



Superlinear Convergence of a Stabilized SQP Method to a Degenerate Solution

STEPHEN J. WRIGHT*

Mathematics and Computer Science Division, Argonne National Laboratory, 9700 South Cass Avenue, Argonne, Illinois 60439, U.S.A.

Received February 13, 1997; Revised November 17, 1997; Accepted November 21, 1997

Abstract. We describe a slight modification of the well-known sequential quadratic programming method for nonlinear programming that attains superlinear convergence to a primal-dual solution even when the Jacobian of the active constraints is rank deficient at the solution. We show that rapid convergence occurs even in the presence of the roundoff errors that are introduced when the algorithm is implemented in floating-point arithmetic.

Keywords: nonlinear programming, sequential quadratic programming, degenerate solutions

1. Introduction

We describe a slight modification of the well-known sequential quadratic programming (SQP) algorithm for the nonlinear programming problem

$$\min \phi(z) \quad \text{subject to } g(z) \leq 0, \tag{1}$$

where $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are smooth functions. Our particular interest is in the case in which the active constraints at the solution z^* —those for which $g_i(z^*) = 0$ —have linearly dependent gradients. This property can interfere with the superlinear convergence rate of SQP, as we demonstrate with a simple example. The advantage of our *stabilized SQP* algorithm is that superlinear convergence is still attainable when this property holds.

The Lagrangian for (1) is

$$\mathcal{L}(z, \lambda) = \phi(z) + \sum_{i=1}^m \lambda_i g_i(z) = \phi(z) + \lambda^T g(z), \tag{2}$$

where $\lambda \in \mathbb{R}^m$ is the vector of Lagrange multipliers. When a constraint qualification holds at z^* (see discussion below), first-order necessary conditions for $z^* \in \mathbb{R}^n$ to be a solution

*This work was supported by the Mathematical, 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.

of (1) are that there exists a vector $\lambda^* \in \mathbb{R}^m$ such that

$$\nabla_z \mathcal{L}(z^*, \lambda^*) = 0, \quad g(z^*) \leq 0, \quad \lambda^* \geq 0, \quad (\lambda^*)^T g(z^*) = 0. \quad (3)$$

These relations are the well-known Karush-Kuhn-Tucker (KKT) conditions. The *active set* at z^* is defined by

$$\mathcal{B} = \{i = 1, 2, \dots, m \mid g_i(z^*) = 0\}, \quad (4)$$

while its complement is

$$\mathcal{N} = \{1, 2, \dots, m\} \setminus \mathcal{B}. \quad (5)$$

We partition the Lagrange multiplier vector λ and the function $g(\cdot)$ according to $\mathcal{B} \cup \mathcal{N}$ as

$$g_{\mathcal{B}}(\cdot) = [g_i(\cdot)]_{i \in \mathcal{B}}, \quad g_{\mathcal{N}}(\cdot) = [g_i(\cdot)]_{i \in \mathcal{N}}, \quad \lambda_{\mathcal{B}} = [\lambda_i]_{i \in \mathcal{B}}, \quad \lambda_{\mathcal{N}} = [\lambda_i]_{i \in \mathcal{N}}.$$

For notational convenience, we often omit transpose notation and write $\lambda = (\lambda_{\mathcal{B}}, \lambda_{\mathcal{N}})$ and $(z, \lambda) = (z, \lambda_{\mathcal{B}}, \lambda_{\mathcal{N}})$.

In this article, we assume that for at least one of the points (z^*, λ^*) satisfying the first-order conditions (3), we have

$$\lambda_{\mathcal{B}}^* > 0; \quad (6)$$

that is, *strict complementarity* holds. Since $g_i(z^*) < 0$ for $i \in \mathcal{N}$ (by the definition (4)), we have from (3) that $\lambda_{\mathcal{N}}^* = 0$ for all λ^* that satisfy (3). We also assume the following second-order sufficient condition: For any λ^* such that (z^*, λ^*) satisfies the KKT conditions (3), the two-sided projection of the Lagrangian Hessian $\nabla_{zz} \mathcal{L}(z^*, \lambda^*)$ onto $\ker \nabla g_{\mathcal{B}}(z^*)^T$ is positive definite. That is, there is a $\sigma > 0$ such that

$$w^T \nabla_{zz} \mathcal{L}(z^*, \lambda^*) w \geq \sigma \|w\|^2, \quad (7)$$

for all λ^* such that (z^*, λ^*) satisfies (3), and all $w \in \ker \nabla g_{\mathcal{B}}(z^*)^T$.

Finally, we assume that the Mangasarian-Fromovitz [9] constraint qualification (MFCQ) holds at z^* . That is,

$$\nabla g_{\mathcal{B}}(z^*)^T d < 0 \quad \text{for some } d \in \mathbb{R}^n. \quad (8)$$

This assumption is weaker than the linear independence constraint qualification, which assumes that $\nabla g_{\mathcal{B}}(z^*)$ has full columns rank and which is frequently used in the local convergence analysis of algorithms for nonlinear programming.

We use \mathcal{S} to denote the primal-dual solution set whose z -component is z^* , that is,

$$\mathcal{S} = \{(z^*, \lambda^*) \mid \lambda^* \text{ satisfies (3)}\}. \quad (9)$$

Since $\nabla_z \mathcal{L}(z^*, \lambda)$ is linear in λ , it follows immediately that \mathcal{S} is convex. We use \mathcal{S}_λ to denote the set of optimal Lagrange multipliers for z^* , that is,

$$\mathcal{S}_\lambda = \{\lambda^* \mid (z^*, \lambda^*) \in \mathcal{S}\}. \quad (10)$$

In the best-known form of the SQP algorithm, the step Δz is obtained by solving the following subproblem:

$$\min_{\Delta z} \Delta z^T \nabla \phi(z) + \frac{1}{2} \Delta z^T \nabla_{zz} \mathcal{L}(z, \lambda) \Delta z, \quad \text{subject to } g(z) + Dg(z) \Delta z \leq 0, \quad (11)$$

where (z, λ) is the current primal-dual iterate. Denoting the Lagrange multipliers for the constraints in (11) by λ^+ , we see that the solution Δz satisfies the following KKT conditions (cf. (3)):

$$\nabla_{zz} \mathcal{L}(z, \lambda) \Delta z + \nabla \phi(z) + Dg(z)^T \lambda^+ = 0, \quad (12a)$$

$$g(z) + Dg(z) \Delta z \leq 0, \quad (12b)$$

$$\lambda^+ \geq 0, \quad (12c)$$

$$(\lambda^+)^T [g(z) + Dg(z) \Delta z] = 0. \quad (12d)$$

We can derive the subproblem (11) and the KKT conditions (12) from the following min-max problem involving the Lagrangian of (11):

$$\min_{\Delta z} \max_{\lambda^+ \geq 0} \Delta z^T \nabla \phi(z) + \frac{1}{2} \Delta z^T \nabla_{zz} \mathcal{L}(z, \lambda) \Delta z + (\lambda^+)^T (g(z) + Dg(z) \Delta z).$$

Most practical implementations of SQP perform a line search either along the primal space in the direction Δz or in the primal-dual space in the direction $(\Delta z, \lambda^+ - \lambda)$, with the aim of improving the value of some merit function. When (z, λ) is sufficiently close to the primal-dual solution set \mathcal{S} , these globalization strategies should allow unit steps to be taken to yield rapid convergence; that is,

$$(z, \lambda) \leftarrow (z + \Delta z, \lambda^+).$$

Under the conditions discussed above, the subproblem (11) has a local solution Δz when (z, λ) is sufficiently close to \mathcal{S} . (For a proof of this statement, see Theorem 4.3 and Section 5 of Robinson [12].) However, the Lagrange multiplier λ^+ for the linear constraints in (11) may not be uniquely determined by (11) because of rank deficiency in the Jacobian $\nabla g_B(z)$, so that the Hessian $\nabla_{zz} \mathcal{L}$ may not be uniquely defined at the *next* iteration of SQP. More important, SQP may no longer yield a “quadratic” decrease in distance to the solution set \mathcal{S} , even from points that are arbitrarily close to this set. We illustrate this fact with the following example.

