

**Choosing the Forcing Terms in an  
Inexact Newton Method**

*Stanley C. Eisenstat*  
*Homer F. Walker*

**CRPC-TR94463**  
**May 1994**

Center for Research on Parallel Computation  
Rice University  
6100 South Main Street  
CRPC - MS 41  
Houston, TX 77005

# CHOOSING THE FORCING TERMS IN AN INEXACT NEWTON METHOD \*

STANLEY C. EISENSTAT<sup>†</sup> AND HOMER F. WALKER<sup>‡</sup>

**Abstract.** An inexact Newton method is a generalization of Newton's method for solving  $F(x) = 0$ ,  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$ , in which, at the  $k$ th iteration, the step  $s_k$  from the current approximate solution  $x_k$  is required to satisfy a condition  $\|F(x_k) + F'(x_k) s_k\| \leq \eta_k \|F(x_k)\|$  for a "forcing term"  $\eta_k \in [0, 1)$ . In typical applications, the choice of the forcing terms is critical to the efficiency of the method and can affect robustness as well. Promising choices of the forcing terms are given, their local convergence properties are analyzed, and their practical performance is shown on a representative set of test problems.

**Key words.** forcing terms, inexact Newton methods, Newton iterative methods, truncated Newton methods, Newton's method, iterative linear algebra methods, GMRES

**AMS(MOS) subject classifications.** 65H10, 65F10

**1. Introduction.** Suppose that  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is continuously differentiable in a neighborhood of  $x_* \in \mathbb{R}^n$  for which  $F(x_*) = 0$  and  $F'(x_*)$  is nonsingular. Suppose further that  $F'$  is Lipschitz continuous at  $x_*$  with constant  $\lambda$ , i.e.,

$$(1.1) \quad \|F'(x) - F'(x_*)\| \leq \lambda \|x - x_*\|$$

for  $x$  near  $x_*$ , where  $\|\cdot\|$  denotes some norm on  $\mathbb{R}^n$  and the induced norm on  $\mathbb{R}^{n \times n}$ .

An *inexact Newton method* (Dembo, Eisenstat, and Steihaug [4]) is an extension of classical Newton's method for approximating  $x_*$  formulated as follows:

**Algorithm IN:** Inexact Newton Method [4]

LET  $x_0$  BE GIVEN.

FOR  $k = 0$  STEP 1 UNTIL "CONVERGENCE" DO:

FIND **some**  $\eta_k \in [0, 1)$  AND  $s_k$  THAT SATISFY

$$(1.2) \quad \|F(x_k) + F'(x_k) s_k\| \leq \eta_k \|F(x_k)\|.$$

SET  $x_{k+1} = x_k + s_k$ .

Note that (1.2) expresses both a certain reduction in the norm of  $F(x_k) + F'(x_k) s_k$ , the local linear model of  $F$ , and a certain accuracy in solving the *Newton equation*  $F'(x_k) s = -F(x_k)$ , the exact solution of which is the *Newton step*. In many applications, notably Newton iterative or truncated Newton methods<sup>1</sup>, each  $\eta_k$  is specified

---

\* Submitted to *SIAM Journal on Scientific Computing*, May, 1994; submitted in revised form February, 1995.

<sup>†</sup> Department of Computer Science, Yale University, P. O. Box 208285, New Haven, CT 06520-8285 (eisenstat-stan@cs.yale.edu). The research of this author was supported in part by U. S. Army Research Office contract DAAL03-91-G-0032.

<sup>‡</sup> Department of Mathematics and Statistics, Utah State University, Logan, UT 84322-3900 (walker@math.usu.edu). The work of this author was supported in part by United States Air Force Office of Scientific Research Grant AFOSR-91-0294, United States Department of Energy Grants DE-FG02-92ER25136 and DE-FG03-94ER25221, and National Science Foundation Grant DMS-9400217, all with Utah State University. It was done in part during a visit to the Computational and Applied Mathematics Department and the Center for Research on Parallel Computation, Rice University.

<sup>1</sup> These are implementations of Newton's method in which iterative linear algebra methods are used to solve the Newton equation approximately.

first, and then an  $s_k$  is determined so that (1.2) holds. The role of  $\eta_k$  is, then, to force  $\|F(x_k) + F'(x_k) s_k\|$  to be small in a particular way; accordingly,  $\eta_k$  is often called a *forcing term*.

The local convergence of an inexact Newton method is controlled by the forcing terms. Some specific illustrative results are the following (see Dembo, Eisenstat, and Steihaug [4]): Under the present assumptions, if  $x_0$  is sufficiently close to  $x_*$  and  $0 \leq \eta_k \leq \eta_{\max} < 1$  for each  $k$ , then  $\{x_k\}$  converges to  $x_*$   $q$ -linearly in the norm  $\|\cdot\|_*$ , defined by  $\|v\|_* \equiv \|F'(x_*)v\|$  for  $v \in \mathbb{R}^n$ , with asymptotic rate constant no greater than  $\eta_{\max}$ . Furthermore, if  $\lim_{k \rightarrow \infty} \eta_k = 0$ , then the convergence is  $q$ -superlinear, and if  $\eta_k = O(\|F(x_k)\|)$ , then the convergence is  $q$ -quadratic.<sup>2</sup>

In addition to controlling local convergence, there is another important issue associated with the forcing terms. Away from a solution,  $F$  and its local linear model may disagree considerably at a step that closely approximates the Newton step. Thus choosing  $\eta_k$  too small may lead to *oversolving* the Newton equation, by which we mean imposing an accuracy on an approximation of the Newton step that leads to significant disagreement between  $F$  and its local linear model. Oversolving may result in little or no decrease in  $\|F\|$  and, therefore, little or no progress toward a solution. Moreover, in applications such as Newton iterative or truncated Newton methods, in which additional accuracy in solving the Newton equation requires additional expense, it may entail pointless costs; a less accurate approximation of the Newton step may be both cheaper and more effective.

Our purpose is to propose choices of the forcing terms that achieve desirably fast local convergence and also tend to avoid oversolving. All of the proposed choices incorporate information about  $F$  but are scale independent in that they do not change if  $F$  is multiplied by a constant.

In §2, we outline the proposed choices and analyze the local convergence of Algorithm IN that results from them; we also note some practical safeguards that improve performance. In §3, we discuss numerical experiments. The algorithm used in the experiments is a special case of Algorithm IN and is outlined in §3.1. The test problems are described in §3.2. An example of oversolving is given in §3.3, with additional observations and examples in §3.4. Summary test results are shown in §3.5. A summary discussion is given in §4.

*Preliminaries.* We define some useful constants and formulate several elementary results. Set  $M \equiv \max\{\|F'(x_*)\|, \|F'(x_*)^{-1}\|\}$ . For  $\delta > 0$ , define

$$N_\delta(x_*) \equiv \{x \in \mathbb{R}^n : \|x - x_*\| < \delta\},$$

and let  $\delta_* > 0$  be sufficiently small that

1.  $F$  is continuously differentiable and  $F'$  is nonsingular on  $N_{\delta_*}(x_*)$ ,
2.  $\|F'(x)^{-1}\| \leq 2M$  for  $x \in N_{\delta_*}(x_*)$ ,
3. inequality (1.1) holds for  $x \in N_{\delta_*}(x_*)$ ,
4.  $\delta_* < 2/(\lambda M)$ .

LEMMA 1.1. *If  $x \in N_{\delta_*}(x_*)$  and if  $s$  is such that  $x_+ \equiv x + s \in N_{\delta_*}(x_*)$ , then*

$$\|F(x_+) - F(x) - F'(x)s\| \leq \lambda \left( 2\|x - x_*\| + \frac{\|s\|}{2} \right) \|s\|.$$

---

<sup>2</sup> See, e.g., Dennis and Schnabel [6, §§2.3 and 3.1] for definitions of the types of convergence referred to throughout this paper.

*Proof.* Setting  $x(t) \equiv x + ts$  for  $0 \leq t \leq 1$ , we have

$$\begin{aligned}
\|F(x_+) - F(x) - F'(x)s\| &= \left\| \int_0^1 F'(x(t))s \, dt - F'(x)s \right\| \\
&\leq \left\| \int_0^1 [F'(x(t)) - F'(x_*)] \, dt - [F'(x) - F'(x_*)] \right\| \|s\| \\
&\leq \left( \int_0^1 \lambda [\|x - x_*\| + t\|s\|] \, dt + \lambda\|x - x_*\| \right) \|s\| \\
&= \lambda \left( 2\|x - x_*\| + \frac{\|s\|}{2} \right) \|s\|.
\end{aligned}$$

□

LEMMA 1.2. *There is a  $\mu > 0$  such that*

$$\frac{1}{\mu}\|x - x_*\| \leq \|F(x)\| \leq \mu\|x - x_*\|.$$

whenever  $x \in N_{\delta_*}(x_*)$ .

