

QUASI-NEWTON METHODS,  
MOTIVATION AND THEORY

J.E. Dennis and J. Moré<sup>†</sup>

TR 74-217

November 1974

Department of Computer Science  
Cornell University  
Ithaca, New York 14853

<sup>†</sup> This research was supported in part by the National Science Foundation under Grant GJ-27528



Quasi-Newton Methods,  
Motivation and Theory

J.E. Dennis and J.J. More'

Department of Computer Science  
Cornell University  
Ithaca, N.Y.

Abstract:

This paper is an attempt to motivate and justify quasi-Newton methods as useful modifications of Newton's method for general and gradient nonlinear systems of equations. References are given to ample numerical justification; here we give an overview of many of the important theoretical results and each is accompanied by sufficient discussion to make the results and hence the methods plausible.

Key Words and Phrases: unconstrained minimization; nonlinear simultaneous equations; update methods; quasi-Newton methods.



## 1. INTRODUCTION

Nonlinear problems in finite dimensions are generally solved by iteration. Davidon (1959), for the minimization problem, and Broyden (1965), for systems of equations, introduced new methods which although iterative in nature, were quite unlike any others in use at the time. These papers together with the very important modification and clarification of Davidon's work by Fletcher and Powell (1963) have sparked a large amount of research in the late sixties and early seventies. This work has led to a new class of algorithms which have been called by the names quasi-Newton, variable metric, variance, secant, update, or modification methods. Whatever one calls them (we will use quasi-Newton), they have proved themselves in dealing with practical problems of the two types mentioned; that is, systems of  $n$  equations in  $n$  unknowns, and the unconstrained minimization of functionals.

A predictable consequence of this research is that there has been a proliferation of quasi-Newton methods for unconstrained minimization. Moreover, the derivation and relationship between these methods has usually been obscured by appealing to certain idealized situations such as exact line searches and quadratic functionals. This has not happened in nonlinear equations since the only quasi-Newton method that has been seriously used is the one proposed by Broyden (1965).

In this paper we show that it is possible to derive all of the known practical quasi-Newton methods from very natural considerations and in such a way that the relationship between these methods is clear. In addition, this paper contains a survey of the theoretical results which yield insight into the behavior of quasi-Newton methods, and in order to motivate these methods, there is also some background material in Sections 2 and 6. In either case, we have only given those proofs which are either new, give insight, or are simpler than those previously published, but references are always given.

In Sections 4 and 7 we derive the various quasi-Newton updates. This is done by taking the point of view that these updates are methods for generating approximations to derivatives -- Jacobians in nonlinear equations and Hessians in unconstrained minimization. This point of view suggests how to use quasi-Newton methods in other areas such as least squares and constrained optimization.

The theoretical results are contained in Sections 5 and 8. These results show, in particular, that there are four quasi-Newton updates which are globally and superlinearly convergent for linear problems (even in the absence of orthogonality assumptions or exact line searches), and locally and superlinearly convergent for nonlinear problems. These updates are Broyden's 1965 update for nonlinear equations, Powell's symmetric form of Broyden's update, the Davidon-Fletcher-Powell update, and the Broyden-Fletcher-Goldfarb-Shanno update. The theoretical results quoted tend to explain why these four updates are the ones most used in practical work.

In addition to the above material there are some rate of convergence results in Section 3. In particular, we emphasize superlinear convergence and its geometric interpretation.

We use  $R^n$  to denote  $n$ -dimensional real Euclidean space with the usual inner product  $\langle x, y \rangle = x^T y$  while  $L(R^n)$  is the linear space of all real matrices of order  $n$ . Moreover,  $||\cdot||$  stands for either the  $\ell_2$  vector norm  $||x|| = \langle x, x \rangle^{1/2}$ , or for any matrix norm which is consistent (or subordinate) to the  $\ell_2$  vector norm in the sense that  $||Ax|| \leq ||A|| ||x||$  for each  $x$  in  $R^n$  and  $A$  in  $L(R^n)$ . In particular, the  $\ell_2$  operator norm and the Frobenius norm are consistent with the  $\ell_2$  vector norm. For future reference we note that the Frobenius norm can be computed by

$$(1.1) \quad ||A||_F^2 = \sum_{i=1}^n ||Av_i||^2 = \text{trace}(A^T A)$$

where  $\{v_1, \dots, v_n\}$  is any orthonormal set in  $R^n$ , and that for any pair  $A, B$  in  $L(R^n)$ ,

$$(1.2) \quad ||AB||_F \leq \min\{||A||_2 ||B||_F, ||A||_F ||B||_2\}$$

In addition to the above matrix norms, we also make use of the weighted norms

$$(1.3) \quad ||A||_{M,2} = ||MAM||_2, \quad ||A||_{M,F} = ||MAM||_F$$

where  $M$  is a nonsingular symmetric matrix in  $L(R^n)$ . These norms do not satisfy the sub-multiplicative property  $||AB|| \leq ||A|| ||B||$  which is usually satisfied by matrix norms, but are very useful because they can be used to measure the relative error of approximations to symmetric and positive definite matrices. To be

specific, suppose that  $A$  is symmetric and positive definite, and let  $A^{-1/2}$  be the symmetric positive definite square root of  $A^{-1}$ . Since

$$\frac{||B-A||}{||A||} \leq ||A^{-1/2}(B-A)A^{-1/2}||$$

for either the  $\ell_2$  operator norm or the Frobenius norm, it is clear that if  $M = A^{-1/2}$  then  $||B-A||_{M,2}$  and  $||B-A||_{M,F}$  measure the relative error of  $B$  as an approximation to  $A$  in the  $\ell_2$  and Frobenius norms, respectively.

## 2. VARIATIONS ON NEWTON'S METHOD FOR NONLINEAR EQUATIONS

Let  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a mapping with domain and range in  $\mathbb{R}^n$  and consider the problem of finding a solution to the system of  $n$  equations in  $n$  unknowns given by

$$f_i(x_1, \dots, x_n) = 0, \quad 1 \leq i \leq n,$$

where  $f_1, \dots, f_n$  are the component functions of  $F$ .

The best known method for attacking this problem is Newton's method, but sometimes it is modified so as to improve its computational efficiency. In this section we examine some of these variations and their corresponding advantages and disadvantages. This will help to motivate the introduction of quasi-Newton methods as variations of Newton's method.

For the purpose of analyzing the algorithms for solving  $F(x) = 0$ , the mapping  $F$  is assumed to have the following properties.



- (a) The mapping  $F$  is continuously differentiable in an open convex set  $D$ .
- (2.1) (b) There is an  $x^*$  in  $D$  such that  $F(x^*) = 0$  and  $F'(x^*)$  is nonsingular.

The notation  $F'(x)$  denotes the Jacobian matrix  $(\partial_j f_i(x))$  evaluated at  $x$  so that (2.1) guarantees that  $x^*$  is a locally unique solution to the equations  $F(x) = 0$ .

In addition to (2.1) sometimes we will need the stronger requirement that  $F'$  satisfies a Lipschitz condition at  $x^*$ : There is a constant  $\kappa$  such that

$$(2.2) \quad \|F'(x) - F'(x^*)\| \leq \kappa \|x - x^*\|, \quad x \in D.$$

Note that if  $D$  is sufficiently small then (2.2) is satisfied if, for example,  $F$  is twice differentiable at  $x^*$ .

Newton's method for nonlinear equations can be derived by assuming that we have an approximation  $x_k$  to  $x^*$  and that in a neighborhood of  $x_k$  the linear mapping

$$L_k(x) = F(x_k) + F'(x_k)(x - x_k)$$

is a good approximation to  $F$ . If this is the case, then a presumably better approximation  $x_{k+1}$  to  $x^*$  can be obtained by solving the linear system  $L_k(x) = 0$ . Thus Newton's method proceeds from an initial approximation  $x_0$  to  $x^*$ , and attempts to successively improve  $x_0$  by the iteration

$$x_{k+1} = x_k - F'(x_k)^{-1} F(x_k), \quad k = 0, 1, \dots$$

Actually, this is the form of Newton's method which is convenient for analysis. The computational form consists of carrying out the following steps for  $k = 0, 1, \dots, m$  where  $m$  is the maximum number of iterations allowed.

- (a) Compute  $F(x_k)$  and if  $x_k$  is acceptable, stop.  
Otherwise, compute  $F'(x_k)$ .  
(2.3) (b) Solve the linear system  $F'(x_k)s_k = -F(x_k)$  for  $s_k$  and  
set  $x_{k+1} = x_k + s_k$ .

The advantages of this algorithm are summarized in the following well-known result.

Theorem 2.1: Let  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  satisfy assumptions (2.1). Then there is an open set  $S$  which contains  $x^*$  such that for any  $x_0 \in S$  the Newton iterates are well-defined, remain in  $S$  and converge to  $x^*$ . Moreover, there is a sequence  $\{\alpha_k\}$  which converges to zero and with

$$(2.4) \quad ||x_{k+1} - x^*|| \leq \alpha_k ||x_k - x^*||, \quad k = 0, \dots$$

If, in addition,  $F$  satisfies (2.2) then there is a constant  $\beta$  such that

$$(2.5) \quad ||x_{k+1} - x^*|| \leq \beta ||x_k - x^*||^2, \quad k = 0, \dots$$

For a proof of this result see, for example, Ortega and Rheinboldt (1970), page 312. However, in Section 5 we will show that if  $F$  satisfies (2.2) then the convergence of Newton's method follows from a much more general result. Moreover, (2.4) and (2.5) will follow from results in Section 3.

Two advantages of Newton's method are expressed by Theorem 2.1. The first one is the existence of a domain of attraction  $S$  for Newton's method. The existence of this domain of attraction implies that if the Newton iterates ever land in  $S$ , then they will remain in  $S$  and eventually converge to  $x^*$ . This insures some measure of stability for the iteration.

The other advantage is expressed by (2.4) and is known as superlinear convergence. Moreover, if (2.2) holds then Theorem 2.1 shows that we obtain (at least) second order or quadratic convergence; that is, (2.5) holds. However, the example

$$f(x) = x + |x|^{1+\alpha}, \quad \alpha \in (0,1)$$

shows that in general (2.5) does not hold. If  $\beta ||x^*||$  is not too large, then an informal interpretation of (2.5) is that eventually each iteration doubles the number of significant digits in  $x_k$  as an approximation to  $x^*$ .

The best known disadvantage of Newton's method is that a particular problem may require a very good initial approximation to  $x^*$  if the iteration is to converge. This is due to the fact that the set  $S$  in Theorem 2.1 can be very small. To overcome this disadvantage, special techniques (e.g. Powell's (1970a)) are needed.

On the other hand, for many problems the most important disadvantage of Newton's method is the requirement that  $F'(x_k)$  be determined for each  $k$ . This involves the evaluation of  $n^2$  scalar functions at each step and for most functions this is a very costly operation. It is usually taken to be equivalent to  $n$  evaluations of  $F$ , but the exact cost varies from problem to

problem. If the Jacobian is relatively easy to obtain, then Newton's method is very attractive. If obtaining the Jacobian is relatively expensive, then this problem can be circumvented in some cases by using a finite difference approximation to the Jacobian matrix.

For example,  $F'(x_k)$  could be replaced in (2.3) by the computation of  $A(x_k, h_k) \in L(R^n)$  where

$$(2.6) \quad [A(x, h)]_{i,j} = [f_i(x + \eta_j e_j) - f_i(x)] / \eta_j,$$

and  $h = (\eta_1, \dots, \eta_n)$  is some suitably chosen vector. Of course, we now solve the system

$$(2.7) \quad A(x_k, h_k) s_k = -F(x_k)$$

for  $s_k$ .

There is a significant amount of theoretical and computational support for this approach. For example, if  $F$  satisfies assumptions (2.1) and (2.2), and at each iteration  $\|h_k\| \leq \gamma \|F(x_k)\|$  for some constant  $\gamma$  then all the conclusions of Theorem 2.1 also hold for the finite-difference Newton's method. However the expense of computing  $n^2$  scalar functions still remains. A popular technique for trying to reduce the overall computational effort of the Newton or the finite-difference Newton's method is to hold the Jacobian fixed for a given number of iterations. This is particularly useful when the Jacobian is not changing very rapidly. However, it is always difficult to decide how long the Jacobian should be held fixed. Brent (1973) has shown that although this

technique decreases the rate of convergence, it can increase a certain measure of efficiency.

Finally, note that all the modifications of Newton's method mentioned in this section require the solution of a system of linear equations and therefore  $O(n^3)$  arithmetic operations per iteration. For some problems, the solution of these linear systems is the most expensive part of the iteration, and in these cases one should consider holding the Jacobian matrix fixed for a given number of iterations since in each such iteration this expense would be reduced to  $O(n^2)$ .

### 3. RATES OF CONVERGENCE

It is very important to understand something about the rate of convergence of different algorithms, since to a certain extent the rate of convergence of a method is as important as the fact that it converges; if it converges very slowly we may never be able to see it converge. Therefore, in this section we will outline certain results which give insight into rates of convergence. In particular we emphasize the notion of superlinear convergence and its geometrical interpretation.

A reasonable algorithm should at least be linearly convergent in the sense that if  $\{x_k\}$  is generated by the algorithm and  $\{x_k\}$  converges to  $x^*$ , then for some norm  $||\cdot||$  there is an  $\alpha \in (0,1)$  and  $k_0 \geq 0$  such that

$$||x_{k+1} - x^*|| \leq \alpha ||x_k - x^*||, \quad k \geq k_0.$$

This guarantees that eventually the error will be decreased by the factor  $\alpha < 1$ .

To be competitive an algorithm should be superlinearly convergent in the sense that (2.4) holds for some sequence  $\{\alpha_k\}$  which converges to zero. As noted by Dennis and Moré (1974) one of the properties of superlinearly convergent methods is that

$$(3.1) \quad \lim_{k \rightarrow +\infty} \frac{\|x_{k+1} - x_k\|}{\|x_k - x^*\|} = 1$$

provided, of course, that  $x_k \neq x^*$  for  $k \geq 0$ . The importance of (3.1) stems from the fact that it provides a very convenient stopping criterion. That (3.1) holds is quite easy to prove and follows from (2.4) and the fact that

$$\left| \|x_{k+1} - x_k\| - \|x_k - x^*\| \right| \leq \|x_{k+1} - x^*\|$$

The following result of Dennis and Moré (1974) shows precisely when an iteration is superlinearly convergent.

Theorem 3.1: Let  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  satisfy assumptions (2.1), and let  $\{B_k\}$  in  $L(\mathbb{R}^n)$  be a sequence of nonsingular matrices. Suppose that for some  $x_0$  in  $D$  the sequence

$$(3.2) \quad x_{k+1} = x_k - B_k^{-1} F(x_k), \quad k = 0, 1, \dots,$$

remains in  $D$ ,  $x_k \neq x^*$  for  $k \geq 0$ , and converges to  $x^*$ . Then  $\{x_k\}$  converges superlinearly to  $x^*$  if and only if

$$(3.3) \quad \lim_{k \rightarrow +\infty} \frac{\|[B_k - F'(x^*)](x_{k+1} - x_k)\|}{\|x_{k+1} - x_k\|} = 0$$

Clearly, if  $\{B_k\}$  converges to  $F'(x^*)$  then (3.3) holds and thus Theorem 3.1 explains why Newton's method and the finite-difference Newton's method with  $\|h_k\| = O(\|F(x_k)\|)$  converges superlinearly. However, (3.3)

only requires that  $\{B_k\}$  converge to  $F'(x^*)$  along the directions  $s_k = x_{k+1} - x_k$  of the iterative method. As pointed out in Sections 5 and 8, this is the case for certain quasi-Newton methods, and yet for these methods  $\{B_k\}$  does not, in general, converge to  $F'(x^*)$ .

An equivalent but more geometric formulation of (3.3) is that it requires  $s_k = x_{k+1} - x_k$  in the iterative method to asymptotically approach the Newton correction  $s_k^N = -F'(x_k)^{-1}F(x_k)$  in both length and direction. To see this note that since  $F(x_k) = -B_k s_k$ ,

$$s_k - s_k^N = s_k + F'(x_k)^{-1}F(x_k) = F'(x_k)^{-1}[F'(x_k) - B_k]s_k,$$

and thus (3.3) is equivalent with

$$(3.4) \quad \lim_{k \rightarrow \infty} \frac{\|s_k - s_k^N\|}{\|s_k\|} = 0,$$

Equation (3.4) shows that the relative error of  $s_k$  as an approximation to  $s_k^N$  approaches zero, and it is fairly easy to prove that this is equivalent to requiring that  $s_k$  approach  $s_k^N$  in both length and direction. For future reference, we state this formally.

