

THE CYCLIC COLORING PROBLEM AND ESTIMATION OF SPARSE HESSIAN MATRICES*

THOMAS F. COLEMAN† AND JIN-YI CAI†

Abstract. Numerical optimization algorithms often require the (symmetric) matrix of second derivatives, $\nabla^2 f(x)$. If the Hessian matrix is large and sparse, then estimation by finite differences can be quite attractive since several schemes allow for estimation in much fewer than n gradient evaluations.

The purpose of this paper is to analyze, from a combinatorial point of view, a class of methods known as substitution methods. We present a concise characterization of such methods in graph-theoretic terms. Using this characterization, we develop a complexity analysis of the general problem and derive a roundoff error bound on the Hessian approximation. Moreover, the graph model immediately reveals procedures to effect the substitution process optimally (i.e. using fewest possible substitutions given the differencing directions) in space proportional to the number of nonzeros in the Hessian matrix.

Key words. graph coloring, estimation of Hessian matrices, sparsity, differentiation, numerical differences, NP-complete problems, unconstrained minimization

AMS(MOS) subject classifications. 65K05, 65K10, 65H10, 68L10

1. Introduction. We are concerned with the estimation of a large sparse symmetric matrix of second derivatives $\nabla^2 f(x)$ for some problem function $f: R^n \rightarrow R^1$. In particular, we note that the product $\nabla^2 f(x) \cdot d$ can be estimated, for example, by forward differences

$$(1.1) \quad \nabla^2 f(x) \cdot d = [\nabla f(x+d) - \nabla f(x)] + o(\|d\|).$$

When the structure of $\nabla^2 f(x)$ is known, then usually a few well chosen differencing directions d_1, \dots, d_p affords the recovery of estimates of all nonzeros of $\nabla^2 f(x)$. Let us denote our estimate by H . We will assume that the sparsity pattern of H is known; the diagonal elements are specified as nonzero; H is symmetric. (Restricting the diagonal to be zero-free is reasonable in many contexts: In particular, a minimizer of f usually possesses a positive definite Hessian matrix.) We will be concerned with methods that use differencing directions d_1, d_2, \dots, d_p that are based on a *partition* of columns C_1, \dots, C_p . In particular, let S_k denote the set of columns in group C_k and let h_i be the steplength associated with column i , $i = 1, \dots, n$. Finally, define

$$(1.2) \quad d_k = \sum_{i \in S_k} h_i e_i$$

for $k = 1, \dots, p$, where e_i is the i th column of the identity.

There has been considerable work recently concerned with this problem, especially with trying to make p as small as possible. Curtis, Powell, and Reid [1974] suggested a method, *CPR*, for the unsymmetric problem. Their idea was to build groups of *structurally independent* columns in a left-to-right greedy fashion. (Two columns (vectors) x, y are *structurally independent* if $x_i * y_i = 0$, for all i .) It is easy to see that such a p -partition allows for the estimation of a matrix with p differencing directions. Specifically, let C_1, \dots, C_p be a partition of the columns of H where each group consists of structurally independent columns. Then, if $[\nabla f(x+d_k) - \nabla f(x)]_i \neq 0$ it follows that there is exactly one column j in group C_k with H_{ij} a designated nonzero

* Received by the editors October 10, 1984, and in revised form April 24, 1985. This work was supported in part by the Applied Mathematical Sciences Research Program (KC-04-02) of the Office of Energy Research of the U.S. Department of Energy under contract DE-AC02-83ER13069.

† Computer Science Department, Cornell University, Ithaca, New York 14853.

and we can assign

$$H_{ij} \leftarrow \frac{[\nabla f(x + d_k) - \nabla f(x)]_i}{h_j}.$$

Coleman and Moré [1983] analyzed and modified this method by taking a combinatorial point of view. In particular, a *column intersection graph* can be formed by associating with each column i of H a node v_i and defining an edge between node v_i and node v_j iff there is an index k such that both H_{ki} and H_{kj} are nonzeros. A p -coloring of this graph is an assignment, ϕ , of "colors" to nodes such that if there is an edge between node v_i and node v_j then $\phi(v_i) \neq \phi(v_j)$. It is not hard to see that a p -coloring of this graph induces a valid partition of structurally independent columns and vice versa.

