a superlinear rate, so the overall convergence rate of the process is superlinear in the number of function/derivative evaluations. In sum, the Newton/log-barrier method—a purely primal method—can achieve a superlinear local convergence rate, just like its primal-dual counterpart.

We use the following notation in the rest of the paper. For related positive quantities α and β , we say $\beta = O(\alpha)$ if there is a constant M such that $\beta \leq M\alpha$ for all α sufficiently small. We

say that $\beta = o(\alpha)$ if $\beta/\alpha \to 0$ as $\alpha \to 0$, $\beta = \Omega(\alpha)$ if $\alpha = O(\beta)$, and $\beta = \Theta(\alpha)$ if $\beta = O(\alpha)$ and $\alpha = O(\beta)$. It follows that the expression $\beta = O(1)$ means that $\beta \leq M$ for some constant M and all values of β in the domain of interest.

2 Assumptions, Notation, and Basic Results

In this section, we specify the optimality conditions for the nonlinear program (1), outline our assumptions on its solution x^* , and state an implicit function theorem to be used in the next section.

The Lagrangian function for (1) is

$$\mathcal{L}(x,\lambda) = f(x) - \lambda^T c(x), \qquad (4)$$

where λ is the vector of Lagrange multipliers. The solution x^* of (1) satisfies the first-order conditions for optimality, which are that there exists a Lagrange multiplier vector λ^* such that

$$c(x^*) \ge 0, \quad \lambda^* \ge 0, \quad (\lambda^*)^T c(x^*) = 0, \quad \nabla f(x^*) = \sum_{i=1}^m \lambda_i^* \nabla c_i(x^*).$$
 (5)

The active constraints are the components of c for which $c_i(x^*) = 0$. Without loss of generality we assume these to be the first q components of c. We also assume that the solution is *nondegenerate*, that is,

$$[\nabla c_1(x^*) | \cdots | \nabla c_q(x^*)] \quad \text{has rank } q.$$
(6)

(Note that nondegeneracy implies uniqueness of λ^* .) We also assume *strict complementarity*, that is,

$$\lambda_i^* + c_i(x^*) > 0, \qquad i = 1, 2, \dots, m.$$
 (7)

Finally, we assume that second-order sufficient conditions for optimality are satisfied at (x^*, λ^*) , that is,

$$y^T \mathcal{L}_{xx}(x^*, \lambda^*) y > 0$$
 for all $y \neq 0$ with $\nabla c_i(x^*)^T y = 0$ for all $i = 1, 2, \dots, q$. (8)

It is easy to see that (x^*, λ^*) is a root of the function $F(x, \lambda)$ defined by

$$F(x,\lambda) = \begin{bmatrix} \mathcal{L}_x(x,\lambda) \\ \Lambda c(x) \end{bmatrix} = \begin{bmatrix} g - A(x)\lambda \\ \Lambda c(x) \end{bmatrix},$$
(9)

where

$$\Lambda = \operatorname{diag}(\lambda_1, \lambda_2, \cdots, \lambda_m).$$

The Jacobian of F is

$$\nabla F(x,\lambda) = \begin{bmatrix} \mathcal{L}_{xx}(x,\lambda) & -A(x) \\ \Lambda A(x)^T & C(x) \end{bmatrix},$$
(10)

where $C(x) = \text{diag}(c_1(x), c_2(x), \dots, c_m(x))$. When the nondegeneracy, strict complementarity, and second-order sufficient conditions hold, $\nabla F(x^*, \lambda^*)$ is nonsingular. It follows from the assumed smoothness of f and c that the Jacobian $\nabla F(x, \lambda)$ is nonsingular for all (x, λ) close to (x^*, λ^*) .

Given any strictly feasible point x and any positive value of the barrier parameter μ in (2), we define a vector of Lagrange multiplier estimates $\lambda(x, \mu)$ by

$$\lambda(x,\mu) = \mu C(x)^{-1} e = \left[\frac{\mu}{c_1(x)}, \dots, \frac{\mu}{c_m(x)}\right]^T.$$
 (11)

If x is the exact minimizer $x(\mu)$ of $P(\cdot;\mu)$, we define

$$\lambda(\mu) \stackrel{\Delta}{=} \lambda(x(\mu), \mu). \tag{12}$$

The derivatives of the barrier function (2) are

$$P_x(x;\mu) = \nabla f(x) - \sum_{i=1}^m \frac{\mu}{c_i(x)} \nabla c_i(x), \qquad (13a)$$

$$P_{xx}(x;\mu) = \nabla^2 f(x) + \mu \sum_{i=1}^m \left[\frac{1}{c_i^2(x)} \nabla c_i(x) \nabla c_i(x)^T - \frac{1}{c_i(x)} \nabla^2 c_i(x) \right].$$
(13b)

By combining (11) with (13a), we obtain

$$\nabla f(x) = A(x)\lambda(x,\mu) + P_x(x;\mu), \qquad (14)$$

while for the case $x = x(\mu)$, we have from (12) that $\nabla f(x(\mu)) = A(x(\mu))\lambda(\mu)$.

The following technical lemma, a consequence of the implicit function theorem, is proved in Wright and Jarre [6, Lemma 3.1].

Lemma 2.1 Let the vector pair $(\tilde{x}(z,\zeta), \tilde{\lambda}(z,\zeta))$ be defined implicitly as the solution of the nonlinear system

$$F(\tilde{x}, \tilde{\lambda}) = \begin{bmatrix} z \\ \zeta e \end{bmatrix}, \tag{15}$$

for given (z,ζ) and F defined as in (9). Then there are positive constants $\epsilon > 0$ and M > 0 such that the following statements hold.

