

# The Convergence Properties of Direct Methods of Energy Minimization with Respect to Linear Coefficients in the LCAO-MO-SCF Approach

B. T. Sutcliffe\*

Kemisk Laboratorium IV, H. C. Ørsted Institutet, 2100 Copenhagen Ø, Denmark

Received August 3, 1973

It is shown that in the LCAO-MO-SCF problem, if the molecular orbital orthonormality constraints are introduced in the manner first suggested by Fletcher, then the Hessian of the problem is singular. It is suggested that this singularity may well account for the slow convergence observed using direct energy minimization methods to solve the SCF problem. Ways of avoiding the consequences of this singularity are discussed.

*Key words:* Direct energy minimization – Convergence of LCAO-MO-SCF calculations

## 1. Introduction

Recently interest has been revived in “direct” methods of minimizing the energy with respect to such parameters as nuclear position, orbital exponents and orbital (linear) coefficients, following the pioneering work of McWeeny [1, 2] using the steepest descent methods. In particular Fletcher [3] showed how it was possible to use one of the more modern conjugate-direction techniques in such direct minimization, and an approach similar to that of Fletcher was later exploited by Kari and Sutcliffe [4] and by Claxton and Smith [5]. Claxton and Smith concentrated on optimizing the linear coefficients in an unrestricted Hartree-Fock (UHF) approach for systems which had proved convergent only with difficulty using more conventional techniques. Though Claxton and Smith were able to obtain convergence using a direct method (in fact the Fletcher-Reeves method [6]), they commented that the method proceeded only very slowly and could not compete with conventional methods when these methods worked. The object of this paper is to try explain why it is that direct methods have so far proved so disappointing for linear coefficients in closed and in unrestricted LCAO-MO-SCF calculations.

## 2. Direct Methods of Minimization

Most modern direct minimization techniques are based on choosing a sequence of directions in the co-ordinate space in which the function  $f(\mathbf{x})$  is to be minimized, and finding a sequence of points by minimizing, or at least de-

\* Permanent address: Department of Chemistry, University of York, York YO1 5DD, England.

creasing, the function value along the chosen direction until a minimum point is found. The most effective of the modern methods are based on the supposition that sufficiently close to the minimum the objective function  $f(\mathbf{x})$  can be expanded in a Taylor series to second order.

$$f(\mathbf{x}) = a + \mathbf{b}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} \quad (2.1)$$

where  $\mathbf{x}$  is a column vector of co-ordinates and the matrix  $\mathbf{H}$  (the Hessian matrix at the minimum) is assumed to be a real symmetric, positive-definite, non-singular matrix. If this is possible the gradient of the function  $g(\mathbf{x})$  may be written

$$\mathbf{g}(\mathbf{x}) = \mathbf{b} + \mathbf{H} \mathbf{x}. \quad (2.2)$$

with  $g_i(\mathbf{x}) = \partial f / \partial x_i$ , and hence the minimum point  $\mathbf{a}^0$  may be found from any arbitrary point  $\mathbf{a}$  as

$$(\mathbf{a}^0 - \mathbf{a}) = -\mathbf{H}^{-1} \mathbf{g}(\mathbf{a}) \quad (2.3)$$

provided that  $\mathbf{H}$  is non-singular.

It can be shown that this problem can be solved in just  $n$  steps without having to invert  $\mathbf{H}$  or indeed without explicit knowledge even of  $\mathbf{H}$ , by constructing a sequence of conjugate-directions, that is a sequence of directions  $\mathbf{p}_i$  such that

$$\mathbf{p}_i^T \mathbf{H} \mathbf{p}_j = 0, \quad i \neq j, \quad (2.4)$$

$$\mathbf{g}_{i+1}^T \mathbf{p}_i = 0, \quad (2.5)$$

and minimizing the function along these directions. That is, if at any point  $\mathbf{x} = \mathbf{a}$ , one knows the direction  $\mathbf{p} \equiv \mathbf{p}(\mathbf{a})$ , then one constructs the function

$$F(\lambda) = f(\mathbf{a} + \lambda \mathbf{p}) \quad (2.6)$$

and finds the value of  $\lambda$ ,  $\alpha$  say, that minimizes  $F(\lambda)$  and then the next point in the descent sequence,  $\hat{\mathbf{a}}$ , is chosen according to

$$\hat{\mathbf{a}} = \mathbf{a} + \alpha \mathbf{p}. \quad (2.7)$$

It is easy to see that at the point  $\hat{\mathbf{a}}$

$$\hat{\mathbf{g}}^T \mathbf{p} = 0 \quad (2.8)$$

where

$$\hat{\mathbf{g}} \equiv \mathbf{g}(\hat{\mathbf{a}}) \quad (2.9)$$

and that therefore

$$f(\hat{\mathbf{a}}) - f(\mathbf{a}) = -(\mathbf{g}^T \mathbf{p})^2 / 2 \mathbf{p}^T \mathbf{H} \mathbf{p} \quad (2.10)$$

and that

$$\alpha = -\mathbf{g}^T \mathbf{p} / \mathbf{p}^T \mathbf{H} \mathbf{p}. \quad (2.11)$$

Many methods are available for choosing such conjugate direction, examples are the method of Powell [7] which does not use the gradient matrix, the method of Fletcher and Reeves [6] which uses the gradient matrix and Fletcher and Powell's modification of Davidon's method [8], which uses the gradient matrix

and also yields an estimate of the inverse Hessian at the minimum. A general discussion of such methods in the case of quadratic functions may be found in Huang [9] see also Dixon [10].