Coleman, Garbow, and Moré [1984] have developed FORTRAN 77 codes based on this work. Such (unsymmetric) methods can be applied to the symmetric problem (McCormick [1983] discusses the complexity of this approach); however it is probably worthwhile using symmetry when it is present.

Powell and Toint [1979] were the first to try to exploit symmetry. They pointed out that symmetry can be used both in a direct and an indirect fashion. A direct method is one in which each unknown of H is determined independently of the others. More specifically, let C_1, \dots, C_p be a partition of the columns of H . Since each off-diagonal nonzero is represented twice, it is no longer necessary that each group consist of structurally independent columns. It is necessary, however, that for each nonzero (i, j) either column i resides in a group C_r such that no other column in this group has a nonzero in row j or column j resides in a group C_s such that no other column in this group has a nonzero in row i . If the latter condition were true then H_{ij} would be determined

$$H_{ij} \leftarrow \frac{[\nabla f(x + d_s) - \nabla f(x)]_i}{h_j}$$

and $H_{ji} \leftarrow H_{ij}$. Clearly a similar (symmetric) rule would hold for the former condition.

Coleman and Moré [1984] analyzed such methods from a combinatorial point of view and produced a simple graph-theoretic characterization of all partitions that can be used to induce a direct symmetry-exploiting determination of H . Let us represent the structure of H by the usual adjacency graph $G(H) = (V(H), E(H))$. That is, if H is a symmetric matrix of order n , then $V(H)$ consists of n vertices v_1, \dots, v_n (associate column i of H with vertex v_i) and $E(H)$ consists of pairs of vertices (edges) where $(v_i, v_j) \in E(H)$ if and only if $H_{ij}(H_{ji})$ is considered a nonzero. A p -partition of the columns of H , C_1, \dots, C_p can be viewed as an assignment of colors, ϕ , to the nodes of G , $\phi: V \rightarrow \{1, \dots, p\}$. This assignment is a p -coloring if $(v, w) \in E \Rightarrow \phi(v) \neq \phi(w)$. A *path p -coloring* is a p -coloring with the additional stipulation that every path in G of length 4 (distinct) vertices uses at least 3 colors. The characterization of direct symmetric methods given by Coleman and Moré is simply

THEOREM 1.1. *The mapping ϕ is a path p -coloring if and only if ϕ induces a partition of the columns of H consistent with direct determination.*

Note: We have changed the notation used by Coleman and Moré [1984]; here we use "path coloring" instead of "symmetric coloring" because in our context the term path coloring is more appropriate.

This characterization led to a deeper understanding of the direct estimation problem on symmetric structures which in turn yielded a complexity analysis and algorithmic possibilities.

Indirect estimation of symmetric matrices may be preferable because fewer groups (i.e. differencing directions) will be needed, in general. Powell and Toint concentrated on substitution methods where directions are chosen so that nonzeros can be determined via a substitution process. (They restricted their attention, as we do, to substitution methods based on a partition of columns.) So in this case there is interdependence of the matrix unknowns (nonzeros) to the degree that an underlying lower triangular system is defined. Powell and Toint proposed an algorithm to determine the differencing directions and then solve for the unknowns (lower triangular substitution method (LTS)). Subsequently, Coleman and Moré [1984] analyzed this process from a combinatorial point of view. This analysis led to a modified and empirically superior procedure (the resulting FORTRAN 77 code is described in Coleman, Garbow, and Moré [1985]). However, a simple insightful characterization, in the vein of Theorem 1.1, was not provided.

The purpose of this paper is to provide such a characterization. This result is as simple as Theorem 1.1 and is clearly the analogous result. This view provides enormous insight into the combinatorial nature of the problem as well as suggesting algorithmic possibilities. Furthermore, the graph theoretic interpretation reveals that if a partition of columns allows for the recovery of H via a substitution process, then it is always possible to do so efficiently. In particular, every unknown can be solved for in (roughly) less than $n/2$ substitutions and the space required to compute H is proportional to the number of nonzeros. This is somewhat surprising since the Powell-Toint procedure relies heavily on a regular matrix structure produced by LTS which is not present for an arbitrary feasible partition. Finally, the graph model allows one to derive a growth of error bound for a general substitution method, which is essentially analogous to the result achieved by Powell and Toint for a specific method, LTS.