Example. Consider the problem

$$\min z_1 \quad \text{subject to} \quad \begin{aligned} (z_1 - 2)^2 + z_2^2 &\leq 4, \\ (z_1 - 4)^2 + z_2^2 &\leq 16, \end{aligned} \tag{13}$$

which has a unique minimizer at $z^* = 0$ at which both constraints are active and MFCQ is satisfied. The optimal multiplier set is defined by

$$\mathcal{S}_\lambda = \{(1/4 - 2\alpha, \alpha) \mid 0 \leq \alpha \leq 1/8\}.$$

Given a primal-dual point (z, λ) , the quantities needed to define the SQP subproblem (11) are as follows:

$$\begin{aligned} \nabla \phi(z) &= \begin{bmatrix} 1 \\ 0 \end{bmatrix}, & \nabla_{zz} \mathcal{L}(z, \lambda) &= 2(\lambda_1 + \lambda_2) \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \\ g(z) &= \begin{bmatrix} (z_1 - 2)^2 + z_2^2 - 4 \\ (z_1 - 4)^2 + z_2^2 - 16 \end{bmatrix}, & Dg(z) &= 2 \begin{bmatrix} (z_1 - 2) & z_2 \\ (z_1 - 4) & z_2 \end{bmatrix}. \end{aligned}$$

Suppose we apply SQP from the point $z = (\epsilon, \epsilon)$, $\lambda \geq 0$, where ϵ is small and positive. Note that $\|z - z^*\| = \sqrt{2}\epsilon$. It can be shown that the solution Δz of (11) satisfies the following linear system:

$$\begin{bmatrix} 2(\lambda_1 + \lambda_2) & 0 & 2(\epsilon - 4) \\ 0 & 2(\lambda_1 + \lambda_2) & 2\epsilon \\ -2(\epsilon - 4) & -2\epsilon & 0 \end{bmatrix} \begin{bmatrix} \Delta z_1 \\ \Delta z_2 \\ \lambda_2^+ \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 2\epsilon^2 - 8\epsilon \end{bmatrix},$$

where λ_2^+ is the Lagrange multiplier for the second linearized constraint in (11). (The first constraint is inactive.) By solving this system we obtain

$$\Delta z_1 = -\epsilon \frac{8 - \epsilon(2 - 1/4(\lambda_1 + \lambda_2)) + O(\epsilon^2)}{8 - 2\epsilon + O(\epsilon^2)} = -\epsilon + O(\epsilon^2),$$

$$\lambda_2^+ = \frac{-1}{2(\epsilon - 4)} + O(\epsilon) = 1/8 + O(\epsilon),$$

$$\Delta z_2 = -\frac{\epsilon \lambda_2^+}{\lambda_1 + \lambda_2} = -\frac{\epsilon}{8(\lambda_1 + \lambda_2)} + O(\epsilon^2).$$

If $\lambda = (1/4, 0)$ —an optimal multiplier for (13)—we have

$$\Delta z = (-\epsilon, -\epsilon/2) + O(\epsilon^2),$$

and therefore

$$\|(z + \Delta z) - z^*\| = \|(O(\epsilon^2), \epsilon/2 + O(\epsilon^2))\| = \epsilon/2 + O(\epsilon^2),$$

giving just a “linear” decrease in the distance to the primal optimum z^* on this iteration, even for ϵ arbitrarily small. The new primal-dual iterate $(z + \Delta z, \lambda^+)$ is also only linearly closer to \mathcal{S} than is (z, λ) . In fact, a linear decrease is obtained when λ is *any* optimal multiplier for (13), unless it happens to be close to the extreme point $(0, 1/8)$ of \mathcal{S}_λ .

We now describe a stabilized variant of SQP for which a quadratic improvement in the error is guaranteed whenever (z, λ) is sufficiently close to a certain large subset of the primal-dual solution set \mathcal{S} —a subset that encompasses most of the relative interior of \mathcal{S} . Iterates of the stabilized SQP algorithm are obtained by solving the following min-max problem for $(\Delta z, \lambda^+)$:

$$\begin{aligned} \min_{\Delta z} \max_{\lambda^+ \geq 0} & \Delta z^T \nabla \phi(z) + \frac{1}{2} \Delta z^T \nabla_{zz} \mathcal{L}(z, \lambda) \Delta z + (\lambda^+)^T (g(z) \\ & + Dg(z) \Delta z) - \frac{1}{2} \mu \|\lambda^+ - \lambda\|^2, \end{aligned} \quad (14)$$

where the parameter μ is defined as

$$\mu = \mu(z, \lambda) \stackrel{\text{def}}{=} \|(\nabla_z \mathcal{L}(z, \lambda), g(z)_+, \lambda^T g(z))\|. \quad (15)$$

Note that (14) differs from the standard min-max formulation only in the inclusion of the proximal penalty term $\frac{1}{2} \mu \|\lambda^+ - \lambda\|^2$. The optimality conditions for a candidate solution $(\Delta z, \lambda^+)$ of (14) are likewise similar to (12), namely,

$$\nabla_{zz} \mathcal{L}(z, \lambda) \Delta z + \nabla \phi(z) + Dg(z)^T \lambda^+ = 0, \quad (16a)$$

$$g(z) + Dg(z) \Delta z - \mu(\lambda^+ - \lambda) \leq 0, \quad (16b)$$

$$\lambda^+ \geq 0, \quad (16c)$$

$$(\lambda^+)^T [g(z) + Dg(z) \Delta z - \mu(\lambda^+ - \lambda)] = 0. \quad (16d)$$

We show later that for (z, λ) sufficiently close to some strictly complementary primal-dual solution (z^*, λ^*) , there is a unique solution $(\Delta z, \lambda^+)$ of (16) for which $\|(\Delta z, \lambda^+ - \lambda)\| = O(\mu)$. This solution satisfies the following linear system:

$$\begin{bmatrix} \nabla_{zz} \mathcal{L}(z, \lambda) & \nabla g_B(z) \\ -\nabla g_B(z)^T & \mu I \end{bmatrix} \begin{bmatrix} \Delta z \\ \lambda_B^+ - \lambda_B \end{bmatrix} = \begin{bmatrix} -\nabla \phi(z) - \nabla g_B(z) \lambda_B \\ g_B(z) \\ \lambda_N^+ = 0. \end{bmatrix}, \quad (17)$$

In Section 3, we show that the norm of the stabilized SQP step is small; in fact, $\|(\Delta z, \lambda^+ - \lambda)\|$ approaches zero at the same rate as μ , while other possible solutions $(\Delta z, \lambda^+)$ to (14), if they exist, cannot satisfy this estimate. These observations hold even if the active constraint Jacobian $\nabla g_B(\cdot)$ is rank deficient at or near the solution z^* . We show in Section 4 that a full step along the direction produces a “quadratic” decrease in μ and in the distance to the solution set and that local quadratic convergence follows as a consequence. Our analysis has much in common with the analysis of Ralph and Wright [10, 11], who deal

with an interior-point algorithm rather than an SQP-based algorithm. In the interior-point approach, the step equations to be solved at each iteration are similar to (17), except that μI is replaced by another diagonal matrix that plays a similar stabilization role. Our focus on the SQP algorithm in this paper has practical significance because of the popularity of this method and because full rank of the active constraint Jacobian is frequently violated (or nearly so) on large-scale nonlinear programming problems.

We consider, too, the effects of finite-precision floating-point arithmetic on the step obtained from (17). While the errors in some components of the stabilized SQP steps may grow quite large near the solution, rapid convergence is still attainable. Our analysis of the floating-point case follows directly from the exact analysis because of our use of linear algebra techniques. The relationship between stabilized SQP approach and the interior-point approach could be used to derive similar finite-precision results for primal-dual interior-point algorithms such as those in [10, 11].

Our conclusions extend to the case in which equality constraints are explicitly present, as we outline in Section 5.

2. Assumptions, preliminary results, and notation

Throughout the remainder of the article, we make the following assumption.

Assumption 1. The vector z^* is a local solution of (1) and the functions $\phi(\cdot)$ and $g(\cdot)$ are twice Lipschitz continuously differentiable in an open neighborhood of z^* . The first-order conditions (3) and the second-order condition (7) are satisfied at z^* , and the strict complementarity condition (6) holds for some vector λ^* for which (3) are satisfied.

Assumption 1 and the MFCQ (8) lead to two preliminary results that have appeared in previous work, as noted below.

Lemma 2.1 (Gauvin [6]). *Suppose that Assumption 1 holds. Then \mathcal{S}_λ is bounded if and only if the MFCQ (8) is satisfied.*

A proof of the following result can be found in Bertsekas [1, Proposition 3.3.2], for example.

Lemma 2.2. *If Assumption 1 holds, then z^* is a locally unique solution of (1).*

Given the definition (9) of \mathcal{S} and Lemma 2.2, we define the constant ξ as follows:

$$\xi = \max_{\lambda^* \in \mathcal{S}_\lambda} \min_{i \in \mathcal{B}} \lambda_i^*. \quad (18)$$

Assumption 1 implies that ξ is positive, while Lemma 2.1 implies that it is bounded. For each $\gamma \in (0, 1)$, we define $\mathcal{N}^\gamma(\epsilon)$ by

$$\begin{aligned} \mathcal{N}^\gamma(\epsilon) &= \{(z, \lambda) \mid \|(z, \lambda) - (z^*, \lambda^*)\| \leq \epsilon, \text{ for some} \\ &\quad \lambda^* \in \mathcal{S}_\lambda \text{ with } \lambda_{\mathcal{B}}^* \geq \gamma \xi e, \text{ and } \lambda \geq 0\}. \end{aligned} \quad (19)$$

From Lemma 2.1, the set \mathcal{S} is compact, so by the definition (18), the set $\{(z^*, \lambda^*) \in \mathcal{S} \mid \lambda_{\mathcal{B}}^* \geq \gamma \xi e\}$ is nonempty and compact also for each $\gamma \in (0, 1)$.

In the remainder of the article, we assume that the following collection of assumptions is satisfied.

Standing Assumptions: We assume that Assumption 1 holds and that the MFCQ (8) is satisfied. We assume, too, that γ used to define (19) is a fixed constant in the range $(0, 1)$.