(i) $(\tilde{x}(z,\zeta), \tilde{\lambda}(z,\zeta))$ is a C^2 function of (z,ζ) in the neighborhood defined by

$$\mathcal{N}_{\epsilon} = \{ (z, \zeta) \mid ||z|| + |\zeta| \le \epsilon \}.$$

(ii) For $\zeta > 0$ and $(z,\zeta) \in \mathcal{N}_{\epsilon}$, we have $\tilde{\lambda}_i(z,\zeta) > 0$ and $c_i(\tilde{x}(z,\zeta)) > 0$ for i = 1, 2, ..., m.

(iii) For (z_1, ζ_1) and (z_2, ζ_2) in \mathcal{N}_{ϵ} , we have

$$\begin{bmatrix} \tilde{x}(z_1,\zeta_1)\\ \tilde{\lambda}(z_1,\zeta_1) \end{bmatrix} - \begin{bmatrix} \tilde{x}(z_2,\zeta_2)\\ \tilde{\lambda}(z_2,\zeta_2) \end{bmatrix} = \nabla F(\tilde{x}(z_1,\zeta_1),\tilde{\lambda}(z_1,\zeta_1))^{-1} \begin{bmatrix} z_1-z_2\\ (\zeta_1-\zeta_2)e \end{bmatrix} + r, \quad (16)$$

where

$$||r|| \le M(||z_1 - z_2|| + |\zeta_1 - \zeta_2|)^2.$$
(17)

Uniform nonsingularity of ∇F in a neighborhood of (x^*, λ^*) implies that the first term in the right-hand side of (16) dominates the second-order term for ϵ sufficiently small; that is,

$$\left\| \begin{bmatrix} \tilde{x}(z_1,\zeta_1) \\ \tilde{\lambda}(z_1,\zeta_1) \end{bmatrix} - \begin{bmatrix} \tilde{x}(z_2,\zeta_2) \\ \tilde{\lambda}(z_2,\zeta_2) \end{bmatrix} \right\| = \Theta(\|z_1 - z_2\| + |\zeta_1 - \zeta_2|).$$
(18)

3 Convergence of Newton's Method to the Log-Barrier Minimizer

We now analyze the local convergence properties of Newton's method to the minimizer $x(\mu)$ of the barrier function $P(x;\mu)$, for a fixed value of μ . It is well known that under the second-order assumptions discussed above, the barrier function $P(x;\mu)$ has a minimizer $x(\mu)$ at which the Hessian is positive definite (though ill conditioned). Moreover, since the objective $f(\cdot)$ and constraint functions $c_i(\cdot)$, $i = 1, 2, \ldots, m$ are twice Lipschitz continuously differentiable, then $P(x;\mu)$ is also twice Lipschitz continuously differentiable near $x(\mu)$. Hence, quadratic convergence follows from a standard result (see, for example, Theorem 5.2.1 of Dennis and Schnabel [2]). If w_1, w_2, w_3, \ldots are the Newton iterates, the standard theory yields the estimate

$$\|w_{t+1} - x(\mu)\| \le L_1(\mu)L_2(\mu)\|w_t - x(\mu)\|^2,$$
(19)

where $L_1(\mu)$ is a Lipschitz constant for $P_{xx}(x;\mu)$ in the vicinity of $x(\mu)$ and $L_2(\mu)$ is a bound on $\|P_{xx}(x;\mu)^{-1}\|$ near $x(\mu)$. The estimates we derive later in this section show that

$$L_1(\mu) = O(\mu^{-2}), \quad L_2(\mu) = O(1),$$

so that (19) reduces to

$$||w_{t+1} - x(\mu)|| = O(\mu^{-2})||w_t - x(\mu)||^2.$$

This expression does not even imply convergence of the iteration sequence unless w_1 is in a very small neighborhood of the solution, specifically,

$$||w_1 - x(\mu)|| = O(\mu^2).$$
⁽²⁰⁾

It appears from this estimate that a number of careful line-search Newton iterations would be needed to move from the approximation to the minimizer $x(\mu_{-})$ at the previous value of μ into the neighborhood of quadratic convergence for the current value of μ .

Wright and Jarre [6] investigated the use of a reformulation of (1) in which the objective function is linear. They show that if the final approximation to $x(\mu_{-})$ obtained at the previous value μ_{-} of the barrier parameter is reasonably accurate, then the Newton step for $P(\cdot;\mu)$ from this point for the new value of μ passes quite close to the new minimizer $x(\mu)$. Even in this case, however, the resulting point will not generally lie in the neighborhood (20), except perhaps when a stringent stopping criterion of the form $||P_x(x;\mu)|| = O(\mu)$ is used at the previous value of μ .

In this section, we show that the expressions (19) and (20) are unduly pessimistic and, in fact, that there exists a constant $\bar{\mu} > 0$ such that quadratic convergence to $x(\mu)$ can be obtained from any point w that satisfies

$$\|w - x(\mu)\| \le C\mu^{\sigma}, \quad \text{for all } \mu \in (0, \bar{\mu}], \tag{21}$$

where C and σ are any given constants satisfying C > 0 and $\sigma > 1$. That is, the domain of quadratic convergence for $P(\cdot; \mu)$ shrinks as $\mu \downarrow 0$, but the rate of shrinkage is not especially severe.