It is clear however from the above discussion that considerable difficulties may arise in utilizing one of these methods if the Hessian matrix is not positive definite. Thus in this case it is possible that the denominators in (2.10) and (2.11) become zero so that the location of a minimum in the direction  $\mathbf{p}$  is just not possible. Even if the methods do not fail overtly because of this, it is the case that proof of quadratic termination for the methods depends on the positive definiteness of  $\mathbf{H}$  so that one might well expect poor convergence even when the method does not fail outright.

We shall now show that the Hessian at the minimum in the closed shell SCF problem is indeed singular and we suggest that this may be the origin of the difficulties experienced by Claxton and Smith and others. The demonstration we use may be generalized immediately to the UHF problem.

### 3. The LCAO-MO Closed Shell Problem

Using the notation of McWeeny and Sutcliffe [11] the energy function in the closed shell problem may be written as

$$E = 2 \operatorname{tr} \mathbf{hR} + \operatorname{tr} \mathbf{G(R)R} \quad (3.1)$$

where  $\mathbf{h}$  is the matrix of one electron integrals and  $\mathbf{G(R)}$  the usual electron interaction matrix, both in the atomic orbital basis. The matrix  $\mathbf{R}$  is defined as

$$\mathbf{R} = \mathbf{T T}^T \quad (3.2)$$

where  $\mathbf{T}$  is the  $m$  by  $n$  matrix that relates the  $n$  doubly-occupied molecular orbitals to the chosen atomic orbital basis ( $\eta$ ).

This function as it stands is not a suitable object for use in a direct minimisation procedure since the variables of the problem, the linear coefficients,  $T_{ir}$ , are constrained by the orthonormality requirements among the molecular orbitals, namely

$$\mathbf{T}^T \mathbf{S T} = \mathbf{I} \quad (3.3)$$

where  $\mathbf{S}$  is the overlap matrix in the atomic orbital basis and  $\mathbf{I}$  is the  $n$  dimensional unit matrix.

These constraints can be incorporated, as was first shown by Fletcher [3], by writing

$$\mathbf{T} = \mathbf{Y U} \quad (3.4)$$

where  $\mathbf{Y}$  is an  $m$  by  $n$  matrix of unconstrained variables and the  $n$  by  $n$  matrix  $\mathbf{U}$  is chosen to supply the constraints. In terms of (3.3) it is seen that  $\mathbf{U}$  must satisfy the equation

$$\mathbf{U U}^T = (\mathbf{Y}^T \mathbf{S Y})^{-1} \quad (3.5)$$

and for the sake of brevity we denote  $(Y^T S Y)$  by  $A$ . It follows therefore that we may write

$$R = Y A^{-1} Y^T \quad (3.6)$$

and since the energy depends only on  $R$ , we see that it is unnecessary to specify  $U$  more closely than by (3.5).

Following Fletcher [3] (see also Kari and Sutcliffe [12]) we may determine the gradient of  $E$  with respect to the variables  $Y_{ir}$  by noting that under the change  $Y \rightarrow Y + \delta Y$  such that  $E \rightarrow E + \delta E$ , then

$$R \rightarrow R + Y A^{-1} \delta Y^T (I - SR) + (I - RS) \delta A^{-1} Y^T, \quad (3.7)$$

$$A^{-1} \rightarrow A^{-1} - A^{-1} \delta A A^{-1}, \quad (3.8)$$

where

$$\delta A = \delta Y^T S Y + Y^T S \delta Y \quad (3.9)$$

and that

$$G(R) \rightarrow G(R) + G(\delta R). \quad (3.10)$$

After a little manipulation it may be shown that for real  $Y$ , that

$$\delta E = 4 \operatorname{tr}(I - SR) f Y A^{-1} \delta Y^T \quad (3.11)$$

where

$$f = h + G(R), \quad (3.12)$$

and hence

$$\frac{\partial E}{\partial Y_{ir}} = 4 \{(I - SR) f Y A^{-1}\}_{ir} \quad (3.13)$$

so that the gradient can in this case be represented by an  $m$  by  $n$  matrix

$$V = 4(I - SR) f Y A^{-1} \quad (3.14)$$

with the row indices labelling the atomic, and the column indices the molecular orbitals.