*Proof.* With Lemma 1.1, we have

$$\begin{aligned}
\|F(x)\| &\leq \|F'(x_*)(x - x_*)\| + \|F(x) - F(x_*) - F'(x_*)(x - x_*)\| \\
&\leq M\|x - x_*\| + \frac{\lambda}{2}\|x - x_*\|^2 \leq \left( M + \frac{\lambda\delta_*}{2} \right) \|x - x_*\|
\end{aligned}$$

and

$$\begin{aligned}
\|F(x)\| &\geq \|F'(x_*)(x - x_*)\| - \|F(x) - F(x_*) - F'(x_*)(x - x_*)\| \\
&\geq \frac{1}{M}\|x - x_*\| - \frac{\lambda}{2}\|x - x_*\|^2 \geq \left( \frac{1}{M} - \frac{\lambda\delta_*}{2} \right) \|x - x_*\|.
\end{aligned}$$

The lemma follows with  $\mu \equiv \max \left\{ M + \lambda\delta_*/2, (1/M - \lambda\delta_*/2)^{-1} \right\}$ . □

LEMMA 1.3. *If  $x \in N_{\delta_*}(x_*)$  and  $\|F(x) + F'(x)s\| \leq \eta\|F(x)\|$  for some  $s$  and  $\eta \in [0, 1)$ , then  $\|s\| \leq 4M\|F(x)\|$ .*

*Proof.* We have

$$\begin{aligned}
\|s\| &\leq \|F'(x)^{-1}\| \|F'(x)s\| \\
&\leq 2M(\|F(x)\| + \|F(x) + F'(x)s\|) \\
&\leq 2M(1 + \eta)\|F(x)\| \leq 4M\|F(x)\|.
\end{aligned}$$

□

LEMMA 1.4. *There is a  $B > 0$  such that if  $x \in N_{\delta_*}(x_*)$ ,  $s$  and  $\eta \in [0, 1)$  are such that  $\|F(x) + F'(x)s\| \leq \eta\|F(x)\|$ , and  $x_+ \equiv x + s \in N_{\delta_*}(x_*)$ , then*

$$\|F(x_+)\| \leq (\eta + B\|F(x)\|)\|F(x)\|.$$

*Proof.* With Lemmas 1.1–1.3, we have that

$$\begin{aligned}
\|F(x_+)\| &\leq \|F(x) + F'(x)s\| + \|F(x_+) - F(x) - F'(x)s\| \\
&\leq \eta\|F(x)\| + \lambda \left( 2\|x - x_*\| + \frac{\|s\|}{2} \right) \|s\| \\
&\leq \eta\|F(x)\| + \lambda(2\mu\|F(x)\| + 2M\|F(x)\|) \cdot 4M\|F(x)\| \\
&= (\eta + B\|F(x)\|)\|F(x)\|,
\end{aligned}$$

where  $B \equiv 8\lambda M(\mu + M)$ . □

**2. The proposed choices.** In the analysis in this section, we use the Lipschitz constant  $\lambda$  in (1.1) and the constants  $M$ ,  $\delta_*$ ,  $\mu$ , and  $B$  introduced in the preliminaries in §1. We also let  $\delta$  be such that  $0 < \delta \leq \delta_*/(1 + 4\mu M)$  and note the following consequence of Lemmas 1.2 and 1.3:

**PROPOSITION 2.1.** *If  $x \in N_\delta(x_*)$  and  $\|F(x) + F'(x)s\| \leq \eta\|F(x)\|$  for some  $s$  and  $\eta \in [0, 1)$ , then  $x + s \in N_{\delta_*}(x_*)$ .*

We assume for convenience that Algorithm IN continues indefinitely without termination and that  $F(x_k) \neq 0$  for all  $k$ . Note that if  $x_k \in N_{\delta_*}(x_*)$ , then  $F'(x_k)$  is nonsingular and, therefore, suitable  $s_k$  and  $x_{k+1}$  exist for any  $\eta_k \in [0, 1)$ . Our standing assumptions on  $F$  and  $x_*$  are those made in the first paragraph of §1.

Our first choice is the following:

*Choice 1:* Given  $\eta_0 \in [0, 1)$ , choose

$$(2.1) \quad \eta_k = \frac{\|F(x_k) - F(x_{k-1}) - F'(x_{k-1})s_{k-1}\|}{\|F(x_{k-1})\|}, \quad k = 1, 2, \dots,$$

or

$$(2.2) \quad \eta_k = \frac{\left| \|F(x_k)\| - \|F(x_{k-1}) + F'(x_{k-1})s_{k-1}\| \right|}{\|F(x_{k-1})\|}, \quad k = 1, 2, \dots$$

Note that  $\eta_k$  given by either (2.1) or (2.2) directly reflects the agreement between  $F$  and its local linear model at the previous step. The choice (2.2) may be more convenient to evaluate than (2.1) in some circumstances. Since it is at least as small, local convergence will be at least as fast as with (2.1); however, if it is significantly smaller, then it may be more difficult to find a suitable step in some applications and perhaps risk greater oversolving as well.

**THEOREM 2.2.** *Under the standing assumptions on  $F$  and  $x_*$ , if  $x_0$  is sufficiently near  $x_*$ , then  $\{x_k\}$  produced by Algorithm IN with  $\{\eta_k\}$  given by Choice 1 remains in  $N_{\delta_*}(x_*)$  and converges to  $x_*$  with*

$$(2.3) \quad \|x_{k+1} - x_*\| \leq \beta \|x_k - x_*\| \|x_{k-1} - x_*\|, \quad k = 1, 2, \dots,$$

for a constant  $\beta$  independent of  $k$ .

*Remark:* It follows immediately from (2.3) that the convergence is  $q$ -superlinear and two-step  $q$ -quadratic. As in the case of the classical secant method, it also follows that the convergence is of  $r$ -order  $(1 + \sqrt{5})/2$ ; see, e.g., Stoer and Bulirsch [14, p. 293] for the argument.

*Proof.* It suffices to prove the theorem with  $\{\eta_k\}$  given by (2.1).

Suppose that  $\eta_0 \in [0, 1)$  is given. Let  $\tau$  be such that  $\eta_0 < \tau < 1$ , and let  $\epsilon > 0$  be sufficiently small that  $\eta_0 + B\epsilon \leq \tau$ ,  $[8\lambda M(\mu + M) + B]\epsilon \leq \tau$ , and  $\epsilon < \delta/\mu$ . Note that if  $x \in N_{\delta_*}(x_*)$  and  $\|F(x)\| \leq \epsilon$ , then  $x \in N_\delta(x_*)$  by Lemma 1.2.

Let  $x_0 \in N_\delta(x_*)$  be sufficiently near  $x_*$  that  $\|F(x_0)\| \leq \epsilon$ . Since  $x_0 \in N_\delta(x_*)$ , we have  $x_1 \in N_{\delta_*}(x_*)$  by Proposition 2.1. Also, by Lemma 1.4,

$$(2.4) \quad \begin{aligned} \|F(x_1)\| &\leq (\eta_0 + B\|F(x_0)\|)\|F(x_0)\| \leq (\eta_0 + B\epsilon)\|F(x_0)\| \\ &\leq \tau\|F(x_0)\| \leq \|F(x_0)\| \leq \epsilon, \end{aligned}$$

and, hence,  $x_1 \in N_\delta(x_*)$ .

As an inductive hypothesis, suppose that, for some  $k \geq 1$ , we have  $x_k \in N_\delta(x_*)$ ,  $x_{k-1} \in N_\delta(x_*)$ ,  $\|F(x_k)\| \leq \epsilon$ , and  $\|F(x_{k-1})\| \leq \epsilon$ . Then  $x_{k+1} \in N_{\delta_*}(x_*)$  by Proposition 2.1, and Lemmas 1.1–1.3 give

$$\begin{aligned} \eta_k &= \frac{\|F(x_k) - F(x_{k-1}) - F'(x_{k-1})s_{k-1}\|}{\|F(x_{k-1})\|} \\ &\leq \frac{\lambda(2\|x_{k-1} - x_*\| + \|s_{k-1}\|/2)\|s_{k-1}\|}{\|F(x_{k-1})\|} \\ &\leq \frac{\lambda(2\mu\|F(x_{k-1})\| + 2M\|F(x_{k-1})\|) \cdot 4M\|F(x_{k-1})\|}{\|F(x_{k-1})\|} \\ &= 8\lambda M(\mu + M)\|F(x_{k-1})\|. \end{aligned}$$

Then Lemma 1.4 implies

$$\begin{aligned} (2.5) \quad \|F(x_{k+1})\| &\leq (\eta_k + B\|F(x_k)\|)\|F(x_k)\| \\ &\leq [8\lambda M(\mu + M)\|F(x_{k-1})\| + B\|F(x_k)\|]\|F(x_k)\| \\ &\leq [8\lambda M(\mu + M) + B]\epsilon\|F(x_k)\| \leq \tau\|F(x_k)\|. \end{aligned}$$

Thus  $\|F(x_{k+1})\| \leq \tau\|F(x_k)\| \leq \epsilon$  and, hence,  $x_{k+1} \in N_\delta(x_*)$ .