We derive most of the results in an informal style, favoring the use of order notation for clarity. Explicit bounds are introduced for certain important estimates, to ensure that the inductive proofs of quadratic convergence at the end of this section are rigorous. Our analysis is based on the following consequence of Taylor's theorem. If w is the current iterate and s is the Newton step defined by $s = -P_{xx}(w;\mu)^{-1}P_x(w;\mu)$, we have that

$$P_{x}(w+s;\mu) = P_{x}(w;\mu) + P_{xx}(w;\mu)s + \int_{0}^{1} \left[P_{xx}(w+\tau s;\mu) - P_{xx}(w;\mu)\right]s \,d\tau$$
$$= \int_{0}^{1} \left[P_{xx}(w+\tau s;\mu) - P_{xx}(w;\mu)\right]s \,d\tau.$$
(22)

In the analysis below, and in the next section, we (implicitly) identify a value of $\bar{\mu}$ such that certain estimates are satisfied by certain functions of the point w satisfying (21). We assume without loss of generality that $\bar{\mu}$ is small enough that the neighborhood (21) excludes local minimizers of $P(\cdot; \mu)$ other than $x(\mu)$. Since $\sigma > 1$ in (21), it follows that the ratio $||w - x(\mu)||/\mu$ approaches zero as $\mu \downarrow 0$.

From the definition (12), an application of the implicit function result (18) with

$$\begin{aligned} &(z_1,\zeta_1) = (0,\mu), \text{ that is, } &(\tilde{x}(z_1,\zeta_1),\lambda(z_1,\zeta_1)) = (x(\mu),\lambda(\mu)), \\ &(z_2,\zeta_2) = (0,0), \text{ that is, } &(\tilde{x}(z_2,\zeta_2),\tilde{\lambda}(z_2,\zeta_2)) = (x^*,\lambda^*), \end{aligned}$$

yields

$$\left\| \left[\begin{array}{c} x(\mu) - x^* \\ \lambda(\mu) - \lambda^* \end{array} \right] \right\| = O(\mu).$$

Hence, by using (11), the strict complementarity assumption, and (21), we have for all active indices $i = 1, 2, \dots, q$ that

$$c_i(w) = c_i(x(\mu)) + O(||w - x(\mu)||) = \frac{\mu}{\lambda_i(\mu)} + O(\mu^{\sigma}) = \frac{\mu}{\lambda_i^*} + O(\mu^{\min(2,\sigma)}) = \Theta(\mu). \quad i = 1, 2, \dots, q,$$
(23)

for all μ sufficiently small.

Next, we examine the structure of $P_{xx}(w; \mu)$. By differentiating (2) twice, and partitioning the sums into active and inactive indices, we obtain

$$P_{xx}(w;\mu) = \sum_{i=1}^{q} \frac{\mu}{c_i^2(w)} \nabla c_i(w) \nabla c_i(w)^T + \left\{ \nabla^2 f(w) - \sum_{i=1}^{q} \frac{\mu}{c_i(w)} \nabla^2 c_i(w) \right\} - \sum_{i=q+1}^{m} \left\{ \frac{\mu}{c_i^2(w)} \nabla c_i(w) \nabla c_i(w)^T - \frac{\mu}{c_i(w)} \nabla^2 c_i(w) \right\}.$$
(24)

We deal with the three terms on the right-hand side in turn. Because of (23), we have

$$\frac{\mu}{c_i^2(w)} = \Omega(\mu^{-1}), \qquad i = 1, 2, \dots, q,$$

so by the nondegeneracy assumption, the first sum in (24) is a rank-q matrix, whose q nonzero eigenvalues are all positive with size $\Omega(\mu^{-1})$. Using a standard factorization, we can write

$$\sum_{i=1}^{q} \frac{\mu}{c_i^2(w)} \nabla c_i(w) \nabla c_i(w)^T = \hat{U}(w) \hat{D}(w) \hat{U}(w)^T,$$
(25)

where $\hat{D}(w)$ is a $q \times q$ diagonal matrix whose diagonal elements all have size $\Omega(\mu^{-1})$ and $\hat{U}(w)$ is an $n \times q$ orthonormal matrix whose columns span the range space of $[\nabla c_i(w)]_{i=1}^q$.

Since $||w - x(\mu)||/\mu = o(1)$ it follows that $\mu/c_i(w) \approx \lambda_i^*$, and the second term in (24) is a small perturbation of the Lagrangian Hessian $\mathcal{L}_{xx}(x^*, \lambda^*)$, which by our second-order assumption (8) is positive definite on the null space of $[\nabla c_i(x^*)]_{i=1}^q$. Hence, if we define $\tilde{U}(w)$ to be an $n \times (n-q)$ orthonormal matrix that spans the nearby null space of $[\nabla c_i(w)]_{i=1}^q$, (so that $[\hat{U}(w) | \tilde{U}(w)]$ is orthogonal), straightforward arguments show that the $(n-q) \times (n-q)$ matrix $G_{22}(w)$ defined by

$$\tilde{G}_{22}(w) \stackrel{\Delta}{=} \tilde{U}(w)^T \left\{ \nabla^2 f(w) - \sum_{i=1}^q \frac{\mu}{c_i(w)} \nabla^2 c_i(w) \right\} \tilde{U}(w) = \tilde{U}(w)^T \mathcal{L}_{xx}(x^*, \lambda^*) \tilde{U}(w) + o(1)$$

is positive definite, with all eigenvalues of size $\Omega(1)$.

Since $c_i(w) = \Omega(1)$ for i = q + 1, ..., m and all μ sufficiently small, the third term in (24) is $O(\mu)$.

