

The Convergence of a Class of Double-rank Minimization Algorithms

2. The New Algorithm

C. G. BROYDEN

*Computing Centre, University of Essex,
Wivenhoe Park, Colchester, Essex*

[Received 18 August 1969 and in revised form 30 September 1969]

This paper presents a new minimization algorithm and discusses theoretically some of its properties when applied to quadratic functions. Results of comparative testing for a set of non-quadratic functions are described and reasons for the observed experimental behaviour are suggested.

1. Introduction

IN THE FIRST PART of this paper (Broyden, 1970) we examined a class of algorithms for minimizing the function $F(x)$ where it was assumed that the gradient of $F(x)$, denoted by $f(x)$, is available as an explicit expression. In these algorithms the function $F(x)$ is minimized at the i th iteration in the direction p_i , where p_i is given by

$$p_i = -H_i f_i \tag{1.1}$$

and where f_i is the value of the gradient at x_i . H_i is a symmetric and preferably positive definite matrix. The matrix H_i is updated at each iteration, and the updating equation involved an arbitrary parameter denoted in Part 1 by β_i . This updating equation is somewhat complicated, but its properties may be analysed more simply if $F(x)$, the function to be minimized, is quadratic. If then

$$F(x) = \frac{1}{2}x^T A x - b^T x + c, \tag{1.2}$$

where A is an n th order positive definite matrix, b an n th order vector and c a constant, we can define a matrix K_i by

$$K_i = B H_i B \tag{1.3}$$

where B is the positive definite matrix that satisfies

$$B^2 = A. \tag{1.4}$$

Now we are concerned in this paper only with changes that occur during a single iteration so we may, in order to simplify notation, omit the subscript i and replace the subscript $i+1$ by 1. With this convention it was shown in Part 1 that K is updated according to the equation

$$K_1 = K + [q \quad q_1] \begin{bmatrix} \omega & \xi \\ \xi & \eta \end{bmatrix} \begin{bmatrix} q^T \\ q_1^T \end{bmatrix} \tag{1.5a}$$

where

$$\omega = 1 - q^T K q, \tag{1.5b}$$

$$\xi = -q_1^T K q \tag{1.5c}$$

and where q and q_1 are uniquely determined orthonormal vectors. The parameter η is essentially arbitrary in that it depends upon β . It was suggested in Part 1 that a suitable choice for η would be zero since if it were negative, or large and positive, the matrix K_1 and hence H_1 might become needlessly badly conditioned. It was noted moreover that choosing η in this way gives rise to a new algorithm.

Of the two algorithms in this class already published, that due to Davidon (1959) and modified by Fletcher & Powell (1963) is obtained by putting β equal to zero and it was shown in Part 1 that this led, in general, to negative values of η . We thus expect the sequence of matrices $\{H_i\}$ obtained by that algorithm to exhibit a tendency to singularity and this tendency has been noted by, among others, Broyden (1967) and Pearson (1969). In a more recent algorithm, due to Greenstadt (1967), if H is positive definite the values of η are even more negative than those occurring in the DFP algorithm. One result of this is that for this algorithm the matrices H cannot, unlike those for the DFP algorithm, be proved to be positive definite and this has serious implications when considering numerical stability.

In this paper we show theoretically that the new algorithm is stable and we prove that it is the only member of the class considered for which a certain matrix error norm is reduced strictly monotonically when minimizing quadratic functions. We discuss the effect of rounding and of poor conditioning of H on the attainable accuracy of the final solution and conclude by presenting the results of a numerical survey in which the performance of the new algorithm for a variety of test problem is compared with that of the DFP algorithm.