Lemma 3.2: Let  $u, v$  belong to  $R^n$  with  $\langle u, v \rangle \neq 0$  and let  $\alpha \in (0, 1)$ . If  $\|u - v\| \leq \alpha \|u\|$  then  $\langle u, v \rangle$  is positive and

$$(3.5) \quad \left| 1 - \frac{\|v\|}{\|u\|} \right| \leq \alpha, \quad 1 - \left( \frac{\langle u, v \rangle}{\|u\| \|v\|} \right)^2 \leq \alpha^2.$$

Conversely, if  $\langle u, v \rangle$  is positive and (3.5) holds then

$$\|u - v\| \leq 3\alpha \|u\|.$$

Proof: Assume first that  $\|u - v\| \leq \alpha \|u\|$ . Then

$$\left| \frac{\|u\| - \|v\|}{\|u\|} \right| \leq \frac{\|u - v\|}{\|u\|} \leq \alpha,$$

and thus the first part of (3.5) holds. For the second part let  $\omega = \langle u, v \rangle / (||u|| ||v||)$  and note that

$$||u-v||^2 = ||u||^2 - 2||u|| ||v|| \omega + ||v||^2 \geq ||u||^2 (1 - \omega^2).$$

This proves (3.5). Now note that if  $\omega \leq 0$  then the equality above shows that  $||u-v|| \geq ||u||$ . Hence,  $\alpha < 1$  implies that  $\langle u, v \rangle$  is positive. For the converse note that

$$||u-v||^2 = (||u|| - ||v||)^2 + 2(1-\omega)||u|| ||v|| \leq \alpha^2 ||u||^2 [1+2(1+\alpha)]$$

and since  $\alpha < 1$ , it certainly follows that  $||u-v|| \leq 3\alpha ||u||$  as desired.

Lemma 3.2 shows that (3.4) is equivalent to

$$\lim_{k \rightarrow \infty} \frac{||s_k^N||}{||s_k||} = \lim_{k \rightarrow \infty} \left\langle \frac{s_k}{||s_k||}, \frac{s_k^N}{||s_k^N||} \right\rangle = 1,$$

and thus an iterative method is superlinearly convergent if and only if its directions asymptotically approach the Newton direction in both length and direction.

We would also like to explore second order convergence and for this we need the following estimate.

Lemma 3.3: Let  $F: R^n \rightarrow R^n$  satisfy assumptions (2.1) and (2.2).

Then

$$(3.6) \quad ||F(v) - F(u) - F'(x^*)(v-u)|| \leq \kappa \max\{||v-x^*||, ||u-x^*||\} ||v-u||$$

for all  $v$  and  $u$  in  $D$ .



The proof of this result follows immediately from Theorem 3.2.5 of Ortega and Rheinboldt (1970); note that the assumption  $F(x^*) = 0$  is not necessary for Lemma 3.3 nor is the invertibility of  $F'(x^*)$ .

Using Lemma 3.3 it is not difficult to modify the proof of Theorem 3.1 as given by Dennis and Moré (1974) and show that if the assumptions of Theorem 3.1 are satisfied and (2.2) holds then there is a constant  $\mu_1$  such that

$$||x_{k+1} - x^*|| \leq \mu_1 ||x_k - x^*||^p, \quad k = 0, 1, \dots$$

for some  $p \in (1, 2]$  if and only if there is a constant  $\mu_2$  such that

$$||[B_k - F'(x^*)](x_{k+1} - x_k)|| \leq \mu_2 ||x_{k+1} - x_k||^p, \quad k = 0, 1, \dots$$

However, we have not found any use for this result. The following well-known result is much easier to prove and is apparently just as useful.

Theorem 3.4: Let  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  satisfy assumptions (2.1) and (2.2), and let  $\{B_k\}$  be a sequence of nonsingular matrices. Assume that for some  $x_0$  in  $D$  the sequence (3.2) remains in  $D$  and converges to  $x^*$ . If

$$(3.7) \quad ||B_k - F'(x^*)|| \leq \eta ||x_k - x^*||, \quad k = 0, 1, \dots,$$

then  $\{x_k\}$  converges quadratically to  $x^*$ .

Proof: Since  $\{x_k\}$  converges to  $x^*$ , inequality (3.7) and the Banach Lemma imply that there is a constant  $\gamma$  such that

$$||B_k^{-1}|| \leq \gamma \text{ for } k \text{ sufficiently large. Since}$$

$$x_{k+1} - x^* = -B_k^{-1} \left\{ [F(x_k) - F(x^*) - F'(x^*)(x_k - x^*)] + (F'(x^*) - B_k)(x_k - x^*) \right\},$$

Lemma 3.3 together with (3.7) show that

$$||x_{k+1} - x^*|| \leq \gamma\{\kappa ||x_k - x^*||^2 + \eta ||x_k - x^*||^2\} ,$$

and it follows that  $\{x_k\}$  converges quadratically to  $x^*$ .

The most natural way to guarantee that (3.7) holds is to require that

$$(3.8) \quad ||B_k - F'(x_k)|| \leq \eta_1 ||F(x_k)|| , \quad k \geq 0$$

If this is the case then

$$||B_k - F'(x^*)|| \leq \eta_1 ||F(x_k)|| + \kappa ||x_k - x^*|| ,$$

and Lemma 3.3 implies that (3.7) holds. Note that Newton's method and the finite-difference Newton's method with  $||h_k|| = O(||F(x_k)||)$  satisfy (3.8).

#### 4. BROYDEN'S METHOD

In Section 2 we saw that two disadvantages of Newton's method were its need for  $n^2 + n$  scalar function evaluations and its use of  $O(n^3)$  arithmetic operations at each iteration. We will now derive Broyden's method (1965) and show how it effects an order of magnitude reduction in each of these expenses. The price paid is a reduction from second order to superlinear convergence.

The idea behind Broyden's 1965 proposal is that it is a method for approximating Jacobian matrices. As pointed out in Section 2, one of the major expenses of Newton's method is the calculation of  $F'(x_k)$ ; let us now show how Broyden derived an approximation  $B_k$  to  $F'(x_k)$  such that  $B_{k+1}$  can be obtained from  $B_k$  in  $O(n^2)$  arithmetic operations per iteration and

evaluating  $F$  at only  $x_k$  and  $x_{k+1}$ .

To derive his method, assume that  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is continuously differentiable in an open convex set  $D$  and that for given  $x$  in  $D$  and  $s \neq 0$ , the vector  $\bar{x} = x + s$  belongs to  $D$ . You should associate  $x$  with  $x_k$  and  $\bar{x}$  with  $x_{k+1}$ , so that what we want is a good approximation to  $F'(\bar{x})$ .

Since  $F'$  is continuous at  $\bar{x}$ , given  $\epsilon > 0$  there is a  $\delta > 0$  such that

$$||F(x) - F(\bar{x}) - F'(\bar{x})(x - \bar{x})|| \leq \epsilon ||x - \bar{x}||$$

provided  $||\bar{x} - x|| < \delta$ . It follows that

$$F(x) \approx F(\bar{x}) + F'(\bar{x})(x - \bar{x})$$

the degree of approximation increasing as  $||x - \bar{x}||$  decreases.

Hence, if  $\bar{B}$  is to denote our approximation to  $F'(\bar{x})$ , it seems reasonable to require that  $\bar{B}$  satisfy the equation

$$F(x) = F(\bar{x}) + \bar{B}(x - \bar{x}).$$

This is generally written

$$(4.1) \quad \bar{B}s = y \equiv F(\bar{x}) - F(x)$$

where  $s = \bar{x} - x$ .

In the case of  $n = 1$ , equation (4.1) completely determines  $\bar{B}$  and the secant method would result from using this approximate derivative in a Newton-like iteration. For  $n > 1$ , we can still argue that the only new information about  $F$  has been gained in the direction determined by  $s$ . Now suppose we had an approximation  $B$  to  $F'(x)$ . Broyden reasoned that there really is no justification

for having  $\bar{B}$  differ from  $B$  on the orthogonal complement of  $s$ . This can be expressed as the requirement

$$(4.2) \quad \bar{B}z = Bz \quad \text{if} \quad \langle z, s \rangle = 0.$$

Clearly (4.1) and (4.2) uniquely determine  $\bar{B}$  from  $B$  and in fact

$$(4.3) \quad \bar{B} = B + \frac{(y - Bs)s^T}{\langle s, s \rangle}.$$

Equation (4.1) is central to the development of quasi-Newton methods, and therefore it has often been called the quasi-Newton equation. In fact, it also plays a role in a second derivation of Broyden's update.

The second derivation again starts from the assumption that any matrix that satisfies the quasi-Newton equation (4.1) is a good candidate for  $\bar{B}$ . However, now it is argued that out of all the matrices that satisfy the quasi-Newton equation,  $\bar{B}$  should be the closest to  $B$ . The next result establishes that this matrix is again given by (4.3) if "closest" is measured by the Frobenius norm.

Theorem 4.1: Given  $B \in L(R^n)$ ,  $y \in R^n$  and some nonzero  $s \in R^n$ , define  $\bar{B}$  by (4.3). Then  $\bar{B}$  is the unique solution to the problem

$$\min\{\|\hat{B} - B\|_F : \hat{B}s = y\}.$$

Proof: To show that  $\bar{B}$  is a solution note that if  $y = \hat{B}s$  then

$$\|\bar{B} - B\|_F = \|(\hat{B} - B) \frac{ss^T}{\langle s, s \rangle}\|_F \leq \|\hat{B} - B\|_F$$

That  $\bar{B}$  is the unique solution follows from the fact that the

mapping  $f: L(R^n) \rightarrow R$  defined by  $f(A) = \|B - A\|_F$  is strictly convex in  $L(R^n)$  and that the set of  $\hat{B} \in L(R^n)$  such that  $\hat{B}s = y$  is convex.

By now it should be clear how (4,3) can be used in an iterative method. For example, in its most basic form Broyden's method is defined by:

$$(4.4) \quad x_{k+1} = x_k - B_k^{-1} F(x_k), \quad k = 0, 1, \dots$$

where the matrices  $B_k \in L(R^n)$  are generated by

$$(4.5) \quad B_{k+1} = B_k + \frac{(y_k - B_k s_k) s_k^T}{\langle s_k, s_k \rangle}, \quad k = 0, 1, \dots$$

with

$$(4.6) \quad y_k = F(x_{k+1}) - F(x_k), \text{ and } s_k = x_{k+1} - x_k.$$

As it stands, it is clear that given  $x_0$  and  $B_0$ , Broyden's method can be carried out with  $n$  scalar function evaluations per iteration. However, (4.4) and (4.5) seem to indicate that the solution of the linear system  $B_k s_k = -F(x_k)$  is required. One way to overcome this difficulty requires the following result which is due to Sherman and Morrison (1948).

Lemma 4.2: Let  $u, v \in R^n$  and assume that  $A \in L(R^n)$  is nonsingular. Then  $A + uv^T$  is nonsingular if and only if  $\sigma = 1 + \langle v, A^{-1}u \rangle \neq 0$ . If  $\sigma \neq 0$  then

$$(4.7) \quad (A + uv^T)^{-1} = A^{-1} - (1/\sigma) A^{-1} uv^T A^{-1}$$

Proof: That  $A + uv^T$  is nonsingular if and only if  $\sigma \neq 0$

follows from Lemma 4.4 which will be proved later. It is easy to verify (4.7) because if the matrix on the right hand side is multiplied by  $A + uv^T$  then the result is the identity matrix.

From Lemma 4.2 it follows that if  $H_k = B_k^{-1}$ , then  $H_{k+1} = B_{k+1}^{-1}$  is defined by

$$(4.8) \quad H_{k+1} = H_k + \frac{(s_k - H_k y_k) s_k^T H_k}{\langle s_k, H_k y_k \rangle}$$

provided  $\langle s_k, H_k y_k \rangle \neq 0$ . Therefore, Broyden's method can also be implemented as

$$x_{k+1} = x_k - H_k F(x_k)$$

where  $\{H_k\}$  is generated by (4.8), and in this form Broyden's method only requires  $n$  scalar function evaluations and  $O(n^2)$  arithmetic operations per iteration.

It is also possible to implement (4.5) and use only  $O(n^2)$  arithmetic operations per iteration. For example, Gill and Murray (1972) describe a method by which if  $B_k = Q_k R_k$  where  $Q_k$  is orthogonal and  $R_k$  is upper triangular, then the corresponding factorization of  $B_{k+1}$  can be obtained in  $O(n^2)$  operations. Of course, if  $B_k = Q_k R_k$  is given then the solution of the linear system  $B_k s_k = -F(x_k)$  only involves  $O(n^2)$  operations. One reason why this approach would be preferable over (4.8) is because in (4.5) there are no matrix-vector multiplications; the term  $B_k s_k$  is just  $-F(x_k)$ . Another reason is that the analysis of Section 5 shows that (4.5) is more stable.

Note that we don't need to choose  $s_k = x_{k+1} - x_k$  in either (4.5) or (4.8). It is entirely reasonable to choose  $s_k$  to be any vector such that  $F$  is defined at  $x_k + s_k$  and then set  $y_k = F(x_k + s_k) - F(x_k)$ . For example, if we set  $s_k = \eta e^j$  for some scalar  $\eta$ , then (4.5) shows that  $B_{k+1}$  only differs from  $B_k$  in the  $j$ -th column, and that this column is now

$$[F(x + \eta e^j) - F(x)]/\eta$$

Of course, if  $s_k \neq x_{k+1} - x_k$  then each iteration requires two function evaluations instead of one.

As theoretical justification for his method, Broyden only offered the fact that for affine functions it is norm-reducing with respect to the  $\ell_2$  operator norm. The following result shows that a slightly stronger result holds in the Frobenius norm.

Theorem 4.3: Let  $A \in L(R^n)$  satisfy  $y = As$  for some nonzero

$s \in \mathbb{R}^n$  and  $y \in \mathbb{R}^n$ . Moreover, given  $B \in L(\mathbb{R}^n)$  define  $\bar{B}$  by (4.3).

Then

$$||\bar{B} - A||_F \leq ||B - A||_F$$

with equality if and only if  $y = Bs$ .

Proof: Let  $\{s/||s||, v_2, \dots, v_n\}$  be an orthonormal set. Since  $(\bar{B} - A)s = 0$ ,

$$||\bar{B} - A||_F^2 = \sum_{i=2}^n ||(\bar{B} - A)v_i||^2.$$

Moreover,  $\bar{B}v_i = Bv_i$  for  $2 \leq i \leq n$  and therefore,

$$||\bar{B} - A||_F^2 = ||B - A||_F^2 - (|| (B - A)s || / ||s||)^2.$$

The result follows from this relationship.

If  $\{x_k\}$  is any sequence, and  $s_k, y_k$  are defined by (4.6), then  $y_k = As_k$  for

$$A = \int_0^1 F'(x_k + \theta s_k) d\theta.$$

Thus, Theorem 4.3 guarantees that in the Frobenius norm,  $B_{k+1}$  is a better approximation than  $B_k$  to the average of  $F'$  on the line segment from  $x_k$  to  $x_{k+1}$ . Of course, if  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is affine then  $A$  is the coefficient matrix, and therefore for affine functions Broyden's method is norm-reducing in the Frobenius norm.

To conclude this section we point out that Broyden's method is sometimes implemented in the form

$$(4.9) \quad \bar{B} = B + \theta \frac{(y - Bs) s^T}{\langle s, s \rangle}$$



where  $\theta$  is chosen so as to avoid singularity in  $\bar{B}$ . The following result can be used to decide how to choose  $\theta$ .

Lemma 4.4: Let  $v, w$  in  $R^n$  be given. Then

$$(4.10) \quad \det(I + vw^T) = 1 + \langle v, w \rangle.$$

Proof: Let  $P = I + vw^T$  and assume that  $v \neq 0$  for otherwise the result is trivial. Then any eigenvector of  $P$  is either orthogonal to  $w$  or a multiple of  $v$ . If the eigenvector is orthogonal to  $w$  then the eigenvalue is unity while if it is parallel to  $v$  then the eigenvalue is  $1 + \langle v, w \rangle$ . Equation (4.10) follows.

## 5. LOCAL CONVERGENCE RESULTS

We now would like to present a local convergence result that is available for Broyden's method and some of its variations. The importance of this result lies in the fact that the techniques used in its proof are applicable to other methods and in particular, to the double-rank updates of Section 7.

In this analysis it is assumed that  $x_0$  and  $B_0$  are sufficiently close to  $x^*$  and  $F'(x^*)$ , respectively, where  $F$  satisfies assumptions (2.1) and (2.2). The convergence follows from a very general theorem due to Broyden, Dennis, and Moré (1973). This result was developed to extend, to other quasi-Newton methods, the analysis given by Dennis (1971) for Broyden's method.

To describe the algorithms that this result handles, we will need the concept of an update function. If  $F: R^n \rightarrow R^n$  is defined on a set  $D$ , an update function  $U$  for  $F$  on  $D$  is a set-valued mapping from  $D \times D_M$  into  $D_M \subset L(R^n)$ . Thus if the domain of  $U$  is denoted by  $\text{dom}U$  then  $U(x, B)$  is a nonempty subset of  $D_M$ .

for each  $(x, B) \in \text{dom}U$ .

Update functions are only a means to denote the various Jacobian approximations which might be used in an iterative process;  $D_M$  and  $\text{dom}U$  depend on the particular algorithm. For the iteration

$$(5.1) \quad \begin{aligned} x_{k+1} &= x_k - B_k^{-1} F(x_k) \\ B_{k+1} &\in U(x_k, B_k), \quad k = 0, 1, \dots, \end{aligned}$$

it is convenient to define  $\text{dom}U$  as the set of all  $(x, B)$  in  $D \times D_M$  such that  $B$  is nonsingular and  $\bar{x} = x - B^{-1}F(x)$  belongs to  $D$  and differs from  $x$ . Moreover, usually  $D_M = L(R^n)$ . Therefore, (5.1) is well-defined if  $(x_k, B_k) \in \text{dom}U$  for  $k = 0, 1, \dots$ . Of course, if  $\bar{x} = x$  then  $F(x) = 0$  and the algorithm stops.