By combining all these observations, we find that

$$P_{xx}(w;\mu) = \begin{bmatrix} \hat{U}(w) & \tilde{U}(w) \end{bmatrix} \begin{bmatrix} G_{11}(w) & G_{12}(w) \\ G_{12}^T(w) & G_{22}(w) \end{bmatrix} \begin{bmatrix} \hat{U}(w)^T \\ \tilde{U}(w)^T \end{bmatrix},$$
(26)

where

$$G_{11}(w) \stackrel{\Delta}{=} \hat{U}(w)^T P_{xx}(w;\mu) \hat{U}(w) = \hat{D}(w) + O(1),$$

$$G_{22}(w) \stackrel{\Delta}{=} \tilde{U}(w)^T P_{xx}(w;\mu) \tilde{U}(w) = \tilde{G}_{22}(w) + O(\mu),$$

$$G_{12}(w) \stackrel{\Delta}{=} \hat{U}(w)^T P_{xx}(w;\mu) \tilde{U}(w) = O(1),$$
(27)

where $\hat{D}(w)$ and $\tilde{G}_{22}(w)$ are defined as above.

To uncover the properties of $P_{xx}(w;\mu)^{-1}$, we use the following technical result about the inverse of a block 2×2 matrix; see, for example, Wright [5, Lemma 3.1].

Lemma 3.1 Let G be a symmetric matrix partitioned as

$$G = \left[\begin{array}{cc} G_{11} & G_{12} \\ G_{12}^T & G_{22} \end{array} \right],$$

where G_{11} and G_{22} are square. Suppose that G_{11} and $G_{22} - G_{12}^T G_{11}^{-1} G_{12}$ are nonsingular. Then G is nonsingular and G^{-1} has the form

$$G^{-1} = \begin{bmatrix} H_{11} & H_{12} \\ H_{12}^T & H_{22} \end{bmatrix},$$
 (28)

where

$$\begin{aligned} H_{11} &= G_{11}^{-1} + G_{11}^{-1} G_{12} (G_{22} - G_{12}^T G_{11}^{-1} G_{12})^{-1} G_{12}^T G_{11}^{-1} \\ H_{12} &= -G_{11}^{-1} G_{12} (G_{22} - G_{12}^T G_{11}^{-1} G_{12})^{-1} \\ H_{22} &= (G_{22} - G_{12}^T G_{11}^{-1} G_{12})^{-1}. \end{aligned}$$

In our case—the 2 × 2 block matrix in (27)—we have $\hat{D}(w)^{-1} = O(\mu)$ and so

$$G_{11}(w)^{-1} = \hat{D}(w)^{-1}(I + O(\mu))^{-1} = O(\mu).$$

By using this estimate, we further obtain

$$\left[G_{22}(w) - G_{12}(w)^T G_{11}(w)^{-1} G_{12}(w)\right]^{-1} = \left[\tilde{G}_{22}(w) + O(\mu)\right]^{-1} = \tilde{G}_{22}(w)^{-1} + O(\mu) = O(1),$$

Hence, from (26) and (28), we obtain

$$P_{xx}(w;\mu)^{-1} = \begin{bmatrix} \hat{U}(w) & \tilde{U}(w) \end{bmatrix} \begin{bmatrix} H_{11}(w) & H_{12}(w) \\ H_{12}^T(w) & H_{22}(w) \end{bmatrix} \begin{bmatrix} \hat{U}(w)^T \\ \tilde{U}(w)^T \end{bmatrix},$$
(29)

where

$$H_{11}(w) = O(\mu), \quad H_{12}(w) = O(\mu), \quad H_{22}(w) = O(1).$$
 (30)

We now examine the structure of $P_x(w;\mu)$. From (21) and (23) we have

$$||w - x(\mu)||/c_i(w) = O(\mu^{\sigma-1}) \ll 1, \qquad i = 1, 2, \dots, q,$$
(31)

for all μ sufficiently small. We can use this expression to estimate the difference between the reciprocals $c_i^{-1}(w)$ and $c_i^{-1}(x(\mu))$. We have

$$c_i^{-1}(x(\mu)) = [c_i(w) + O(||w - x(\mu)||)]^{-1}$$

= $c_i^{-1}(w) \left[1 + O(c_i(w)^{-1}||w - x(\mu)||) \right]^{-1}$
= $c_i^{-1}(w) + c_i^{-2}(w)O(||w - x(\mu)||).$ (32)

From (2), noting that $P_x(x(\mu);\mu) = 0$ and partitioning active and inactive indices, we have that

$$P_{x}(w;\mu) = P_{x}(w;\mu) - P_{x}(x(\mu);\mu)$$

= $\sum_{i=1}^{q} \left[\frac{\mu}{c_{i}(w)} \nabla c_{i}(w) - \frac{\mu}{c_{i}(x(\mu))} \nabla c_{i}(x(\mu)) \right] + \left[\nabla f(w) - \nabla f(x(\mu)) \right]$ (33)
+ $\sum_{i=q+1}^{m} \left[\frac{\mu}{c_{i}(w)} \nabla c_{i}(w) - \frac{\mu}{c_{i}(x(\mu))} \nabla c_{i}(x(\mu)) \right] .$

For the second term on the right-hand side of (33), we have by smoothness of f that $\nabla f(w) - \nabla f(x(\mu)) = O(||w - x(\mu)||)$. In the third term, we have for each index $i = q + 1, \ldots, m$ that $c_i(w) = \Omega(1)$ and $c_i(x(\mu)) = \Omega(1)$. Hence, by smoothness of ∇c_i , this term has size $O(\mu ||w - x(\mu)||)$.