It would thus seem that choice of the elements of  $Y$  as the variables of minimization according to (3.4) is an extremely good choice since one is able to express the gradient of the energy in a compact manner in terms of them. Furthermore, as they are peculiarly suitable for a direct minimization procedure precisely because they are unconstrained variables and so do not need to be modified to satisfy an ancillary condition at each iteration. As was pointed out by Fletcher [3] if one chooses a variable set subject to an ancillary condition which one needs to restore at the end of an iteration (for example if one chose  $T$ ) and restored orthonormality one cannot use a direct minimization process because the information from the previous iteration is "spoiled" by the restoration of constraints and so the advantages of direct minimization are lost. However the advantages of many direct minimization procedures may well be lost if the Hessian at the minimum turns out to be singular, and as we shall now show, unfortunately on the basis provided by the elements of  $Y$  the Hessian at the minimum is indeed singular.

To determine the Hessian of the problem we must find the second variation in  $E$ , and it is easy to see that under a variation  $Y \rightarrow Y + \delta Y$  we get  $V \rightarrow V + V^1$

where:

$$V^1 = 4\{-S\delta RfYA^{-1} + (I - SR)(f\delta YA^{-1} - fYA^{-1}\delta AA^{-1}) + (I - SR)G(\delta R)YA^{-1}\} \quad (3.15)$$

$$= 4\left\{(I - SR)(f\delta YA^{-1} - fYA^{-1}\delta AA^{-1}), -S(YA^{-1}\delta Y(I - SR) + (I - RS)\delta YA^{-1}Y)fYA^{-1} + \sum_{j=1}^m \sum_{r=1}^n X^{jr} \delta Y_{jr}\right\} \quad (3.16)$$

where

$$X_{is}^{jr} = \sum_{klqp=1} (I - SR)_{iq}(I - SR)_{jp}(B_{pl,qk} + B_{lp,qk})(YA^{-1})_{lr}(YA^{-1})_{ks} \quad (3.17)$$

where  $B_{pl,qk}$  is the electron-repulsion supermatrix with elements

$$B_{pl,qk} = 2\langle qp|g|kl\rangle - \langle qp|g|lk\rangle \quad (3.18)$$

where the integral notation is that of McWeeny and Sutcliffe [11].

Since we are interested only in the Hessian at the minimum, we can use the fact that at the minimum  $V = \mathbf{0}$ , to simplify (3.16) somewhat, and it is easily seen that at the minimum the second and third terms in (3.16) vanish to give

$$V^1 = 4\left\{(I - SR)(f\delta YA^{-1} - S\delta YA^{-1}YfYA^{-1}) + \sum_{j=1}^m \sum_{r=1}^n X^{jr} \delta Y_{jr}\right\} \quad (3.19)$$

that is

$$\frac{\partial^2 E}{\partial Y_{jr} \partial Y_{is}} = 4\{((I - SR)f)_{ij}A_{rs}^{-1} - ((I - SR)S)_{ij} \cdot (A^{-1}Y^T fYA^{-1})_{rs} + X_{is}^{jr}\} \quad (3.20)$$

so that the Hessian at the minimum has elements given by (3.20) where it is understood that all quantities dependent on  $Y$  in (3.20) are given in terms of the minimizing  $Y$ , though this is not explicitly indicated in the equation.

From (3.14) it is easy to see that at the minimum

$$fRS - SRfRS = \mathbf{0} \quad (3.21)$$

and that

$$SRf - SRfRS = \mathbf{0} \quad (3.22)$$

so that the matrix  $(I - SR)f$  is symmetric at the minimum and hence the Hessian itself is symmetric, as required. We can therefore write the Hessian as a partitioned matrix of dimension  $n$  by  $n$  in blocks of dimension  $m$  by  $m$ , the MO's labelling the blocks and the AO's labelling the runs and columns within a block.

The  $r, s$  block clearly has the structure

$$\begin{aligned} H^{rs} &= 4A_{rs}^{-1}(\mathbf{I} - \mathbf{SR})\mathbf{f} - 4\bar{A}_{rs}^{-1}(\mathbf{I} - \mathbf{SR})\mathbf{S} \\ &\quad + 4\mathbf{Z}^{rs} \end{aligned} \quad (3.23)$$

where

$$\bar{A} = \mathbf{A}(\mathbf{Y}^T \mathbf{f} \mathbf{Y})^{-1} \mathbf{A} \quad (3.24)$$

and

$$\mathbf{Z}_{ji}^{rs} = X_{is}^{jr}. \quad (3.25)$$

Now let us suppose that we have found a matrix  $\mathbf{T}$ , that minimises  $E$  by satisfying the usual equations

$$\mathbf{fT} = \mathbf{ST}\varepsilon, \quad (3.26)$$

$$\mathbf{T}^T \mathbf{ST} = \mathbf{I}. \quad (3.27)$$

Then we know that we can write the minimizing  $\mathbf{R}$  as  $\mathbf{TT}^T$  and the minimum  $\mathbf{f}$  as the one obtained from (3.26). Consider now the  $mn$  by 1 column matrix  $\mathbf{t}$  whose first  $m$  rows are  $\mathbf{T}_1$  whose second  $m$  rows are  $\mathbf{T}_2$ , and so on, where  $\mathbf{T}_r$  is the  $r$ 'th column of  $\mathbf{T}$ . We can then construct