To illustrate these concepts note that for Newton's method  $U(x, B) = \{F'(\bar{x})\}$ , while for Broyden's method  $U(x, B) = \{\bar{B}\}$  where  $\bar{B}$  is defined by (4.3) with  $y = F(\bar{x}) - F(x)$  and  $s = \bar{x} - x$ . Also note that the finite difference form of Newton's method defined by (2.4) and (2.5) can be described by  $U(x, B) = \{A(x, h) : 0 \leq \|h\| \leq \gamma \|F(x)\|\}$  where  $\gamma$  is a fixed non-negative constant. This description has the advantage of not requiring a precise specification of the choice of  $h$ . Another illustration of the ease of description furnished by update functions is the following. Let  $U$  be given and for  $(x, B) \in \text{dom}U$ , set  $\hat{U}(x, B) = U(x, B) \cup \{B\}$ . Then  $\hat{U}$  defines the modification to (5.1) in which  $B_k$  is not necessarily changed at each iteration.

Update functions also apply to the minimization algorithms

of Sections 6 and 7. These algorithms are of the form (5.1), at least in a neighborhood of a local minimizer, where  $U$  is an update function for a gradient mapping. In this case  $D_M$  is usually the set of all symmetric matrices in  $L(R^n)$ .

For the following result recall that if  $U$  is an update function for (5.1) then  $\text{dom}U$  is the set of all  $(x, B)$  in  $D \times D_M$  such that  $B$  is nonsingular and  $\bar{x} = x - B^{-1}F(x)$  belongs to  $D$  and differs from  $x$ .

**Theorem 5.1:** Let  $F: R^n \rightarrow R^n$  satisfy assumptions (2.1) and (2.2), and let  $U$  be an update function for  $F$  such that for all  $(x, B) \in \text{dom}U$  and  $\bar{B} \in U(x, B)$ ,

$$(5.2) \quad ||\bar{B} - F'(x^*)|| \leq [1 + \alpha_1 \sigma(x, \bar{x})] ||B - F'(x^*)|| + \alpha_2 \sigma(x, \bar{x})$$

for some constants  $\alpha_1$  and  $\alpha_2$  where  $\bar{x} = x - B^{-1}F(x)$  and

$$(5.3) \quad \sigma(x, \bar{x}) = \max\{||\bar{x} - x^*||, ||x - x^*||\}.$$

Then there are positive constants  $\epsilon$  and  $\delta$  such that if  $x_0 \in D$  and  $B_0 \in D_M$  satisfy  $||x_0 - x^*|| < \epsilon$  and  $||B_0 - F'(x^*)|| < \delta$  then iteration (5.1) is well-defined and converges linearly to  $x^*$ .

By definition iteration (5.1) is locally convergent at  $x^*$  if there is an  $\epsilon > 0$  and a  $\delta > 0$  such that whenever  $x_0 \in D$  and  $B_0 \in D_M$  satisfy  $||x_0 - x^*|| < \epsilon$  and  $||B_0 - F'(x^*)|| < \delta$  then  $\{x_k\}$  is well-defined and converges to  $x^*$ . Thus Theorem 5.1 guarantees the local and linear convergence of (5.1). Note that local convergence depends on  $D_M$  but since  $D_M$  is usually  $L(R^n)$  or the set of symmetric matrices,  $D_M$  is large enough to make Theorem 5.1 meaningful. Also note that since all matrix norms are equivalent there is no restriction on the matrix norm (5.2).

Now obviously Theorem 5.1 cannot guarantee better than linear convergence since the stationary iteration  $U(x, B) = \{B\}$  satisfies (5.2) with  $\alpha_1 = \alpha_2 = 0$ . The usual procedure is to use this theorem to prove the existence and convergence of  $\{x_k\}$  and then apply Theorem 3.1 or Theorem 3.4 to make a more precise statement about the rate of convergence. We illustrate this below.

If  $F$  satisfies (2.2) then for Newton's method,  $U(x, B) = \{F'(\bar{x})\}$  satisfies (5.2) with  $\alpha_1 = 0$ ,  $\alpha_2 = \kappa$  and  $D_M = L(R^n)$ . This proves the local convergence of Newton's method. The quadratic convergence follows from Theorem 3.4.

The proof of Theorem 2.1 that we have just given generalizes quite readily to the finite difference Newton's method defined by  $x_{k+1} = x_k + s_k$  where  $s_k$  satisfies (2.6) and (2.7) with  $\|h\| \leq \gamma \|F(x)\|$  for some constant  $\gamma$ . We now turn to the application of Theorem 5.1 to Broyden's method.

Theorem 5.2: Let  $F: R^n \rightarrow R^n$  satisfy assumptions (2.1) and (2.2), and consider Broyden's method as defined by equations (4.4), (4.5) and (4.6). Then Broyden's method is locally and superlinearly convergent at  $x^*$ .

Proof: We will prove that Broyden's method is locally convergent at  $x^*$  by showing that (5.2) is satisfied with  $D_M = L(R^n)$ . For this note that (4.5) implies that

$$(5.4) \quad \bar{B} - F'(x^*) = [B - F'(x^*)] \left[ I - \frac{ss^T}{\langle s, s \rangle} \right] + \frac{(y - F'(x^*)s)s^T}{\langle s, s \rangle}$$

In particular,

$$\|\bar{B} - F'(x^*)\| \leq \|B - F'(x^*)\| + \frac{\|y - F'(x^*)s\|}{\|s\|}$$

where the matrix norm is either the  $\ell_2$  operator norm or the Frobenius norm. Therefore Lemma 3.3 implies that (5.2) is satisfied with  $\alpha_1 = 0$  and  $\alpha_2 = \kappa$ . This proves the linear convergence of Broyden's method.

Like Newton's method, the more precise rate of convergence requires further work.. In fact, we will show that (3.3) holds. For this, note that direct computation using  $||A||_F^2 = \text{trace}(A^T A)$  shows that

$$(5.5) \quad ||E[I - \frac{ss^T}{\langle s, s \rangle}]||_F^2 = ||E||_F^2 - (\frac{||Es||}{||s||})^2$$

for any  $E \in L(R^n)$ , and since  $(\alpha^2 - \beta^2)^{1/2} \leq \alpha - (2\alpha)^{-1}\beta^2$ ,

$$||E[I - \frac{ss^T}{\langle s, s \rangle}]||_F \leq ||E||_F - (2||E||_F)^{-1} (\frac{||Es||}{||s||})^2.$$

Now define  $\eta_k = ||B_k - F'(x^*)||_F$  and use the above inequality and Lemma 3.3 in (5.4) to obtain that

$$\eta_{k+1} \leq [1 - (2\eta_k^2)^{-1}\psi_k^2]\eta_k + \kappa\sigma_k$$

where  $\sigma_k = \max\{||x_{k+1} - x^*||, ||x_k - x^*||\}$  and

$$\psi_k = \frac{||[B_k - F'(x^*)]s_k||}{||s_k||}.$$

Since  $\eta_{k+1} \leq \eta_k + \kappa\sigma_k$  and  $\{x_k\}$  is linearly convergent, it follows that  $\{\eta_k\}$  is bounded, and if  $\eta$  is an upper bound then

$$(2\eta)^{-1}\psi_k^2 \leq \eta_k - \eta_{k+1} + \kappa\sigma_k.$$

Thus

$$(2\eta)^{-1} \sum_{k=0}^{\infty} \psi_k^2 \leq \eta_0 + \kappa \sum_{k=0}^{\infty} \sigma_k,$$

forcing  $\{\psi_k\}$  to converge to zero. Hence, (3.3) holds and this concludes the proof.

There are several interesting points about the proof of Theorem 5.2. The first is that although (3.3) holds, it does not necessarily follow that  $\{B_k\}$  converges to  $F'(x^*)$ .

Example 5.3: Let  $F: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be defined by  $x = (\xi_1, \xi_2)^T$  and  $F(x) = (\xi_1, \xi_2 + \xi_2^3)^T$ , and consider Broyden's method with  $x_0 = (0, \varepsilon)^T$  and

$$B_0 = \begin{pmatrix} 1 + \delta & 0 \\ 0 & 1 \end{pmatrix}$$

It is easy to verify that the (1,1) element of  $B_k$  is always  $1 + \delta$  and thus  $\{B_k\}$  does not converge to  $F'(x^*)$ .

Another point of interest about this proof is that it generalizes to the modification of Broyden's method given by (4.9). Thus Moré and Trangenstein (1974) prove that a parameter  $\theta_k$  can be chosen so that if (4.5) is replaced by

$$(5.6) \quad (a) \quad B_{k+1} = B_k + \theta_k \frac{(y_k - B_k s_k) s_k^T}{\langle s_k, s_k \rangle}$$

$$(b) \quad B_{k+1} \text{ nonsingular, } |\theta_k - 1| \leq \hat{\theta} \quad \text{and} \quad \hat{\theta} \in (0, 1),$$

then Theorem 5.2 holds. They also noted that if  $F$  is affine then for this modification the  $\varepsilon$  and  $\delta$  in Theorem 5.1 are infinite.

Theorem 5.4: Let  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  be defined by  $F(x) = Ax - b$  where  $A \in L(\mathbb{R}^n)$  is nonsingular and  $b \in \mathbb{R}^n$ , and consider Broyden's

method as defined by equations (4.4), (4.6) and (5.6). Then Broyden's method is globally and superlinearly convergent to  $A^{-1}b$ .

There are other variations of Broyden's method for which Theorem 5.2 holds. For example, if we decide that  $s_k = x_{k+1} - x_k$  is not a suitable direction, we can use (4.4), (4.5) but replace (4.6) by  $y_k = F(x_k + s_k) - F(x_k)$  where  $s_k$  is any nonzero vector such that

$$||s_k|| \leq \eta \max\{||x_{k+1} - x^*||, ||x_k - x^*||\}$$

for some constant  $\eta$ . For example, the choice  $s_k = ||F(x_{k+1})||e_j$  is suitable for each  $j$ . Of course, if  $s_k \neq x_{k+1} - x_k$  then the computation of  $y_k$  involves two evaluations of  $F$ .

There is a variation of Broyden's method which is of interest in the case that  $F'(x)$  is sparse. In this variation equations (4.4), (4.5) and (4.6) are used to define  $B_{k+1}$  from  $B_k$  but before it is used,  $B_{k+1}$  is forced to have the same sparsity pattern as  $F'(x)$ . That Theorem 5.2 holds follows from the observation that forcing  $B_k$  to have the same sparsity pattern as  $F'(x)$  decreases  $||B_k - F'(x^*)||$ . Schubert (1970) has proposed an algorithm along these lines and Broyden (1971a) has shown that it is locally convergent. However, it is not known whether Schubert's algorithm is superlinearly convergent.

We conclude this section by discussing two important variations of Theorem 5.2. The following variation arises because for some algorithms it is more natural to think of them as generating approximations to the inverse of the Jacobian. In this case

$\text{dom}U$  will be the set of all  $(x, H)$  in  $D \times D_M$  such that  $\bar{x} = x - HF(x)$  belongs to  $D$  and differs from  $x$ .

Theorem 5.5: Let  $F: R^n \rightarrow R^n$  satisfy assumptions (2.1) and (2.2), and let  $U$  be an update function for  $F$  such that for all  $(x, H) \in \text{dom}U$  and  $\bar{H} \in U(x, H)$ ,

$$(5.7) \quad ||\bar{H} - F'(x^*)^{-1}|| \leq [1 + \alpha_1 \sigma(x, \bar{x})] ||H - F'(x^*)^{-1}|| + \alpha_2 \sigma(x, \bar{x})$$

for some constants  $\alpha_1$  and  $\alpha_2$  where  $\bar{x} = x - HF(x)$  and  $\sigma(x, \bar{x})$  is defined by (5.3). Then there are  $\epsilon > 0$  and  $\delta > 0$  such that if  $x_0 \in D$  and  $H_0 \in D_M$  satisfy  $||x_0 - x^*|| < \epsilon$  and  $||H_0 - F'(x^*)^{-1}|| < \delta$  then the iteration

$$(5.8) \quad \begin{aligned} x_{k+1} &= x_k - H_k F(x_k) \\ H_{k+1} &\in U(x_k, H_k), \quad k = 0, 1, \dots \end{aligned}$$

is well-defined and converges linearly to  $x^*$ .

The same remarks that we made after Theorem 5.2 for iteration (5.1) also apply, with suitable modifications, to (5.8). In particular, if (5.8) satisfies the conclusions of Theorem 5.5, then by definition (5.8) is locally and, of course, linearly convergent at  $x^*$ .

We also note that although Theorems 5.1 and 5.5 as well as their proofs are very similar, the two results are independent of each other. In fact, in Section 8 we will discuss two important algorithms and show that the local convergence of one of these algorithms follows from Theorem 5.1 while the other needs Theorem 5.5.



Finally we note that it is possible to generalize both these theorems by showing that the conclusions still hold if instead of (5.1) and (5.8) we consider the sequence

$$x_{k+1} = x_k - \lambda_k B_k^{-1} F(x_k) \equiv x_k - \lambda_k H_k F(x_k)$$

provided the sequence  $\{\lambda_k\}$  satisfies  $|\lambda_k - 1| < \hat{\lambda}$  for some  $\hat{\lambda} \in (0,1)$ .

## 6. VARIATIONS OF NEWTON'S METHOD FOR UNCONSTRAINED MINIMIZATION

Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be a functional defined on an open set  $D$  and consider the problem of finding a  $z$  in  $D$  such that  $f(z) \leq f(x)$  for each  $x$  in  $D$ . In this case  $z$  is a global minimizer of  $f$  and even if it is known to exist, finding it is usually an intractable task. Generally, one seeks  $z$  among the local minimizers of  $f$ ; that is, find  $x^*$  in  $D$  such that for some  $\delta > 0$ ,

$$(6.1) \quad f(x^*) \leq f(x), \quad ||x - x^*|| \leq \delta, \quad x \in D.$$

In this section we provide some background material and outline some of the methods that are used to solve (6.1). In particular, we stress the differences and analogies between the methods considered here and those in previous sections. This will help to motivate the introduction of quasi-Newton methods for unconstrained minimization.

We only consider the solution of (6.1) if  $f$  is differentiable. In this case (6.1) is usually attacked by trying to find a zero of  $\nabla f$  - the gradient of  $f$ . This approach is based on the fact that if  $x^*$  is a local minimizer of  $f$  in the open set  $D$  and  $f$  is

differentiable at  $x^*$  then  $\nabla f(x^*) = 0$ . Moreover, in this section only descent methods are considered.

A descent method for solving (6.1) generates for each iterate  $x_k$  a direction  $p_k$  of local descent in the sense that there is a  $\lambda_k$  such that  $f(x_k + \lambda p_k) < f(x_k)$  for  $\lambda \in (0, \lambda_k)$ . The next iterate is of the form  $x_{k+1} = x_k + \lambda_k p_k$  where the parameter  $\lambda_k$  is chosen so that  $f(x_{k+1}) < f(x_k)$ . The directions  $p_k$  and the parameters should be chosen in such a way that  $\{\nabla f(x_k)\}$  converges to zero. If  $\|\nabla f(x_k)\|$  is small then usually  $x_k$  is near a zero of  $\nabla f$  while the fact that  $\{f(x_k)\}$  is decreasing indicates that this zero of  $\nabla f$  is probably a local minimizer of  $f$ .

The simplest example of a descent method is the method of steepest descent. In this method we ask for the vector  $\hat{p}$  of unit length (in the  $\ell_2$  norm) such that for some  $\hat{\lambda} > 0$ ,

$$f(x + \lambda \hat{p}) < f(x + \lambda p), \quad \lambda \in (0, \hat{\lambda})$$

for all  $\|p\| = 1$ . It is not difficult to show that if  $\nabla f(x) \neq 0$  then  $\hat{p} = -\nabla f(x) / \|\nabla f(x)\|$ . Therefore, the method of steepest descent is given by

$$(6.2) \quad x_{k+1} = x_k - \lambda_k \nabla f(x_k), \quad k = 0, 1, \dots,$$

where the parameter  $\lambda_k$  is needed to guarantee that  $f(x_{k+1}) < f(x_k)$ ; that such a parameter exists is a consequence of the following simple result.

Lemma 6.1: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be defined in an open set  $D$  and differentiable at  $x$  in  $D$ . If  $\langle \nabla f(x), p \rangle < 0$  for some  $p$  in

$\mathbb{R}^n$  then there is a  $\lambda^* = \lambda^*(x, p)$  such that  $\lambda^* > 0$  and

$$f(x + \lambda p) < f(x), \quad \lambda \in (0, \lambda^*) .$$

The proof of this result is quite easy and follows from the fact that

$$\lim_{\lambda \rightarrow 0^+} [f(x + \lambda p) - f(x)]/\lambda = \langle \nabla f(x), p \rangle .$$

Lemma 6.1 guarantees, in particular, that the parameter  $\lambda_k$  in the steepest descent method can be chosen so that  $f(x_{k+1}) < f(x_k)$ . This is not sufficient to show that  $\{x_k\}$  gets close to a zero of  $\nabla f$  since  $\lambda_k$  may be arbitrarily small. In fact,  $\lambda_k > 0$  can be chosen so that  $\|x_{k+1} - x_k\| \leq \epsilon/2^k$  and therefore  $\{x_k\}$  converges to a point  $x^*$  with  $\|x_0 - x^*\| \leq 2\epsilon$ . If  $\nabla f(x_0) \neq 0$  and  $\nabla f$  is continuous at  $x_0$ , then  $\epsilon$  can be chosen so that  $\nabla f(x^*) \neq 0$ . At the end of this section we discuss a specific method for choosing  $\lambda_k$  which avoids this problem, and note that if  $\lambda_k$  is chosen appropriately, then the following result holds.