2. The New Algorithm

The new algorithm is obtained by setting η equal to zero and it then follows from Part 1, equations (6.5) and (6.10a), that in order to achieve this β must be chosen to satisfy

$$\beta t z^T K z = 1. \tag{2.1}$$

This equation becomes, from Part 1, equations (2.5), (3.1), (3.4) and (3.5),

$$\beta = -1/(t p^T t) \tag{2.2}$$

which is equivalent, from equations (2.8) and (3.2b) of Part 1, to

$$\beta = 1/(t p^T y). \tag{2.3}$$

If this value of β is substituted into the general matrix updating equation (Part 1, equation (3.2)) we obtain the updating equation of the new algorithm in a form suited to computation,

$$H_1 = H + \frac{1}{p^T y} (p p^T - p y^T H - H y p^T) \tag{2.4a}$$

where

$$\rho = t + \frac{y^T H y}{p^T y}. \tag{2.4b}$$

We now prove that the new algorithm is stable in the sense of Broyden (1967).
THEOREM 1. *If t is chosen to minimize $F(x)$, H is positive definite and H_1 is given by equation (2.4) then H_1 is positive definite.*

Proof. Since β is given by equation (2.2) it follows from equation (1.1) that

$$\beta = 1/(t^T H f). \quad (2.5)$$

But since H is assumed to be positive definite both $f^T H f$ and t are positive (for proof of the latter assertion see, e.g. Broyden, 1967). Thus $\beta > 0$ and the proof follows from Broyden (1967), Theorem 8.

It follows from this theorem that the matrix update of the new algorithm is always well-defined, and hence suitable for automatic computation. Since the algorithm in addition minimizes a quadratic function in at most n iterations it possesses in theory all the properties that made the DFP algorithm so successful. Our final theorem shows that the new algorithm possesses a property not possessed by the DFP algorithm, and although it is strictly relevant only to quadratic functions it is likely, since it involves matrix error norms, to have a bearing on the more general case. Expressed crudely the final theorem states that of all the double rank algorithms considered in Part 1, the new algorithm approximates most closely to Newton's method. We first prove a lemma.

LEMMA 1. If $M = [m_{ij}]$ is a 2×2 matrix, $X = [x_1, x_2]$ and $Y = [y_1, y_2]$ where x_1, x_2, y_1 and y_2 are n th order vectors, then

$$\text{Tr}(XMY^T) = m_{11}y_1^T x_1 + m_{12}y_2^T x_1 + m_{21}y_1^T x_2 + m_{22}y_2^T x_2.$$

Proof. On expanding the matrix product we obtain

$$\text{Tr}(XMY^T) = \text{Tr}(x_1 m_{11} y_1^T + x_1 m_{12} y_2^T + x_2 m_{21} y_1^T + x_2 m_{22} y_2^T)$$

and the lemma follows from the identities

$$\text{Tr}(uv^T) \equiv v^T u \quad (2.6)$$

and

$$\text{Tr}(A+B) \equiv \text{Tr}(A) + \text{Tr}(B). \quad (2.7)$$

THEOREM 2. Let the generalized algorithm as defined in Part 1 be applied to the quadratic function (1.2) and let E be given by

$$E = K - I, \quad (2.8)$$

where K is defined by equation (1.3) and (1.4). Then at least one of the following two possibilities occurs:

$$(a) \quad e_1 = 0, \quad (2.9a)$$

where e is the difference between the solution and the i th approximation, and

$$(b) \quad \|E_1\|^2 < \|E\|^2 + \eta(\eta + 2q_1^T E q_1), \quad (2.9b)$$

where $\|\cdot\|$ denotes the Euclidean matrix norm.

Proof. Equations (1.5) and (2.8) give

$$E_1 = E - QMQ^T \quad (2.10)$$

where

$$Q = [q, q_1] \quad (2.11a)$$

and

$$M = [m_{ij}] = \begin{bmatrix} q^T E q & q^T E q_1 \\ q_1^T E q & -\eta \end{bmatrix}. \quad (2.11b)$$

Now $\|E_1\|^2 = \text{Tr}(E_1^T E_1)$ and it then follows from symmetry and the orthonormality of q and q_1 that

$$\|E_1\|^2 = \text{Tr}(E^2 - EQMQ^T - QMQ^T E + QM^2 Q^T). \quad (2.12)$$

Applying the lemma to equation (2.12) then gives, since M is symmetric,

$$\|E_1\|^2 = \|E\|^2 - 2(m_{11}q^T E q + m_{12}q_1^T E q + m_{21}q^T E q + m_{22}q_1^T E q_1) + m_{11}^2 + 2m_{12}^2 + m_{22}^2$$

which reduces from equation (2.11b) to

$$\|E_1\|^2 = \|E\|^2 - (q^T E q)^2 - 2(q^T E q_1)^2 + \eta(\eta + 2q_1^T E q_1). \quad (2.13)$$