$$\hat{\mathbf{t}} = \mathbf{Ht} \quad (3.28)$$

where  $\hat{\mathbf{t}}$  is a column matrix whose first  $m$  rows are

$$\hat{\mathbf{t}}^1 = \sum_{s=1}^n \mathbf{H}^{1s} \mathbf{T}_s \quad (3.29)$$

and so on. If we write out the expression for  $\mathbf{t}^r$  explicitly we get

$$\mathbf{t}^r = 4 \sum_{s=1}^n A_{rs}^{-1}(\mathbf{I} - \mathbf{SR})\mathbf{fT} - A_{rs}^{-1}(\mathbf{I} - \mathbf{SR})\mathbf{ST} + \mathbf{Z}^{rs} \mathbf{T}_s. \quad (3.30)$$

But

$$\begin{aligned} (\mathbf{I} - \mathbf{SR})\mathbf{fT}_s &= \mathbf{fT}_s - \mathbf{STT}^+ \mathbf{fT}_s \\ &= \mathbf{fT}_s - \varepsilon_s \mathbf{ST}_s \\ &= \mathbf{0}, \end{aligned} \quad (3.31)$$

$$\begin{aligned} (\mathbf{I} - \mathbf{SR})\mathbf{ST}_s &= \mathbf{ST}_s - \mathbf{STTST}_s \\ &= \mathbf{ST}_s - \mathbf{ST}_s \\ &= \mathbf{0}, \end{aligned} \quad (3.32)$$

and

$$\begin{aligned} (\mathbf{Z}^{rs} \mathbf{T}_s)_j &= \sum_{i=1}^m \mathbf{Z}_{ji}^{rs} \mathbf{T}_{is} \\ &= \sum_{i,p=1}^m X_{rs,jp} (\mathbf{I} - \mathbf{SR})_{ip} \mathbf{T}_{is} \end{aligned} \quad (3.33)$$

where

$$X_{rs,jp} = \sum_{k,l,q=1}^m (\mathbf{I} - \mathbf{SR})_{jq} (B_{pl,qk} + B_{pl,kq}) (\mathbf{YA}^{-1})_{lr} (\mathbf{YA}^{-1})_{ks} \quad (3.34)$$

so that

$$\begin{aligned}
 (\mathbf{Z}^{rs} \mathbf{T}_s)_j &= \sum_{i,p=1}^m X_{rs,jp} \left( \delta_{ip} T_{is} - \sum_1 R_{pl} S_{li} T_{is} \right) \\
 &= \sum_{p=1}^m X_{rs,jp} \left( T_{ps} - \sum_{il=1}^m \sum_{u=1}^n T_{pu} T_{lu} S_{li} T_{is} \right) \\
 &= \sum_{p=1}^m X_{rs,jp} \left( T_{ps} - \sum_{u=1}^n \delta_{us} T_{pu} \right) \\
 &= 0.
 \end{aligned} \tag{3.35}$$

We therefore conclude that  $\hat{\mathbf{t}}' = \mathbf{0}$ , so that we can write:

$$\mathbf{H}\mathbf{t} = \mathbf{0}\mathbf{t} \tag{3.36}$$

and hence we conclude that  $\mathbf{t}$  is an eigenvector of  $\mathbf{H}$  with zero eigenvalue so that  $\mathbf{H}$  is singular and not positive definite.

In fact the above demonstration may easily be extended to show that the Hessian at the minimum has precisely  $n^2$  zero roots, by the following means. We can regard the Hessian as being defined in an  $mn$  dimensional vector space and we can define a basis in this space by choosing  $n^2$  vectors  $\mathbf{t}_p$ ,  $p = 1, 2, \dots, mn$ , according to the following specification. Select one column  $\mathbf{T}_q$  from the  $n$  possible columns of  $\mathbf{T}$ . Let  $\mathbf{t}_p$  be the vector that has  $\mathbf{T}_q$  as its first  $m$  rows and is null elsewhere, let  $\mathbf{t}_{p+1}$  have  $\mathbf{T}_q$  as the second  $m$  rows and be null elsewhere and so on. It is easy to see that the  $n^2$  vectors so chosen are linearly independent since they are orthonormal in a metric specified by the matrix partitioned as is  $\mathbf{H}$  but with  $\mathbf{S}$  forming the diagonal blocks and with null blocks elsewhere. It follows at once from (3.29) to (3.34) that these  $n^2$  vectors are eigen-vectors of  $\mathbf{H}$  with null eigenvalues. The vector  $\mathbf{t}$  that we chose in equation (3.28) is of course just the linear combination of degenerate eigenvectors

$$\mathbf{t} = \mathbf{t}_1 + \mathbf{t}_{n+2} + \mathbf{t}_{2n+3} + \dots + \mathbf{t}_{n^2}. \tag{3.37}$$