It follows from this induction that  $\{x_k\} \subset N_\delta(x_*) \subset N_{\delta_*}(x_*)$ . Furthermore, (2.4) and (2.5) give  $\|F(x_{k+1})\| \leq \tau\|F(x_k)\|$  for each  $k \geq 0$ ; hence,  $F(x_k) \rightarrow 0$  and, by Lemma 1.2,  $x_k \rightarrow x_*$  as well.

To show (2.3), we note that (2.4) and (2.5) give, for  $k \geq 1$ ,  $\|F(x_k)\| \leq \|F(x_{k-1})\|$  and

$$\begin{aligned} \|F(x_{k+1})\| &\leq [8\lambda M(\mu + M)\|F(x_{k-1})\| + B\|F(x_k)\|]\|F(x_k)\| \\ &\leq [8\lambda M(\mu + M) + B]\|F(x_{k-1})\|\|F(x_k)\|. \end{aligned}$$

With Lemma 1.2, this implies (2.3) with  $\beta \equiv \mu^3[8\lambda M(\mu + M) + B]$ .  $\square$

One possible way to obtain faster local convergence while retaining the potential advantages of (2.1) and (2.2) is to raise those expressions to powers greater than one. A particular possibility that we considered in our numerical experiments is squaring those expressions. We note without proof that this leads to local convergence with

$$\|x_{k+1} - x_*\| \leq \beta \max \left\{ \|x_{k-1} - x_*\|^2, \|x_k - x_*\| \right\} \|x_k - x_*\|, \quad k = 1, 2, \dots,$$

which implies that  $x_k \rightarrow x_*$   $r$ -quadratically. However, this possibility was not as successful in our experiments as the other choices proposed here, and we do not consider it further.

Our second choice is the following:

*Choice 2:* Given  $\gamma \in [0, 1]$ ,  $\alpha \in (1, 2]$ , and  $\eta_0 \in [0, 1)$ , choose

$$(2.6) \quad \eta_k = \gamma \left( \frac{\|F(x_k)\|}{\|F(x_{k-1})\|} \right)^\alpha, \quad k = 1, 2, \dots$$

The choice (2.6) does not directly reflect the agreement between  $F$  and its local linear model, as does Choice 1. However, the experiments in §3 show that it results in

little oversolving in practice, and the following theorem shows that it offers attractive local convergence.

**THEOREM 2.3.** *Under the standing assumptions on  $F$  and  $x_*$ , if  $x_0$  is sufficiently near  $x_*$ , then  $\{x_k\}$  produced by Algorithm IN with  $\{\eta_k\}$  given by Choice 2 remains in  $N_{\delta_*}(x_*)$  and converges to  $x_*$ . If  $\gamma < 1$ , then the convergence is of  $q$ -order  $\alpha$ . If  $\gamma = 1$ , then the convergence is of  $r$ -order  $\alpha$  and of  $q$ -order  $p$  for every  $p \in [1, \alpha)$ .*

*Proof.* Suppose that  $\eta_0 \in [0, 1)$  is given and let  $\epsilon > 0$  be sufficiently small that  $\eta_0 + B\epsilon \leq \eta_0^{1/\alpha}$  and  $\epsilon < \delta/\mu$ . Note that if  $x \in N_{\delta_*}(x_*)$  and  $\|F(x)\| \leq \epsilon$ , then  $x \in N_\delta(x_*)$  by Lemma 1.2.

Let  $x_0 \in N_\delta(x_*)$  be sufficiently near  $x_*$  that  $\|F(x_0)\| \leq \epsilon$ . As an inductive hypothesis, suppose that, for some  $k \geq 0$ , we have  $x_k \in N_\delta(x_*)$ ,  $\|F(x_k)\| \leq \epsilon$ , and  $\eta_k \leq \eta_0$ . Since  $x_k \in N_\delta(x_*)$ , we have  $x_{k+1} \in N_{\delta_*}(x_*)$  by Proposition 2.1. Also, by Lemma 1.4,

$$(2.7) \quad \begin{aligned} \|F(x_{k+1})\| &\leq (\eta_k + B\|F(x_k)\|)\|F(x_k)\| \\ &\leq (\eta_0 + B\epsilon)\|F(x_k)\| \leq \eta_0^{1/\alpha} \|F(x_k)\|. \end{aligned}$$

Then  $\|F(x_{k+1})\| \leq \eta_0^{1/\alpha} \epsilon \leq \epsilon$ , and it follows that  $x_{k+1} \in N_\delta(x_*)$ . Furthermore, (2.7) gives

$$\eta_{k+1} = \gamma(\|F(x_{k+1})\|/\|F(x_k)\|)^\alpha \leq \gamma\eta_0 \leq \eta_0.$$

It follows from this induction that  $\{x_k\} \subset N_\delta(x_*) \subset N_{\delta_*}(x_*)$ . Furthermore, (2.7) gives  $\|F(x_{k+1})\| \leq \eta_0^{1/\alpha} \|F(x_k)\|$  for each  $k \geq 0$ ; hence,  $F(x_k) \rightarrow 0$  and, by Lemma 1.2,  $x_k \rightarrow x_*$  as well.

It remains to show the desired rates of convergence. Note that, for  $k > 0$ , (2.7) and (2.6) give

$$(2.8) \quad \|F(x_{k+1})\| \leq \left[ \gamma \left( \frac{\|F(x_k)\|}{\|F(x_{k-1})\|} \right)^\alpha + B\|F(x_k)\| \right] \|F(x_k)\|.$$

First, suppose that  $\gamma < 1$  and set  $\rho_k \equiv \|F(x_k)\|/\|F(x_{k-1})\|^\alpha$  for  $k > 0$ . From (2.8) and (2.7), we have  $\rho_{k+1} \leq \gamma\rho_k + B\|F(x_k)\|^{2-\alpha} \leq \gamma\rho_k + B\|F(x_0)\|^{2-\alpha}$  for  $k > 0$ , and it follows inductively that

$$\rho_{k+1} \leq \gamma^k \rho_1 + \left( \sum_{j=0}^{k-1} \gamma^j \right) B\|F(x_0)\|^{2-\alpha} \leq \rho_1 + \frac{B}{1-\gamma} \|F(x_0)\|^{2-\alpha}.$$

Thus  $\{\rho_k\}$  is uniformly bounded. Consequently,  $F(x_k) \rightarrow 0$  with  $q$ -order  $\alpha$ , and it follows from Lemma 1.2 that  $x_k \rightarrow x_*$  with  $q$ -order  $\alpha$  as well.

Now, suppose that  $\gamma = 1$ . We first show that the convergence is of  $q$ -order  $p$  for  $p \in [1, \alpha)$ . For  $k > 0$ , (2.8) gives

$$(2.9) \quad \begin{aligned} \|F(x_{k+1})\| &\leq \left[ \left( \frac{\|F(x_k)\|}{\|F(x_{k-1})\|} \right)^\alpha + B\|F(x_k)\| \right] \|F(x_k)\| \\ &= \left[ \left( \frac{\|F(x_k)\|}{\|F(x_{k-1})\|} \right)^{\alpha-p} \frac{\|F(x_k)\|}{\|F(x_{k-1})\|^p} + B\|F(x_k)\|^{2-p} \right] \|F(x_k)\|^p. \end{aligned}$$

For each  $k > 0$ , set  $\sigma_k \equiv \|F(x_k)\|/\|F(x_{k-1})\|^p$  and recall that (2.7) gives  $\|F(x_k)\| \leq \eta_0^{1/\alpha} \|F(x_{k-1})\|$ , whence  $\|F(x_k)\| \leq \eta_0^{k/\alpha} \|F(x_0)\|$ . Then for  $k > 0$ , (2.9) implies

$$\sigma_{k+1} \leq \eta_0^{1-p/\alpha} \sigma_k + B\eta_0^{k(2-p)/\alpha} \|F(x_0)\|^{2-p} \leq \xi \sigma_k + \xi^k C,$$

where  $\xi \equiv \eta_0^{1-p/\alpha}$  and  $C \equiv B\|F(x_0)\|^{2-p}$ . It follows inductively that

$$\sigma_{k+1} \leq \xi^k (\sigma_1 + kC),$$

and, hence,

$$\|F(x_{k+1})\| \leq \xi^k (\sigma_1 + kC) \|F(x_k)\|^p.$$

Since  $\xi^k (\sigma_1 + kC) \rightarrow 0$  as  $k \rightarrow \infty$ , we conclude that  $F(x_k) \rightarrow 0$  with  $q$ -order  $p$  and, by Lemma 1.2,  $x_k \rightarrow x_*$  with  $q$ -order  $p$  as well.