We now consider the two possibilities:

$$(1) \quad Kq = q\lambda \quad (2.14)$$

where λ is some scalar, non-zero since K is non-singular. This implies, from the definition of q , that

$$K^2 z = Kz\lambda. \quad (2.15)$$

Now Part 1, equation (3.6), states that

$$z_1 = z - Kz t$$

but it follows from Part 1, equation (3.8) and equation (2.15) (above) that $t = \lambda^{-1}$, so that

$$z_1 = z - Kz\lambda^{-1}. \quad (2.16)$$

But since K is non-singular equation (2.15) may be pre-multiplied by K^{-1} and this gives, from equation (2.16), $z_1 = 0$. The first alternative of the theorem thus occurs.

$$(2) \quad Kq \neq q\lambda. \quad (2.17)$$

We prove that, in this case, $q^T E q_1 \neq 0$. It was shown in Part 1 (equation (4.4)) that

$$K_{i+1} z_{i+1} = K_i z_i \phi_i - K_i^2 z_i \psi_i,$$

where ϕ_i and ψ_i are scalars, and it was moreover established that $\psi_i \neq 0$.

From the definition of q_i and the subscript convention this equation may then be written

$$Kq = q\alpha + q_1\beta \quad (2.18)$$

and it follows immediately from inequality (2.17) that $\beta \neq 0$. Pre-multiplication of equation (2.18) by q_1^T gives, from equation (2.8),

$$q_1^T E q = \beta$$

so that $q_1^T E q \neq 0$. The theorem then follows immediately from equation (2.13) and the symmetry of E .

Corollary. In the new algorithm at least one of the following two possibilities occurs:

$$(a) \quad e_1 = 0,$$

$$(b) \quad \|E_1\| < \|E\|.$$

Proof. Put $\eta = 0$ in inequality (2.9b).

The relevance of this theorem and its corollary is due to the fact that for quadratic functions the norms of the successive errors e_i are bounded in terms of the spectral norm of E (Part 1, Theorem 8), and although for quadratic functions the vector errors are independent of the successive values of β_i , for non-quadratic functions it may well be advisable to choose these parameters with a view to the greatest possible reduction

of the vector error norms. This may hopefully be achieved if the matrix error norms $\|E_i\|$ are kept as small as possible, and though the spectral norms may increase when using the new algorithm the Euclidean norms are reduced strictly monotonically. It is of course possible that, for some problems, algorithms for which $\eta \neq 0$ may reduce the matrix error norms by a greater amount than the new algorithm, and we might expect the new algorithm to be inferior in these cases. On the other hand, if $\eta \neq 0$, it is possible that a very large increase in the matrix error norm might occur and in this case we would expect the use of the new algorithm to result in a substantial improvement. Since reducing the matrix error norm makes the iteration matrix look more like the inverse Jacobian we are justified to some extent in claiming that of the general class of double-rank methods, the new one is that which most resembles Newton's method. We would thus expect the performance of the new algorithm to reflect that of Newton's method, so that it might perform comparatively badly if the Jacobian matrix is singular at the solution. This did, indeed, occur in the one example attempted of this type of problem.

3. Numerical Examples

One of the principal initial difficulties in carrying out a programme of comparative testing was finding problems sufficiently difficult to reveal any significant differences between the performance of the new algorithm and the DFP algorithm. The tests commenced with four standard problems, namely Rosenbrock's, Helical Valley, Powell in four variables (for all these see Fletcher & Powell, 1963) and a problem