We conclude this section with some items of notation to be used in subsequent sections. We define integers \bar{m} and \check{m} by

$$\bar{m} = |\mathcal{B}|, \quad \check{m} = \text{rank } \nabla g_{\mathcal{B}}(z^*), \quad (20)$$

so that $0 \leq \check{m} \leq \bar{m} \leq m$. Since Lemma 2.1 implies boundedness of \mathcal{S}_{λ} , there is a constant $\bar{\sigma} > 0$ such that

$$\|\nabla_{zz} \mathcal{L}(z^*, \lambda^*)\| \leq \bar{\sigma}, \quad \text{for all } \lambda^* \in \mathcal{S}_{\lambda}. \quad (21)$$

For convenience, we assume that the problem is scaled so that $\bar{\sigma}$ and $\|Dg(z^*)\|$ are not too much different from 1.

The notation $\delta(z, \lambda)$ denotes the Euclidean distance from the point $(z, \lambda) \in \mathbb{R}^n \times \mathbb{R}^m$ to the primal-dual solution set \mathcal{S} defined by (9), that is,

$$\delta(z, \lambda) = \inf\{\|(z^*, \lambda^*) - (z, \lambda)\| \mid (z^*, \lambda^*) \in \mathcal{S}\}.$$

We use $P(\cdot)$ to denote projection onto the set of optimal Lagrange multipliers, that is,

$$P(\lambda) = \arg \min_{\lambda^* \in \mathcal{S}_{\lambda}} \|\lambda^* - \lambda\|.$$

We use order notation in the following way. If a scalar quantity β is a function of another scalar quantity α , we write $\beta = O(\alpha)$ if there is a constant M such that $|\beta| \leq M|\alpha|$ for all α sufficiently small. We write $\beta = \Theta(\alpha)$ if M can be defined so that $M^{-1}|\alpha| \leq |\beta| \leq M|\alpha|$ for all α sufficiently small.

We use $\|\cdot\|$ to denote the Euclidean norm of a matrix or vector, and $\kappa(A) = \|A\| \|A^{-1}\|$ to denote the condition number of a nonsingular matrix with respect to this norm. We use $\ker A$ to denote the kernel (null space) of the matrix A .

Finally, we mention that when functions such as $g_{\mathcal{B}}$, $\nabla_z \mathcal{L}$, and $\nabla_{zz} \mathcal{L}$ appear without specific arguments, the arguments are understood to be the current points z or (z, λ) , as appropriate.

3. Step size estimates

In this section, we show that the step calculated from any point $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$ via (17) in exact arithmetic satisfies the estimate

$$(\Delta z, \Delta \lambda_B) \stackrel{\text{def}}{=} (\Delta z, \lambda_B^+ - \lambda_B) = O(\mu), \quad (22)$$

while any other local solution of (16) cannot satisfy this estimate. We also discuss the effect of finite-precision floating-point arithmetic on this estimate.

Our first result shows that μ defined in (15) is closely related to the distance from the current point to the solution set \mathcal{S} .

Lemma 3.1. *Suppose that the standing assumptions hold. Then there is a constant $\epsilon > 0$ such that for all $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$ we have*

$$\delta(z, \lambda) = \Theta(\mu). \quad (23)$$

Proof: We show first that $\mu = O(\delta(z, \lambda))$. Let $(z^*, P(\lambda))$ be the projection of (z, λ) onto the (compact) solution set \mathcal{S} , so that

$$\|(z, \lambda) - (z^*, P(\lambda))\| = \delta(z, \lambda).$$

By the optimality condition (3) and Assumption 1, we have

$$\|\nabla_z \mathcal{L}(z, \lambda)\| = \|\nabla_z \mathcal{L}(z, \lambda) - \nabla_z \mathcal{L}(z^*, P(\lambda))\| = O(\delta(z, \lambda)). \quad (24)$$

Similarly, we have

$$\|g(z)_+\| = \|g(z)_+ - g(z^*)_+\| \leq \|g(z) - g(z^*)\| = O(\delta(z, \lambda)). \quad (25)$$

By boundedness of \mathcal{S} , we have that $P(\lambda)$ is bounded by a constant independent of λ .

Hence, we can write

$$\begin{aligned} \lambda^T g(z) &= \lambda^T g(z) - P(\lambda)^T g(z^*) \\ &= [\lambda - P(\lambda)]^T [g(z) - g(z^*)] + P(\lambda)^T [g(z) - g(z^*)] + [\lambda - P(\lambda)]^T g(z^*) \\ &= O(\delta(z, \lambda)). \end{aligned} \quad (26)$$

We obtain the result by substituting (24), (25), and (26) into (15).

The proof of reverse estimate— $\delta(z, \lambda) = O(\mu)$ —follows the proof of [11, Lemma 5.5] closely. The condition $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$ is important in deriving this estimate; that is, the λ component should be close to a “sufficiently strictly complementary” point in \mathcal{S}_λ rather than near some extreme point of this set. We omit the details and refer the reader to the earlier paper. \square

To analyze the step $(\Delta z, \Delta \lambda_B)$, we decompose it to conform with the singular value decomposition (svd) of the optimal Jacobian $\nabla g_B(z^*)$. From the definitions (4) and (20), this matrix has dimensions $n \times \bar{m}$ and rank \bar{m} . We write its svd as

$$\nabla g_B(z^*) = [\hat{U} \quad \hat{V}] \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U^T \\ V^T \end{bmatrix}, \quad (27)$$

where S is diagonal with diagonal elements $\sigma_1 \geq \dots \geq \sigma_{\bar{m}} > 0$, U is $\bar{m} \times \bar{m}$, V is $\bar{m} \times (\bar{m} - \bar{m})$, \hat{U} is $n \times \bar{m}$, and \hat{V} is $n \times (n - \bar{m})$, and $[U \quad V]$ and $[\hat{U} \quad \hat{V}]$ are orthogonal.

Note, in particular, that the columns of \hat{V} constitute an orthonormal basis for $\ker \nabla g_B(z^*)^T$.

3.1. Exact arithmetic

We first analyze the step $(\Delta z, \Delta \lambda_B)$ obtained by solving (17), assuming exact arithmetic. We decompose this step as

$$\Delta z = \hat{U} y_{\hat{U}} + \hat{V} y_{\hat{V}}, \quad (28a)$$

$$\Delta \lambda_B = \lambda_B^+ - \lambda_B = U w_U + V w_V. \quad (28b)$$

By substituting (28) into (17) and premultiplying the blocks of this system by the matrices \hat{U}^T , \hat{V}^T , U^T , and V^T , we obtain

$$\begin{bmatrix} \hat{U}^T(\nabla_{zz}\mathcal{L})\hat{U} & \hat{U}^T(\nabla_{zz}\mathcal{L})\hat{V} & \hat{U}^T(\nabla g_B)U & \hat{U}^T(\nabla g_B)V \\ \hat{V}^T(\nabla_{zz}\mathcal{L})\hat{U} & \hat{V}^T(\nabla_{zz}\mathcal{L})\hat{V} & \hat{V}^T(\nabla g_B)U & \hat{V}^T(\nabla g_B)V \\ -U^T(\nabla g_B)^T\hat{U} & -U^T(\nabla g_B)^T\hat{V} & \mu I & 0 \\ -V^T(\nabla g_B)^T\hat{U} & -V^T(\nabla g_B)^T\hat{V} & 0 & \mu I \end{bmatrix} \begin{bmatrix} y_{\hat{U}} \\ y_{\hat{V}} \\ w_U \\ w_V \end{bmatrix} = \begin{bmatrix} r_{\hat{U}} \\ r_{\hat{V}} \\ r_U \\ r_V \end{bmatrix}, \quad (29)$$

where the right-hand side is

$$\begin{bmatrix} r_{\hat{U}} \\ r_{\hat{V}} \\ r_U \\ r_V \end{bmatrix} = \begin{bmatrix} -\hat{U}^T(\nabla \phi + \nabla g_B \lambda_B) \\ -\hat{V}^T(\nabla \phi + \nabla g_B \lambda_B) \\ U^T g_B \\ V^T g_B \end{bmatrix}. \quad (30)$$

Since

$$\nabla g_B(z) - \nabla g_B(z^*) = O(\|z - z^*\|) = O(\mu),$$

by Assumption 1 and Lemma 3.1, we have from (27) that

$$\hat{U}^T \nabla g_B(z) U = S + O(\mu), \quad (31a)$$

$$\hat{U}^T \nabla g_{\mathcal{B}}(z) V = O(\mu), \quad (31b)$$

$$\hat{V}^T \nabla g_{\mathcal{B}}(z) U = O(\mu), \quad (31c)$$

$$\hat{V}^T \nabla g_{\mathcal{B}}(z) V = O(\mu). \quad (31d)$$

Meanwhile, by the second-order condition (7) and orthonormality of \hat{V} , we have that the matrix $\hat{V}^T \nabla_{zz} \mathcal{L}(z^*, \lambda^*) \hat{V}$ satisfies

$$v^T \hat{V}^T \nabla_{zz} \mathcal{L}(z^*, \lambda^*) \hat{V} v \geq \sigma \|v\|^2, \quad (32)$$

for all $v \in \mathbb{R}^{n-\bar{m}}$, all $\lambda^* \in \mathcal{S}_\lambda$, and some $\sigma > 0$ independent of v and λ^* .