#### 4. The Origin of the Zero Roots in the Hessian

Let us suppose for the moment that we had formulated our energy expression originally in terms of a set of non-orthonormal orbitals related to the atomic basis  $\eta$  by the matrix  $\mathbf{Y}$ . Then it is easy to see, using the formulae for matrix elements between determinants of non-orthogonal orbitals (see e.g. [11], p. 49–51) that the energy expression obtained is just (3.1) but now with  $\mathbf{R}$  defined directly by (3.6). Thus had we worked without any constraints at all, we would have obtained precisely the same equations as we have already for the gradient and for the Hessian and would, in consequence, have encountered precisely the same difficulties. In the light of this it is perhaps misleading to regard  $\mathbf{U}$  in Eq. (3.4) as a constraint supplying matrix, but rather as a constraint removing matrix. One can therefore regard the minimization problem we have so far formulated as the one of determining the non-orthonormal molecular orbital at any stage

and then, using the freedom that one has on the one-determinant approximation, performing a linear transformation among them to produce an orthonormal set. It would seem likely therefore that it is precisely because we have, even at the minimum, this freedom to perform an arbitrary  $n$  by  $n$  linear transformation among the solution vectors that we have a Hessian with  $n^2$  degenerate zero roots. It further seems likely that if instead of removing the constraints we had used the constraints to remove variables from the problem and hence effectively to remove the freedom to perform an arbitrary linear transformation among the solution vectors, then we should not have a singular Hessian for the problem. That this removal is, in certain circumstances, possible and does indeed have the required results is easy to see from the following example.

Let us suppose that we are working, not in an arbitrary basis of AO's, but in a basis of exact canonical MO's, numbered in order of increasing orbital energy  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_r, \dots, \varepsilon_n$ , and let us further suppose that we require the orbitals at any stage to be orthonormal. At the minimum we shall have a matrix

$$T = \begin{matrix} n \\ m-n \end{matrix} \begin{pmatrix} Q^T \\ \mathbf{0} \end{pmatrix} = YU \quad (4.1)$$

where  $Q$  is an arbitrary orthogonal matrix and if the minimization process happens to yield the canonical molecular orbitals then  $Q$  will be the unit matrix.

In the canonical MO basis we shall have:

$$S = I$$

$$f_{ij} = \varepsilon_i \delta_{ij}$$

and at the minimum

$$R_{ij} = \begin{cases} \delta_{ij}; & i, j \leq n \\ 0; & i, j > n \end{cases} \quad (4.2)$$

and after a little manipulation we can re-write (3.20) as

$$\begin{aligned} \frac{\partial^2 E}{\partial Y_{jr} \partial Y_{is}} &= 0; \quad i \text{ or } j \leq n, \quad \text{all } r, s \\ &4\delta_{ij}(\varepsilon_j A_{rs}^{-1}) - (\bar{U} \varepsilon^0 U^T)_{rs} \\ &+ \sum_{vt=1}^n (B_{jt; jv} + B_{tj; jv}) \bar{U}_{vs} \bar{U}_{ts}; \quad \text{with } i \text{ and } j > n \end{aligned} \quad (4.3)$$

where

$$\bar{U} = UQ \quad (4.4)$$

and where  $\varepsilon^0$  is the diagonal matrix composed of  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ .

If it so happens that the minimization process yields the canonical orbitals with  $U = I$  also, that is of  $Y_{ir} = T_{ir}$ , then (4.3) simplifies further so that the Hessian has elements:

$$\begin{aligned} \frac{\partial^2 E}{\partial Y_{jr} \partial Y_{is}} &= 0; \quad i, j \leq n \\ &= [4(\varepsilon_j - \varepsilon_r) \delta_{ij} \delta_{rs} \\ &\quad + 8\langle ij|g|sr\rangle - 4\langle ij|g|rs\rangle]; \quad i, j > n \\ &\quad + 8\langle ir|g|sj\rangle - 4\langle ir|g|js\rangle]. \end{aligned} \quad (4.5)$$



This matrix may easily be transformed by elementary operations to be a matrix null except for one diagonal block of side  $mn - n^2$ , hence confirming that the Hessian has just  $n^2$  zero roots.

Now let us partition  $Y$  as  $T$  is partitioned in (4.1) viz:

$$Y = \begin{matrix} n \\ m-n \end{matrix} \begin{pmatrix} Y_t \\ Y_b \end{pmatrix} \quad (4.6)$$

thus from (4.1), at the minimum we must have

$$Y_t U = Q^T, \quad (4.7)$$

$$Y_b U = 0. \quad (4.8)$$

Assuming that  $U$  is non-singular (4.8) implies that at the minimum  $Y_b$  is a null matrix, and therefore at the minimum the orthogonality constraint (3.3) becomes