In the first term, we have for each active index *i*, using the smoothness of ∇c_i together with (23), (31), and (32), that

$$\frac{\mu}{c_i(x(\mu))} \nabla c_i(x(\mu)) = \left[\frac{\mu}{c_i(w)} + \frac{\mu}{c_i^2(w)} O(\|w - x(\mu)\|) \right] \left[\nabla c_i(w) + O(\|w - x(\mu)\|) \right] \\
= \frac{\mu}{c_i(w)} \nabla c_i(w) + O(\mu^{-1} \|w - x(\mu)\|) \nabla c_i(w) + O(\|w - x(\mu)\|).$$
(34)

By collecting these observations and substituting into (33), we obtain

$$P_x(w;\mu) = \sum_{i=1}^q O(\mu^{-1} ||w - x(\mu)||) \nabla c_i(w) + O(||w - x(\mu)||).$$
(35)

From the definitions of $\hat{U}(w)$ and $\tilde{U}(w)$, it follows immediately that

$$\hat{U}(w)^T P_x(w;\mu) = O(\mu^{-1} ||w - x(\mu)||), \qquad (36a)$$

$$\tilde{U}(w)^T P_x(w;\mu) = O(||w - x(\mu)||).$$
 (36b)

Meanwhile from (30), we have for the Newton step from w that

$$s = -P_{xx}(w;\mu)^{-1}P_{x}(w;\mu) = \begin{bmatrix} \hat{U}(w) & \tilde{U}(w) \end{bmatrix} \begin{bmatrix} H_{11}(w) & H_{12}(w) \\ H_{12}^{T}(w) & H_{22}(w) \end{bmatrix} \begin{bmatrix} \hat{U}(w)^{T}P_{x}(w;\mu) \\ \tilde{U}(w)^{T}P_{x}(w;\mu) \end{bmatrix} = \begin{bmatrix} \hat{U}(w) & \tilde{U}(w) \end{bmatrix} \begin{bmatrix} O(\mu) & O(\mu) \\ O(\mu) & O(1) \end{bmatrix} \begin{bmatrix} \hat{U}(w)^{T}P_{x}(w;\mu) \\ \tilde{U}(w)^{T}P_{x}(w;\mu) \end{bmatrix}.$$

Hence, by combining with (36), we find that we can choose $\bar{\mu}$ in (21) such that for all w and μ satisfying (21), we have

$$\|s\| \leq O(\mu \| \hat{U}(w)^T P_x(w; \mu) \| + \| \tilde{U}(w)^T P_x(w; \mu) \|)$$

$$\leq C_3 \left[\mu \| \hat{U}(w)^T P_x(w; \mu) \| + \| \tilde{U}(w)^T P_x(w; \mu) \| \right]$$
 (37a)

$$\leq C_1 \| w - x(\mu) \| = O(\mu^{\sigma}),$$
 (37b)

for some positive numbers C_1 and C_3 independent of μ .

Note that the naive estimate of ||s|| obtained by ignoring the structure of $P_x(w;\mu)$ and $P_{xx}(w;\mu)$ would be simply $||s|| \leq ||P_{xx}(w;\mu)^{-1}|| ||P_x(w;\mu)|| = O(\mu^{\sigma-1})$, which is too pessimistic for our purposes.

We now examine the integrand in (22), partitioning it into the subspaces defined by the active constraint gradients at the next Newton iterate w + s. We start by partitioning the integrand as follows:

$$[P_{xx}(w+\tau s;\mu) - P_{xx}(w;\mu)]s = q_1 + q_2 + q_3 + q_4,$$
(38)

where $\tau \in [0, 1]$ and

$$\begin{split} q_{1} &= \sum_{i=1}^{q} \left\{ \frac{\mu}{c_{i}^{2}(w+\tau s)} \nabla c_{i}(w+\tau s) \nabla c_{i}(w+\tau s)^{T}s - \frac{\mu}{c_{i}^{2}(w)} \nabla c_{i}(w) \nabla c_{i}(w)^{T}s \right\} \\ q_{2} &= -\sum_{i=1}^{q} \left\{ \frac{\mu}{c_{i}(w+\tau s)} \nabla^{2}c_{i}(w+\tau s)s - \frac{\mu}{c_{i}(w)} \nabla^{2}c_{i}(w)s \right\} \\ q_{3} &= \left[\nabla^{2}f(w+\tau s) - \nabla^{2}f(w) \right]s, \\ q_{4} &= \sum_{i=q+1}^{m} \left\{ \frac{\mu}{c_{i}^{2}(w+\tau s)} \nabla c_{i}(w+\tau s) \nabla c_{i}(w+\tau s)^{T}s - \frac{\mu}{c_{i}^{2}(w)} \nabla c_{i}(w) \nabla c_{i}(w)^{T}s - \frac{\mu}{c_{i}(w+\tau s)} \nabla^{2}c_{i}(w+\tau s)s + \frac{\mu}{c_{i}(w)} \nabla^{2}c_{i}(w)s \right\}. \end{split}$$

To estimate the first term q_1 , we note from (37b) that $\mu^{-1} ||s|| = O(\mu^{\sigma-1}) \ll 1$. Hence, as in (32), we have

$$\frac{1}{c_i(w+\tau s)} = \frac{1}{c_i(w)} + \frac{1}{c_i^2(w)}O(\|\tau s\|) = \frac{1}{c_i(w)} + O(\mu^{-2}\|s\|), \quad i = 1, 2, \dots, q,$$
(39a)

$$\frac{1}{c_i^2(w+\tau s)} = \frac{1}{c_i^2(w)} \left[1 + \frac{1}{c_i(w)} O(\|\tau s\|) \right]^{-2} = \frac{1}{c_i^2(w)} + O(\mu^{-3} \|s\|), \quad i = 1, 2, \dots, q.$$
(39b)