By substituting the estimates (31) into (29), we obtain

$$\begin{bmatrix} \hat{U}^T (\nabla_{zz} \mathcal{L}) \hat{U} & \hat{U}^T (\nabla_{zz} \mathcal{L}) \hat{V} & S + O(\mu) & O(\mu) \\ \hat{V}^T (\nabla_{zz} \mathcal{L}) \hat{U} & \hat{V}^T (\nabla_{zz} \mathcal{L}) \hat{V} & O(\mu) & O(\mu) \\ -S + O(\mu) & O(\mu) & \mu I & 0 \\ O(\mu) & O(\mu) & 0 & \mu I \end{bmatrix} \begin{bmatrix} y_{\hat{U}} \\ y_{\hat{V}} \\ w_U \\ w_V \end{bmatrix} = \begin{bmatrix} r_{\hat{U}} \\ r_{\hat{V}} \\ r_U \\ r_V \end{bmatrix}. \quad (33)$$

By eliminating w_V from this system and rearranging the resulting block 3×3 coefficient matrix, we obtain

$$w_V = \mu^{-1} r_V + O(\|y_{\hat{U}}\|) + O(\|y_{\hat{V}}\|) \quad (34)$$

and

$$\begin{aligned} & \begin{bmatrix} S + O(\mu) & \hat{U}^T (\nabla_{zz} \mathcal{L}) \hat{V} + O(\mu) & \hat{U}^T (\nabla_{zz} \mathcal{L}) \hat{U} + O(\mu) \\ O(\mu) & \hat{V}^T (\nabla_{zz} \mathcal{L}) \hat{V} + O(\mu) & \hat{V}^T (\nabla_{zz} \mathcal{L}) \hat{U} + O(\mu) \\ \mu I & O(\mu) & -S + O(\mu) \end{bmatrix} \begin{bmatrix} w_U \\ y_{\hat{V}} \\ y_{\hat{U}} \end{bmatrix} \\ &= \begin{bmatrix} r_{\hat{U}} + O(\|r_V\|) \\ r_{\hat{V}} + O(\|r_V\|) \\ r_U \end{bmatrix}. \end{aligned} \quad (35)$$

(Note that elimination of w_V has introduced $O(\mu)$ perturbations into the $\nabla_{zz} \mathcal{L}$ blocks.) The coefficient matrix in (35) is an $O(\mu)$ perturbation of the block upper triangular matrix $M(z, \lambda)$ defined by

$$M(z, \lambda) = \begin{bmatrix} S & \hat{U}^T (\nabla_{zz} \mathcal{L}) \hat{V} & \hat{U}^T (\nabla_{zz} \mathcal{L}) \hat{U} \\ 0 & \hat{V}^T (\nabla_{zz} \mathcal{L}) \hat{V} & \hat{V}^T (\nabla_{zz} \mathcal{L}) \hat{U} \\ 0 & 0 & -S \end{bmatrix}, \quad (36)$$

whose condition number can be bounded independently of (z, λ) for all $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$ for ϵ sufficiently small. To verify this claim, note first that the above-diagonal blocks are

bounded in norm by $\|\nabla_{zz}\mathcal{L}\|$ which, by (21), Assumption 1, and Lemma 3.1, satisfies the bound

$$\begin{aligned} \|\nabla_{zz}\mathcal{L}(z, \lambda)\| &\leq \|\nabla_{zz}\mathcal{L}(z^*, P(\lambda))\| + O(\|(z, \lambda) - (z^*, P(\lambda))\|) \\ &\leq \bar{\sigma} + O(\mu) \\ &\leq 2\bar{\sigma}, \end{aligned} \quad (37)$$

for a sufficiently small choice of ϵ . By decreasing ϵ further if necessary, we have by applying (32), Assumption 1, and Lemma 3.1 that

$$v^T \hat{V}^T \nabla_{zz}\mathcal{L}(z, \lambda) \hat{V} v = v^T \hat{V}^T \nabla_{zz}\mathcal{L}(z^*, P(\lambda)) \hat{V} v + O(\mu) \|v\|^2 \geq (\sigma/2) \|v\|^2,$$

for all $v \in \mathbb{R}^{n \times (n-\check{m})}$. Hence, using (37) again, we obtain

$$\kappa(\hat{V}^T \nabla_{zz}\mathcal{L} \hat{V}) = \|\hat{V}^T \nabla_{zz}\mathcal{L} \hat{V}\| \|[\hat{V}^T \nabla_{zz}\mathcal{L} \hat{V}]^{-1}\| \leq 4(\bar{\sigma}/\sigma). \quad (38)$$

The other diagonal blocks in (36) are also well conditioned, since by the definition of S we have

$$\kappa(S) = \|S\| \|S^{-1}\| = \sigma_1/\sigma_{\check{m}}. \quad (39)$$

Since

$$M(z, \lambda)^{-1} = \begin{bmatrix} S^{-1} & M_{12} & M_{13} \\ 0 & [\hat{V}^T \nabla_{zz}\mathcal{L} \hat{V}]^{-1} & M_{23} \\ 0 & 0 & -S^{-1} \end{bmatrix},$$

where

$$\begin{aligned} M_{12} &= -S^{-1}(\hat{U}^T \nabla_{zz}\mathcal{L} \hat{V})(\hat{V}^T \nabla_{zz}\mathcal{L} \hat{V})^{-1}, \\ M_{23} &= (\hat{V}^T \nabla_{zz}\mathcal{L} \hat{V})^{-1}(\hat{V}^T \nabla_{zz}\mathcal{L} \hat{U})S^{-1}, \\ M_{13} &= S^{-1}[-(\hat{U}^T \nabla_{zz}\mathcal{L} \hat{V})M_{23} + (\hat{U}^T \nabla_{zz}\mathcal{L} \hat{U})S^{-1}], \end{aligned}$$

it is easy to see from (37), (38), and (39) that, for a sufficiently small choice of ϵ , the quantities $\|M(z, \lambda)^{-1}\|$ and $\kappa(M(z, \lambda))$ are bounded for all $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$. That is, we can define a constant C_M such that

$$\|M(z, \lambda)^{-1}\| \leq C_M, \quad \kappa(M(z, \lambda)) \leq C_M, \quad \text{for all } (z, \lambda) \in \mathcal{N}^\gamma(\epsilon).$$

Using (36), we rewrite (35) as

$$[M(z, \lambda) + O(\mu)] \begin{bmatrix} w_U \\ y_{\hat{V}} \\ y_{\hat{U}} \end{bmatrix} = \begin{bmatrix} r_{\hat{U}} + O(\|r_V\|) \\ r_{\hat{U}} + O(\|r_V\|) \\ r_U \end{bmatrix}. \quad (40)$$

By decreasing ϵ if necessary, we can ensure nonsingularity of the coefficient matrix and in fact that

$$\|[M(z, \lambda) + O(\mu)]^{-1}\| \leq 2\|M(z, \lambda)^{-1}\| \leq 2C_M.$$

Hence, we have immediately from (40) that

$$\|(w_U, y_{\hat{V}}, y_{\hat{U}})\| = O(\|(r_{\hat{U}}, r_{\hat{V}}, r_U, r_V)\|). \quad (41)$$

It follows by substitution in (34) that

$$\|w_V\| = O(\mu^{-1})\|r_V\| + O(\|(r_{\hat{U}}, r_{\hat{V}}, r_U)\|). \quad (42)$$

We obtain the estimate (22) for $(\Delta z, \Delta \lambda_B)$ by using the definition of the right-hand side in (30). We have

$$\|(r_{\hat{U}}, r_{\hat{V}})\| = \|\nabla\phi(z) + \nabla g_B(z)\lambda_B\| \leq \|\nabla_z \mathcal{L}(z, \lambda)\| + \|\nabla g_N(z)\| \|\lambda_N\|.$$

From the definition (15), we have directly that $\|\nabla_z \mathcal{L}(z, \lambda)\| \leq \mu$. By compactness of $\mathcal{N}^\gamma(\epsilon)$ and smoothness of $\nabla g(\cdot)$, we have that $\|\nabla g_N(z)\|$ is bounded above, while for $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$ we have from Lemma 3.1 that $\|\lambda_N\| \leq \delta(z, \lambda) = O(\mu)$. We conclude that

$$\|(r_{\hat{U}}, r_{\hat{V}})\| = O(\mu). \quad (43)$$

Since $g_B(z^*) = 0$, we have from Assumption 1 and Lemma 3.1 that

$$\|r_U\| = \|U^T g_B(z)\| \leq \|g_B(z) - g_B(z^*)\| = O(\|z - z^*\|) = O(\mu). \quad (44)$$

The remaining estimate is slightly more refined. By Assumption 1 and Lemma 3.1, we have

$$\begin{aligned} r_V &= V^T g_B(z) \\ &= V^T [g_B(z^*) + \nabla g_B(z^*)^T (z - z^*)] + O(\|z - z^*\|^2) \\ &= O(\mu^2), \end{aligned} \quad (45)$$

since $g_B(z^*) = 0$ and $V^T \nabla g_B(z^*)^T = 0$ by (27). Substitution of (43), (44), and (45) into (41) and (42) yields

$$\|(y_{\hat{U}}, y_{\hat{V}}, w_U, w_V)\| = O(\mu).$$

The estimate (22) follows immediately from (28).