$$U^T Y_t^T Y_t U = (Y_t U)^T (Y_t U) = I \quad (4.9)$$

thus at the minimum the matrix  $Y_t U$  must, by virtue of the orthogonality constraint alone, be an orthogonal matrix. This means that we can make an arbitrary initial choice of  $Y_t$ , and keep this choice fixed throughout the minimization as long as  $U$  remains defined and a minimum can be found with this choice of  $Y_t$ . It is obvious in this case that a minimum can be found if we choose initially  $Y_t$  as the unit matrix (or indeed as any orthogonal matrix) and make an initial estimate of  $Y_b$  such that  $Y_b^T Y_b$  has all its eigenvalues greater than  $-1$ . In fixing  $Y_t$  we have removed  $n^2$  variables from the problem and thus the gradient matrix now has  $mn - n^2$  elements and the Hessian matrix consists of blocks  $H^{rs}$  composed simply of the non-zero parts of (4.3). It is easy to see that with this choice all the singularities in the Hessian have been eliminated. Indeed it is fairly clear that one need not (though one may on this basis) fix all  $n^2$  elements of  $Y_t$ . It is sufficient to fix  $\frac{1}{2}n(n+1)$  of them, say by choosing  $Y_t$  initially to be a unit upper triangle, and the singularities in the Hessian will still be removed and a solution will still be possible.

In the case where we work in an arbitrary basis the situation is not so clear however precisely because we cannot say which elements of  $Y$  we may fix and still obtain a minimizing solution. Indeed it is not clear even in the very simple case in which  $Y$  consists of a single column only. In this case the problem is very similar to the problem of finding the lowest eigenvalue of a real symmetric matrix by minimizing the Rayleigh quotient, a problem which is discussed by Fletcher and Bradbury [13]. Here one can obviously avoid the singularity by fixing one element of the column at 1. This effectively confines the possible solutions of the problem to those lying on the faces of hypercube and the face of the hypercube along which minimization is to occur is determined by the choice of element set equal to 1. However it is not known in advance whether a minimum exists on this face. If a minimum does not exist on the chosen face then this should be shown up by an element other than that chosen becoming greater than 1, in which case the vector must be renormalized so that the emerging element is now 1, this

element must then be fixed and the minimization process started afresh along the newly chosen face.

In the general case the equivalent tactic would be to fix  $\frac{1}{2}n(n+1)$  values of  $Y$  and then to attempt to minimize the energy. Presumably the equivalent behaviour in the event of no minimum existing would be for the matrix  $A$  to become singular. If this occurred then one would simply have to choose a new set of fixed values that avoided this. Alternatively, of course, there is the possibility that fixing elements of  $Y$  makes the closest stationary point of the function, one other than the lowest minimum sought (an "excited state"). In this case also one would have simply to restart the process from with a new choice of fixed elements.

It is clearly not possible to specify any general strategy in respect of choosing the elements of  $Y$  to be fixed, beyond saying that of course no more than  $n$  elements and no less than one element should be fixed in any one column. It seems a case where numerical trials alone can decide whether a general strategy is possible for the problem. In this context it is worthwhile remembering that numerical experience up to now does indicate that the Fletcher-Reeves method, for example, converges quite well, even allowing all the elements of  $Y$  to vary, until one is fairly close to the minimum. A sensible strategy might therefore be to treat the full problem for a number of iterations and then to inspect the elements of the resulting  $T$  matrix and fix  $\frac{1}{2}n(n+1)$  of these in some plausible manner to define a starting point for the minimization in the space of reduced dimension.

It should also be noticed here that it is of course possible to reformulate our problem completely in terms of  $mn - \frac{1}{2}n(n+1)$  independent variables and a way of doing this has been developed, for example by Raffinetti and Ruedenberg [14]. If such a course is adopted then it seems clear that the Hessian in the basis of independent variables should be free of singularities and thus a direct method cast in terms of these variables, should work efficiently. However it would seem to be the case that it is very difficult to get analytic expressions for the gradient matrix (and even more difficult to obtain the Hessian) in such bases. This difficulty effectively confines the choice of a direct method to a non-gradient method (such as Powell's conjugate direction method). Such methods are generally found to be slower than gradient methods, so that it is not clear that the formulation in terms of independent variables should lead to a superior convergence rate in practice and whether it is more effective to use a  $Y$  matrix formulation or not therefore remains an open question. Our analysis so far however does suggest another method of tackling this problem in which we can actually utilize the presence of zero roots in the Hessian.

### 5. A Direct Method Which Avoids the Zero Root Problem

Suppose that at any stage of a minimization process we have a matrix

$$T_1 = YU, \quad T_1^T S T_1 = I_n \quad (5.1)$$

and from the matrix  $f$  evaluated at  $T_1$  we construct a matrix  $\bar{f}$  such that

$$\bar{f} = T_1^T f T_1 \quad (5.2)$$

and find a matrix  $Q_1$  which diagonalizes  $\bar{f}$  such that

$$\bar{f}Q_1 = Q_1\epsilon_1, \quad Q_1^T Q_1 = I_n \tag{5.3}$$

and define a matrix  $Z_1 = T_1 Q_1$ . Let us also invent an  $m$  by  $m - n$  matrix  $Z_2$  and define a new basis

$$\phi = (\phi_1 | \phi_2) = \eta(Z_1 | Z_2) \tag{5.4}$$

such that in this basis

$$f \rightarrow \bar{f} = \begin{pmatrix} \epsilon_1 & | & \bar{f}_{12} \\ \hline \bar{f}_{21} & | & \epsilon_2 \end{pmatrix}, \quad \bar{f}_{ij} = Z_i^T f Z_j \tag{5.5}$$

where  $\epsilon_2$  is diagonal

$$S \rightarrow \bar{S} = I_m \tag{5.6}$$

and in this basis clearly

$$T_1 \rightarrow \bar{T} = \begin{pmatrix} I_n & \\ & \mathbf{0} \end{pmatrix}. \tag{5.7}$$