By smoothness of ∇c_i , we obtain

$$\frac{\mu}{c_i^2(w+\tau s)} \nabla c_i(w+\tau s) \nabla c_i(w+\tau s)^T s
= \left[\frac{\mu}{c_i^2(w)} + O(\mu^{-2} \|s\|)\right] \left[\nabla c_i(w) + O(\|s\|)\right] \left[\nabla c_i(w)^T s + O(\|s\|^2)\right]
= \frac{\mu}{c_i^2(w)} \nabla c_i(w) \nabla c_i(w)^T s + O(\mu^{-2} \|s\|^2) \nabla c_i(w) + O(\mu^{-1} \|s\|^2),$$

where we have used the fact that $\mu^{-1}||s|| \ll 1$ (see (37b)) to absorb higher-order terms. By substituting into the definition of q_1 , we obtain

$$q_1 = \sum_{i=1}^q O(\mu^{-2} \|s\|^2) \nabla c_i(w) + O(\mu^{-1} \|s\|^2) = \sum_{i=1}^q O(\mu^{-2} \|s\|^2) \nabla c_i(w+s) + O(\mu^{-1} \|s\|^2), \quad (40)$$

where the change of argument from w to w + s in the first term causes a perturbation that can be absorbed in the second term.

For the second term q_2 , we have from (39a) that

$$\frac{\mu}{c_i(w+\tau s)} \nabla^2 c_i(w+\tau s)s = \left[\frac{\mu}{c_i(w)} + O(\mu^{-1} \|s\|)\right] \left[\nabla^2 c_i(w)s + O(\|s\|^2)\right]$$
$$= \frac{\mu}{c_i(w)} \nabla^2 c_i(w)s + O(\mu^{-1} \|s\|^2).$$

The remaining terms q_3 and q_4 are less significant. By Lipschitz continuity of ∇f , we have $q_3 = O(||s||^2)$. In q_4 , the denominators all have size $\Omega(1)$, so it is easy to show that $q_4 = O(\mu ||s||^2)$.

By collecting all these estimates into (38), performing the integration, and substituting into (22), we obtain

$$P_x(w+s;\mu) = \sum_{i=1}^q O(\mu^{-2} ||s||^2) \nabla c_i(w+s) + O(\mu^{-1} ||s||^2).$$

Hence, after a possible adjustment in $\overline{\mu}$ in (21), we have that there is a positive number C_2 independent of μ such that

$$\hat{U}(w+s)^T P_x(w+s;\mu) \leq C_2 \mu^{-2} \|s\|^2,$$
(41a)

$$\tilde{U}(w+s)^T P_x(w+s;\mu) \leq C_2 \mu^{-1} ||s||^2,$$
(41b)

for any orthonormal matrix $\hat{U}(w+s)$ that spans the column space of $[\nabla c_i(x+s)]_{i=1}^q$, and for any orthonormal matrix $\tilde{U}(w+s)$ such that $[\hat{U}(w+s) | \tilde{U}(w+s)]$ is orthogonal. (For future reference, we assume without loss of generality that $C_2 \geq 1$.)

At this point, we have identified a threshold $\bar{\mu}$ and constants C_1 , C_2 , and C_3 such that if w lies in the neighborhood defined by (21) for given values of C and σ , the Newton step s satisfies the important relationships (37a), (37b), and (41). An important corollary of these relationships is that if the next Newton iterate w + s also lies in the neighborhood (21), then we have from (37a) and (41) that the next Newton step s_+ satisfies

$$\|s_{+}\| \le 2C_2 C_3 \mu^{-1} \|s\|^2.$$
(42)

We now use all these estimates to show that if we choose the starting point w_1 for the Newton iteration in a slightly more restrictive neighborhood than (21), then *all* Newton iterates will remain inside the full neighborhood (21) and quadratic convergence of the newton sequence to $x(\mu)$ will be observed. We state the result formally as a theorem.

Theorem 3.2 Let the constants C > 0 and $\sigma > 1$ be given, and let C_1 , C_2 , C_3 , and $\overline{\mu}$ be defined as above, in such a way that the relationships (37a), (37b), (41), and (42) hold. Let the constants $C_0 > 0$ and $\overline{\mu}_0$ be chosen in such a way that the following inequalities are satisfied:

$$(1+2C_1)C_0 \le C, \qquad 2C_0C_1C_2C_3\bar{\mu}_0^{\sigma-1} \le 1/4.$$
 (43)

Then if $\mu \in (0, \overline{\mu}_0]$ and w_1 is any point that satisfies

$$\|w_1 - x(\mu)\| \le C_0 \mu^{\sigma},\tag{44}$$

then Newton's method with full steps, applied to the function $P(\cdot; \mu)$ and starting from w_1 , converges Q-quadratically to $x(\mu)$.

Proof. By (37b), we have that the first Newton step s_1 satisfies

$$||s_1|| \le C_1 ||w_1 - x(\mu)||, \tag{45}$$

and so, because of the definition of C_0 in (43), the next iterate $w_2 = w_1 + s_1$ satisfies

$$||w_2 - x(\mu)|| \le ||w_1 - x(\mu)|| + ||s_1|| \le (1 + C_1)||w_1 - x(\mu)|| \le C_0(1 + C_1)\mu^{\sigma} < C\mu^{\sigma}.$$

Hence, w_2 also lies in the neighborhood (21), so we can apply (42) to obtain the following estimate for the next Newton step s_2 :

$$\|s_2\| \le 2C_2 C_3 \mu^{-1} \|s_1\|^2.$$
(46)

From (44) and (45), we have that

$$u^{-1} \|s_1\| \le C_0 C_1 \mu^{\sigma - 1},$$

so by substituting (46) and using the definition (43), we obtain

$$||s_2|| \le 2C_0C_1C_2C_3\mu^{\sigma-1}||s_1|| \le (1/4)||s_1||.$$

Hence, for the next Newton iterate $w_3 = w_2 + s_2$, we have

$$\|w_3 - x(\mu)\| \le \|w_1 - x(\mu)\| + \|s_1\| + \|s_2\| \le \|w_1 - x(\mu)\| + (5/4)\|s_1\| \le C_0(1 + (5/4)C_1)\mu^{\sigma} < C\mu^{\sigma},$$

so that w_3 also lies in the neighborhood defined by (21).