We summarize this result as a theorem.

Theorem 3.2. Suppose that the standing assumptions hold. Then there is a constant $\epsilon > 0$ such that for all $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$, the solution $(\Delta z, \lambda^+)$ of (17) satisfies the estimate

$$\|(\Delta z, \lambda^+ - \lambda)\| \leq C' \mu, \quad (46)$$

for some $C' > 0$ that depends on γ and ϵ but not on μ or (z, λ) .

Proof: We have shown above that (22) holds. It remains only to examine the \mathcal{N} -components of λ^+ and λ . Since $\lambda_{\mathcal{N}}^* = 0$ for all $\lambda^* \in \mathcal{S}_\lambda$, we have

$$\|\lambda_{\mathcal{N}} - \lambda_{\mathcal{N}}^+\| = \|\lambda_{\mathcal{N}}\| \leq \delta(z, \lambda) = O(\mu). \quad \square$$

It is not difficult to show that the solution $(\Delta z, \lambda^+)$ of (17) is the only local solution of (16) satisfying the estimate (46). Let $(\widehat{\Delta z}, \widehat{\lambda}^+)$ be any solution of (16) that satisfies (46). If there is an index i for which $i \notin \mathcal{N}$ but $\widehat{\lambda}_i^+ = 0$, we have by the definition (19) and $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$ that $\lambda_i \geq \gamma\xi - \epsilon$. Hence, by choosing ϵ smaller than $\gamma\xi/2$, we have that

$$\|\lambda - \widehat{\lambda}^+\| \geq |\lambda_i - \widehat{\lambda}_i| \geq \gamma\xi/2,$$

for this particular index i . Hence, each $i \notin \mathcal{N}$ must have $\widehat{\lambda}_i^+ > 0$.

If, on the other hand, there is an index $i \in \mathcal{N}$ for which $\widehat{\lambda}_i^+ > 0$, by complementarity and Eq. (12b) we have

$$g_i(z) + Dg_i(z)\widehat{\Delta z} - \mu(\widehat{\lambda}_i^+ - \lambda_i) = 0.$$

If the estimate $\|(\widehat{\Delta z}, \widehat{\lambda}_i^+ - \lambda_i)\| = O(\mu)$ holds (as assumed), this equation yields $g_i(z) = O(\mu)$. However, Lemma 3.1 and $g_i(z^*) < 0$ imply that $g_i(z)$ is bounded away from zero for μ sufficiently small, giving a contradiction.

We conclude from these two cases that $i \in \mathcal{N}$ if and only if $\widehat{\lambda}_i^+ = 0$, so that the index partition for the solution $(\widehat{\Delta z}, \widehat{\lambda}^+)$ is simply $\mathcal{B} \cup \mathcal{N}$. Since the coefficient matrix in (17) is nonsingular, we must have $(\widehat{\Delta z}, \widehat{\lambda}^+) = (\Delta z, \lambda^+)$.

3.2. Finite-precision arithmetic

We now examine the effect of finite-precision floating-point arithmetic on the step calculated from (17). In our discussion of floating-point arithmetic, we use \mathbf{u} to denote unit roundoff, which we define by the following statement: When x and y are any two floating-point numbers, op denotes $+, -, \times$, or $/$, and $f(z)$ denotes the floating-point approximation of any real number z , we have

$$f(x \text{ op } y) = (x \text{ op } y)(1 + \zeta), \quad |\zeta| \leq \mathbf{u}.$$

We also introduce the unfamiliar notation $\delta_{\mathbf{u}}$ to represent a scalar quantity that is a modest multiple of \mathbf{u} .

In solving (17), errors due to floating-point arithmetic arise from two sources:

- (a) Error incurred during the evaluation of the components of the matrix and the right-hand side, and
- (b) Error incurred during factorization of the matrix in (17) and the subsequent triangular substitutions.

We consider (a) first. The quantities on the right-hand side— $\nabla\phi(z) + \nabla g_B(z)\lambda_B$ and $g_B(z)$ —are $O(\mu)$ in exact arithmetic. However, they are typically evaluated by adding and subtracting quantities whose size is independent of μ , and hence they contain evaluation errors of size δ_u . Of all floating-point errors, these are the most significant in their effect on the accuracy of the step and on the algorithm's convergence behavior. We assume that

$$\mu \gg u, \quad (47)$$

since, if not, the perturbed right-hand side may bear no relation to the exact version, so we could not expect any similarity between the perturbed solution and its exact counterpart.

Evaluation errors may also appear in the blocks of the coefficient matrix in (17), except for the $(2, 2)$ block μI . Since we have assumed for convenience that $\|\nabla_{zz}\mathcal{L}(z^*, \lambda^*)\|$ and $\nabla g_B(z^*)$ are not too different from 1, these errors can be accounted for by introducing perturbations of size δ_u into the $(1, 1)$, $(1, 2)$, and $(2, 1)$ blocks of the matrix in (17).

In assessing the errors that arise from cause (b), we assume that the matrix is factored by a stable procedure, one in which the elements of the submatrices that arise during the factorization are not too large relative to the norm of the original matrix. Suitable algorithms could include Gaussian elimination with pivoting or a Bunch-Parlett or Bunch-Kaufman algorithm applied to a symmetric indefinite reformulation of the problem. (See, for example, Fourer and Mehrotra [5] and Wright [13] for a discussion of algorithms for symmetric indefinite matrices and their stability properties.) A standard backward error analysis applied to the general square system $Mx = r$ shows that the approximate solution \hat{x} computed by stable factorization and triangular substitutions satisfies

$$(M + E_M)\hat{x} = r, \quad \text{where } \|E_M\| \leq \delta_u \|M\|$$

(see for example Golub and Van Loan [7, Chapter 3]). The effects of the errors of type (b) can be accounted for by introducing perturbations of size δ_u into all the elements of the matrix in (17).

Collating these errors, we find that the computed approximation $(\tilde{\Delta}z, \tilde{\Delta}\lambda_B)$ to the step $(\Delta z, \Delta \lambda_B)$ satisfies the equation

$$\left\{ \begin{bmatrix} \nabla_{zz}\mathcal{L}(z, \lambda) & \nabla g_B(z) \\ -\nabla g_B(z)^T & \mu I \end{bmatrix} + \bar{E} \right\} \begin{bmatrix} \tilde{\Delta}z \\ \tilde{\Delta}\lambda_B \end{bmatrix} = \begin{bmatrix} -\nabla_z\mathcal{L}(z, \lambda) \\ g_B(z) \end{bmatrix} + \bar{e}, \quad (48)$$

where

$$\|\bar{E}\| \leq \delta_u, \quad \|\bar{e}\| \leq \delta_u. \quad (49)$$

By decomposing the approximate step as

$$\tilde{z} = \hat{U}\tilde{y}_{\hat{U}} + \hat{V}\tilde{y}_{\hat{V}}, \quad \Delta\tilde{\lambda}_B = U\tilde{w}_U + V\tilde{w}_V, \quad (50)$$

and partitioning the system as in (33), we obtain

$$\begin{aligned} & \left\{ \begin{bmatrix} \hat{U}^T(\nabla_{zz}\mathcal{L})\hat{U} & \hat{U}^T(\nabla_{zz}\mathcal{L})\hat{V} & S + O(\mu) & O(\mu) \\ \hat{V}^T(\nabla_{zz}\mathcal{L})\hat{U} & \hat{V}^T(\nabla_{zz}\mathcal{L})\hat{V} & O(\mu) & O(\mu) \\ -S + O(\mu) & O(\mu) & \mu I & 0 \\ O(\mu) & O(\mu) & 0 & \mu I \end{bmatrix} + E \right\} \begin{bmatrix} \tilde{y}_{\hat{U}} \\ \tilde{y}_{\hat{V}} \\ \tilde{w}_U \\ \tilde{w}_V \end{bmatrix} \\ &= \begin{bmatrix} r_{\hat{U}} + e_{\hat{U}} \\ r_{\hat{V}} + e_{\hat{V}} \\ r_U + e_U \\ r_V + e_V \end{bmatrix}, \end{aligned} \quad (51)$$

where the norms of E , $e_{\hat{U}}$, $e_{\hat{V}}$, e_U , and e_V all have size δ_u . From the last block row we obtain

$$\tilde{w}_V = [\mu I + E_{VV}]^{-1}[r_V + e_V + (O(\mu) + \delta_u)(\|\tilde{y}_{\hat{U}}\| + \|\tilde{y}_{\hat{V}}\|) + \delta_u\|\tilde{w}_U\|], \quad (52)$$

where E_{VV} is the lower right block of E . We use the condition (47) and the estimate $\|E\| = \delta_u$ to deduce that

$$[\mu I + E_{VV}]^{-1} = [I + \mu^{-1}E_{VV}]^{-1}\mu^{-1} = O(\mu^{-1}). \quad (53)$$