Still assuming  $\gamma = 1$ , we now show that  $x_k \rightarrow x_*$  with  $r$ -order  $\alpha$ . By Lemma 1.2, it suffices to show that  $\|F(x_k)\| \rightarrow 0$  with  $r$ -order  $\alpha$ ; we shall prove the somewhat stronger result that  $\tau_k \equiv \|F(x_k)\|/\|F(x_{k-1})\| \rightarrow 0$  with  $r$ -order  $\alpha$ .

It follows from the results above that  $\tau_k \rightarrow 0$ . Then there is a  $k_0$  such that  $(2\tau_{k_0+1})^{(\alpha-1)} + 2B\|F(x_{k_0})\| \leq 1$ . For convenience, we re-index if necessary so that  $k_0 = 0$ . Then  $(2\tau_1)^{(\alpha-1)} + 2B\|F(x_0)\| \leq 1$ , which implies  $D \equiv 1/(2\tau_1) > 1$ . Set  $\beta_k \equiv D\tau_k$  for  $k \geq 0$ . Note that  $\beta_1 = 1/2$ . It suffices to show that  $\beta_k \rightarrow 0$  with  $r$ -order  $\alpha$ .

We claim that  $\beta_k \leq \beta_1^{\alpha^{k-1}}$  for  $k = 1, 2, \dots$ , from which it follows that  $\beta_k \rightarrow 0$  with  $r$ -order  $\alpha$ . The claim clearly holds for  $k = 1$ . Suppose that it holds up to some  $k \geq 1$ . Then Lemma 1.4 implies

$$\|F(x_{k+1})\| \leq (\tau_k^\alpha + B\|F(x_k)\|) \|F(x_k)\|,$$

whence

$$\tau_{k+1} \leq \tau_k^\alpha + B\tau_k \dots \tau_1 \|F(x_0)\|.$$

From this, we obtain

$$\begin{aligned} \beta_{k+1} &\leq D^{1-\alpha} \beta_k^\alpha + \frac{B\|F(x_0)\|}{D^{k-1}} \beta_k \dots \beta_1 \\ &\leq D^{1-\alpha} \left(\beta_1^{\alpha^{k-1}}\right)^\alpha + B\|F(x_0)\| \beta_1^{(\alpha^{k-1} + \dots + 1)} \\ &\leq \left(D^{1-\alpha} + B\|F(x_0)\|/\beta_1\right) \beta_1^{\alpha^k} \\ &= \left[(2\tau_1)^{\alpha-1} + 2B\|F(x_0)\|\right] \beta_1^{\alpha^k} \leq \beta_1^{\alpha^k}, \end{aligned}$$

and the proof is complete.  $\square$

It is possible to show local convergence for Algorithm IN when  $\{\eta_k\}$  is given by Choice 2 with  $\gamma > 1$ , provided  $\eta_0$  is sufficiently small. However, Choice 2 with  $\gamma > 1$  was not competitive in our experiments.



**2.1. Practical safeguards.** Although the forcing term choices given above are usually effective in avoiding oversolving, we have observed in experiments that they occasionally become too small far away from a solution. There is a particular danger of the Choice 1 forcing terms becoming too small; indeed, an  $\eta_k$  given by (2.1) or (2.2) can be undesirably small because of either a very small step or coincidental very good agreement between  $F$  and its local linear model. In our experiments, we observed relatively few occasions on which the Choice 2 forcing terms became undesirably small; however, this did occur.

We introduce safeguards here that are intended to prevent the forcing terms from becoming too small too quickly. The rationale is that if large forcing terms are appropriate at some point, then subsequent forcing terms should not be allowed to become much smaller until this has been justified over several iterations. These are not claimed to be the most effective safeguards that might be devised for general use or even for the test problems used in our experiments. However, they were consistently effective in our tests, more so than several other possibilities that we tried, and they serve to demonstrate the usefulness of safeguards.

For each choice, we restrict  $\eta_k$  to be no less than a certain minimum value, but only if that minimum value is above a threshold. The minimum value is determined by raising  $\eta_{k-1}$  to a power associated with the rate of convergence of the (unsafeguarded) choice. The threshold that we use here is .1; this is clearly somewhat arbitrary but was effective in our experiments. Note that, in each case, the minimum value eventually drops below the threshold whenever there is convergence to a solution. Thus the safeguards eventually become inactive whenever there is convergence, and the asymptotic convergence is that for the unsafeguarded choice given by the theorems above.

For Choice 1, the safeguard is the following:

*Choice 1 safeguard:* Modify  $\eta_k$  by  $\eta_k \leftarrow \max\{\eta_k, \eta_{k-1}^{(1+\sqrt{5})/2}\}$  whenever  $\eta_{k-1}^{(1+\sqrt{5})/2} > .1$ .

For perspective, recall from the remark after Theorem 2.2 that the convergence of (2.3) implies convergence of  $r$ -order  $(1 + \sqrt{5})/2$ . For Choice 2, the safeguard is the following:

*Choice 2 safeguard:* Modify  $\eta_k$  by  $\eta_k \leftarrow \max\{\eta_k, \gamma \eta_{k-1}^\alpha\}$  whenever  $\gamma \eta_{k-1}^\alpha > .1$ .

Finally, we note that, away from a solution, it may be possible for each of the proposed choices to be greater than one. Accordingly, it may be necessary in practice to impose an additional safeguard to make sure that  $\eta_k \in [0, 1)$  for each  $k$ , as in the algorithm in §3.1 below that was used in our experiments.

**3. Numerical experiments.** In this section, we report on numerical experiments with the forcing term choices outlined in §2, modified with the safeguards given in §2.1. In the experiments, for computational convenience, we always used  $\eta_k$  given by (2.2) for Choice 1. For Choice 2, we used  $\gamma = 1, .9, .5$  and  $\alpha = 2, (1 + \sqrt{5})/2$ . The latter value of  $\alpha$  results in an order of convergence roughly comparable to that for Choice 1; see Theorem 2.3 and the remark after Theorem 2.2. For a broader comparison, we also included the following representative forcing term choices:

1. The choice  $\eta_k = 10^{-1}$ , which requires modestly accurate approximations of Newton steps and results in local  $q$ -linear convergence in the norm  $\|\cdot\|_*$ .
2. The choice  $\eta_k = 10^{-4}$  used by Cai, Gropp, Keyes, and Tidiri [3], which requires uniformly close approximations of Newton steps for all  $k$  and results

- in fast local  $q$ -linear convergence in the norm  $\|\cdot\|_*$ .
3. The choice  $\eta_k = 1/2^{k+1}$  of Brown and Saad [2]. This choice results in local  $q$ -superlinear convergence and allows relatively inaccurate approximations of Newton steps for small  $k$ , when  $x_k$  may not be near  $x_*$ ; however, it incorporates no information about  $F$ .
  4. The choice  $\eta_k = \min\{1/(k+2), \|F(x_k)\|\}$  of Dembo and Steihaug [5]. This choice results in  $q$ -quadratic local convergence and also may allow relatively inaccurate approximations of Newton steps for small  $k$ . It incorporates some information about  $F$ ; however, it does not reflect the agreement of  $F$  and its local linear model and, in addition, depends on the scale of  $F$ .

**3.1. The algorithm.** A globalized inexact Newton algorithm was necessary because initial approximate solutions were not always near a solution. We used Algorithm INB of Eisenstat and Walker [7, §6]. This is an inexact Newton method globalized by backtracking, which we write here as follows:

**Algorithm INB:** Inexact Newton Backtracking Method [7]

LET  $x_0, \eta_{\max} \in [0, 1), t \in (0, 1),$  AND  $0 < \theta_{\min} < \theta_{\max} < 1$  BE GIVEN.

FOR  $k = 0$  STEP 1 UNTIL “CONVERGENCE” DO:

CHOOSE AN **initial**  $\eta_k \in [0, \eta_{\max}]$  AND  $s_k$  SUCH THAT

$$\|F(x_k) + F'(x_k) s_k\| \leq \eta_k \|F(x_k)\|.$$

WHILE  $\|F(x_k + s_k)\| > [1 - t(1 - \eta_k)] \|F(x_k)\|$  DO:

CHOOSE  $\theta \in [\theta_{\min}, \theta_{\max}]$ .