TABLE 1

Problem	n	N_0	ε	$m(\text{NM})$	$M(\text{DFP})$	$s(\text{NM})$	$s(\text{DFP})$
1	2	2.3_{10}^2	10^{-6}	19	23	188	239
2	3	1.9_{10}^3	10^{-6}	21	21	167	167
3	4	3.6_{10}^3	10^{-6}	26	18	231	182
4	2	1.3_{10}^2	10^{-6}	15	14	152	157

attributed to Beale (Shah, Buehler & Kempthorne, 1964). We denote these as Problems 1, 2, 3, and 4 respectively. Each problem was deemed to be solved when $\|f\| < \varepsilon$, where ε is some arbitrary tolerance, and the results of the calculations are summarized in Table 1. In the tables we denote the number of independent variables by n , the total number of gradient evaluations by m and the total number of function evaluations by s . This latter figure represents the amount of labour involved in linear minimization, a process that was carried out using a well-tried quadratic interpolation algorithm to be described more fully by Fielding (to appear). The initial value of $\|f\|$ is denoted by N_0 and the final value, as stated, is less than ε . We choose to consider $\|f\|$ rather than the function value since the former is zero at the solution whereas in general no such statement may be made about the latter.

A glance at Table 1 reveals no significant differences between the methods except for

Problem 3, where the new method is markedly inferior. In this problem the Jacobian is singular at the solution and the remarks of the previous section apply.

Problems 5-10 were cases of the trigonometric function of Fletcher & Powell (1963). The maximum number of independent variables, 45, was dictated by the size of the available computer and it is seen from Table 2 that for these problems there is nothing to choose between the two methods.

TABLE 2

Problem	n	N_0	ε	$m(\text{NM})$	$m(\text{DFP})$	$s(\text{NM})$	$s(\text{DFP})$
5	5	4.89_{10}^8	2.24_{10}^{-5}	15	15	86	83
6	10	1.86_{10}^9	3.16_{10}^{-5}	21	21	151	145
7	20	1.95_{10}^{10}	4.47_{10}^{-5}	29	29	217	225
8	30	3.21_{10}^{10}	5.48_{10}^{-5}	46	46	350	350
9	40	2.62_{10}^{11}	6.32_{10}^{-5}	53	54	382	403
10	45	1.63_{10}^{11}	6.71_{10}^{-5}	63	63	480	499

The remaining tests were attempts to fit data by a sum of exponentials, it being thought that this would combine maximum scope for testing with minimum extra programming. If (x_i, y_i) , $i = 1, 2, \dots, p$ represents p data points we define a function S by

$$S = \sum_{i=1}^p \sigma_i^2 \quad (3.1)$$

where

$$\sigma_i = y_i - \sum_{j=1}^q \alpha_j e^{-\beta_j x_i} \quad (3.2)$$

Explicit expressions for $\partial S / \partial \alpha_j$ and $\partial S / \partial \beta_j$ were then found and these formed the elements of the vector f , of order $2q$. The q values of α_j and q values of β_j formed the vector of independent variables with respect to which S is minimized. Data for these tests was obtained in three ways. For Problems 11 and 12 the values of y_i were the sum of three exponentials evaluated at 13 values of x_i but for problems 13 and 14 a freehand curve of faintly exponential character was drawn, and 17 points taken from it. The values of q for Problems 11 to 14 were 1, 3, 1 and 2 respectively, the choice of three for Problem 12 being governed by the knowledge that the data did in fact represent three exponentials. The data for Problem 15 was provided by a user, who required a sum of six exponentials fitting to 54 data points. The values of y_i were obtained experimentally. The results of the five exponential problems are summarized in Table 3.

The behaviour of the DFP algorithm in these examples was extremely interesting. It appeared to get reasonably close to the solution in only a few more iterations than required by the new algorithm, and it then proceeded to "mark time" for perhaps twenty iterations or so. A characteristic example was Problem 12. After 33 iterations $\|f\|$ had been reduced to approximately 10^{-3} , and it then hovered around this value

until iteration 60 when it was reduced to about 10^{-4} . Subsequent iterations then reduced $\|f\|$ steadily until at the 65th iteration it fell below 10^{-6} and the program terminated. This behaviour, which was repeated in a more complex manner in Problem 14, was in marked contrast to the new algorithm which converged extremely rapidly when close to the solution.