It then follows immediately from (3.14) (5.5), (5.6), and (5.7) that the gradient matrix with respect to the elements of  $T$  is just

$$V = 4 \begin{pmatrix} Q \\ \hline \bar{f}_{21} \end{pmatrix}. \tag{5.8}$$

We can find the Hessian at the point  $\bar{T}$  (that is, *not* at the minimum) by noting from (3.15) that in general we must add to (3.20) a term for the  $(jr, is)$  element

$$-(SYA^{-1})_{ir} V_{js} + (SYA^{-1})_{js} V_{ir} + (VY^T S)_{ij}(A^{-1})_{rs} \tag{5.9}$$

where  $V$  is given by (3.14).

After some manipulation it can be shown that the Hessian with respect to the elements of  $\bar{T}$  at  $\bar{T}$  is blocked [cf. Eq. (3.23)] with blocks of the form

$$H^{rs} = 4\delta_{rs} \begin{pmatrix} \mathbf{0} & | & B_{rr}^{(1)} \\ \hline B_{rr}^{(2)} & | & H_r \end{pmatrix} + 4 \begin{pmatrix} \mathbf{0} & | & B_{rs}^{(1)} \\ \hline B_{rs}^{(2)} & | & H_{rs} \end{pmatrix}. \tag{5.10}$$

Where  $H_r$  consists of the non-zero diagonal terms from (4.5):

$$(H_r)_{ii} = (\epsilon_{n+i} - \epsilon_r) + 3\langle (n+i)r | g | r(n+i) \rangle - \langle (n+i)r | g | (n+i)r \rangle \tag{5.11}$$

and the matrix  $H_{rs}$  consists of the equivalent off-diagonal terms.

The matrix  $B_{rs}^{(1)}$  is null but for its  $r$ 'th row which is the  $s$ 'th row of  $\bar{f}_{12}$  and  $B_{rs}^{(2)}$  is similarly null except for its  $s$ 'th column which is the  $r$ 'th column of  $\bar{f}_{21}$ .

Now we have seen in Eq. (2.10) that it is always possible to decrease the value of locally quadratic function if we choose a direction  $p$  such that  $g^T p \neq 0$  and  $p^T H p > 0$ , even if  $H$  is not positive definite. Now if we choose our direction vector in this problem so that its  $r$ 'th group of  $m$  rows is given by

$$p_r = -\frac{1}{4} \begin{pmatrix} \mathbf{0} & | & \mathbf{0} \\ \hline \mathbf{0} & | & H_r^{-1} \end{pmatrix} V_r \tag{5.12}$$

where  $V_r$  is the  $r$ 'th column of  $V$ , then it follows at once that

$$\mathbf{g}^T \mathbf{p} = \mathbf{p}^T \mathbf{H} \mathbf{p} = -4 \sum_{r=1}^m \sum_{j=n+1}^m \bar{f}_{jr} (\mathbf{H}_r^{-1}) \bar{f}_{jr} \quad (5.13)$$

and since the elements of  $\mathbf{H}_r$  are practically certain to be positive we have achieved the desired result. It is interesting to note as an aside that (5.13) also implies that the optimum step length,  $\alpha$ , along  $\mathbf{p}$  is unity and that at a true minimum even though  $\bar{f}_{12} = 0$ , no divergence is encountered in the descent formulae (2.10) and (2.11).

We can re-write (5.12) in terms of a rectangular  $m$  by  $r$  matrix  $\mathbf{P}$  with elements

$$P_{ir} = \begin{cases} 0, & i \leq n \\ -(\bar{f}_{21})_{i-n,r} / [(\varepsilon_i - \varepsilon_r) + 3\langle ir | g | ri \rangle - \langle ir | g | ir \rangle] \end{cases} \quad (5.14)$$

and then the next point in the descent is found by constructing

$$\hat{\mathbf{T}} = \bar{\mathbf{T}} + \lambda \mathbf{P}$$

and minimizing the energy with respect to  $\lambda$ . Using (5.4) we can write the change in terms of the  $\mathbf{Z}_i$  as

$$\hat{\mathbf{Z}}_1 = \mathbf{Z}_1 + \lambda \mathbf{Z}_2 \mathbf{P}_b \quad (5.15)$$

where  $\mathbf{P}_b$  is just the non-zero part of  $\mathbf{P}$  written as an  $(m-n)$  by  $n$  matrix.