UPDATE  $s_k \leftarrow \theta s_k$  AND  $\eta_k \leftarrow 1 - \theta(1 - \eta_k)$ .

SET  $x_{k+1} = x_k + s_k$ .

Note that Algorithm INB requires  $\eta_k \in [0, \eta_{\max}]$  for each initial  $\eta_k$ . For the safeguarded choices in §2, this necessitates the additional safeguard  $\eta_k \leftarrow \min\{\eta_k, \eta_{\max}\}$ .

Theorem 6.1 of Eisenstat and Walker [7] states that if  $\{x_k\}$  generated by Algorithm INB has a limit point  $x_*$  such that  $F'(x_*)$  is invertible, then  $F(x_*) = 0$  and  $x_k \rightarrow x_*$ . Furthermore, in this case, the initial  $\eta_k$  and  $s_k$  are accepted without modification for all sufficiently large  $k$ ; it follows in particular that the asymptotic convergence to  $x_*$  is determined by the initial  $\eta_k$ 's.

In implementing Algorithm INB, we first chose each initial  $\eta_k$  (with  $\eta_0 = 1/2$  for Choices 1 and 2) and then determined an initial  $s_k$  by approximately solving the Newton equation using GMRES( $m$ ), the restarted GMRES method of Saad and Schultz [12], with restart value  $m = 20$ . Products of  $F'(x_k)$  with vectors were evaluated analytically in some cases and approximated by finite differences of  $F$ -values in others; see §3.2. When finite-difference approximations were used, a second-order central difference was used to evaluate the initial residual at the beginning of each cycle of 20 GMRES steps, and subsequently first-order forward differences were used within the cycle. This selective second-order differencing gave essentially the same accuracy as if central differences had been used throughout, but at much lower cost (see Turner and Walker [16]).

The parameters used were  $\eta_{\max} = .9, t = 10^{-4}, \theta_{\min} = 1/10,$  and  $\theta_{\max} = 1/2$ . The norm was the Euclidean norm  $\|\cdot\|_2$ . In the while-loop, each  $\theta$  was chosen to minimize over  $[\theta_{\min}, \theta_{\max}]$  the quadratic  $p(\theta)$  for which  $p(0) = g(0), p'(0) = g'(0),$  and  $p(1) = g(1),$  where  $g(\theta) \equiv \|F(x_k + \theta s_k)\|_2^2$ . Convergence was declared when either

$\|F(x_k)\|_2 \leq 10^{-12}\|F(x_0)\|_2$  or  $\|s_k\|_2 \leq 10^{-12}$ . These tight stopping tolerances allowed asymptotic convergence behavior to become evident.<sup>3</sup> Failure was declared when one of the following occurred: (1)  $k$  reached 200 without convergence, (2) an initial  $s_k$  was not found in 1000 GMRES(20) iterations, or (3) ten iterations of the while-loop failed to produce an acceptable step. All computing was done in double precision on Sun Microsystems workstations using the Sun Fortran compiler.

**3.2. The test problems.** The test set consisted of four PDE problems and two integral equation problems. The PDE problems are all elliptic boundary value problems posed on  $\Omega \equiv [0, 1] \times [0, 1] \subseteq \mathbb{R}^2$ .

**3.2.1. A PDE problem.** The problem is

$$\Delta u + u^3 = 0 \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega.$$

This problem has multiple solutions, but only one that is positive everywhere (McKenna [10], Schaaf [13]). These properties appear to be shared by the discretized problem, and finding the everywhere-positive solution can be difficult without a good initial approximate solution. Discretization was by the usual centered differences on a  $100 \times 100$  uniform grid, so that  $n = 10^4$ . The discretized problem was preconditioned on the right using a fast Poisson solver from FISHPACK (Swartztrauber and Sweet [15]). Products of  $F'$  with vectors were evaluated analytically. The initial approximate solution was a discretization of  $u_0(x) \equiv \kappa x_1(1 - x_1)x_2(1 - x_2)$ , which should lead to the everywhere-positive solution for large  $\kappa$ . Two test cases were considered:  $\kappa = 100$  and  $\kappa = 1000$ . For the latter value, the initial approximate solution is farther from the solution and the problem is harder.

**3.2.2. The (modified) Bratu problem.** The problem is

$$\Delta u + \kappa \frac{\partial u}{\partial x_1} + \lambda e^u = 0 \text{ in } \Omega, \quad u = 0 \text{ on } \partial\Omega.$$

The actual Bratu (or Gelfand) problem has  $\kappa = 0$ ; see, e.g., Glowinski, Keller, and Reinhart [8] or the description by Glowinski and Keller in the collection of nonlinear model problems assembled by Moré [11, pp. 733-737]. As  $\kappa$  and  $\lambda$  grow, solving the Newton equations for the discretized problem becomes harder for GMRES(20). Discretization and preconditioning were as in §3.2.1. Products of  $F'$  with vectors were evaluated analytically. The initial approximate solution was zero. Two test cases were considered:  $\kappa = \lambda = 10$  and  $\kappa = \lambda = 20$ .

**3.2.3. The driven cavity problem.** The problem is

$$(1/Re)\Delta^2\psi + \frac{\partial\psi}{\partial x_1}\frac{\partial}{\partial x_2}\Delta\psi - \frac{\partial\psi}{\partial x_2}\frac{\partial}{\partial x_1}\Delta\psi = 0 \text{ in } \Omega,$$

$$\psi = 0 \quad \text{and} \quad \frac{\partial\psi}{\partial n} = g \quad \text{on } \partial\Omega,$$

---

<sup>3</sup> In some applications, less stringent convergence tolerances are commonly used. As a result, asymptotic convergence behavior may not be very important, and it may be appropriate to use forcing terms that are not asymptotically increasingly demanding, such as constant forcing terms that give adequately fast  $q$ -linear convergence.

where  $g(x_1, x_2) = 1$  if  $x_2 = 1$  and  $g(x_1, x_2) = 0$  if  $0 \leq x_2 < 1$ . This is a widely used test problem; see, e.g., Brown and Saad [2] or Glowinski, Keller, and Reinhart [8]. The numerical problem becomes harder as the Reynolds number  $Re$  increases. Discretization was by piecewise-linear finite elements on a uniform  $63 \times 63$  grid<sup>4</sup>, so that  $n = 3969$ . The discretized problem was preconditioned on the right using a fast biharmonic solver of Bjørstad [1]. Products of  $F'$  with vectors were approximated with finite differences. The initial approximate solution was zero. Two test cases were considered:  $Re = 100$  and  $Re = 500$ .

**3.2.4. The porous medium equation.** The problem considered here is

$$\Delta(u^2) + d \frac{\partial}{\partial x_1}(u^3) + f = 0 \text{ in } \Omega,$$

with  $u = 1$  on the bottom and left sides of  $\Omega$  and  $u = 0$  on the top and right sides. This is more or less a steady-state special case of a general problem considered by van Duijn and de Graaf [17]. Discretization was by the usual centered differences on a  $64 \times 64$  uniform grid, so that  $n = 4096$ . The discretized problem was preconditioned on the right using the tridiagonal part of the Jacobian. Products of  $F'$  with vectors were evaluated analytically. The function  $f$  was a point source of magnitude 50 at the lower left grid point. The initial approximate solution was a discretization of  $u_0(x) \equiv 1 - x_1 x_2$  on the interior grid points, which tended to require more backtracking for negative  $d$  and to cause more oversolving for positive  $d$ . Two test cases were considered:  $d = 50$  and  $d = -50$ .

**3.2.5. An integral equation.** The problem, from Kelley and Northrup [9], is

$$cu(x)^2 - \frac{1}{2} \int_0^1 \cos(yu(x))u(y) dy + \frac{1}{2} \sin 1 - c = 0, \quad x \in [0, 1].$$

Clearly,  $u(x) \equiv 1$  is always a solution, and there exist other solutions for at least some values of  $c$ . The discretized problem was determined by approximating integrals using 20-point Gaussian quadrature<sup>5</sup> over 20 subintervals of  $[0, 1]$ , so that  $n = 400$ . No preconditioning was necessary. Products of  $F'$  with vectors were approximated with finite differences. The initial approximate solution was a discretization of  $u_0(x) \equiv 1 + \kappa \cos 9\pi x$ . One test case was considered:  $c = \kappa = 1.25$ .

**3.2.6. The Chandrasekhar H-equation.** The problem is

$$u(x) - \frac{1}{1 - Lu(x)} = 0, \quad x \in [0, 1],$$

where

$$Lu(x) \equiv \frac{c}{2} \int_0^1 \frac{xu(\xi)}{x + \xi} d\xi.$$

This problem arises in radiative transfer problems; see, e.g., the description by Kelley in the Moré problem collection [11, pp. 737-739]. The continuous problem is singular at  $c = 1$ , and so is the discretized problem considered here with discretization as in

<sup>4</sup> We thank P. N. Brown for providing the code for this.