**Theorem 6.2:** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable and bounded below on  $\mathbb{R}^n$ , and assume that  $x_0$  is such that  $\nabla f$  is uniformly continuous on the level set

$$L(x_0) = \{x \in \mathbb{R}^n \mid f(x) \leq f(x_0)\} .$$

Then there is a sequence  $\{\lambda_k\}$  such that the steepest descent sequence (6.2) is well defined,  $\{f(x_k)\}$  is decreasing, and  $\{\nabla f(x_k)\}$  converges to zero.

If  $f$  is continuously differentiable on  $R^n$  and  $L(x_0)$  is compact, then the rest of the assumptions of Theorem 6.2 are automatically satisfied and in addition,  $f$  has a global minimizer and  $\{\nabla f(x_k)\}$  converges to zero. However, not even in this case is the steepest descent sequence guaranteed to converge to a local minimizer of  $f$ . An example reported by Wolfe (1971) shows that the steepest descent sequence may converge to a saddle point of  $f$ . Nevertheless, Theorem 6.2 is quite a strong convergence result. The fact that  $\{\nabla f(x_k)\}$  converges to zero implies that any limit point of  $\{x_k\}$  is a zero of  $\nabla f$  and that for any  $\epsilon > 0$  the stopping criterion  $\|\nabla f(x_k)\| < \epsilon$  will be satisfied in a finite number of steps. Unfortunately, steepest descent usually converges linearly.

The slow rate of convergence of steepest descent can be improved by switching to a faster method in a neighborhood of a zero of  $\nabla f$ . Since  $F = \nabla f$  is a mapping from  $R^n$  to  $R^n$ , any of the methods discussed in Sections 2 and 4 could be used. For example, if  $f$  is twice differentiable then Newton's method is given by

$$(6.3) \quad x_{k+1} = x_k - \nabla^2 f(x_k)^{-1} \nabla f(x_k), \quad k = 0, 1, \dots$$

where  $\nabla^2 f(x)$  is the Hessian matrix of  $f$  at  $x$ ; that is,  $\nabla^2 f(x)$  is just the Jacobian matrix of  $\nabla f$ . It should be clear that Theorem 2.1 applies to (6.3) with  $F = \nabla f$ , and that under the appropriate conditions we obtain local and quadratic convergence of (6.3) to a zero of  $\nabla f$ .

In view of the global convergence of steepest descent and the fast local convergence of Newton's method, it would be desirable to have a method that behaves like Newton's method near a local minimizer but like steepest descent far from a local minimizer. Most descent methods of this type are of the form

$$(6.4) \quad x_{k+1} = x_k - \lambda_k B_k^{-1} \nabla f(x_k), \quad k = 0, 1, \dots$$

where  $B_k$  is a symmetric, positive definite matrix which resembles  $\nabla^2 f(x_k)$ , at least in a neighborhood of a local minimizer.

As an example of such a method, Goldfeldt, Quandt and Trotter (1966) suggested the iteration

$$(6.5) \quad x_{k+1} = x_k - \lambda_k (\nabla^2 f(x_k) + \mu_k I)^{-1} \nabla f(x_k), \quad k = 0, 1, \dots$$

where the scalar  $\mu_k \geq 0$  is chosen so that  $\nabla^2 f(x_k) + \mu_k I$  is positive definite. To justify the claim that (6.5) behaves like Newton's method in a neighborhood of a local minimizer, recall that if  $f$  is differentiable in an open set  $D$  and twice differentiable at a local minimizer  $x^*$  of  $f$  in  $D$  then  $\nabla^2 f(x^*)$  is positive semidefinite. Therefore if  $x_k$  is in a neighborhood of a local minimizer, then very small values of  $\mu_k$  will suffice to make  $\nabla^2 f(x_k) + \mu_k I$  positive definite. Also note that if

$$s(\mu) = -(\nabla^2 f(x) + \mu I)^{-1} \nabla f(x),$$

then  $s(0)$  is the Newton direction while as  $\mu \rightarrow +\infty$  the angle between  $s(\mu)$  and  $-\nabla f(x)$  decreases monotonically to zero. Thus for large  $\mu$  iteration (6.5) behaves like steepest descent.

In order to preserve, in (6.5), the good local properties of Newton's method, one has to choose  $\mu_k$  and  $\lambda_k$  with some care. It is easy to see from Theorem 3.1 that as long as  $\{\mu_k\}$  and  $\{\lambda_k\}$  converge to zero and unity, respectively, iteration (6.5) is superlinearly convergent. Moreover, Theorem 3.4 shows that if  $\mu_k \leq \eta \|\nabla f(x_k)\|$  for some constant  $\eta$  and  $\lambda_k = 1$  for all sufficiently large  $k$ , then (6.5) converges quadratically. Unfortunately, these results do not indicate how to choose  $\{\mu_k\}$  globally, and in fact, this has turned out to be a hard problem.

There is a method of the form (6.4) which avoids the problem of choosing  $\mu_k$  in (6.5) and yet resembles (6.5). In this method we try to obtain a Cholesky decomposition of  $\nabla^2 f(x_k)$ ; that is, we try to find a nonsingular, lower triangular matrix  $L_k$  such that  $\nabla^2 f(x_k) = L_k L_k^T$ . Of course, if  $\nabla^2 f(x_k)$  is not positive definite then this decomposition does not even exist, but the idea is that as the decomposition proceeds it is possible to add to the diagonal of  $\nabla^2 f(x_k)$  and ensure that we obtain the Cholesky decomposition of a well-conditioned, positive definite matrix which differs from  $\nabla^2 f(x_k)$  in some minimal way. In particular, if  $\nabla^2 f(x_k)$  is a well-conditioned positive definite matrix then  $\nabla^2 f(x_k) = L_k^T L_k$ . The details are given by Murray (1972), page 64.

In the remainder of this section we describe some of the selection rules for  $\lambda_k$  which are used in methods of the form (6.4) and more generally, in any descent method of the form

$$(6.6) \quad x_{k+1} = x_k + \lambda_k p_k, \quad k = 0, 1, \dots$$

where  $\langle \nabla f(x_k), p_k \rangle < 0$ . The development of these particular rules are due to the initial work of Goldstein (1965) and Armijo (1966).

In a descent method  $\lambda_k$  should satisfy  $f(x_{k+1}) < f(x_k)$  but we have already noted that this requirement can be satisfied by arbitrarily small  $\lambda_k$  and then  $\{x_k\}$  may converge to a point at which  $\nabla f$  is not zero. A more reasonable requirement is that

$$(6.7) \quad f(x_k + \lambda_k p_k) \leq f(x_k) + \alpha \lambda_k \langle \nabla f(x_k), p_k \rangle, \quad \alpha \in (0, 1/2).$$

The reason for choosing  $\alpha < 1/2$  is that with this choice, Theorem 6.4 shows that if  $\{x_k\}$  converges to a local minimizer of  $f$  at which  $\nabla^2 f(x^*)$  is positive definite, and  $\{p_k\}$  converges to the Newton step  $-\nabla^2 f(x_k)^{-1} \nabla f(x_k)$  in both length and direction, then  $\lambda_k = 1$  will satisfy (6.7) for all sufficiently large  $k$ .

Since the right hand side of (6.7) is a straight line in  $\lambda$  which interpolates  $f(x_k + \lambda p_k)$  at  $\lambda = 0$  and whose slope is larger than  $\langle \nabla f(x_k), p_k \rangle$ , it is clear that there is a  $\lambda_k$  which satisfies (6.7). If  $\alpha$  is close to zero then (6.7) is not a very stringent requirement, and  $\alpha$  is generally chosen in this way with  $[10^{-4}, 10^{-1}]$  being the usual range. However, it is not a good idea to fix  $\lambda_k$  by just requiring that it satisfy (6.7) since, for instance,  $\lambda_k = 0$  is then admissible. In general, unreasonably small  $\lambda_k$  are ruled out by the numerical search procedure but theoretically we need to impose another requirement. One such requirement is that

$$(6.8) \quad \langle \nabla f(x_k + \lambda_k p_k), p_k \rangle \geq \beta \langle \nabla f(x_k), p_k \rangle, \quad \beta \in (\alpha, 1).$$

To show that there are  $\lambda_k$  which satisfy (6.7) and (6.8) assume that  $\bar{f}$  is defined on  $R^n$  and  $f(x_k + \lambda p_k)$  is bounded below for  $\lambda \geq 0$ . It is then geometrically obvious that there are  $\lambda_k > 0$  for which equality holds in (6.7). If  $\lambda_k$  is the first such  $\lambda_k$  then the mean value theorem implies that

$$f(x_k + \hat{\lambda}_k p_k) - f(x_k) = \langle \nabla f(x_k + \theta_k \hat{\lambda}_k p_k), p_k \rangle = \alpha \langle \nabla f(x_k), p_k \rangle$$

for some  $\theta_k \in (0,1)$ , and since  $\alpha < \beta$ ,

$$\langle \nabla f(x_k + \theta_k \hat{\lambda}_k p_k), p_k \rangle \geq \beta \langle \nabla f(x_k), p_k \rangle.$$

Thus  $\lambda_k = \theta_k \hat{\lambda}_k$  satisfies (6.7) and (6.8). However, we emphasize that a search routine for  $\lambda$  should not necessarily try to satisfy (6.7) and (6.8). In fact, the intervals which satisfy these two conditions can be quite small and therefore difficult to find. Moreover, to test whether or not (6.8) is satisfied requires the evaluation of  $\nabla f$ . Instead, the search routine should produce a  $\lambda_k$  which satisfies (6.7) and not be too small; (6.8) just guarantees that  $\lambda_k$  is not too small.

Theorem 6.3: Let  $f: R^n \rightarrow R$  satisfy the assumptions of Theorem 6.2, and consider an iteration of the form (6.6) where the search directions  $p_k$  satisfy  $\langle \nabla f(x_k), p_k \rangle < 0$ . Then there is a sequence  $\{\lambda_k\}$  which satisfies (6.7) and (6.8) and

$$(6.9) \quad \lim_{k \rightarrow +\infty} \langle \nabla f(x_k), \frac{p_k}{\|p_k\|} \rangle = 0.$$

Theorem 6.3 is due to Wolfe (1969) who also pointed out that for many iterations (6.9) implies that  $\{\|\nabla f(x_k)\|\}$  converges to zero; it is only necessary to verify that the angle between  $p_k$



and  $\nabla f(x_k)$  stays bounded away from ninety degrees. For example, if  $p_k = -\nabla f(x_k)$ , or more generally, if  $p_k = -B_k^{-1} \nabla f(x_k)$  where  $\{B_k\}$  is a sequence of symmetric, positive definite matrices with uniformly bounded condition numbers, then

$$-\langle \nabla f(x_k), \frac{p_k}{\|p_k\|} \rangle \geq \mu \|\nabla f(x_k)\|$$

where  $\mu$  is an upper bound on the condition number of  $B_k$ . Hence, (6.9) ensures that  $\{\|\nabla f(x_k)\|\}$  converges to zero.

To conclude this section we assume that the vectors  $p_k$  converge in direction and length to the Newton step and show that  $\lambda_k = 1$  will eventually satisfy (6.7) and (6.8).

**Theorem 6.4:** Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be twice continuously differentiable in an open set  $D$  and consider iteration (6.6) where  $\langle \nabla f(x_k), p_k \rangle < 0$  and  $\lambda_k$  is chosen to satisfy (6.7) and (6.8). If  $\{x_k\}$  converges to a point  $x^*$  in  $D$  at which  $\nabla^2 f(x^*)$  is positive definite and

$$(6.10) \quad \lim_{k \rightarrow \infty} \frac{\|\nabla f(x_k) + \nabla^2 f(x_k) p_k\|}{\|p_k\|} = 0,$$

then there is an index  $k_0 \geq 0$  such that  $\lambda_k = 1$  is admissible for  $k \geq k_0$ . Moreover,  $\nabla f(x^*) = 0$  and  $\{x_k\}$  converges super-linearly to  $x^*$ .

**Proof:** As a first step note that a consequence of (6.10) is that there is an  $\eta > 0$  such that

$$(6.11) \quad -\langle \nabla f(x_k), p_k \rangle \geq \eta \|p_k\|^2$$

for all  $k$  large enough. This follows since

$$-\langle \nabla f(x_k), p_k \rangle = \langle \nabla^2 f(x_k) p_k, p_k \rangle - \langle \nabla^2 f(x_k) p_k + \nabla f(x_k), p_k \rangle,$$

so that (6.11) follows from (6.10) and the fact that  $\nabla^2 f(x)$  is positive definite for all  $x$  close enough to  $x^*$ .

To show that (6.7) is eventually satisfied by  $\lambda_k = 1$  use the mean value theorem to obtain  $u_k$  in the line segment from  $x_k$  to  $x_k + p_k$  such that

$$f(x_k + p_k) - f(x_k) - 1/2 \langle \nabla f(x_k), p_k \rangle = 1/2 \langle \nabla^2 f(u_k) p_k, p_k \rangle + \langle \nabla f(x_k), p_k \rangle.$$

Now (6.9) and (6.11) show that  $\{p_k\}$  converges to zero; therefore (6.10) implies that for all  $k$  sufficiently large

$$(6.12) \quad f(x_k + p_k) - f(x_k) - 1/2 \langle \nabla f(x_k), p_k \rangle \leq (1/2 - \alpha) \eta \|p_k\|^2,$$

and thus (6.11) and (6.12) show that (6.7) is satisfied by  $\lambda_k = 1$ . To prove that (6.8) is also eventually satisfied by  $\lambda_k = 1$  we again use the mean value theorem to show that there is a  $v_k$  such that

$$\langle \nabla f(x_k + p_k), p_k \rangle = \langle \nabla f(x_k) + \nabla^2 f(v_k) p_k, p_k \rangle.$$

Thus (6.10) and (6.11) imply that for all  $k$  large enough ,

$$\langle \nabla f(x_k + p_k), p_k \rangle \leq \eta \beta \|p_k\|^2 \leq -\beta \langle \nabla f(x_k), p_k \rangle.$$

Hence  $\lambda_k = 1$  satisfies (6.8) and this concludes the first part of the proof. For the remainder, note that since  $\{p_k\}$  converges to zero, (6.10) shows that  $\nabla f(x^*) = 0$ . The superlinear convergence of  $\{x_k\}$  follows from Theorem 3.1.

## 7. QUASI-NEWTON METHODS FOR UNCONSTRAINED MINIMIZATION

The derivation of updates suitable for unconstrained optimization proceeds along lines similar to the development in Section 4. For nonlinear equations only Broyden's method appears to be satisfactory, but here some notable differences, motivated by the discussion in Section 6, will lead us to single out four reasonable update formulas.

One important consideration is the desire to have the quasi-Newton step  $-B_k^{-1} \nabla f(x_k)$  define a descent direction. In fact, the most widespread use of these methods is in conjunction with iterations of the form (6.4). In this context the update formula should generate a sequence of symmetric positive definite matrices  $\{B_k\}$  such that  $B_k$  resembles  $\nabla^2 f(x_k)$ , at least when  $x_k$  is near a local minimizer of  $f$ . We will examine these updates in Section 7.2.

In Section 7.1 we examine quasi-Newton methods which can be used to approximate the Hessian in such a way that the direction  $p_k = -B_k^{-1} \nabla f(x_k)$  resembles the true Newton direction. In this case  $p_k$  may not be a descent direction, so that the direction is usually further modified. For example, it may be modified by adding to  $B_k$  a suitable multiple of the identity matrix as in iteration (6.5).

It is also possible to look at the updates of Sections 7.1 and 7.2 from an "inverse" point of view in which we try to generate approximations to the inverse of the Hessian. It turns out that this gives rise to at least one other important update. These inverse updates and their relationship to the updates of Section 7.1

and 7.2 are examined in Section 7.3.

Throughout this section we assume  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  to be twice differentiable in the open convex set  $D$ , and that we have an approximation  $B$  to  $\nabla^2 f(x)$  for some  $x$  in  $D$ , and a derivation  $s$  such that  $x + s$  belongs to  $D$ . We now want to obtain a good approximation  $\bar{B}$  to  $\nabla^2 f(\bar{x})$  where  $\bar{x} = x + s$ .

### 7.1 Symmetry and the Quasi-Newton Equation

In view of the above discussion, and since the Hessian is symmetric, we want the update formula to have the property of hereditary symmetry; that is,

$$(7.1) \quad B \text{ symmetric implies } \bar{B} \text{ symmetric.}$$

Moreover, because of our desire to approximate the Hessian, arguments similar to those in Section 4 lead us to require that  $\bar{B}$  satisfy

$$(7.2) \quad \bar{B}s = y \equiv \nabla f(\bar{x}) - \nabla f(x),$$

Note that (7.2) is just the quasi-Newton equation (4.1) for  $F = \nabla f$ .