TABLE 3

Problem	n	N_0	ϵ	$m(\text{NM})$	$m(\text{DFP})$	$s(\text{NM})$	$s(\text{DFP})$
11	2	$9.9 \cdot 10^1$	10^{-6}	9	9	90	90
12	6	$5.7 \cdot 10^1$	10^{-6}	33	65	404	886
13	2	$7.9 \cdot 10^3$	10^{-6}	11	—	148	—
14	4	2.1	10^{-6}	20	61	267	912
15	12	$3 \cdot 10^3$	10^{-4}	56	more than 150	783	not recorded

One more feature of Table 3 calls for comment, the fact that the DFP algorithm did not solve Problems 13 and 15. This was due in the latter case to exceeding an arbitrary limit on the number of iterations and in the former to a value of β_j becoming negative. This resulted, for one of the larger values of x_i , in exponential overflow when evaluating σ_i .

4. Discussion

It was noted in the previous section that the performance of the new algorithm was substantially the same as that of the DFP algorithm in the initial stages of the solution of a problem, but that the characteristics of the algorithms during the final stages were markedly different. That this behaviour is not unreasonable may be inferred from a consideration of the values of β for the two algorithms. We have, in fact,

$$\beta = 0 \quad (4.1)$$

for the DFP algorithm and

$$\beta = 1/(f^T H f) \quad (4.2)$$

for the new one. At the beginning of the iteration, when the gradients are usually large, equation (4.2) implies that provided t is not too small then β may well approach zero, so that the two algorithms become effectively identical. In practice it has been found that initially β has usually been very close to zero (values of 10^{-4} have been recorded), as for example in Problem 9, where the values of x_i obtained by both algorithms were identical (to four significant figures) for the first five or so iterations.

On the other hand, as the solution is approached, β for the new method becomes extremely large (a value of 10^4 has been monitored) and the maximum discrepancy between the two methods occurs. The other eventuality that could give rise to a large value of β is severe ill-conditioning of the Hessian of $F(x)$. If this occurred in a region where $F(x)$ was strongly non-quadratic it would be possible, despite a large value of $\|f\|$, for both $\|Hf\|$ and t to be small so that in this case also the DFP algorithm would

differ markedly from the new one. Since, as was shown in Part 1, this would cause the DFP algorithm to yield a new value of H that would be much more badly conditioned than that given by the new algorithm, and since the occurrence of a near singular H would explain the observed poor performance of the DFP algorithm for badly conditioned problems, we regard the discrepancy between the two algorithms in this case as highly encouraging.

We turn now to the effect of rounding error in computing the gradient. Suppose that f is the true value of the gradient but the computed value is $f + \Delta f$. We then have, ignoring the possibility of further error,

$$p = -Hf \quad (4.3)$$

and

$$p + \Delta p = -H(f + \Delta f) \quad (4.4)$$

where p is the correct step vector, and $p + \Delta p$ is the computed one. Thus, from equation (4.3) and (4.4),

$$\Delta p = -H\Delta f. \quad (4.5)$$

Equations (4.3) and (4.5) now yield

$$\|f\| \leq \|H^{-1}\| \|p\|$$

and

$$\|\Delta p\| \leq \|H\| \|\Delta f\|$$

so that

$$\|\Delta p\| / \|p\| \leq k(H) \|\Delta f\| / \|f\|. \quad (4.6)$$

Thus the relative error in the computed step vector is bounded by the product of the relative error in the computed gradient multiplied by the condition number of H . Since $\|\Delta f\| / \|f\|$ becomes very large as the solution is approached equation (4.5) gives a rough practical guide as to the attainable accuracy. The quantity $\|\Delta f\|$ may be assessed roughly from the knowledge of the details of the computation of f and the word-length of the computer, and f of course is known. It remains therefore to estimate the condition number of H , and a lower bound of this may be determined very simply. Consider the function ϕ defined by

$$\phi = \frac{\|p\| \|f\|}{|p^T f|}$$

which since $p = -Hf$, may be written

$$\phi = \frac{\|Hf\| \|f\|}{|f^T H f|}.$$

Now since H is symmetric and positive definite its spectral norm is equal to its largest eigenvalue, say λ_{\max} . Thus

$$\|Hf\| \|f\| \leq \lambda_{\max} \|f\|^2.$$

But $f^T H f \geq \lambda_{\min} \|f\|^2$, where λ_{\min} is the smallest eigenvalue of H , and combining these inequalities gives

$$\phi \leq \lambda_{\max} / \lambda_{\min}. \quad (4.7)$$

But $\lambda_{\max} / \lambda_{\min}$ is the condition number of H , so that a lower bound of the condition number may readily be computed. The use of this bound, coupled with the inequality (4.6), is sufficient to give an idea of the attainable accuracy in a given case. It would be

of use, for instance, in deciding whether failure to converge was due to some peculiarity of the function or to an unreasonably small value of ϵ .