The matrix  $\mathbf{Z}_1$  does not of course satisfy the orthonormality constraints, so it should be regarded as a next estimate of  $\mathbf{Y}$  and treated accordingly to determine  $\lambda$ . It violates the constraints only by terms of order  $(P_{ir})^2$  which vanish as the minimum is reached.

The up-dating Eq. (5.15) is seen at once to be of precisely the same form as that proposed by Hillier and Saunders [15] and indeed our up-dating matrix  $\mathbf{P}_b$  differs from the up-dating matrix  $\mathbf{B}$  proposed by these authors only in the presence in  $\mathbf{P}_b$  of electron interaction terms in the denominator. Since these terms are probably in most instances small compared with the orbital energy differences, it is perhaps legitimate to neglect them and we can regard this method as being essentially equivalent to that of Hillier and Saunders, but obtained from different considerations. The implementation scheme that we would propose for the method is exactly that proposed by Hillier and Saunders except that we would recommend a numerical search to minimize the energy with respect to  $\lambda$  rather than the analytical method described by these authors since by this means orthogonality among the orbitals can be vigorously maintained at every stage.

Finally it is interesting to note that one can make a much more simple assumption about the structure of  $\mathbf{P}$  by replacing  $\mathbf{H}_r^{-1}$  in (5.12) by a multiple of the unit matrix and this still preserves the local descent properties of the method. The method in this case becomes a variant of steepest descents and in this case (5.18) is replaced by

$$\hat{\mathbf{T}} = \bar{\mathbf{T}} - \lambda \bar{\mathbf{V}} \quad (5.16)$$

and (5.19) becomes

$$\begin{aligned}\hat{\mathbf{Z}}_1 &= \mathbf{Z}_1 - 4\lambda \mathbf{Z}_2 \bar{f}_{21} \\ &= \mathbf{Z}_1 - \lambda \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{f} \mathbf{Z}_1\end{aligned}\quad (5.17)$$

and using the resolution of the identity we see that

$$\mathbf{Z}_2 \mathbf{Z}_2^T = \mathbf{S}^{-1} - \mathbf{R} \quad (5.18)$$

and therefore

$$\mathbf{Z}_1 = \mathbf{Z}_1 - 4\lambda [\mathbf{S}^{-1} - \mathbf{R}] \mathbf{f} \mathbf{Z}_1. \quad (5.19)$$

Obviously (5.19) is a much simpler formula to implement than (5.15) and as such might well be worth some numerical investigation even though steepest descent methods are not generally considered to be very effective.

## 6. Summary

We have demonstrated in this article that the Hessian matrix at the minimum with respect to the linear coefficients can be singular in the closed shell LCAO-MO-SCF problem, for a particular and highly convenient way of choosing the coefficients. We suggest that this is the reason why only slow convergence has been achieved using such methods as the Fletcher-Reeves method, in minimizing the LCAO-MO energy. The generalization of our results to the UHF problem is immediate and obvious.

We have suggested a method by which the singularities may be avoided by constraining the coefficients and which may well therefore have superior convergence properties. We have also established the status of the method of Hillier and Saunders among direct methods as a quasi-Newton method and suggested a new steepest descent method in the spirit of Hillier and Saunders' method.

It would seem likely from our analysis that the singularities in the Hessian arise because of the freedom one has in the closed shell case to perform a linear transformation among the occupied orbitals while leaving the energy invariant. One may perhaps conjecture that such singularities may well arise in all problems where one has such a freedom but not in problems where this freedom is absent.

## References

1. McWeeny, R.: Proc. Roy. Soc. A **235**, 496 (1956)
2. McWeeny, R.: Rev. Phys. **32**, 335 (1960)
3. Fletcher, R.: Mol. Phys. **19**, 55 (1970)
4. Kari, R., Sutcliffe, B. T.: Chem. Phys. Letters **7**, 149 (1970)
5. Claxton, T., Smith, W.: Theoret. Chim. Acta (Berl.) **22**, 399 (1971)
6. Fletcher, R., Reeves, C. M.: Comput. J. **7**, 149 (1964)
7. Powell, M. J. D.: Comput. J. **7**, 155 (1964)
8. Fletcher, R., Powell, M. J. D.: Comput. J. **6**, 163 (1963)
9. Huang, H. Y.: J. Opt. Th: Appl. **5**, 405 (1970)
10. Dixon, L. C. W.: Math. Programming **2**, 383 (1972)

11. McWeeny, R., Sutcliffe, B. T.: *Methods of molecular quantum mechanics*. London, New York: Academic Press 1969
12. Kari, R., Sutcliffe, B. T.: *Int. J. Quant. Chem.* **7**, 459 (1973)
13. Bradbury, W., Fletcher, R.: *Num. Math.* **9**, 259 (1966)
14. Raffanetti, R., Ruedenberg, K.: *Int. J. Quant. Chem.* **35**, 625 (1971)
15. Hillier, I., Saunders, V.: *Proc. Roy. Soc. A* **320**, 161 (1970/71)

Dr. B. T. Sutcliffe  
Department of Chemistry  
University of York  
York YO1 5DD/England