<sup>5</sup> We thank C. T. Kelley for providing the code for this.

§3.2.5. The discretized problem becomes more difficult to solve as  $c \rightarrow 1$  but is still tractable at  $c = 1$ . As in §3.2.5, no preconditioning was necessary. Products of  $F'$  with vectors were approximated with finite differences. The initial approximate solution was zero. Three test cases were considered:  $c = .5$ ,  $c = .999$ , and  $c = 1$ .

**3.3. An example of oversolving.** Algorithm INB with the Dembo–Steihaug [5] choice  $\eta_k = \min\{1/(k + 2), \|F(x_k)\|_2\}$  was applied to the driven cavity problem with  $Re = 500$ . The results are shown in Figure 3.1, in which the logarithms of the norms of  $F$  and its local linear model are plotted as dotted and solid curves, respectively, versus the numbers of GMRES(20) iterations. (Most of the  $F$ -values used for Figures 3.1–3.4 would not normally be available but were computed for these illustrations.) Triangles indicate the start of new inexact Newton steps. In this example,  $\eta_k = \|F(x_k)\|_2$  for each  $k > 0$ ; the safeguard value  $\eta_k = 1/(k + 2)$  was never invoked for  $k > 0$ .

In Figure 3.1, gaps between the solid and dotted curves indicate oversolving. Note that once oversolving begins, there is virtually no further reduction in  $\|F\|_2$  until the beginning of the next inexact Newton step; thus further GMRES(20) iterations represent wasted effort. Note also the vertical discontinuity in the dotted curve at the end of the fourth inexact Newton step (after 45 GMRES(20) iterations); this indicates a reduction of the initial inexact Newton step through backtracking.

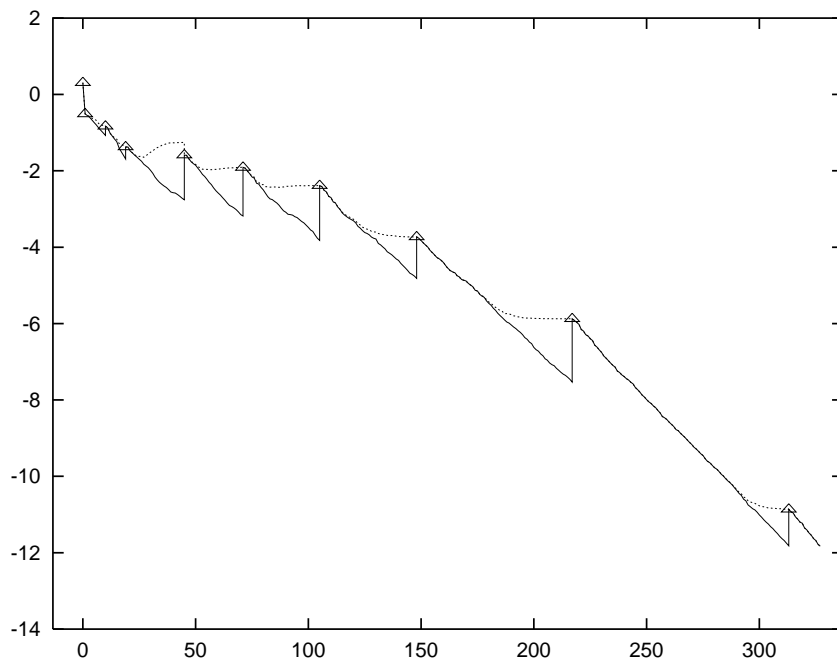


FIG. 3.1. Illustration of oversolving with  $\eta_k = \min\{1/(k + 2), \|F(x_k)\|_2\}$  on the driven cavity problem with  $Re = 500$ . The horizontal axis indicates the number of GMRES(20) iterations. The solid curve is  $\log_{10} \|F + F's\|_2$ ; the dotted curve is  $\log_{10} \|F\|_2$ . Triangles indicate new inexact Newton steps.

To show the benefits gained by reducing oversolving, we applied Algorithm INB with  $\eta_k$  given by the safeguarded Choice 1 to the same problem. The results are shown in Figure 3.2. Note that oversolving is almost eliminated and there are no step reductions through backtracking. Also, the total number of GMRES(20) iterations is 221, compared to 327 in the previous case. However, the number of inexact Newton steps is 12, compared to 10 previously.

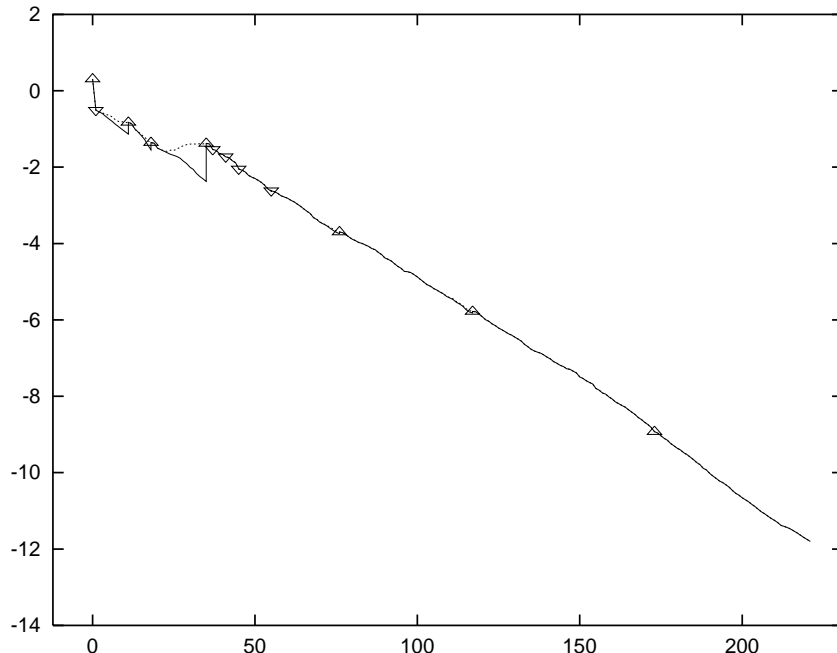


FIG. 3.2. Illustration of reduction of oversolving with the safeguarded Choice 1 forcing terms on the driven cavity problem with  $Re = 500$ . The horizontal axis indicates the number of GMRES(20) iterations. The solid curve is  $\log_{10} \|F + F' s\|_2$ ; the dotted curve is  $\log_{10} \|F\|_2$ . Triangles indicate new inexact Newton steps: “ $\Delta$ ” indicates  $\eta_k$  given by (2.2); “ $\nabla$ ” indicates the safeguard value.

**3.4. Additional observations and examples.** In an algorithm such as the implementation of Algorithm INB used here, choosing a very small forcing term may risk more than needless expense in obtaining an unnecessarily accurate solution of the Newton equation. First, if oversolving results, then disagreement between  $F$  and its local linear model may require significant work from the globalization procedure or even cause it to fail. In the example in §3.3, the choice  $\eta_k = \min\{1/(k+2), \|F(x_k)\|_2\}$  required one backtracking, while the safeguarded Choice 1 did not. We observed a more dramatic example involving the PDE problem of §3.2.1 with  $\kappa = 1000$ . With the safeguarded Choice 1, the iterates from Algorithm INB converged to the everywhere-positive solution in 40 GMRES(20) iterations; two backtracks were required. With the choice  $\eta_k = \min\{1/(k+2), \|F(x_k)\|_2\}$ , 164 GMRES(20) iterations and 11 backtracks were necessary; furthermore, convergence was to a solution other than the everywhere-positive solution. Such convergence to a “wrong” solution may or may not be undesirable per se, but it does indicate the potentially serious effects of disagreement between  $F$  and its local linear model.

Second, unless special care is taken, a very small forcing term may require more residual reduction than an iterative linear solver such as GMRES can accurately deliver, especially when products of  $F'$  with vectors are approximated with finite differences. Recall from §3.1 that our implementation of Algorithm INB uses selective second-order differencing to obtain essentially the same accuracy as if second-order differences were used throughout. Using the safeguarded Choice 2 forcing terms with  $\alpha = 2$  and  $\gamma = .9$ , we applied this implementation to the driven cavity problem with  $Re = 500$ ; the results are shown in Figure 3.3. There is no evidence of inaccuracy in GMRES(20), and 218 iterations were required for successful termination.