Hence, from (52) we have

$$\|\tilde{w}_V\| = O(\mu^{-1})\|r_V + e_V\| + O(1)[\|\tilde{y}_{\hat{U}}\| + \|\tilde{y}_{\hat{V}}\|] + O(\mu^{-1}\delta_u)\|\tilde{w}_U\|. \quad (54)$$

By eliminating \tilde{w}_V from (51) and using the estimate (53), we obtain

$$\begin{aligned} & \left\{ \begin{bmatrix} S & \hat{U}^T(\nabla_{zz}\mathcal{L})\hat{V} & \hat{U}^T(\nabla_{zz}\mathcal{L})\hat{U} \\ 0 & \hat{V}^T(\nabla_{zz}\mathcal{L})\hat{V} & \hat{V}^T(\nabla_{zz}\mathcal{L})\hat{U} \\ 0 & 0 & -S \end{bmatrix} + \tilde{E} \right\} \begin{bmatrix} \tilde{w}_U \\ \tilde{y}_{\hat{V}} \\ \tilde{y}_{\hat{U}} \end{bmatrix} \\ &= \begin{bmatrix} r_{\hat{U}} + e_{\hat{U}} \\ r_{\hat{V}} + e_{\hat{V}} \\ r_U + e_U \end{bmatrix} + \begin{bmatrix} O(\|r_V + e_V\|) \\ O(\|r_V + e_V\|) \\ O(\mu^{-1}\delta_u\|r_V + e_V\|) \end{bmatrix}, \end{aligned} \quad (55)$$

where

$$\|\tilde{E}\| \leq O(\mu) + \delta_u.$$

As before, we choose the neighborhood radius ϵ small enough to ensure that the error-free part of the coefficient matrix in (55) dominates any $O(\mu)$ perturbations. We assume, too, that \mathbf{u} is small enough that perturbations of size $\delta_{\mathbf{u}}$ are also dominated by the error-free part. By applying the logic that follows equation (35), we obtain as in (41) that $\|(\tilde{w}_U, \tilde{y}_{\hat{V}}, \tilde{y}_{\hat{U}})\|$ is of the order of the right-hand side in (55), that is,

$$\begin{aligned}\|(\tilde{w}_U, \tilde{y}_{\hat{V}}, \tilde{y}_{\hat{U}})\| &= O(\|(r_{\hat{U}} + e_{\hat{U}}, r_{\hat{V}} + e_{\hat{V}}, r_U + e_U, r_V + e_V)\|) \\ &= O(\mu) + \delta_{\mathbf{u}}.\end{aligned}\quad (56)$$

By using some simple manipulation involving (35) and (55), we can show further that

$$\|(\tilde{w}_U, \tilde{y}_{\hat{V}}, \tilde{y}_{\hat{U}}) - (w_U, y_{\hat{V}}, y_{\hat{U}})\| = O(\mu^2) + \delta_{\mathbf{u}}, \quad (57)$$

so that the relative accuracy of the computed step components $(\tilde{w}_U, \tilde{y}_{\hat{V}}, \tilde{y}_{\hat{U}})$ remains high. (Note that these estimates hold even when μ is similar in size to \mathbf{u} , provided that the \tilde{E} term is small enough in the sense described above.)

Returning to the component \tilde{w}_V , we have from (54), (45), (49), and (56) that

$$\|\tilde{w}_V\| = O(\mu + \mu^{-1} \delta_{\mathbf{u}}). \quad (58)$$

Note that if we allow $\mu \approx \mathbf{u}$, the matrix in (55) may be singular and the norm of \tilde{w}_V can be arbitrarily large.

From the estimates (56) and (58), we conclude that finite-precision arithmetic has little effect on the $(w_U, y_{\hat{V}}, y_{\hat{U}})$ step components, while it has a potentially significant effect on the w_V component. In fact, for $\mu < \sqrt{\mathbf{u}}$, the estimate (58) indicates that the error in \tilde{w}_V can dominate the “exact” contribution. Fortunately, as we see below, the potentially large error in this component has little effect on the local convergence properties of the algorithm.

By substituting the estimates (56) and (58) into (50), we obtain

$$\|\widetilde{\Delta z}\| = O(\mu) + \delta_{\mathbf{u}}, \quad \|\widetilde{\Delta \lambda}_{\mathcal{B}}\| = O(\mu + \mu^{-1} \delta_{\mathbf{u}}). \quad (59)$$

4. Local convergence

We now examine the effect of the exact and inexact steps on the decrease in the optimality measure μ . We show that the exact step yields a “quadratic” decrease in μ , indicating a quadratic rate of convergence of the iterates (z, λ) to the primal-dual solution set \mathcal{S} . In finite precision, this convergence behavior is affected less severely than one might expect, but we show that reduction of μ below the level of \mathbf{u} cannot be achieved in general and is, in any case, undesirable.

Since we are interested in the asymptotic behavior, we assume that μ (and hence $\delta(z, \lambda)$) is small enough that the new iterate is obtained from (17), even when floating-point arithmetic is used.

First, we analyze the case of exact steps. Using the definition (2) of the Lagrangian \mathcal{L} , Taylor's theorem, Assumption 1, and $\lambda_{\mathcal{N}}^+ = 0$, we find that

$$\begin{aligned}
 \nabla_z \mathcal{L}(z + \Delta z, \lambda^+) &= \nabla \phi(z + \Delta z) + \nabla g(z + \Delta z)^T \lambda^+ \\
 &= \nabla \phi(z) + \nabla^2 \phi(z) \Delta z + O(\|\Delta z\|^2) + \nabla g(z) \lambda^+ \\
 &\quad + \sum_{i=1}^m \lambda_i^+ \nabla^2 g_i(z) \Delta z + O(\|\lambda + \Delta \lambda\| \|\Delta z\|^2) \\
 &= \nabla_{zz} \mathcal{L}(z, \lambda) \Delta z + \nabla g_{\mathcal{B}}(z) \lambda_{\mathcal{B}}^+ + \nabla \phi(z) \\
 &\quad + \sum_{i=1}^m \Delta \lambda_i \nabla^2 g_i(z) \Delta z + O(\|\Delta z\|^2).
 \end{aligned} \tag{60}$$

From the first block row of (17), the first three terms in this expression sum to zero. Moreover, since $\|\lambda_{\mathcal{N}}\| \leq \delta(z, \lambda)$, we have from (60) and (22) that

$$\nabla_z \mathcal{L}(z + \Delta z, \lambda^+) = O((\|\Delta \lambda_{\mathcal{B}}\| + \|\Delta \lambda_{\mathcal{N}}\|) \|\Delta z\|) + O(\|\Delta z\|^2) = O(\mu^2). \tag{61}$$

For the \mathcal{B} components of $g(\cdot)$, we have from the second block row in (17) and (22) that

$$\begin{aligned}
 g_{\mathcal{B}}(z + \Delta z) &= g_{\mathcal{B}}(z) + \nabla g_{\mathcal{B}}(z) \Delta z + O(\|\Delta z\|^2) \\
 &= \mu \Delta \lambda_{\mathcal{B}} + O(\|\Delta z\|^2) \\
 &= O(\mu^2).
 \end{aligned} \tag{62}$$

For the \mathcal{N} components, since $g_{\mathcal{N}}(z)$ is negative and bounded away from zero for $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$, we have

$$g_{\mathcal{N}}(z + \Delta z) < 0,$$

provided that ϵ is chosen to be sufficiently small. By combining the last two expressions we obtain

$$\|g(z + \Delta z)_+\| = \|g_{\mathcal{B}}(z + \Delta z)_+\| \leq \|g_{\mathcal{B}}(z + \Delta z)\| = O(\mu^2). \tag{63}$$

For the third component in the definition (15) of μ , we obtain from (17), (22), $\lambda_{\mathcal{N}}^+ = 0$, and boundedness of \mathcal{S}_λ that

$$\begin{aligned}
 (\lambda^+)^T g(z + \Delta z) &= (\lambda_{\mathcal{B}} + \Delta \lambda_{\mathcal{B}})^T g_{\mathcal{B}}(z + \Delta z) \\
 &= (\lambda_{\mathcal{B}}^+)^T (g_{\mathcal{B}}(z) + \nabla g_{\mathcal{B}}(z) \Delta z + O(\|\Delta z\|^2)) \\
 &= (\lambda_{\mathcal{B}}^+)^T (\mu \Delta \lambda_{\mathcal{B}} + O(\|\Delta z\|^2)) \\
 &= O(\mu \|\Delta \lambda_{\mathcal{B}}\|) + O(\|\Delta z\|^2) \\
 &= O(\mu^2).
 \end{aligned} \tag{64}$$

By combining the estimates above, we have the following result concerning the decrease in μ produced by a unit step.

Theorem 4.1. *Suppose that the standing assumptions hold. Then there is a constant $\epsilon > 0$ such that for all $(z, \lambda) \in \mathcal{N}^\gamma(\epsilon)$, the solution $(\Delta z, \lambda^+)$ of (17) yields*

$$\mu_+ \stackrel{\text{def}}{=} \mu(z + \Delta z, \lambda^+) \leq C\mu^2, \quad (65)$$

for some $C > 0$ that depends on γ and ϵ but not on μ or (z, λ) .