It is natural to ask whether it is possible to satisfy (7.1) and (7.2) with a rank one update formula. To see whether this can be done, first note that the general single-rank update that satisfies the quasi-Newton equation (7.2) is given by

$$(7.3) \quad \bar{B} = B + \frac{(y - Bs)c^T}{\langle c, s \rangle}$$

for  $c \in \mathbb{R}^n$  with  $\langle c, s \rangle \neq 0$ . If  $\bar{B}$  is to satisfy (7.1) then it is easy to show that

$$(7.4) \quad \bar{B} = B + \frac{(y - Bs)(y - Bs)^T}{\langle y - Bs, s \rangle}$$

is the only solution provided  $\langle y - Bs, s \rangle \neq 0$ . If  $y = Bs$  then  $\bar{B} = B$  is the solution while if  $y \neq Bs$  but  $\langle y - Bs, s \rangle = 0$  then there is no solution.

This update is known as the symmetric single-rank formula. It seems to have been first published by Broyden (1967) although several independent discoveries of (7.4) apparently occurred at about the same time. If  $H = B^{-1}$  and  $\bar{H} = \bar{B}^{-1}$  both exist and  $B$  is symmetric then the inverse relation

$$(7.5) \quad \bar{H} = H + \frac{(s - Hy)(s - Hy)^T}{\langle s - Hy, y \rangle}$$

holds. The following theorem, essentially due to Fiacco and McCormick (1968), shows that this method has very interesting behavior when it is applied to a quadratic functional.

Theorem 7.1: Let  $A \in L(R^n)$  be a nonsingular symmetric matrix, and set  $y_k = As_k$  for  $0 \leq k \leq m$  where  $\{s_0, \dots, s_m\}$  spans  $R^n$ . Let  $H_0$  be symmetric and for  $k = 0, \dots, m$  generate the matrices

$$(7.6) \quad H_{k+1} = H_k + \frac{(s_k - H_k y_k)(s_k - H_k y_k)^T}{\langle s_k - H_k y_k, y_k \rangle}$$

where it is assumed that

$$(7.7) \quad \langle s_k - H_k y_k, y_k \rangle \neq 0.$$

Then  $H_{m+1} = A^{-1}$ .

The proof of this result consists of verifying, by induction, that

$$H_k y_j = s_j, \quad 0 \leq j < k, \quad \text{for } k = 1, \dots, m+1.$$

Once this is done,

$$H_{m+1}y_j = H_{m+1}As_j = s_j, \quad 0 \leq j \leq m,$$

and the result follows from the assumption that  $\{s_0, \dots, s_m\}$  spans  $R^n$ .

The gist of Theorem 7.1 lies in the fact that if we have an iteration of the form  $x_{k+1} = x_k + s_k$  and (7.7) holds, then the use of (7.6) allows one to minimize a quadratic functional in a finite number of steps. Unfortunately, there is no guarantee that (7.7) will hold although it is not difficult to show that if  $A^{-1} - H_0$  is semidefinite (positive or negative) and if  $\{H_k\}$  is generated by (7.6) when (7.7) holds, and  $H_{k+1} = H_k$  otherwise, then  $H_{m+1} = A^{-1}$ .

The fact that the vectors  $s - Hy$  and  $y$  can be orthogonal forces a certain amount of numerical instability on the symmetric single-rank method. In particular, update (7.4) does not satisfy (5.2) or (5.7) and is therefore not of bounded deterioration. These difficulties have led to several improvements in the basic algorithm, and in its modified form the method has been quite successful. See, for example, the numerical results of Dixon (1972b).

The numerical difficulties with the symmetric single-rank method have led to a whole class of updates which satisfy (7.1) and (7.2). The technique used to derive this class is due to Powell (1970 d) who used it to obtain a double-rank version of Broyden's method. Dennis (1972) then showed that Powell's technique could be used to derive most of the well-known quasi-Newton updates.

In this derivation we begin with a symmetric  $B \in L(R^n)$  and consider

$$C_1 = B + \frac{(y - Bs)c^T}{\langle c, s \rangle}$$

as a possible candidate for  $\bar{B}$ . In general  $C_1$  is not symmetric, so consider

$$C_2 = (C_1 + C_1^T)/2$$

However, since  $C_2$  does not satisfy the quasi-Newton equation, we repeat the process. In this way a sequence  $\{C_k\}$  is generated by

$$(7.8) \quad C_{2k+1} = C_{2k} + \frac{(y - C_{2k}s)c^T}{\langle c, s \rangle},$$

$$C_{2k+2} = (C_{2k+1} + C_{2k+1}^T)/2, \quad k = 0, 1, \dots,$$

where  $C_0 = B$ . It turns out that  $\{C_k\}$  has a limit  $\bar{B}$  given by

$$(7.9) \quad \bar{B} = B + \frac{(y - Bs)c^T + c(y - Bs)^T}{\langle c, s \rangle} - \frac{\langle y - Bs, s \rangle}{\langle c, s \rangle^2} cc^T$$

and it is clear that this update satisfies (7.1) and (7.2).

Lemma 7.2: Let  $B \in L(R^n)$  be symmetric and let  $c, s$ , and  $y$  be in  $R^n$  with  $\langle c, s \rangle \neq 0$ . If the sequence  $\{C_k\}$  is defined by (7.8) with  $C_0 = B$  then  $\{C_k\}$  converges to  $\bar{B}$  as defined by (7.9).

Proof: We only need to prove that the sequence  $\{C_{2k}\}$  converges.

If  $G_k = C_{2k}$ , then (7.8) shows that

$$(7.10) \quad G_{k+1} = G_k + (1/2) [w_k C_k^T + C_k w_k^T] / \langle c, s \rangle$$

where  $w_k = y - G_k s$ . In particular,

$$w_{k+1} = P w_k, \quad P = (1/2) [I - \frac{cs^T}{\langle c, s \rangle}].$$

It is clear that  $P$  has one zero eigenvalue and all other eigenvalues equal to  $1/2$ , so that the Neumann Lemma (e.g. Ortega and Rheinboldt (1973), page 45) implies that

$$(7.11) \quad \sum_{k=0}^{\infty} w_k = \sum_{k=0}^{\infty} P^k (y - Bs) = (I - P)^{-1} (y - Bs)$$

Since

$$\lim_{k \rightarrow \infty} G_k = B + \sum_{k=0}^{\infty} (G_{k+1} - G_k),$$

it follows from (7.0) and (7.11) that  $\{G_k\}$  converges, and since

$$(I - P)^{-1} = 2[I - (1/2) \frac{cs^T}{\langle c, s \rangle}] ,$$

equations (7.10) and (7.11) also imply that the limit of  $\{G_k\}$  is  $\bar{B}$  as defined by (7.9).

Once  $c$  is chosen, (7.9) is a rank two update which satisfies (7.1) and (7.2). Before looking at special cases of (7.9), we show that this update solves a problem similar to the one specified in Theorem 3.1.

Theorem 7.3: Let  $B \in L(R^n)$  be symmetric, and let  $c, s$ , and  $y$  be in  $R^n$  with  $\langle c, s \rangle > 0$ . Assume that  $M \in L(R^n)$  is any nonsingular, symmetric matrix such that  $Mc = M^{-1}s$ . Then  $\bar{B}$  as defined by (7.9) is the unique solution to the problem

$$(7.12) \quad \min\{\|\hat{B} - B\|_{M,F} : \hat{B} \text{ symmetric}, \hat{B}s = y\}$$

where  $\|\cdot\|_{M,F}$  is defined by (1.3).

Proof: Let  $\hat{B}$  be any symmetric matrix such that  $y = \hat{B}s$ , and pre- and post-multiply (7.9) by  $M$ . If  $My = M^{-1}s = z$  it follows that



$$\bar{E} = \frac{Ezz^T + zz^TE}{\langle z, z \rangle} = \frac{\langle Ez, z \rangle}{\langle z, z \rangle^2} zz^T$$

where  $E = M(\hat{B} - B)M$  and  $\bar{E} = M(\bar{B} - B)M$ . Now it is clear that  $\|\bar{E}z\| = \|Ez\|$ , and that if  $v$  is orthogonal to  $z$  then  $\|\bar{E}v\| \leq \|Ev\|$ . Thus  $\|\bar{E}\|_F \leq \|E\|_F$  as desired. To show uniqueness just note that the mapping  $f: L(R^n) \rightarrow R$  defined by  $f(A) = \|B - A\|_{M,F}$  is strictly convex on the convex set of symmetric  $\hat{B}$  such that  $\hat{B}s = y$ .

Theorem 7.3 was inspired and is closely related to the results of Greenstadt (1970) and Goldfarb (1970) and it shows that the updates obtained by Greenstadt (1970) could also have been obtained by the symmetrization argument of Lemma 7.2.

Powell (1970d) used the argument of Lemma 7.2 to obtain formula (7.9) in the case  $c = s$ . Since in this case the underlying single-rank method is often called the Powell symmetric Broyden update, or the PSB update:

$$(7.13) \quad \bar{B}_{PSB} = B + \frac{(y - Bs)s^T + s(y - Bs)^T}{\langle s, s \rangle} - \frac{\langle y - Bs, s \rangle ss^T}{\langle s, s \rangle^2}$$

Theorem 7.3 implies that  $\bar{B}_{PSB}$  is the unique solution to the problem

$$\min\{\|\hat{B} - B\|_F : \hat{B} \text{ symmetric, } \hat{B}s = y\}$$

and this property is reminiscent of Theorem 4.1. In fact, it can be shown that if neither  $B$  nor  $\bar{B}$  are required to be symmetric then the unique solution to (7.12) is given by (7.13). Theorem 7.3 also leads us to believe that  $\bar{B}_{PSB}$  is a good approximation to the Hessian. To justify this claim note that (7.13) implies that for any symmetric  $A$  and  $B$  in  $L(R^n)$ ,

$$\bar{B}_{PSB} - A = P^T(B-A)P + [(y-As)s^T + s(y-As)^TP]/\langle s, s \rangle$$

with  $P = I - (ss^T/\langle s, s \rangle)$ . Therefore (1.2) shows that

$$||\bar{B}_{PSB} - A||_F \leq ||B - A||_F + 2 \frac{||y - As||}{||s||}$$

If  $A = \nabla^2 f(x)$  and  $\nabla^2 f$  is Lipschitz continuous (with constant  $k$ ) in the open convex set  $D$  then

$$||\bar{B}_{PSB} - \nabla^2 f(\bar{x})||_F \leq ||B - \nabla^2 f(x)||_F + 3k||s||$$

whenever  $x$  and  $\bar{x}$  lie in  $D$ . This relationship shows that the absolute error of  $B_k$  as an approximation to  $\nabla^2 f(x_k)$  grows linearly with  $||s_k||$ , and that this holds independent of the position of  $x$  in  $D$ .

## 7.2 Positive Definiteness

We now turn to updates which in addition to satisfying (7.1) and (7.2) generate positive definite matrices. For this, let us investigate the property of hereditary positive definiteness; that is,

$$(7.14) \quad B \text{ positive definite implies } \bar{B} \text{ positive definite.}$$

Note that if an update satisfies (7.2) and (7.14), then  $y = \bar{B}s$  and therefore  $\langle y, s \rangle > 0$  whenever  $B$  is positive definite. This imposes a restriction on the angle between  $y$  and  $s$ , which although not severe, must be kept in mind. In fact, if  $\langle \nabla f(x), s \rangle < 0$  then  $\langle y, s \rangle > 0$  is equivalent to the existence of a  $\beta \in (0, 1)$  such that  $\langle \nabla f(\bar{x}), s \rangle \geq \beta \langle \nabla f(x), s \rangle$ . For this reason the requirement (6.8) is very natural for quasi-Newton methods.

To investigate the property of hereditary positive definiteness, we need a result from the perturbation theory of symmetric matrices, e.g. Wilkinson (1965), pages 95-98:

Lemma 7.4: Let  $A \in L(R^n)$  be symmetric with eigenvalues

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n ,$$

and let  $A^* = A + \sigma uu^T$  for some  $u \in R^n$ . If  $\sigma \geq 0$  then  $A^*$  has eigenvalues  $\lambda_i^*$  such that

$$\lambda_1 \leq \lambda_1^* \leq \lambda_2 \leq \dots \leq \lambda_n \leq \lambda_n^* ,$$

while if  $\sigma \leq 0$  then the eigenvalues of  $A^*$  can be arranged so that

$$\lambda_1^* \leq \lambda_1 \leq \lambda_2^* \leq \dots \leq \lambda_n^* \leq \lambda_n .$$

Lemma 7.4 and the next two results will lead us to a choice of  $c$  in (7.9) which naturally satisfies (7.14). This development is a bit long, but it gives a lot of insight.

Theorem 7.5: Let  $B \in L(R^n)$  be symmetric and positive definite, and let  $c, s$ , and  $y$  be in  $R^n$  with  $\langle c, s \rangle \neq 0$ . Then  $\bar{B}$  as defined by (7.9) is positive definite if and only if  $\det \bar{B} > 0$ .

Proof: If  $\bar{B}$  is positive definite then clearly  $\det \bar{B} > 0$ . For the converse first note that we can write

$$\bar{B} = B + vw^T + wv^T$$

where  $w = c$  and

$$v = \frac{y - Bs}{\langle c, s \rangle} - (1/2) \frac{\langle y - Bs, s \rangle}{\langle c, s \rangle^2} c$$

Therefore,

$$\bar{B} = B + (1/2) [(v+w)(v+w)^T - (v-w)(v-w)^T] ,$$

and thus we have written  $\bar{B}$  as  $B$  plus the sum of two symmetric

rank one matrices. If  $B$  is positive definite then Lemma 7.4 implies that  $\bar{B}$  can have at most one non-positive eigenvalue. Therefore, if  $\det \bar{B} > 0$  then all the eigenvalues must be positive and thus  $\bar{B}$  is positive definite.

In view of Theorem 7.5, conditions (7.1) and (7.14) for the updates defined by (7.9) require that if  $B$  is symmetric and positive definite then  $\det \bar{B} > 0$ . To find out what choices of  $c$  satisfy this requirement we need an expression for  $\det \bar{B}$ .

Lemma 7.6: Let  $u_i \in \mathbb{R}^n$  for  $i = 1, 2, 3, 4$ . Then

$$\det(I + u_1 u_2^T + u_3 u_4^T) = (1 + \langle u_1, u_2 \rangle)(1 + \langle u_3, u_4 \rangle) - \langle u_1, u_4 \rangle \langle u_2, u_3 \rangle$$

Proof: A proof of this result can be found in Pearson (1969); the following is an alternate argument.

Assume for the moment that  $\langle u_1, u_2 \rangle \neq -1$ . Then  $I + u_1 u_2^T$  is nonsingular and

$$I + u_1 u_2^T + u_3 u_4^T = (I + u_1 u_2^T)(I + (I + u_1 u_2^T)^{-1} u_3 u_4^T).$$

The result now follows by using Lemmas 4.2 and 4.4. Since it holds for  $\langle u_1, u_2 \rangle \neq -1$ , a continuity argument shows that it holds in general.

Now apply Lemma 7.6 to (7.9). After some algebra it follows that

$$(7.15) \quad \det \bar{B} = \det B [(\langle c, Hy \rangle^2 - \langle c, Hc \rangle \langle y, Hy \rangle + \langle c, Hc \rangle \langle y, s \rangle) / \langle c, s \rangle^2]$$

where  $H = B^{-1}$ . If we assume that  $B$  is positive definite and let  $v = H^{1/2} y$  and  $w = H^{1/2} c$  then

$$(7.16) \quad \det \bar{B} = \det B [(\langle v, w \rangle^2 - \|v\|^2 \|w\|^2 + \|w\|^2 \langle y, s \rangle) / \langle c, s \rangle^2] ,$$

and Theorem 7.5 implies that  $\bar{B}$  is positive definite if and only if

$$(7.17) \quad \|v\|^2 \langle y, s \rangle > \|v\|^2 \|w\|^2 - \langle v, w \rangle^2 .$$

It is now apparent that the most natural way to satisfy (7.17) is to choose  $w$  to be a multiple of  $v$  so that (7.17) only requires that  $\langle y, s \rangle$  be positive. In this case  $c$  is a multiple of  $y$  and then (7.9) reduces to an update introduced by Davidon (1959), and later clarified and improved by Fletcher and Powell (1963). The DFP update is then given by

$$(7.18) \quad \bar{B}_{DFP} = B + \frac{(y - Bs)y^T + y(y - Bs)^T}{\langle y, s \rangle} - \frac{\langle y - Bs, s \rangle yy^T}{\langle y, s \rangle^2} \\ = (I - \frac{ys^T}{\langle y, s \rangle}) B (I - \frac{sy^T}{\langle y, s \rangle}) + \frac{yy^T}{\langle y, s \rangle}$$

Some of its properties are given in the following result, but first we note that the underlying single-rank formula (7.3) where  $c$  is a multiple of  $y$  is an update due to Pearson (1969).

Theorem 7.7: Let  $B \in L(R^n)$  be a nonsingular, symmetric matrix and define  $\bar{B}_{DFP} \in L(R^n)$  by (7.18) for any vectors  $y$  and  $s$  in  $R^n$  with  $\langle y, s \rangle \neq 0$ . Then  $\bar{B}_{DFP}$  is nonsingular if and only if  $\langle y, Hy \rangle \neq 0$  where  $H = B^{-1}$ . If  $\bar{B}_{DFP}$  is nonsingular then  $H_{DFP}^{-1} = \bar{B}_{DFP}$  can be expressed as

$$(7.19) \quad \bar{H}_{DFP} = H + \frac{ss^T}{\langle s, y \rangle} - \frac{Hy y^T H}{\langle y, Hy \rangle} .$$

Furthermore, if  $B$  is positive definite then  $\bar{B}_{DFP}$  is positive definite if and only if  $\langle y, s \rangle > 0$ .