The argument continues inductively. We find in general that for all t = 1, 2, 3..., we have

$$\|s_{t+1}\| \le 2C_2C_3\mu^{-1}\|s_t\|^2 \le (2C_2C_3\mu^{-1}\|s_1\|)\|s_t\| \le (1/4)\|s_t\|,$$
(47)

and that

$$\|w_{t+1} - x(\mu)\| \leq \|w_1 - x(\mu)\| + \sum_{j=1}^t \|s_j\|$$

$$\leq \|w_1 - x(\mu)\| + \sum_{j=1}^t 4^{-(j-1)} \|s_1\|$$

$$\leq C_0 (1 + (4/3)C_1) \mu^{\sigma} < C \mu^{\sigma},$$

so that all Newton iterates w_1, w_2, w_3, \ldots belong to the neighborhood (21).

From (47), we have that $||s_t||$, t = 1, 2, 3... decreases geometrically (in fact, quadratically) to zero. Therefore, $\{w_t\}$ is a Cauchy sequence, so it converges, say to a point $w_*(\mu)$. It follows from (41) that this limit point must satisfy

$$P_x(w_*(\mu);\mu) = 0.$$

Moreover, by the second-order condition (8), we have by the choice of $\bar{\mu}$ and the discussion about the Hessian $P_{xx}(\cdot;\mu)$ and its inverse that $P_{xx}(w;\mu)$ is positive definite for all w satisfying (21). Hence, $w_*(\mu)$ is a local minimizer of $P(\cdot;\mu)$. Since $x(\mu)$ is the only local minimizer of this function in the neighborhood (21) by assumption, we must have $w_*(\mu) = x(\mu)$.

To prove that the convergence of $\{w_t\}$ to $x(\mu)$ is quadratic, we estimate the error $||w_t - x(\mu)||$ in terms of $||s_t||$. By using (47), we have for all t = 1, 2, 3, ... that

$$\|w_t - x(\mu)\| = \left\|\sum_{j=t}^{\infty} s_j\right\| \le \sum_{j=t}^{\infty} \|s_j\| \le \sum_{j=t}^{\infty} 4^{-(j-1)} \|s_t\| \le (4/3) \|s_t\|.$$

Similarly, we have that

$$||w_t - x(\mu)|| \ge ||s_t|| - \sum_{j=t+1}^{\infty} ||s_j|| \ge (2/3) ||s_t||.$$

Hence, from (42), we have

$$\|w_{t+1} - x(\mu)\| \le (4/3) \|s_{t+1}\| \le (8/3) C_2 C_3 \mu^{-1} \|s_t\|^2 \le (128/27) C_2 C_3 \mu^{-1} \|w_t - x(\mu)\|^2,$$

indicating that the convergence is Q-quadratic, as claimed.

Note from (44), (45), and (47) that we have

$$||s_{t+1}|| \le [2C_0C_1C_2C_3]^{2^t} (2C_2C_3)^{-1} \mu^{2^t(\sigma-1)+1}, \quad t = 0, 1, 2, \dots,$$

and therefore from (41) we have

$$\|P_x(w_{t+1};\mu)\| \le 2C_2\mu^{-2}\|s_t\|^2 \le \left[2C_0C_1C_2C_3\right]^{2^t}\left(2C_2C_3^2\right)^{-1}\mu^{2^t(\sigma-1)}, \quad t=1,2,3,\dots$$
(48)

4 Superlinear Convergence When the Objective Function Is Linear

The final result deals with the special case in which the objective function $f(\cdot)$ in (1) is linear. For background, we state a result due to Wright and Jarre [6], which shows that when we have a sufficiently good approximation to the log-barrier minimizer $x(\mu_{-})$ at the current barrier value μ_{-} , and we then reduce μ_{-} to μ (but not too rapidly), then the Newton step for $P_x(\cdot;\mu)$ passes quite close to the new minimizer $x(\mu)$. Moreover, the steplength $\alpha = \mu/\mu_{-}$ is asymptotically optimal. The result is a combination of Theorems 3.2 and 3.3 from [6].

Theorem 4.1 Suppose that f is linear and that the barrier parameter values μ_{-} and μ satisfy the condition

$$\mu \in [\rho_{\min}\mu_{-}^{\tilde{\sigma}}, \rho_{\max}\mu_{-}], \tag{49}$$

where $\rho_{\min} > 0$, $\rho_{\max} \in (0,1)$, and $\tilde{\sigma} \in (1,2]$ are constants. Suppose too that the bound

$$\|P_x(x;\mu_-)\| \le \mu_-^{\hat{\sigma}/2} \tag{50}$$

is satisfied at the current value of x, where $\hat{\sigma} \in (\tilde{\sigma}, 2]$. Then if s is the Newton direction for $P(\cdot; \mu)$ from x, the line segment $x + \tau(\mu/\mu_{-})s$, $\tau \in [0, 1]$ is strictly feasible. Moreover the function

$$\phi(\tau) \stackrel{\triangle}{=} P(x + \tau(\mu/\mu_{-})s;\mu)$$

has a local minimizer τ^* such that

$$1 - \tau^* = O(\mu^{\hat{\sigma} - 1}).$$

Finally, for the value $\tau = 1$, we have

$$\|x + (\mu/\mu_{-})s - x(\mu)\| = O(\|P_x(x;\mu_{-})\|^2 + \mu_{-}^2) = O(\mu_{-}^{\hat{\sigma}}).$$
(51)