Proof: We obtain (65) by substituting (61), (63), and (64) into the definition (15). \square

The local convergence result follows as a simple corollary of Theorems 3.2 and 4.1.

Corollary 4.2. *Suppose that the standing assumptions hold, and let ϵ be small enough that Theorems 3.2 and 4.1 both apply. Then if (z^0, λ^0) is a point that satisfies*

$$(z^0, \lambda^0) \in \mathcal{N}^\gamma(\epsilon/2), \quad C'\mu_0 \leq \epsilon/4, \quad C\mu_0 \leq 1/2,$$

the stabilized SQP method with unit steps converges Q-quadratically to a point $(z^, \lambda^*) \in \mathcal{S}$.*

Proof: The first step of the stabilized SQP satisfies

$$\|(\Delta z^0, \lambda^1 - \lambda^0)\| \leq C'\mu_0 \leq \epsilon/4,$$

so from the definition of $\mathcal{N}^\gamma(\epsilon)$ (19), we have

$$(z^1, \lambda^1) \in \mathcal{N}^\gamma(\epsilon/2 + \epsilon/4) = \mathcal{N}^\gamma(3\epsilon/4). \quad (66)$$

Theorem 4.1 applies, and therefore $\mu_1 \leq C\mu_0^2 \leq \mu_0/2$ from (65). Because of (66), (z^1, λ^1) satisfies the hypotheses of Theorem 3.2, so we have

$$\|(\Delta z^1, \lambda^2 - \lambda^1)\| \leq C'\mu_1 \leq C'\mu_0/2 \leq \epsilon/8,$$

and therefore

$$(z^2, \lambda^2) \in \mathcal{N}^\gamma(7\epsilon/8) \quad \text{and} \quad \mu_2 \leq C\mu_1^2 \leq \mu_1/2.$$

By repeating this argument, we find that $(z^k, \lambda^k) \in \mathcal{N}^\gamma(\epsilon)$ for all k . Moreover, $\{\mu_k\}$ converges Q-quadratically to zero, and therefore, by Lemma 3.1, $\{\delta(z^k, \lambda^k)\}$ also converges Q-quadratically to zero. For any indices k and l with $l > k$, we have

$$\|(z^k, \lambda^k) - (z^l, \lambda^l)\| \leq \sum_{j=k}^{l-1} \|(\Delta z^j, \lambda^{j+1} - \lambda^j)\| \leq C' \sum_{j=k}^{l-1} \mu_j \leq 2C'\mu_k \rightarrow 0, \quad (67)$$

as k and l approach ∞ , so the sequence $\{(z^k, \lambda^k)\}$ is Cauchy. Hence, this sequence has a single limit point $(z^*, \lambda^*) \in \mathcal{S}$. We prove Q-quadratic convergence to this point by using Lemma 3.1 and an argument like (67). We have

$$\begin{aligned} & \| (z^{k+1}, \lambda^{k+1}) - (z^*, \lambda^*) \| \\ & \leq 2C' \mu_{k+1} \leq 2C' C \mu_k^2 = O(\delta(z^k, \lambda^k)^2) = O(\|(z^k, \lambda^k) - (z^*, \lambda^*)\|^2). \end{aligned} \quad \square$$

We turn now to the case of inexact arithmetic and continue to assume that μ and \mathbf{u} satisfy $\mu \gg \mathbf{u}$ so that the estimates (59) apply. As in (48), we have that

$$\begin{aligned} \nabla_z \mathcal{L}(z + \widetilde{\Delta z}, \lambda + \widetilde{\Delta \lambda}) &= \nabla_z \mathcal{L}(z, \lambda) + \nabla g_B(z) \widetilde{\Delta \lambda}_B + \nabla_{zz} \mathcal{L}(z, \lambda) \widetilde{\Delta z} \\ &\quad + O(\|\widetilde{\Delta \lambda}_B\| \|\widetilde{\Delta z}\|) + O(\|\widetilde{\Delta z}\|^2). \end{aligned} \quad (68)$$

From (48), (49), and (59), and the estimate $\mu \gg \mathbf{u}$, we obtain

$$\begin{aligned} \|\nabla_z \mathcal{L}(z, \lambda) + \nabla g_B(z) \Delta \widetilde{\lambda}_B + \nabla_{zz} \mathcal{L}(z, \lambda) \widetilde{\Delta z}\| &\leq \delta_{\mathbf{u}} \|\Delta \widetilde{\lambda}_B\| + \delta_{\mathbf{u}} \|\widetilde{\Delta z}\| + \delta_{\mathbf{u}} \\ &= \delta_{\mathbf{u}}^2 + O(\mu) \delta_{\mathbf{u}} + O(\mu^{-1}) \delta_{\mathbf{u}}^2 + \delta_{\mathbf{u}} \\ &= \delta_{\mathbf{u}}. \end{aligned}$$

Hence, by substituting into (68) and using (59) again, we obtain

$$\|\nabla_z \mathcal{L}(z + \widetilde{\Delta z}, \lambda + \widetilde{\Delta \lambda})\| = \delta_{\mathbf{u}} + O(\mu^{-1}) \delta_{\mathbf{u}}^2 + O(\mu^2) = \delta_{\mathbf{u}} + O(\mu^2). \quad (69)$$

For $g_B(\cdot)$, we have as in (62), by using (48), (49), and (59), that

$$\begin{aligned} \|g_B(z + \widetilde{\Delta z})\| &= \|g_B(z) + \nabla g_B(z)^T \widetilde{\Delta z}\| + O(\|\widetilde{\Delta z}\|^2) \\ &\leq \mu \|\Delta \widetilde{\lambda}_B\| + \delta_{\mathbf{u}} \|\widetilde{\Delta z}\| + \delta_{\mathbf{u}} \|\Delta \widetilde{\lambda}_B\| + O(\|\widetilde{\Delta z}\|^2) \\ &= O(\mu^2) + \delta_{\mathbf{u}}. \end{aligned}$$

Since $\|\widetilde{\Delta z}\|$ is small, we have by the usual argument that $g_N(z + \widetilde{\Delta z}) < 0$. Hence, as in (63), we obtain

$$\|g(z + \widetilde{\Delta z})_+\| = O(\mu^2) + \delta_{\mathbf{u}}. \quad (70)$$

For the third component of μ , we have as in (64), using (48), (49), and (59), that

$$\begin{aligned} (\lambda + \widetilde{\Delta \lambda})^T g(z + \widetilde{\Delta z}) &= (\lambda_B + \widetilde{\Delta \lambda}_B)^T [\mu \Delta \widetilde{\lambda}_B + \delta_{\mathbf{u}} \|\Delta \widetilde{\lambda}_B\| + \delta_{\mathbf{u}} \|\Delta \widetilde{\lambda}_B\| + \delta_{\mathbf{u}}] \\ &= O(\mu^2) + \delta_{\mathbf{u}}. \end{aligned} \quad (71)$$

By combining (69), (70), and (71), we obtain from the definition (15) that

$$\mu(z + \widetilde{\Delta z}, \lambda + \widetilde{\Delta \lambda}) = O(\mu^2) + \delta_{\mathbf{u}}.$$

This expression suggests that the effect of finite-precision arithmetic on the convergence of the algorithm is not really evident until μ reaches the level of $\sqrt{\mathbf{u}}$. Below this threshold, the algorithm will continue to run and to converge rapidly until μ reaches the level of \mathbf{u} . (Typically, just one or two iterations suffice to reduce μ from $\sqrt{\mathbf{u}}$ to \mathbf{u} .) However, we cannot in general reduce μ below the level of \mathbf{u} . We would not want to do so in any case because, by our assumptions on evaluation error, even an exact solution (z^*, λ^*) may yield a μ value of size $\delta_{\mathbf{u}}$. Our pleasing (and slightly surprising) conclusion is that rapid local convergence to a nearly exact solution occurs even though the stabilized SQP step contains large errors.