Proof: Recall that for the DFP update  $c$  is a multiple of  $y$  so that (7.16) reduces to

$$(7.20) \quad \det \bar{B}_{DFP} = \det B \left[ \frac{\langle y, Hy \rangle}{\langle y, s \rangle} \right]$$

Thus  $\bar{B}_{DFP}$  is nonsingular if and only if  $\langle y, Hy \rangle \neq 0$ . To verify that  $\bar{H}_{DFP}$  is given by (7.19) one can either show that  $\bar{H}_{DFP} \bar{B}_{DFP} = I$  or one can use Lemma 4.2 twice on (7.18). In either case the proof is straightforward but tedious and is therefore omitted. Finally, assume that  $B$  is positive definite. If  $\langle y, s \rangle$  is positive then (7.20) shows that  $\det \bar{B}_{DFP}$  is also positive and thus Theorem 7.5 implies that  $\bar{B}_{DFP}$  is positive definite. Conversely, if  $\bar{B}_{DFP}$  is positive definite then

$$\langle y, s \rangle = \langle \bar{B}_{DFP} s, s \rangle > 0$$

which is the desired result.

One way to use the DFP update to generate a quasi-Newton direction and only use  $O(n^2)$  arithmetic operations per iteration would be to generate  $B_k^{-1} = H_k$  via equation (7.19). Another approach is based on the fact that if  $A$  is positive definite and  $A = LL^T$  for some lower triangular matrix, then the corresponding decomposition of

$$\bar{A} = A + \alpha z z^T$$

can be obtained in  $O(n^2)$  operations provided  $\bar{A}$  is positive definite. Methods for doing this are surveyed by Gill, Golub, Murray and Saunders (1974). That these techniques apply to (7.18)

follows from the proof of Theorem 7.5 which shows that (7.18) can be written as

$$\bar{B}_{DFP} = B + \alpha_1 z_1 z_1^T + \alpha_2 z_2 z_2^T$$

where  $\alpha_1 > 0$ ,  $\alpha_2 > 0$  and  $z_1, z_2$  are linear combinations of  $Bs$  and  $y$ . If the DFP update is used in a method of the form (6.4) then an advantage of the latter approach is that (7.18) requires no matrix-vector products.

Finally we remark that the matrices generated by the DFP formula are good approximations to the Hessian. In fact in Section 8 (see (8.16)) we will derive a general result which can be interpreted as follows: If  $\|s\|$  is small then the relative error (as measured in Section 1) of  $\bar{B}_{DFP}$  as an approximation to a positive definite  $\nabla^2 f(x)$  cannot be much larger than the corresponding error in  $B$ . Moreover the possible increase in this error is governed by a relative measure of how much  $f$  differs from a quadratic on  $D$ .

### 7.3 Inverse Updates

So far we have been thinking in terms of obtaining an approximation to the Hessian, but it is perhaps equally reasonable to try to obtain an approximation to the inverse of the Hessian. In particular, it should be clear that it is possible to use the techniques that we have been discussing to develop updating formulas for the inverse. These updates are sometimes called inverse updates while the updates developed in Sections 7.1 and 7.2 could be called direct updates.

To develop inverse updates, assume that we have an approximation  $H$  to  $\nabla^2 f(x)^{-1}$  and try to obtain a good approximation  $\bar{H}$  to  $\nabla^2 f(\bar{x})^{-1}$  where  $\bar{x} = x + s$ . For inverse updates the analogue of the quasi-Newton equation is

$$(7.21) \quad \bar{H}y = s ,$$

and therefore, the general single rank formula which satisfies (7.21) is

$$(7.22) \quad \bar{H} = H + \frac{(s - Hy)d^T}{\langle d, y \rangle}$$

for any  $d \in R^n$  with  $\langle d, y \rangle \neq 0$ .

It is important to realize the relationship between (7.3) and (7.22). If Lemma 4.2 is applied to (7.3) we obtain

$$\bar{B}^{-1} = B^{-1} + \frac{(s - B^{-1}y)c^T B^{-1}}{\langle c, B^{-1}y \rangle} .$$

Therefore, (7.3) and (7.22) represent the same update if  $c = B^T d$ .

Just as in Section 7.1, it is possible to study the property of hereditary symmetry, which for inverse updates is

$$(7.23) \quad H \text{ symmetric implies } \bar{H} \text{ symmetric.}$$

It is easy to verify that the only single rank formula which satisfies the quasi-Newton equation (7.21) and the hereditary symmetric property (7.23) is again given by the symmetric single rank formula (7.5).

To obtain other inverse updates which satisfy (7.21) and (7.23) we carry out the symmetrization argument of Lemma 7.2 on (7.22) to obtain



$$(7.24) \quad \bar{H} = H + \frac{(s - Hy)d^T + d(s - Hy)^T}{\langle d, y \rangle} - \frac{\langle s - Hy, y \rangle}{\langle d, y \rangle^2} dd^T$$

This result is due to Dennis (1972) who also noted that if  $\bar{B}$  and  $\bar{H}$  are defined by (7.9) and (7.24), respectively, then in general  $\bar{B}\bar{H} \neq I$  even if  $B$  is symmetric,  $BH = I$  and  $c = Bd$ . At first this is surprising because under these assumptions (7.3) and (7.22) represent the same update; however, in the argument of Lemma 7.3 we used the symmetrization operation  $(A + A^T)/2$ , and in general, the symmetrization and inversion operations do not commute.

It is also possible to prove an analogue of Theorem 7.3 for updates (7.24). In particular, if  $H$  is symmetric, then the unique solution to the problem

$$\min\{\|\hat{H} - H\|_F : \hat{H} \text{ symmetric, } \hat{H}y = s\}$$

is given by (7.24) with  $d = y$ . This update was proposed by Greenstadt (1970), but it has not received any more attention in the literature since it does not perform as well as the PSB update (7.13). It is interesting that the underlying single rank method was obtained by Broyden (1965), but that this update has also been neglected because of its poor numerical performance.

The most important instance of (7.24) was given by Broyden (1969), (1970), and independently by Fletcher (1970), Goldfarb (1970) and Shanno (1970). This update can be obtained by asking for the update of the general form (7.24) which "naturally" has the property of hereditary positive definiteness for inverse updates; that is,  $H$  positive definite implies  $\bar{H}$  positive definite. It should be

clear from the development in Section 7.2 that this update corresponds to choosing  $d = s$  in (7.24) and therefore the Broyden-Fletcher-Goldfarb-Shanno or BFGS update can be written in the form

$$(7.25) \quad \bar{H}_{\text{BFGS}} = (I - \frac{sy^T}{\langle y, s \rangle}) H (I - \frac{ys^T}{\langle y, s \rangle}) + \frac{ss^T}{\langle y, s \rangle}$$

At this point we note that the BFGS update is sometimes called the complementary DFP update and that the underlying single rank method (7.22) in which  $d = s$  was proposed by G. McCormick (see Pearson (1969)).

There is growing evidence that the BFGS is the best current update formula for use in unconstrained minimization. For example, see the results of Dixon (1972b). For this reason, and for future reference we state the following analogue of Theorem 7.7.

Theorem 7.8: Let  $H \in L(R^n)$  be a nonsingular symmetric matrix, and define  $\bar{H}_{\text{BFGS}} \in L(R^n)$  by (7.25) for any vectors  $y$  and  $s$  in  $R^n$  with  $\langle y, s \rangle \neq 0$ . Then  $\bar{H}_{\text{BFGS}}$  is nonsingular if and only if  $\langle s, Bs \rangle \neq 0$  where  $B = H^{-1}$ . If  $\bar{H}_{\text{BFGS}}$  is nonsingular then

$\bar{B}_{\text{BFGS}} = \bar{H}_{\text{BFGS}}^{-1}$  can be expressed as

$$\bar{B}_{\text{BFGS}} = B + \frac{yy^T}{\langle y, s \rangle} - \frac{Bss^TB}{\langle s, Bs \rangle}$$

Furthermore, if  $H$  is positive definite then  $\bar{H}_{\text{BFGS}}$  is positive definite if and only if  $\langle y, s \rangle > 0$ .

The remark at the end of Section 7.2 about the behavior of  $\bar{B}_{\text{BFGS}}$  as a relative approximation to the Hessian holds for  $\bar{H}_{\text{BFGS}}$  as a relative approximation to the inverse Hessian. (See equation (8.18)) Also note

that there is a close relationship between the matrices generated by the DFP and BFGS updates for it is easy to verify that if  $H$  is positive definite then

$$(7.26) \quad \bar{H}_{\text{BFGS}} = \bar{H}_{\text{DFP}} + vv^T$$

where  $v$  is the vector

$$(7.27) \quad v = \langle y, Hy \rangle^{-1/2} \left[ \frac{s}{\langle s, y \rangle} - \frac{Hy}{\langle y, Hy \rangle} \right],$$

while if  $B$  is positive definite then

$$(7.28) \quad \bar{B}_{\text{DFP}} = \bar{B}_{\text{BFGS}} + ww^T$$

where  $w$  is the vector

$$w = \langle s, Bs \rangle^{-1/2} \left[ \frac{y}{\langle s, y \rangle} - \frac{Bs}{\langle s, Bs \rangle} \right].$$

By virtue of Lemma 7.4, relations (7.26) and (7.28) imply that the eigenvalues of  $\bar{H}_{\text{BFGS}}$  ( $\bar{B}_{\text{BFGS}}$ ) are larger (smaller) than the eigenvalues of  $\bar{H}_{\text{DFP}}$  ( $\bar{B}_{\text{DFP}}$ ). However, there does not seem to be any relationship between the condition number of  $\bar{H}_{\text{BFGS}}$  and the condition number of  $\bar{H}_{\text{DFP}}$ .

From a purely algebraic point of view, the developments of Sections 7.1 and 7.2 are identical to those in Section 7.3. This follows from the fact that (7.22) and (7.24) can be obtained from (7.3) and (7.9), respectively, by interchanging  $y$  and  $s$ , replacing  $B$ 's by  $H$ 's and  $c$  by  $d$ . In particular Theorem 7.7 and 7.8 are identical since both of them follow from a more general result which relates  $\bar{A}$  and  $A$  where

$$\bar{A} = \left( I - \frac{uv^T}{\langle u, v \rangle} \right) A \left( I - \frac{vu^T}{\langle u, v \rangle} \right) + \frac{uu^T}{\langle u, v \rangle}$$

and  $\langle u, v \rangle \neq 0$ . In spite of these remarks we have opted for a separate development for expository purposes. Nevertheless, it is useful to note that the DFP and BFGS are related by the transformation

$$(7.29) \quad s \leftrightarrow y, \quad B \leftrightarrow H, \quad \bar{B} \leftrightarrow \bar{H}.$$

In fact, Fletcher's (1970) derivation of the BFGS update was through this transformation.

Finally, we note that if a direct and inverse update are related by the transformation (7.28) then these updates are sometimes called "dual" or "complementary" updates, and this is the reason why the BFGS is also called the complementary DFP formula.

## 8. CONVERGENCE RESULTS FOR RANK-TWO QUASI-NEWTON METHODS

Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be continuously differentiable in an open set  $D$  and consider a method of the form

$$(8.1) \quad x_{k+1} = x_k - \lambda_k H_k \nabla f(x_k), \quad k = 0, 1, \dots,$$

where the matrices  $H_k$  are generated by one of the methods of Section 7 and  $\lambda_k$  is suitably chosen. In this section we examine some of the convergence and rate of convergence results that are available for (8.1).

In a lot of theoretical work sufficient conditions are assumed so that  $\lambda_k$  can be chosen by an exact line search. This usually means that either

$$(8.2) \quad \lambda_k = \min\{\lambda > 0: \langle \nabla f(x_k + \lambda p_k), p_k \rangle = 0\}$$

where  $p_k = -H_k \nabla f(x_k)$ , or that  $\lambda_k$  is the first local minimizer of

$f(x_k + \lambda p_k)$  for  $\lambda \geq 0$ . Either choice is unrealistic as usually it is not possible to find  $\lambda_k$  to much accuracy in a reasonable amount of time unless, for example,  $f$  is a quadratic, positive definite functional. In this case

$$(8.3) \quad f(x) = (1/2)\langle x, Ax \rangle - \langle x, b \rangle + c$$

for some symmetric positive definite  $A \in L(\mathbb{R}^n)$ , and the  $\lambda_k$  which satisfies (8.2) is given by

$$(8.4) \quad \lambda_k = - \langle Ax_k - b, p_k \rangle / \langle Ap_k, p_k \rangle .$$

The earlier convergence results for quasi-Newton methods were given for  $f$  defined by (8.3) and  $\lambda_k$  chosen by (8.4). It was shown that if  $\{x_k\}$  is generated by (8.1) and  $H_k$  correspond to, say the DFP or BFGS updates, then  $x_\ell = A^{-1}b$  for some  $0 \leq \ell \leq n$ , and if  $\ell = n$  then  $H_n = A^{-1}$ . This type of finite termination property has sometimes been called quadratic termination. The relevance of the quadratic termination property to the general nonlinear problem was originally based on the assumption that if a method terminates in a finite number of steps for a quadratic then this implies superlinear convergence for nonlinear functionals. There has never been any theoretical or numerical support for this belief. (See, however, the discussion following Theorem 8.9). Nevertheless, quadratic termination seems to be a desirable property although as Broyden's method shows, it is not indispensable for superlinear convergence.

In order to describe the quadratic termination properties for symmetric rank two quasi-Newton methods, consider the following class of updates:

$$(8.5) \quad \bar{H}_\phi = (1 - \phi)\bar{H}_{\text{DFP}} + \phi\bar{H}_{\text{BFGS}}$$

where  $\phi$  is a parameter which may depend on  $s$ ,  $y$ ,  $H$  and the iteration counter. This class of updates was introduced by Broyden (1967) although not in the form (8.5). It was Fletcher (1970) who showed that Broyden's class, which had been given in terms of a parameter  $\beta$ , could be written in the form (8.5) and that the relationship between  $\phi$  and  $\beta$  is that  $\phi = \beta \langle y, s \rangle$ . Fletcher also noted that if  $H$  is positive definite then equation (7.26) implies that update (8.5) can be written as

$$\bar{H}_\phi = \bar{H}_{DFP} + \phi v v^T$$

where the vector  $v$  is defined by (7.27). It is immediate from this expression that if  $\phi \geq 0$  then  $\bar{H}_\phi$  shares the property of hereditary positive definiteness with  $\bar{H}_{DFP}$ . For future reference and to state the quadratic termination properties of (8.5), note that Broyden's class is generated by

$$H_{k+1} = H_k + \frac{s_k s_k^T}{\langle s_k, y_k \rangle} - \frac{H_k y_k y_k^T H_k}{\langle y_k, H_k y_k \rangle} + \phi_k v_k v_k^T, \quad (8.6) (a)$$

$$v_k = \langle y_k, H_k y_k \rangle^{1/2} \left[ \frac{s_k}{\langle s_k, y_k \rangle} - \frac{H_k y_k}{\langle y_k, H_k y_k \rangle} \right],$$

and where the vectors  $s_k$  and  $y_k$  are usually defined by

$$(8.6) (b) \quad s_k = x_{k+1} - x_k, \quad y_k = \nabla f(x_{k+1}) - \nabla f(x_k).$$

Theorem 8.1: Assume that  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is the positive definite quadratic functional (8.3) and that  $H_0 \in L(\mathbb{R}^n)$  is symmetric and positive definite. For any given  $x_0 \in \mathbb{R}^n$ , let  $\{x_k\}$  be generated by (8.1) where  $\lambda_k$ ,  $H_k$  satisfy (8.4) and (8.6), respectively, and  $\phi_k$  may depend on  $s_k$ ,  $y_k$  and  $H_k$ . If  $\phi_k \geq 0$  then there is an integer  $0 \leq \ell \leq n$  such that  $x_\ell = A^{-1}b$  and if  $\ell = n$  then  $H_n = A^{-1}$ .

A typical proof of Theorem 8.1 proceeds by induction to show that the directions  $s_k$  are A-conjugate in the sense that

$$\langle s_i, A s_j \rangle = 0, \quad i > j,$$

and that also

$$H_i y_j = s_j, \quad i > j.$$

This was the argument used by Broyden (1967); it shows that  $x_{k+1}$  minimizes  $f$  in the hyperplane  $x_0 + L$  where  $L$  is the linear span of  $s_0, \dots, s_k$ . Broyden (1971b) and Powell (1972b, 1973) have extended and refined Theorem 8.1; in particular, Powell (1972b) shows that  $A^{1/2} H_k A^{1/2}$  has at least  $k$  unit eigenvalues. However, in all these results finite termination depends on choosing  $\lambda_k$  by (8.4). Also note that if

$$\phi_k = \frac{\langle y_k, s_k \rangle}{\langle s_k - H_k y_k, y_k \rangle}$$

then (8.6) (a) reduces to the symmetric rank one formula (7.6) but Theorems 7.1 and 8.1 are not comparable.

In a certain sense, Theorem 8.1 generalizes to nonlinear functionals. The relevant result is due to Powell (1971, 1972a).