However, when the implementation was changed to use only first-order forward differences throughout, we obtained the results in Figure 3.4. Note the increase in the linear residual norm curve (the solid curve) just after iteration 200. The linear residual norm values used for this curve were evaluated directly at the beginning of each GMRES(20) cycle and then maintained recursively within the cycle; the observed increase occurs after the direct evaluation at iteration 200 and indicates that the recursively maintained values have become inaccurate. We note also that the number of GMRES(20) iterations required for termination has increased to 232.

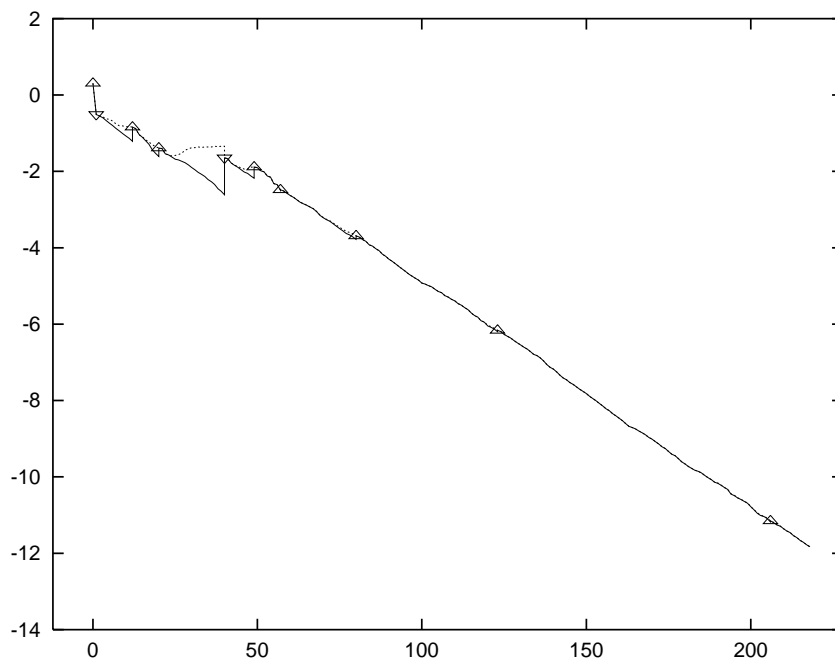


FIG. 3.3. Illustration of the performance of Algorithm INB with selective second-order differencing and safeguarded Choice 2 forcing terms,  $\alpha = 2$ ,  $\gamma = .9$ , on the driven cavity problem with  $Re = 500$ . The horizontal axis indicates the number of GMRES(20) iterations. The solid curve is  $\log_{10} \|F + F's\|_2$ ; the dotted curve is  $\log_{10} \|F\|_2$ . Triangles indicate new inexact Newton steps: “ $\Delta$ ” indicates  $\eta_k$  given by (2.6); “ $\nabla$ ” indicates the safeguard value.

**3.5. Summary test results.** In Table 3.1, we summarize the results of applying Algorithm INB to all test problem cases described in §3.2. In Table 3.2, we summarize the results over the PDE problem cases only. The results for the PDE problems are broken out in a separate table not only because these problems constitute an important problem class but also because the characteristic performance of Algorithm INB on these problems differed from that on the integral equations. On the integral equations, and on the H-equation in particular, GMRES(20) was so effective that the effects of different forcing term choices tended to be obscured. In most cases, only one to three GMRES(20) iterations were required for each inexact Newton step, and the linear residual norm was often reduced by several orders of magnitude in a single iteration. On the PDE problems, many more GMRES(20) iterations were typically required for each inexact Newton step, with only modest linear residual norm reduction per GMRES(20) iteration. Thus the PDE problems gave a somewhat more refined view of the effects of different forcing term choices.

The first three columns of Tables 3.1 and 3.2 give geometric means of the numbers

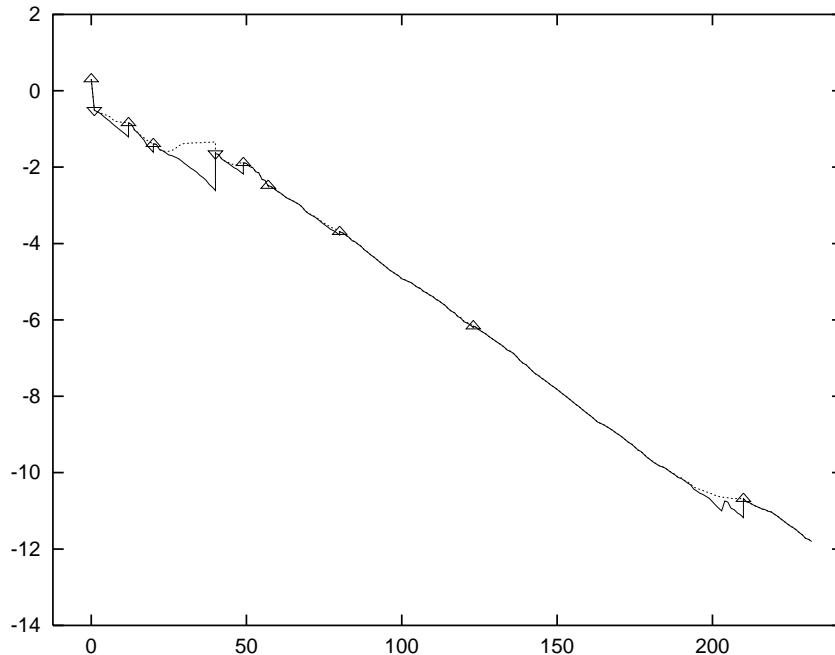


FIG. 3.4. *Illustration of the performance of Algorithm INB with first-order differencing throughout and safeguarded Choice 2 forcing terms,  $\alpha = 2$ ,  $\gamma = .9$ , on the driven cavity problem with  $Re = 500$ . The horizontal axis indicates the number of GMRES(20) iterations. The solid curve is  $\log_{10} \|F + F's\|_2$ ; the dotted curve is  $\log_{10} \|F\|_2$ . Triangles indicate new inexact Newton steps: “ $\Delta$ ” indicates  $\eta_k$  given by (2.6); “ $\nabla$ ” indicates the safeguard value.*

of linear iterations (GMRES(20) iterations), inexact Newton steps, and “function evaluation equivalents”, where, for each test case, we define the number of “function evaluation equivalents” to be the sum of the numbers of linear iterations, backtracks, and inexact Newton steps. The number of linear iterations is the same as the number of products of  $F'$  with vectors; if these products were always approximated by first-order forward differences, then the number of “function evaluation equivalents” would be equal to the number of function evaluations. This number provides a rough relative measure of overall work for the test problems used here. It would be a less suitable measure, e.g., if there were additional costs associated with beginning a new inexact Newton step, such as initializing a new preconditioner. The fourth column gives numbers of backtracks over all test cases, i.e., numbers of step-reductions in the while-loop in Algorithm INB. The fifth column gives numbers of instances of convergence to a “wrong” solution, i.e., convergence to a solution other than the everywhere-positive solution in the PDE problem of §3.2.1 or to a solution other than  $u \equiv 1$  in the integral equation problem of §3.2.5. As noted previously, convergence to a “wrong” solution illustrates the potentially serious effects of disagreement between  $F$  and its local linear model. The sixth column gives the number of failures over all test cases. If failure occurred in a test case, then that case was not included in the statistics for columns 1–5; consequently, those statistics are not fully comparable to those for which all runs were successful.

One sees from Table 3.1 and 3.2 that the best overall performances were from Choice 1 and from Choice 2 with  $\gamma = .9$  and  $\gamma = 1$ . Taking  $\gamma = .5$  in Choice 2 resulted in significantly less efficiency with  $\alpha = 2$ ; in addition, it led to increased numbers of



TABLE 3.1

Summary test results over all problems. GMLI, GMINS, and GMFEE are geometric means of the numbers of linear iterations, inexact Newton steps, and “function evaluation equivalents”, respectively. NB, NW, and NFAIL are the total numbers of backtracks, instances of convergence to a “wrong” solution, and failures, respectively. Results marked “\*” were over successful runs only.