5. Equality constraints

The algorithm can be modified easily to handle the case in which equality constraints are present explicitly in the formulation; that is, we have

$$\min \phi(z) \quad \text{subject to } g(z) \leq 0, h(z) = 0,$$

where $h : \mathbb{R}^n \rightarrow \mathbb{R}^p$ is also smooth. The Lagrangian \mathcal{L} is redefined as

$$\mathcal{L}(z, \lambda, \eta) = \phi(z) + \lambda^T g(z) + \eta^T h(z),$$

where $\eta \in \mathbb{R}^p$ is the vector of Lagrange multipliers for the equality constraints. The extension of the KKT conditions to this case is well known, while the second-order condition now requires positive definiteness of the Hessian $\nabla_{zz}\mathcal{L}(z^*, \lambda^*, \eta^*)$ on the subspace

$$\ker \begin{bmatrix} \nabla g_B(z^*)^T \\ \nabla h(z^*)^T \end{bmatrix}.$$

The appropriate extension of the Mangasarian-Fromovitz condition (8) is that $Dh(z^*)$ has full row rank and that

$$\nabla h(z^*)^T d = 0, \quad \nabla g_B(z^*)^T d < 0, \quad \text{for some } d \in \mathbb{R}^n.$$

To extend the algorithm, we first redefine μ as

$$\mu = \mu(z, \lambda) \stackrel{\text{def}}{=} \|(\nabla_z \mathcal{L}(z, \lambda), g(z)_+, h(z), \lambda^T g(z))\|. \quad (72)$$

The min-max subproblem (14) becomes

$$\begin{aligned} \min_{\Delta z} \max_{\lambda^+ \geq 0, \eta^+} & \Delta z^T \nabla \phi(z) + \frac{1}{2} \Delta z^T \nabla_{zz} \mathcal{L}(z, \lambda) \Delta z + (\eta^+)^T (h(z) + \nabla h(z)^T \Delta z) \\ & + (\lambda^+)^T (g(z) + \nabla g(z)^T \Delta z) - \frac{1}{2} \mu \|\lambda^+ - \lambda\|^2. \end{aligned} \quad (73)$$

(No stabilization with respect to the Lagrange multipliers of the equality constraints is needed.) Analogously to (17), we can show that for (z, λ, η) sufficiently close to a strictly

complementary primal-dual solution (z^*, λ^*, η^*) , the solution of (73) satisfies the following system:

$$\begin{aligned} & \begin{bmatrix} \nabla_{zz} \mathcal{L}(z, \lambda, \eta) & \nabla g_{\mathcal{B}}(z) & \nabla h(z)^T \\ -\nabla g_{\mathcal{B}}(z)^T & \mu I & 0 \\ -\nabla h(z)^T & 0 & 0 \end{bmatrix} \begin{bmatrix} \Delta z \\ \lambda_{\mathcal{B}}^+ - \lambda_{\mathcal{B}} \\ \eta^+ - \eta \end{bmatrix} \\ &= \begin{bmatrix} -\nabla \phi(z) - \nabla g_{\mathcal{B}}(z)\lambda_{\mathcal{B}} - \nabla h(z)\eta \\ g_{\mathcal{B}}(z) \\ h(z) \end{bmatrix}, \quad \lambda_{\mathcal{N}}^+ = 0. \end{aligned} \quad (74)$$

The proof of Lemma 3.1 can be extended to show that μ in (72) remains a good measure of the distance to the solution set, for points that are sufficiently strictly complementary. To prove that the solution of (74) satisfies the estimate

$$(\Delta z, \lambda^+ - \lambda, \eta^+ - \eta) = O(\mu), \quad (75)$$

we extend the analysis of Section 3.1 by redefining the svd in (27) as

$$\begin{bmatrix} \nabla g_{\mathcal{B}}(z^*) \\ \nabla h(z^*) \end{bmatrix} = [\hat{U} \quad \hat{V}] \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U^T \\ V^T \end{bmatrix}.$$

Because we assume full row rank of $Dh(z^*)$, we can show that V has the form

$$V = \begin{bmatrix} V_1 \\ 0 \end{bmatrix}, \quad V_1 \in \mathbb{R}^{\bar{m} \times (\bar{m} + p - \bar{m})},$$

that is, the last p rows of V are zero. In rewriting the system (74) analogously to (29), then, we obtain a block μI in the diagonal position corresponding to the w_V component—an important feature in proving the estimate $\|w_V\| = O(\mu)$. Estimates for the other components, and hence for the full step (75), proceed as in Section 3.1. It is also easy to show, as in Section 4, that a unit step produces a quadratic reduction in μ . As before, finite-precision computations have little effect on the convergence behavior.

6. Discussion

Our superlinear convergence result rests on two assumptions that deserve further comment. The first is the MFCQ, which is used to obtain boundedness of the optimal multiplier set (Lemma 2.1) and therefore to show that μ is a valid estimate of $\delta(z, \lambda)$. We can replace MFCQ with an assumption that the algorithm eventually generates an iterate (z^k, λ^k) that lies inside a neighborhood of some bounded subset of \mathcal{S} , where the radius ϵ of the neighborhood depends on the bound on this subset. A step estimate like Theorem 3.2 and a local capture and convergence result like Corollary 4.2 could then be proved. Some kind of

constraint qualification is still required, however, because without one the KKT conditions are not even necessary, as demonstrated by the following well-known example:

$$\min z_1 \quad \text{subject to } z_2 \leq z_1^3, z_1 \geq 0.$$

The second assumption concerns the requirement that the starting guess for λ is sufficiently strictly complementary, that is, not too close to the boundary of the set \mathcal{S}_λ . Without this assumption, it is not possible to prove, at least using techniques like those in Lemma 3.1, that μ is a measure of the distance to optimality $\delta(z, \lambda)$. Nor is it possible to apply the linear-algebra-based analysis of Section 3 without significant modification, since it is no longer true that all indices in \mathcal{B} are active at the solution of the stabilized SQP subproblem (14).

We can accommodate this restriction with the help of an active set identification technique such as that of Facchinei et al. [3]. Our algorithm can periodically make an estimate of the active constraint set and then solve a subproblem to modify the current value of λ to make it “more strictly complementary” without increasing μ by too much. If $\delta(z, \lambda)$ is small enough that \mathcal{B} is identified correctly, a single λ adjustment should suffice. Subsequent iterates are captured within a neighborhood such as $\mathcal{N}^\gamma(\epsilon)$ and the superlinear convergence follows.

Still, it would be more satisfactory to know that the algorithm exhibited the desired behavior without this identification/adjustment step. Hager [8] describes the behavior of a variant of stabilized SQP in which the coefficient μ in the stabilization term is replaced by a value that is bounded below by a sufficiently large multiple of $\delta(z, \lambda)$. His analysis, which rests partly on perturbation results developed by Dontchev and Hager [2], does not require the starting λ to be sufficiently far from the boundary of \mathcal{S}_λ , nor does it require existence of a strictly complementary multiplier. By using a suitably extended second-order sufficient condition (which requires positive definiteness of $\nabla_{zz}\mathcal{L}(z^*, \lambda^*)$ on a larger subspace than that of (7) when λ^* is at a boundary of \mathcal{S}_λ), Hager proves that superlinear convergence of stabilized SQP is still attainable.

Simultaneously with the work described in this paper, Fischer [4] developed an algorithm in which a special procedure for choosing the Lagrange multiplier estimate is inserted between standard SQP iterations. He proves superlinear convergence under a different set of assumptions from ours. Existence of a strictly complementary multiplier is not required, but an alternative condition referred to as “weak complementarity” is needed, along with a constant-rank assumption on the Jacobian $\nabla g_{\mathcal{B}}(z)$ of active constraints in a neighborhood of z^* . We note that this constant-rank assumption is not satisfied by the example of Section 1, while our Assumption 1 is satisfied.

Acknowledgments

I am grateful to Bill Hager for helpful discussions and pointers during the preparation of this work.

References

1. D.P. Bertsekas, Nonlinear Programming, Athena Scientific: Belmont, MA, 1995.
2. A.L. Dontchev and W.W. Hager, "Lipschitzian stability for state constrained nonlinear optimal control," SIAM J. Control Optim., vol. 36, pp. 696–718, 1998.
3. F. Facchinei, A. Fischer, and C. Kanzow, "On the accurate identification of active constraints," DIS Technical Report 22.96, Universita' di Roma "La Sapienza", 1996.
4. A. Fischer, "Modified Wilson method for nonlinear programs with nonunique multipliers," Technical Report MATH-NM-04-1997, Technische Universität Dresden, Institute of Numerical Mathematics, Technische Universität Dresden, D-01062, Dresden, Germany, Feb. 1997.
5. R. Fourer and S. Mehrotra, "Solving symmetric indefinite systems in an interior-point method for linear programming," Mathematical Programming, vol. 62, pp.15–39, 1993.
6. J. Gauvin, "A necessary and sufficient regularity condition to have bounded multipliers in nonconvex programming," Mathematical Programming, vol. 12, pp.136–138, 1977.
7. G.H. Golub and C.F. Van Loan, Matrix Computations, 2nd edition, The Johns Hopkins University Press: Baltimore, 1989.
8. W.W. Hager, "Convergence of Wright's stabilized SQP algorithm," Technical Report, Department of Mathematics, University of Florida, Gainesville, Fl., Jan. 1997.
9. O.L. Mangasarian and S. Fromovitz, "The Fritz-John necessary optimality conditions in the presence of equality and inequality constraints," Journal of Mathematical Analysis and Applications, vol. 17, pp. 37–47, 1967.
10. D. Ralph and S.J. Wright, "Superlinear convergence of an interior-point method despite dependent constraints," Preprint ANL.MCS-P622–1196, Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, Ill., Nov. 1996.
11. D. Ralph and S.J. Wright, "Superlinear convergence of an interior-point method for monotone variational inequalities," in Complementarity and Variational Problems: State of the Art, M.C. Ferris and J. Pang (Eds.), SIAM Publications, 1997, pp. 345–385.
12. S.M. Robinson, "Local epi-continuity and local optimization," Mathematical Programming, vol. 19, pp. 208–222, 1987.
13. S.J. Wright, "Stability of augmented system factorizations in interior-point methods," Preprint MCS-P446-0694, Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, Ill., June 1994. To appear in SIAM Journal of Matrix Analysis and Applications.