Theorem 8.2: Let  $f: R^n \rightarrow R$  be twice continuously differentiable and convex on  $R^n$  and assume that for a given  $x_0 \in R^n$  the level set  $L(x_0)$  is bounded. Suppose that  $\{x_k\}$  is generated by (8.1) and that  $\lambda_k, H_k$  are chosen by an exact line search and the DFP update, respectively. Then for any symmetric, positive definite  $H_0 \in L(R^n)$  and  $\epsilon > 0$  there is an index  $k$  such that  $\|\nabla f(x_k)\| < \epsilon$ .

It is possible to show that if  $\lambda_k$  is chosen by an exact line search, then the conclusion of Theorem 6.3 still holds; that is,

$$\lim_{k \rightarrow +\infty} \langle \nabla f(x_k), \frac{p_k}{\|p_k\|} \rangle = 0 ,$$

where  $p_k = -H_k \nabla f(x_k)$ . Therefore, an obvious but so far unsuccessful approach to proving Theorem 8.2 would be to show that some subsequence of  $\{H_k\}$  has uniformly bounded condition numbers. Instead, the proof consists of assuming that  $\|\nabla f(x_k)\| \geq \epsilon$  holds for all  $k \geq 0$ , and then reaching a contradiction. Moreover, the techniques used in arriving at this contradiction depend very heavily on the use of exact line searches.

It would be very interesting to show that Theorem 8.2 still holds if the convexity assumption is relaxed or the choice of  $\lambda_k$  is more realistic. In doing this, the choice of  $\lambda_k$  should guarantee that the matrices  $H_k$  remain positive definite. A choice of  $\lambda_k$  which satisfies this requirement is given in Theorem 6.3 since in this case

$$\langle y_k, s_k \rangle \geq \lambda_k (\beta - 1) \langle \nabla f(x_k), p_k \rangle > 0 .$$

Note that exact line searches also satisfy this requirement.

That Theorem 8.2 extends to other methods in the Broyden class (8.5) follows from the following remarkable result of Dixon (1972a).

Theorem 8.3: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be differentiable on  $\mathbb{R}^n$ , and assume that for a given  $x_0 \in \mathbb{R}^n$  the level set  $L(x_0)$  is bounded. Given a symmetric positive definite  $H_0 \in L(\mathbb{R}^n)$  suppose that  $\{x_k\}$  is generated by (8.1) where  $\lambda_k, H_k$  are chosen according to (8.2) and (8.6), respectively. If  $\phi_k \geq 0$  then the sequence  $\{x_k\}$  is independent of  $\{\phi_k\}$ .



Dixon's result is actually more general than Theorem 8.3 since it allows negative values of  $\phi_k$ . However, the above formulation suffices for our purposes, and moreover, the more general formulation requires additional assumptions on  $\{H_k\}$ .

All the results presented so far on the convergence of rank-two quasi-Newton methods depend on exact line searches. Moreover, none of these results give any indication of how to choose  $\phi_k \in [0, \infty)$  when exact searches are not used. The following result of Fletcher (1970) shows that for quadratic functionals the updates with  $\phi_k \in [0, 1]$  have a very desirable property which does not depend on exact line searches.

Theorem 8.4: Assume that  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is the positive definite quadratic functional (8.3) and that  $H_0 \in L(\mathbb{R}^n)$  is symmetric and positive definite. Let  $\{s_k\}$  be any sequence of nonzero vectors, let  $\{H_k\}$  be generated by (8.6)(a) with  $y_k = As_k$ , and let  $\lambda_i^{(k)}$ ,  $i = 1, \dots, n$  be the eigenvalues of  $A^{1/2}H_kA^{1/2}$ . If  $\phi_k \in [0, 1]$  then

$$\min\{\lambda_i^{(k)}, 1\} \leq \lambda_i^{(k+1)} \leq \max\{\lambda_i^{(k)}, 1\}.$$

Since Fletcher (1970) showed that Theorem 8.4 fails if  $\phi_k \notin [0, 1]$ , this result indicates that the most reasonable updates in Broyden's class (8.5) correspond to  $\phi \in [0, 1]$ . In fact, numerical results of Dixon (1972b) suggest that  $\phi = 1$  is to be preferred. Of course, Theorem 8.4 does not say anything about the PSB update since it is not of the form (8.5).

As noted by Fletcher (1970), Theorem 8.4 implies that the sequence  $H_k$  has uniformly bounded condition numbers and therefore it can be used to give a global convergence result for strictly

convex quadratic functionals without exact line searches (for example if  $\lambda_k$  is chosen according to Theorem 6.3 with  $\phi_k = -H_k \nabla f(x_k)$ ). This result indicates that Theorem 8.2 might hold with a more realistic line search.

We have surveyed the global convergence results for rank-two quasi-Newton methods; since the analysis of the asymptotic rate of convergence is of major importance, we now investigate this topic as well as the local convergence properties of the BFGS and DFP methods.

For the local convergence of these updates we show how to choose a norm so that (4.2) and (4.7) hold, respectively. First consider the DFP method and recall that  $\bar{B}_{DFP}$  and  $B$  are related by (7.18). It follows that for any symmetric  $A$  and  $B$  in  $L(R^n)$

$$(8.7) \quad \bar{B}_{DFP} - A = P^T(B - A)P + [(y - As)y^T + y(y - As)^T P] / \langle y, s \rangle$$

where

$$(8.8) \quad P = I - \frac{sy^T}{\langle y, s \rangle}.$$

A similar relationship holds between  $H$  and  $\bar{H}_{BFGS}$  for the BFGS update. In this case  $\bar{H}_{BFGS}$  and  $H$  are related by (7.25) so that if  $A$  and  $H$  are symmetric, and  $A$  is nonsingular, then

$$(8.9) \quad \bar{H}_{BFGS} - A^{-1} = P(H - A^{-1})P^T + [(s - A^{-1}y)s^T + s(s - A^{-1}y)^T P^T] / \langle y, s \rangle$$

where  $P$  is again defined by (8.8). In order to show that (8.7) satisfies (4.2), and (8.9) satisfies (4.7), we need the following result of Broyden (1970).

Lemma 8.5: If  $Q \in L(R^n)$  is defined by

$$(8.10) \quad Q = I - \frac{uv^T}{\langle u, v \rangle}$$

with  $u, v$  in  $R^n$  and  $\langle u, v \rangle \neq 0$ , then

$$||Q||_2 = \frac{||u|| ||v||}{|<u,v>|}.$$

Proof: The most straightforward way to verify this result is to recall that  $||Q||_2^2$  is the largest eigenvalue of  $Q^T Q$  and to calculate the eigenvalues of  $Q^T Q$  with Lemma 7.6.

Lemma 8.5 shows that  $||P||_2$  is the secant of the angle between  $y$  and  $s$ , and since  $y$  and  $s$  are not in general parallel,  $||P||_2$  may be arbitrarily large. Therefore the  $l_2$ -norm does not seem to be suitable for estimating (8.7) or (8.9). However, near  $x^*$  we do have that  $A^{-1/2}y$  and  $A^{1/2}s$  are nearly parallel if  $A = \nabla^2 f(x^*)$  and this suggests the use of a weighted norm. For the DFP method an appropriate norm is defined by

$$(8.11) \quad ||E||_{DFP} = ||A^{-1/2}EA^{-1/2}||_F.$$

Then Lemma 8.5 and (1.2) imply that

$$(8.12) \quad ||P^T(B-A)P||_{DFP} \leq ||A^{1/2}PA^{-1/2}||_2^2 ||B-A||_{DFP} \equiv \frac{1}{\omega} ||B-A||_{DFP}$$

where

$$(8.13) \quad \omega = \frac{<y,s>}{||A^{-1/2}y|| ||A^{1/2}s||} = \frac{<A^{-1/2}y, A^{1/2}s>}{||A^{-1/2}y|| ||A^{1/2}s||}.$$

Similar estimates of the other terms in (8.7) yield

$$(8.14) \quad ||\frac{y(y-As)^T P}{<y,s>}||_{DFP} \leq \frac{1}{\omega^2} \frac{||A^{-1/2}y - A^{1/2}s||}{||A^{1/2}s||},$$

$$(8.15) \quad ||\frac{(y-As)y^T}{<y,s>}||_{DFP} \leq \frac{1}{\omega} \frac{||A^{-1/2}y - A^{1/2}s||}{||A^{1/2}s||}.$$

Now place (8.12), (8.14) and (8.15) together to obtain

$$(8.16) \quad ||\bar{B}_{DFP} - A||_{DFP} \leq \frac{1}{\omega^2} ||B-A||_{DFP} + \frac{2}{\omega^2} \frac{||A^{-1/2}y - A^{1/2}s||}{||A^{1/2}s||}.$$

An analogous relationship holds for the BFGS update in this case the appropriate norm is defined by

$$(8.17) \quad ||E||_{\text{BFGS}} = ||A^{1/2}EA^{1/2}||_F,$$

and it is not difficult to verify that the analogue of (8.16) is

$$(8.18) \quad ||\bar{H}_{\text{BFGS}}^{-1}||_{\text{BFGS}} \leq \frac{1}{\omega^2} ||H-A^{-1}||_{\text{BFGS}} + \frac{2}{\omega^2} \frac{||A^{1/2}s-A^{-1/2}y||}{||A^{-1/2}y||}.$$

As noted in Section 7, an interpretation of (8.18) is that if  $A = \nabla^2 f(x)$  is positive definite and  $||s||$  is small, then the relative error of  $\bar{H}_{\text{BFGS}}$  as an approximation to  $\nabla^2 f(x)^{-1}$  is not too much larger than the relative error of  $H$  as an approximation to  $\nabla^2 f(x)^{-1}$ . Furthermore, the possible growth in this relative error is determined by how much  $f$  differs on the points  $x$  and  $\bar{x}$  from the quadratic whose Hessian is  $A$ . This difference is measured in two ways but both have to do with how well  $A^{-1/2}y$  is approximated by  $A^{1/2}s$ ; there is an additive term which is the relative error in this approximation and a multiplicative term which is the square of the secant of the angle between these two vectors. Of course, we easily see that the additive term does not exceed the product of the square root of the condition number of  $A$  and the relative error in the approximation of  $y$  by  $As$ . An analogous discussion holds for (8.16).

Another consequence of (8.16), (8.18) is the local convergence of the DFP and BFGS methods as given by Broyden, Dennis and Moré (1973).

Theorem 8.6: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be twice continuously differentiable in an open convex set  $D$ , and assume that  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive definite for some  $x^*$  in  $D$ . Suppose, in addition, that

$$(8.19) \quad ||\nabla^2 f(x) - \nabla^2 f(x^*)|| \leq \kappa ||x - x^*||, \quad x \in D,$$

and consider the DFP and BFGS methods as defined by (8.1) with  $\lambda_k \equiv 1$ . Then the DFP and BFGS methods are locally and superlinearly convergent at  $x^*$ .

Proof: To prove that these methods are locally and linearly convergent at  $x^*$ , we only need to show that (4.2) and (4.7) are satisfied when  $D_M$  is the set of all symmetric matrices and  $F = \nabla f$ . For the DFP method, first note that (8.13) with  $A = \nabla^2 f(x^*)$  and Lemmas 3.2 and 3.3 imply that

$$1 - \omega^2 \leq \left[ \mu \frac{\|y - \nabla^2 f(x^*)s\|}{\|s\|} \right]^2 \leq [\mu \kappa \sigma(x, \bar{x})]^2$$

where  $\sigma(x, \bar{x})$  is defined by (4.3) and  $\mu = \|\nabla^2 f(x^*)^{-1}\|$ . Thus if  $x$  and  $\bar{x}$  lie in a neighborhood  $N_1$  of  $x^*$  such that  $\sigma(x, \bar{x}) \leq (2\mu\kappa)^{-1}$  then  $\omega^2 \geq 1/2$ . In particular,

$$(8.20) \quad \frac{1}{\omega^2} = 1 + \frac{1-\omega^2}{\omega^2} \leq 1 + \kappa\mu\sigma(x, \bar{x}).$$

Therefore (8.16) with  $A = \nabla^2 f(x^*)$  and Lemma 3.2 imply that (5.2) holds with  $\alpha_1 = \kappa\mu$  and  $\alpha_2 = 4\kappa\mu$  and where  $\|\cdot\|$  is the matrix norm defined by (8.11). This proves the local convergence of the DFP method. For the BFGS first let  $\varepsilon$  be a positive lower bound for the eigenvalues of  $\nabla^2 f(x)$  in a neighborhood  $N_2$  of  $x^*$  so that  $\|y\| \geq \varepsilon\|s\|$  provided  $x$  and  $\bar{x}$  lie in  $N_2$ . If  $N_2 \subset N_1$  then (8.18) with  $A = \nabla^2 f(x^*)^{-1}$ , (8.20) and Lemma 3.3 imply that (5.7) holds with  $\alpha_1 = \kappa\mu$  and  $\alpha_2 = 4(\mu\rho)^{1/2}\kappa/\varepsilon$ . Here  $\rho = \|\nabla^2 f(x^*)\|$  and  $\|\cdot\|$  is the norm defined by (8.17).

To prove that the DFP and BFGS methods are superlinearly convergent is more difficult and requires careful estimation of the terms  $\|P^T(B-A)P\|_{DFP}$  and  $\|P(H-A^{-1})P^T\|_{BFGS}$  respectively. These estimates were obtained by Broyden, Dennis and Moré (1973) and then used by Dennis and Moré (1974) to prove superlinear convergence for various choices of  $\lambda_k$  in (8.1).

Theorem 8.7: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  satisfy the assumptions of Theorem 8.6, and suppose that  $\{x_k\}$  is a sequence in  $D$  such that

$$(8.21) \quad \sum_{k=0}^{\infty} \|x_k - x^*\| < +\infty.$$

If the sequence  $\{H_k\}$  is defined by (8.6) with either  $\phi_k \equiv 0$  or  $\phi_k \equiv 1$  and  $\langle y_k, s_k \rangle$  is positive for  $k \geq 0$ , then for any symmetric positive definite  $H_0$  in  $L(\mathbb{R}^n)$  the matrices  $H_k$  are well-defined and positive definite with uniformly bounded condition numbers.

Moreover, if  $B_k = H_k^{-1}$  then

$$(8.22) \quad \lim_{k \rightarrow \infty} \frac{\|[B_k - \nabla^2 f(x^*)]s_k\|}{\|s_k\|} = 0.$$

Since  $\langle y_k, s_k \rangle$  is positive for  $k \geq 0$ , Theorems 7.7 and 7.8 imply that in either case  $H_k$  is well-defined and positive definite. The remainder of the proof is somewhat long, so we only outline it. First it is shown that  $\{\|B_k\|\}$  is bounded and (8.22) holds for the DFP method. A similar argument for BFGS shows that  $\{\|H_k\|\}$  is bounded and instead of (8.22),

$$(8.23) \quad \lim_{k \rightarrow \infty} \frac{\|[H_k - \nabla^2 f(x^*)^{-1}]y_k\|}{\|y_k\|} = 0.$$

However, if  $\{\|B_k\|\}$  is bounded then

$$[B_k - \nabla^2 f(x^*)]s_k = [I - B_k \nabla^2 f(x^*)^{-1}](y_k - \nabla^2 f(x^*)s_k) + B_k [H_k - \nabla^2 f(x^*)^{-1}]y_k$$

shows that (8.23) implies (8.22). Hence, the final step in the proof consists of using the techniques of Powell (1971, pages 31-32) to prove that  $\{\|B_k\|\}$  and  $\{\|H_k\|\}$  are bounded for the BFGS and DFP methods respectively.

Dennis and Moré (1974) elaborate on Theorem 8.7 and give examples which show that (8.24) does not necessarily imply that

$\{B_k\}$  converges to  $\nabla^2 f(x^*)$ . Also note that Theorem 8.7 implies that the DFP and BFGS methods of Theorem 8.6 are superlinearly convergent. This is also a consequence of the following more general result.

Theorem 8.8: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  satisfy the assumptions of Theorem 8.6 and consider the DFP and BFGS methods as defined by (8.1) with  $\lambda_k$  determined by any strategy such that (8.22) implies that  $\{\lambda_k\}$  converges to unity. If the sequence  $\{x_k\}$  generated by the DFP and BFGS method satisfies (8.21) then  $\{x_k\}$  converges superlinearly to  $x^*$ .

Proof: Theorem 3.1 implies that  $\{x_k\}$  converges superlinearly to  $x^*$  if

$$(8.24) \quad \lim_{k \rightarrow \infty} \frac{||[\lambda_k^{-1} B_k - \nabla^2 f(x^*)] s_k||}{||s_k||} = 0.$$

On the other hand, Theorem 8.8 and our assumptions show that (8.22) holds and hence  $\{\lambda_k\}$  converges to unity. Thus (8.24) also holds and hence  $\{x_k\}$  converges superlinearly to  $x^*$ .

The remark at the end of Section 5 shows that assumption (8.21) is not needed if  $(x_0, H_0)$  is sufficiently close to  $(x^*, \nabla^2 f(x^*)^{-1})$ . In a similar vein, it would be interesting to prove Theorem 8.8 assuming only that  $\{x_k\}$  converges to  $x^*$  instead of (8.21). However, as it stands Theorem 8.8 shows that either the DFP and BFGS methods converge superlinearly or they converge sublinearly in the sense that

$$\limsup_{k \rightarrow \infty} ||x_k - x^*||^{1/k} = 1.$$

In practice, sublinear convergence is essentially equivalent to non-convergence, so Theorem 8.8 covers all the computationally interesting cases.