$\eta_k$ choice	GMLI	GMINS	GMFEE	NB	NW	NFAIL
$10^{-1}$	65.5*	12.00*	82.3*	2*	1*	1
$10^{-4}$	90.2*	7.21*	103.3*	1*	0*	2
$1/2^{k+1}$	70.3*	9.24*	85.4*	6*	1*	1
$\min\{1/(k+2), \ F(x_k)\ _2\}$	72.2	8.72	86.5	18	2	0
Choice 1	51.7	9.14	65.3	5	0	0
Choice 2, $\alpha = 2, \gamma = 1$	51.8	8.38	64.3	6	0	0
Choice 2, $\alpha = 2, \gamma = .9$	52.5	7.89	64.7	8	0	0
Choice 2, $\alpha = 2, \gamma = .5$	66.8	7.93	79.4	13	1	0
Choice 2, $\alpha = \frac{1+\sqrt{5}}{2}, \gamma = 1$	50.0	9.05	63.2	4	0	0
Choice 2, $\alpha = \frac{1+\sqrt{5}}{2}, \gamma = .9$	51.5	8.91	64.9	6	0	0
Choice 2, $\alpha = \frac{1+\sqrt{5}}{2}, \gamma = .5$	59.4*	7.67*	70.9*	4*	1*	1

backtracks with  $\alpha = 2$  and to one failure and one instance of convergence to a “wrong” solution with  $\alpha = (1 + \sqrt{5})/2$ , which suggest less robustness when  $\gamma$  is as small as .5. The other choices included in the tests were notably less effective.

Among Choice 1 and Choice 2 with  $\gamma = .9$  and  $\gamma = 1$ , Choice 2 with  $\gamma = 1$  and  $\alpha = (1 + \sqrt{5})/2$  placed first in every category except mean numbers of inexact Newton steps; thus this choice might be judged the winner. However, its margin of superiority was slight: for example, in “function evaluation equivalents”, the best and worst means for these choices differ by less than 4% over all problems and by less than 5% over the PDE problems. Furthermore, there was considerable variance in the relative performance and ranking of these choices among the individual test cases.

The results for Choice 2 illustrate that more aggressive choices of the forcing terms, i.e., choices that are smaller or result in faster asymptotic convergence, may decrease the number of inexact Newton steps up to a point but, through oversolving, may also lead to more linear iterations, more backtracking, and less robustness. Less aggressive choices, on the other hand, may reduce the number of linear iterations up to a point and improve robustness but may also result in increased numbers of inexact Newton steps.

**4. Summary discussion.** We have outlined forcing term choices that result in desirably fast local convergence and also tend to avoid oversolving the Newton equation, i.e., imposing an accuracy on an approximation of the Newton step that leads to significant disagreement between  $F$  and its local linear model. The choices, along with theoretical support and practical safeguards, are given in §2. Practical performance on a representative set of test problems is discussed in §3.

Choice 1 directly reflects the agreement between  $F$  and its local linear model at the previous step. It results in a certain  $q$ -superlinear local convergence; see Theorem 2.2 and the following remark for precise statements. Choice 2 does not directly reflect the agreement between  $F$  and its local linear model; however, it performed effectively in our tests. Furthermore, it can give up to  $q$ -quadratic local convergence

TABLE 3.2

Summary test results over the PDE problems. GMLI, GMINS, and GMFEE are geometric means of the numbers of linear iterations, inexact Newton steps, and “function evaluation equivalents”, respectively. NB, NW, and NFAIL are the total numbers of backtracks, instances of convergence to a “wrong” solution, and failures, respectively. Results marked “\*” were over successful runs only.

$\eta_k$ choice	GMLI	GMINS	GMFEE	NB	NW	NFAIL
$10^{-1}$	102.2*	11.89*	117.8*	0*	0*	1
$10^{-4}$	152.4*	6.68*	163.7*	1*	0*	1
$1/2^{k+1}$	104.2*	8.95*	118.4*	3*	0*	1
$\min\{1/(k+2), \ F(x_k)\ _2\}$	117.6	8.22	130.3	15	1	0
Choice 1	83.5	8.94	96.4	3	0	0
Choice 2, $\alpha = 2, \gamma = 1$	81.7	8.18	93.8	4	0	0
Choice 2, $\alpha = 2, \gamma = .9$	83.3	7.57	95.2	6	0	0
Choice 2, $\alpha = 2, \gamma = .5$	98.4	7.57	110.4	10	0	0
Choice 2, $\alpha = \frac{1+\sqrt{5}}{2}, \gamma = 1$	79.6	8.80	91.9	2	0	0
Choice 2, $\alpha = \frac{1+\sqrt{5}}{2}, \gamma = .9$	83.0	8.70	95.9	4	0	0
Choice 2, $\alpha = \frac{1+\sqrt{5}}{2}, \gamma = .5$	91.9*	6.98*	101.3*	0*	0*	1

(see Theorem 2.3), and the parameters  $\alpha$  and  $\gamma$  appearing in it allow flexibility that may be useful in applications.

The best performances in our tests were from Choice 1 and from Choice 2 with  $\gamma = .9$  and  $\gamma = 1$ . (With Choice 2, the values  $\alpha = 2$  and  $\alpha = (1 + \sqrt{5})/2$  were used in the tests. The latter value was chosen to give convergence roughly comparable to that for Choice 1.) Of these choices, Choice 2 with  $\gamma = 1$  and  $\alpha = (1 + \sqrt{5})/2$  could be considered most effective in these tests, but only by a small margin; any of these choices might be best for a particular application.

The numerical results in §3 illustrate that, in a globalized Newton iterative or truncated Newton method such as the implementation of Algorithm INB used here, oversolving resulting from inappropriately small forcing terms not only may incur unnecessary expense in solving the Newton equation but also may place significant demands on the globalization and even cause it to fail. In addition, unless special care is taken, very small forcing terms may call for more residual reduction than the iterative linear solver can accurately obtain, especially when finite differences are used to approximate products of  $F'$  with vectors. Conversely, choosing larger forcing terms may reduce oversolving and avoid inaccuracy in the iterative linear solver but increase the number of the inexact Newton steps required for convergence.

## REFERENCES

- [1] P. BJØRSTAD, *Fast numerical solution of the biharmonic Dirichlet problem on rectangles*, SIAM J. Numer. Anal., 20 (1983), pp. 59–71.
- [2] P. N. BROWN AND Y. SAAD, *Hybrid Krylov methods for nonlinear systems of equations*, SIAM J. Sci. Stat. Comput., 11 (1990), pp. 450–481.
- [3] X.-C. CAI, W. D. GROPP, D. E. KEYES, AND M. D. TIDRIRI, *Newton–Krylov–Schwarz methods in CFD*, in Proceedings of the International Workshop on the Navier–Stokes Equations, R. Rannacher, ed., Notes in Numerical Fluid Mechanics, Braunschweig, 1994 (to appear), Vieweg Verlag.
- [4] R. S. DEMBO, S. C. EISENSTAT, AND T. STEihaug, *Inexact Newton methods*, SIAM J. Numer.

- Anal., 19 (1982), pp. 400–408.
- [5] R. S. DEMBO AND T. STEihaug, *Truncated Newton algorithms for large-scale optimization*, Math. Prog., 26 (1983), pp. 190–212.
  - [6] J. E. DENNIS, JR. AND R. B. SCHINABEL, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Series in Automatic Computation, Prentice-Hall, Englewood Cliffs, NJ, 1983.
  - [7] S. C. EISENSTAT AND H. F. WALKER, *Globally convergent inexact Newton methods*, SIAM J. Optimization, 4 (1994), pp. 393–422.
  - [8] R. GLOWINSKI, H. B. KELLER, AND L. REINHART, *Continuation-conjugate gradient methods for the least squares solution of nonlinear boundary value problems*, SIAM J. Sci. Stat. Comput., 6 (1985), pp. 793–832.
  - [9] C. T. KELLEY AND J. I. NORTHRUP, *A pointwise quasi-Newton method for integral equations*, SIAM J. Numer. Anal., 25 (1988), pp. 1138–1155.
  - [10] P. J. MCKENNA, 1992. Private communication.
  - [11] J. J. MORÉ, *A collection of nonlinear model problems*, in Computational Solution of Nonlinear Systems of Equations, E. L. Allgower and K. Georg, eds., Lectures in Applied Mathematics Vol. 26, American Mathematical Society, 1990, pp. 723–762.
  - [12] Y. SAAD AND M. H. SCHULTZ, *GMRES: a generalized minimal residual method for solving nonsymmetric linear systems*, SIAM J. Sci. Stat. Comput., 7 (1986), pp. 856–869.
  - [13] R. SCHAAF, 1994. Private communication.
  - [14] J. STOER AND R. BULIRSCH, *Introduction to Numerical Analysis*, Springer Verlag, New York, 1980.
  - [15] P. N. SWARTZTRAUBER AND R. A. SWEET, *Algorithm 541: Efficient Fortran subprograms for the solution of separable elliptic partial differential equations*, ACM Trans. Math. Software, 5 (1979), pp. 352–364.
  - [16] K. TURNER AND H. F. WALKER, *Efficient high accuracy solutions with GMRES(m)*, SIAM J. Sci. Stat. Comput., 13 (1992), pp. 815–825.
  - [17] C. J. VAN DUJIN AND J. M. DE GRAAF, *Large time behaviour of solutions of the porous medium equation with convection*, J. Differential Equations, 84 (1990), pp. 183–203.