We do not propose to discuss the technical details of the programs here beyond saying that apart from the matrix updating routines the programs for both DFP and the new algorithm were identical. It is hoped that an ALGOL procedure embodying the new algorithm will be published in due course (Fielding, to appear) and it seems more appropriate to defer publication of the relevant computational details until then.

5. Conclusions

The experimental results quoted in Section 3 above support the view that the DFP algorithm exhibits a tendency to "mark time" for certain problems. The success of the new algorithm in overcoming this tendency in the cases tried suggests that the explanation of this behaviour given in Part 1, namely that in the DFP algorithm the matrices H_i are predisposed towards singularity, is probably correct. It is also clear that the strategy of choosing β to eliminate this tendency appears to have been largely successful.

Of the 15 cases attempted, in only one was the DFP algorithm significantly superior to the new algorithm and this case was special in that the Jacobian at the solution was singular. For all other cases the number of iterations required was either comparable, or substantially favoured the new algorithm. Indeed the DFP algorithm required more than three times as many iterations than the new algorithm to solve Problem 14, and this ratio could have been exceeded in Problem 15 where the DFP test was terminated after 150 iterations. The choice of ϵ in Problem 15 was perhaps a little unfortunate for the DFP algorithm since the value of $\|f\|$ achieved by this algorithm hovered at only slightly above ϵ for some 85 iterations. This example did, though, confirm the apparent inability of the DFP algorithm to administer the *coup de grâce* to a difficult problem.

Further perusal of the tables shows that not only is the new algorithm on the whole superior to the DFP algorithm in terms of number of iterations but it is also, if the number of evaluations of $F(x)$ per iteration is taken to be the yardstick, slightly better in terms of work done during each iteration. Indeed for Problem 4 the number of iterations for the new method is one more than required by the DFP method but the total number of evaluations of $F(x)$ is five fewer, and taking an average over the 13 problems for which full results are available we see that the mean ratio of function to gradient evaluations is 8.55 for the new method and 9.88 for the DFP method. We therefore take as a measure of the relative effectiveness of the two algorithms the total number of gradient evaluations and on this criterion the new algorithm would be regarded as better than the DFP algorithm for Problems 1, 12, 13, 14, and 15 and comparable for Problems 2 and 4-11.

It must be borne in mind that the above results represent only a limited amount of numerical experience applied to a restricted set of problems, and to this extent will not necessarily reflect the overall merit of the two algorithms. They in no way constitute a case for modifying existing programs using the DFP algorithm, simple though this may be. Nor are they substantiated by any form of convergence proof apart from those based upon the erroneous hypothesis that the function to be minimized is

quadratic, and to this extent they are unsatisfactory and incomplete. They do, though, indicate that the new algorithm may repay further consideration, especially in those cases where the problems are known to be difficult or where convergence using existing methods has been less than rapid.

The author is extremely grateful to Mr K. Fielding, of the University of Essex, for carrying out the comparative testing of the algorithms, and to Mr W. Temple of the same University for providing one of the test problems.

REFERENCES

- BROYDEN, C. G. 1967 *Maths Comput.* **21**, 368.
 BROYDEN, C. G. 1970 *J. Inst. Maths Applics* **6**, 76.
 DAVIDON, W. C. 1959 *A.E.C. Research and Development Report ANL-5990*.
 FLETCHER, R. & POWELL, M. J. D. 1963 *Comput. J.* **6**, 163.
 GREENSTADT, J. 1967 I.B.M. Scientific Center Report No. 320-2901.
 PEARSON, J. D. 1969 *Comput. J.* **12**, 171.
 SHAH, B. V., BUEHLER, R. J. & KEMPTHORNE, O. 1964 *J. Soc. ind. appl. Math.* **12**, 74.