As pointed out in Section 3, condition (8.22) is equivalent to requiring the vectors  $p_k = -B_k^{-1} \nabla f(x_k)$  to approach the Newton step in both direction and length, so most algorithms for finding  $\lambda_k$  satisfy the assumptions of Theorem 8.8. For example, Theorem 6.4 shows that the strategy of Theorem 6.3 satisfies Theorem 8.8 while Dennis and Moré (1974) have shown that this is also the case for exact line searches. This indicates that the DFP method with exact line searches may be superlinearly convergent; in fact, more is known.

Theorem 8.9: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  and  $x_0$  satisfy the assumptions of Theorem 8.2 and let  $\{x_k\}$  be generated by the DFP method with perfect line searches. If  $\{x_k\}$  converges to a point  $x^*$  at which  $\nabla^2 f(x^*)$  is positive definite and (8.19) holds, then  $\nabla f(x^*) = 0$  and  $\{x_k\}$  converges superlinearly to  $x^*$ . In addition there is an  $\eta > 0$  such that

$$(8.25) \quad \|x_{k+n} - x^*\| \leq \eta \|x_k - x^*\|^2, \quad k \geq 0.$$

That  $\{x_k\}$  converges superlinearly to  $x^*$  is due to Powell (1971), but (8.25) -- which is known as n-step quadratic convergence -- is due to Burmeister (1973). The proofs of these two results are completely different; Burmeister's proof depends on the finite termination property of DFP while, as pointed out by Dennis and Moré (1974), Powell's result can be proved by showing that (8.21) holds and then applying Theorem 8.8. It should also be appreciated that a sequence may converge n-step quadratically but not be superlinearly convergent and conversely. However, n-step quadratic convergence does imply that

$$\lim_{k \rightarrow \infty} \|x_k - x^*\|^{1/k} = 0,$$



and thus  $\{x_k\}$  is R-superlinearly convergent in the terminology of Ortega and Rheinboldt (1970); Q-superlinear convergence corresponds to the notion used in this paper.

Theorem 8.10: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be the strictly convex quadratic functional (8.3). Then for  $\lambda_k \equiv 1$  the DFP and BFGS methods converge globally and superlinearly to  $A^{-1}b$ .

Proof: It is clear that the iterations are well defined. To prove the result for the DFP method note that since  $y=As$ , equation (8.7) implies that

$$(8.26) \quad ||\bar{B}_{DFP} - A||_{DFP} = ||Q^T [A^{-1/2} (B-A) A^{-1/2}] Q||_F$$

where  $z = A^{1/2}s$  and

$$Q = I - \frac{zz^T}{\langle z, z \rangle}.$$

However, in the proof of Theorem 5.2 we showed that for any  $E \in L(\mathbb{R}^n)$  and  $z \in \mathbb{R}^n$ ,

$$||EQ||_F \leq ||E||_F - (2||E||_F)^{-1} \left( \frac{||Ez||_2}{||z||} \right)^2.$$

Thus, if we let  $\eta_k = ||B_k - A||_{DFP}$ , and use (1.2) and the above estimate in (8.26) then

$$\eta_{k+1} \leq [1 - (2\eta_k^2)^{-1} \psi_k^2] \eta_k$$

where

$$\psi_k = \frac{||A^{-1/2} (B_k - A) s_k||}{||A^{1/2} s_k||}.$$

It is now clear that  $\{\eta_k\}$  is monotone decreasing and hence convergent.

If  $\eta$  is an upper bound for  $\{\eta_k\}$  then

$$(2\eta)^{-1} \psi_k^2 \leq \eta_k - \eta_{k+1},$$

and since  $\{\eta_k\}$  is converging it follows that  $\{\psi_k\}$  tends to zero.

Consequently,

$$(8.27) \quad \lim_{k \rightarrow \infty} \frac{|| (B_k - A)s_k ||}{|| s_k ||} = 0 .$$

Moreover,

$$A(x_{k+1} - x^*) = \nabla f(x_{k+1}) = \nabla f(x_k) + As_k = (A - B_k)s_k$$

implies that

$$\frac{|| x_{k+1} - x^* ||}{|| s_k ||} \leq || A^{-1} || \frac{|| (B_k - A)s_k ||}{|| s_k ||} ,$$

and thus (8.27) shows that  $\{x_k\}$  converges superlinearly to  $x^*$ .

For the BFGS method similar calculations with (8.9) yield

$$(8.28) \quad \lim_{k \rightarrow \infty} \frac{|| [H_k - A^{-1}]y_k ||}{|| y_k ||} = 0 .$$

Moreover, Theorem 8.4 shows that  $\{|| B_k ||\}$  is bounded.. Thus, as noted after Theorem 8.7, (8.28) implies that (8.27) holds and now the proof is completed as before.

Theorem 8.10 seems to be just a curiosity since if  $\lambda_k$  is chosen by an exact line search, then convergence will take place in at most  $n$  steps. However, it does give an indication of the stability of the DFP and BFGS updates without exact line searches.

We have now finished our study of the asymptotic behavior of the DFP and BFGS methods. It is also possible to study the PSB update, but since it does not generate positive definite matrices, results like Theorem 8.8 have to be modified for the PSB update.

The PSB update is not generally used in a descent implementation, but Powell (1970c, 1970d) has described and analyzed a quite competitive algorithm which uses the PSB algorithm in a hybrid

implementation, and has shown that if certain "special iterations" are taken then the algorithm converges globally and superlinearly. These special iterations guarantee that the directions used by the PSB update are uniformly linearly independent and therefore, that the sequence  $\{B_k\}$  generated by the PSB update converges to the Jacobian evaluated at the solution -- for a discussion of the concept of uniform linear independence and its relationship to the Broyden and PSB update see More' and Trangenstein (1974). Powell (1974) later proved that in theory the algorithm converged globally and superlinearly even if these special iterations are not used. In practice however they cannot be taken away from the algorithm, without a significant loss in efficiency.

The above results of Powell deserve further investigation. In fact, the whole question of how to globalize an algorithm is very important and represents an open field of research.

In its simplest form the PSB method is given by

$$(8.29) \quad x_{k+1} = x_k - B_k^{-1} \nabla f(x_k)$$

where  $\{B_k\}$  is generated by

$$(8.30) \quad B_{k+1} = B_k + \frac{(y_k - B_k s_k) s_k^T + s_k (y_k - B_k s_k)^T}{\langle s_k, s_k \rangle} - \frac{\langle y_k - B_k s_k, s_k \rangle}{\langle s_k, s_k \rangle^2} s_k s_k^T,$$

and as usual

$$(8.31) \quad y_k = \nabla f(x_{k+1}) - \nabla f(x_k), \quad s_k = x_{k+1} - x_k.$$

Note that since the matrices  $\{B_k\}$  are not necessarily positive definite it is not possible to carry out the above iteration by

updating an  $LDL^T$  decomposition of  $B_k$ . To avoid  $O(n^3)$  operations per step it is usual to generate  $H_k = B_k^{-1}$ . An alternative would be to update a factorization of the form  $QTQ^T$ , where  $Q$  is orthogonal and  $T$  is tridiagonal, but this approach has not been investigated.

The following result of Broyden, Dennis and Moré (1973) covers the above iteration.

Theorem 8.11: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  satisfy the assumptions of Theorem 8.6 except that now  $\nabla^2 f(x^*)$  is not required to be positive definite, and consider the PSB method as defined by (8.29), (8.30) and (8.31). Then the PSB method is locally and superlinearly convergent.

Since the proof of this result is so similar to that of Theorem 4.2 we omit it. It is also worthwhile noting that the remarks made about Broyden's method after Theorem 4.2 apply, with obvious modifications, to the PSB method, and that Dennis (1971, 1972) has given Kantorovich theorems for the Broyden and PSB methods.

To conclude this section we point out that the PSB method, if properly modified, is globally and superlinearly convergent for the quadratic functional (8.3) if  $A$  is any nonsingular, symmetric matrix. It is only necessary to modify  $B_{k+1}$  so that it is nonsingular. For example, if instead of (8.30) we define

$$(8.32) \quad B_{k+1} = B_k + \theta_k \frac{(y_k - B_k s_k) s_k^T + s_k (y_k - B_k s_k)^T}{\langle s_k, s_k \rangle} - \theta_k^2 \frac{\langle y_k - B_k s_k, s_k \rangle}{\langle s_k, s_k \rangle^2} s_k s_k^T$$

and  $B_k$  is nonsingular then it is possible to choose  $\theta_k$  so that

$$(8.33) \quad B_{k+1} \text{ is nonsingular, } |\theta_k - 1| \leq \hat{\theta} \text{ for some } \hat{\theta} \in (0, 1).$$

Moré and Trangenstein (1974) elaborate on how this can be done, and also prove the following result.

Theorem 8.12: Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be given by (8.3) where  $A \in L(\mathbb{R}^n)$  is any nonsingular, symmetric matrix and consider the PSB method (8.29) where  $\{B_k\}$  is generated by (8.32), (8.31) and  $\theta_k$  satisfies (8.33). Then the PSB method is globally and superlinearly convergent to  $A^{-1}b$ .

More' and Trangenstein (1974) also point out that Theorem 8.11 holds if the PSB method (8.29) is defined by (8.32), (8.31) and  $\theta_k$  satisfies (8.33).

## 9. CONCLUDING REMARKS

We have tried to write this survey in such a way that the important problems suggest themselves, so instead of ending with remarks about directions for future research, we end with an admission of certain omissions.

Although we have indicated several approaches to the computation of the updates, all these approaches are based on an additive correction of rank at most two. Other approaches are possible; Brodlie, Gourlay and Greenstadt (1973) discuss multiplicative corrections so that their direct updates are of the form

$$\bar{B} = (I + uv^T)B(I + vu^T) ,$$

and show that the DFP and BFGS can be written in this factored form.

We have not mentioned any particular implementations because there are a number of very promising algorithms now being tested and such remarks would likely be out of date before their publication. See, however, the paper of Fletcher (1972), which discusses several of the currently available algorithms.

Finally, we apologize for the omission of several excellent papers which only deal with quasi-Newton methods as applied to

strictly convex quadratic functionals. In particular, the paper of Huang (1970) introduces a class of updates which has many of the properties of the Broyden class. We have restricted our attention to the Broyden class since it is that subclass of the Huang class which satisfies the quasi-Newton equation and has the hereditary symmetry property.

## References

- Armijo, L.(1966), Minimization of functions having Lipschitz continuous first partial derivatives, Pacific J. Math. 16, 1-3.
- Brent, R.P.(1973), Some efficient algorithms for solving systems of nonlinear equations, SIAM J. Numer. Anal. 10, 327-344.
- Brodlić, K.W., A.R. Gourlay and J. Greenstadt(1973), Rank-one and rank-two corrections to positive definite matrices expressed in product form, J. Inst. Math. Appl. 11, 73-82.
- Broyden, C.G.(1965), A class of methods for solving nonlinear simultaneous equations, Math. Comp. 19, 577-593.
- Broyden, C.G.(1967), Quasi-Newton methods and their application to function minimization, Math. Comp. 21, 368-381.
- Broyden, C.G.(1969), A new double-rank minimization algorithm, AMS Notices 16, 670.
- Broyden, C.G.(1970), The convergence of single-rank quasi-Newton methods, Math. Comp. 24, 365-382.
- Broyden, C.G.(1971a), The convergence of an algorithm for solving sparse nonlinear systems, Math. Comp. 25, 285-294.
- Broyden, C.G.(1971b), The convergence of a class of double-rank minimization algorithms, Parts I and II, J. Inst. Math. Appl. 7, 76-90, 222-236.
- Broyden, C.G., J.E. Dennis and J.J. Moré(1973), On the local and superlinear convergence of quasi-Newton methods, J. Inst. Math. Appl. 12, 223-246.
- Burmeister, W.(1973), Die Konvergenzordnung des Fletcher-Powell Algorithmus, ZAMM 53, 693-699.
- Davidon, W.C.(1959), Variable metric method for minimization, Argonne Nat. Labs. report ANL-5990 Rev.

- Dennis, J.E.(1971), On the convergence of Broyden's method for nonlinear systems of equations, Math. Comp. 25, 559-567.
- Dennis, J.E.(1972), On some methods based on Broyden's secant approximation to the Hessian, in Numerical Methods for Non-linear Optimization, edited by F.A. Lootsma, Academic Press, London.
- Dennis, J. E. and J.J. Moré(1974), A characterization of superlinear convergence and its application to quasi-Newton methods, Math. Comp. 28, 549-560.
- Dixon, L.C.W.(1972a), All the quasi-Newton family generate identical points, J.O.T.A. 10, 34-40.
- Dixon, L.C.W.(1972b), Variable metric algorithms: Necessary and sufficient conditions for identical behavior on non-quadratic functions, in Numerical Methods for Nonlinear Optimization, edited by F.A. Lootsma, Academic Press, London.
- Fiacco, A.V. and G.P. McCormick(1968), Nonlinear Programming.. Sequential Unconstrained Minimization Techniques, John Wiley, New York.
- Fletcher, R. and M.J.D. Powell(1963), A rapidly convergent descent method for minimization, Comput. J. 6, 163-168.
- Fletcher, R.(1970), A new approach to variable metric algorithms, Comput. J. 13, 317-322.
- Fletcher, R.(1972), A survey of algorithms for unconstrained optimization, in Numerical Methods for Unconstrained Optimization, edited by W. Murray, Academic Press, London.
- Gill, P.E., G. Golub, W. Murray and M.A. Saunders(1974), Methods for modifying matrix factorizations, Math. Comp. 28, 505-536.



- Gill, P.E. and W. Murray(1972), Quasi-Newton methods for unconstrained minimization, J. Inst. Math. Appl. 9, 91-108.
- Goldfarb, D.(1970), A family of variable-metric methods derived by variational means, Math. Comp. 24, 23-26.
- Goldfeldt, S.M., R.E. Quandt, and H.F. Trotter(1966), Maximization by quadratic hill-climbing, Econometrics 34, 541-551.
- Goldstein, A.A.(1965), On steepest descent, SIAM J. Control 3, 147-151.
- Greenstadt, J.(1970), Variations on variable-metric methods, Math. Comp. 24, 1-18.
- Huang, H.Y.(1970), Unified approach to quadratically convergent algorithms for function minimization, J.O.T.A. 5, 405-423.
- More', J.J. and J.A. Trangenstein, On the global convergence of Broyden's method, Cornell University Computer Science Technical Report 74-216.
- Murray, W.(1972), Second derivative methods, in Numerical Methods for Unconstrained Optimization, edited by W. Murray, Academic Press, London.
- Ortega, J.M. and W.C. Rheinboldt(1970), Iterative Solution of Non-linear Equations in Several Variables, Academic Press, New York.
- Pearson, J.D.(1969), Variable metric methods of minimization, Comput. J. 12, 171-178.
- Powell, M.J.D.(1970a), A hybrid method for nonlinear equations, in Numerical Methods for Non-linear Algebraic Equations, . P. Rabinowitz, ed., Gordon and Breach, London.
- Powell, M.J.D., (1970b), A FORTRAN subroutine for solving systems of nonlinear algebraic equations, in Numerical Methods for Non-linear Algebraic Equations, P. Rabinowitz, Ed., Gordon and Breach, London.

- Powell, M.J.D.(1970c), A FORTRAN subroutine for unconstrained minimization, requiring first derivatives of the objective function, A.E.R.E., Harwell, R.6469.
- Powell, M.J.D.(1970d), A new algorithm for unconstrained optimization, in Nonlinear Programming, edited by J.B. Rosen, O.L. Mangasarian, K. Ritter, Academic Press, New York.
- Powell, M.J.D.(1971), On the convergence of the variable metric algorithm, J. Inst. Math. Appl. 7, 21-36.
- Powell, M.J.D.(1972a), Some properties of the variable metric method, in Numerical Methods for Non-linear Optimization, edited by F.A. Lootsma, Academic Press, London.
- Powell, M.J.D.(1972b), Unconstrained minimization and extensions for constraints, A.E.R.E., Harwell, T.P.495.
- Powell, M.J.D.(1973), Quadratic termination properties of minimization algorithms Parts I and II, J. Inst. Math. Appl. 10, 333-342, 343-357.
- Powell, M.J.D.(1974), Convergence properties of a class of minimization algorithms, A.E.R.E., Harwell T.R. C.S.S.8
- Schubert, L.K.(1970), Modification of a quasi-Newton method for nonlinear equations with a sparse Jacobian, Math. Comp. 24, 27-30.
- Shanno, D.F.(1970), Conditioning of quasi-Newton methods for function minimization, Math. Comp. 24, 647-656.
- Sherman, J. and W.J. Morrison(1949), Adjustment of an inverse matrix corresponding to changes in the elements of a given column or a given row of the original matrix, Ann. Math. Statist. 20, 621.
- Wilkinson, J.H.(1965), The Algebraic Eigenvalue Problem, Oxford University Press, London.

Wolfe, P.(1969), Convergence conditions for ascent methods, SIAM Review 11, 226-235.

Wolfe, P.(1971), Convergence conditions for ascent methods II: Some corrections, SIAM REVIEW 13, 185-188.