If we define

$$w_1 \stackrel{\triangle}{=} x + (\mu/\mu_-)s$$

to be the point obtained by taking this truncated Newton step, we have from (49) and (51) that

$$\|w_1 - x(\mu)\| = O(\mu^{\hat{\sigma}/\tilde{\sigma}}).$$
(52)

Choosing σ in the range $(1, \hat{\sigma}/\tilde{\sigma})$, we find that there is a threshold $\bar{\mu}_1 \in (0, \bar{\mu}_0]$ such that

$$\|w_1 - x(\mu)\| = O(\mu^{\hat{\sigma}/\tilde{\sigma} - \sigma})\mu^{\sigma} \le C_0\mu^{\sigma} \quad \text{for all } \mu \in (0, \bar{\mu}_1].$$

Hence, the results of Theorem 3.2, and in particular the relation (48), apply when μ is sufficiently small in this sense. If we use a convergence criterion of the form (50) at this iteration too, that is,

$$\|P_x(w_{t+1};\mu)\| \le \mu^{\hat{\sigma}/2},\tag{53}$$

we can prove the following theorem about the number of Newton iterations needed to satisfy this bound.

Theorem 4.2 Suppose that the assumptions of Theorem 4.1 hold and that, given some choice of σ in the range $(1, \hat{\sigma}/\tilde{\sigma})$, the assumptions and notation of Theorem 3.2 hold as well. Assume without loss of generality that $2C_2C_3^2 \ge 1$. Define the constant $\bar{\mu}_2$ by

$$\bar{\mu}_2 \stackrel{\Delta}{=} \min\left(1, \bar{\mu}_1, \bar{\mu}_0^2\right).$$

Then for $\mu \in (0, \overline{\mu}_2]$, the criterion (53) is satisfied for all t with

$$t \ge \frac{\log(\hat{\sigma}/(\sigma-1))}{\log 2}.$$
(54)

Proof. Because $\log \mu \leq 2 \log \overline{\mu}_0$, we have from (43) that

$$\log(2C_0C_1C_2C_3) + \frac{\sigma - 1}{2}\log\mu \le \log(2C_0C_1C_2C_3) + (\sigma - 1)\log\bar{\mu}_0 < 0.$$

Hence, for the log of the right-hand side of (48), and using $2C_2C_3^2 \ge 1$, we have that

$$2^{t} \log(2C_{0}C_{1}C_{2}C_{3}) - \log(2C_{2}C_{3}^{2}) + 2^{t}(\sigma - 1) \log \mu \le 2^{t} \frac{\sigma - 1}{2} \log \mu.$$

Hence, we see that (53) is satisfied if

$$2^t \frac{\sigma - 1}{2} \log \mu \le \frac{\hat{\sigma}}{2} \log \mu,$$

which, since $\log \mu < 0$, is equivalent to

$$2^t(\sigma - 1) \ge \hat{\sigma}.$$

The result follows immediately.

For reasonable values of $\tilde{\sigma}$ and $\hat{\sigma}$, the required values of t are quite small. We give two examples:

- (i) $\hat{\sigma} = 2$ and $\tilde{\sigma} = 1.5$, giving a convergence tolerance of $||P_x(x;\mu)|| \le \mu$. We choose $\sigma = 1.3$ to lie in the range (1, 2/1.5) = (1, 4/3). Then (54) yields the bound $t \ge 3$.
- (ii) $\hat{\sigma} = 1.5$ and $\tilde{\sigma} = 1.25$, giving a convergence tolerance of $||P_x(x;\mu)|| \le \mu^{.75}$. We choose $\sigma = 1.15$ to lie in the range (1, 1.5/1.25) = (1, 1.2). Then (54) yields $t \ge 4$.

In both cases, we can take a "superlinear" decrease in μ —with $\mu = O(\mu_{-}^{1.5})$ and $\mu = O(\mu_{-}^{1.25})$, respectively—and take at most four or five Newton steps to move from an approximate minimizer of $P(\cdot; \mu_{-})$ to an approximate minimizer of $P(\cdot; \mu)$. The number of Newton steps per value of μ is bounded by a constant, so the overall rate of convergence of the Newton/log-barrier process to x^* is superlinear.

Acknowledgment

I am extremely grateful to Florian Jarre for his advice and many helpful comments on earlier drafts of this paper.

References

- A. R. CONN, N. I. M. GOULD, AND P. L. TOINT, A note on using alternative second-order models for the subproblems arising in barrier function methods for minimization, Numerische Mathematik, 68 (1994), pp. 17–33.
- [2] J. E. DENNIS AND R. B. SCHNABEL, Numerical Methods for Unconstrained Optimization, Prentice-Hall, Englewood Cliffs, NJ, 1983.
- [3] A. V. FIACCO AND G. P. MCCORMICK, Nonlinear Programming: Sequential Unconstrained Minimization Techniques, Wiley, New York, 1968. reprinted by SIAM Publications, 1990.
- [4] M. H. WRIGHT, Why a pure primal Newton barrier step may be infeasible, SIAM Journal on Optimization, 5 (1995), pp. 1–12.
- [5] S. J. WRIGHT, Stability of linear equations solvers in interior-point methods, SIAM Journal on Matrix Analysis and Applications, 16 (1994), pp. 1287–1307.
- [6] S. J. WRIGHT AND F. JARRE, The role of linear objective functions in barrier methods, Preprint MCS-P485-1294, Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, Ill., 1994. Revised, August 1997.