Section 2 will provide the characterization of substitution methods followed by a roundoff error discussion in § 3. In § 4 we establish the complexity of the problem and discuss its combinatorial relationship to the symmetric direct problem (path coloring). Section 5 deals with algorithms for effecting the substitution process in space proportional to $|E|$ (i.e. the number of nonzeros of H). Finally, observations on parallelism are provided in § 6.

2. Substitution methods and cyclic coloring. A partition of columns of a symmetric matrix induces a substitution method if there is an ordering of the matrix unknowns such that all unknowns can be solved for, in that order, using symmetry and previously solved elements. This notion is fully general (subject to the partition restriction) but seems to be a difficult one to work with. There is, however, a very elegant and simple graph theoretic interpretation. The major purpose of this section is to present this characterization.

First it is necessary to formalize the concept of a substitution method in matrix terms. Let U be the set of indices of matrix unknowns (identify (i, j) with (j, i)) and suppose that U is ordered: $U = \{(i_k, j_k)\}$. Let the columns of H be partitioned $\{C_1, \dots, C_p\}$ and define

$$(*) \quad S_0 = \emptyset, \quad S_k = S_{k-1} \cup \{(i_k, j_k)\}, \quad 1 \leq k \leq |U|.$$

The ordering induces a substitution method iff

either j_k belongs to a group C , say, and if l is any other column in C with a nonzero in row i_k then $(i_k, l) \in S_{k-1}$ or i_k belongs to a group C' , say, and if l' is another column in C' with a nonzero in row j_k then $(j_k, l') \in S_{k-1}$.

The essence of this statement is that, at the k th step, it is possible to solve for element (i_k, j_k) or, equivalently (j_k, i_k) , by substitution. We call a partition, for which there exists such an ordering, *substitutable*. For example, if H is a tridiagonal matrix, then it is easy to verify that the partition $(\{1, 3, \dots\}, \{2, 4, \dots\})$ is substitutable.

Obviously there are substitutable partitions for any symmetric matrix. For example, every partition consistent with a path coloring is substitutable. Alternatively, a partition that induces a "lower triangular substitution method" is substitutable. (Coleman and Moré [1984] and Powell and Toint [1979] discussed lower triangular substitution methods.) However, here we are interested in minimizing the number of groups in a general substitutable partition. The above 2 examples are restrictive in that they consider only particular classes of substitutable partitions. The general problem is

Partition problem. Obtain a substitutable partition of the columns of a given symmetric matrix H with the fewest groups.

How difficult is the partition problem? This is a hard question to answer considering the rather clumsy matrix formalization of a substitution method. Fortunately a substitutable partition has a simple expression in the language of graphs.

DEFINITION. A mapping $\phi: V \rightarrow \{1, 2, \dots, p\}$ is a *cyclic p -coloring* of G if ϕ is a p -coloring and if ϕ uses at least 3 colors in every cycle of G .

As the following theorem indicates, we now have a simple characterization of a substitutable partition.

THEOREM 2.1. Let H be a symmetric matrix with a nonzero diagonal. The mapping ϕ induces a substitution method if and only if ϕ is a cyclic coloring of $G(H)$.

Before providing the proof, let us consider an informal argument based on the following example. Let the adjacency graph of H , $G(H)$ be as shown in Fig. 1. Both assignments of the colors r, s, t are valid colorings but assignment 1 is *not a valid cyclic coloring*: the cycle v_1, v_2, v_3, v_4 uses only 2 colors. Assignment 2 is a valid cyclic coloring. The edges (off-diagonal nonzeros) can be determined by considering each pair of colors in turn. For example, consider the subgraph, $F_{r,s}$, induced by the nodes colored r or s as shown in Fig. 2. Edges $(1, 8)$, $(3, 7)$, and $(3, 9)$ can all be determined immediately. Consider for example edge $(3, 7)$. Column 3 has a nonzero in row 7 and resides in group C_r . There is no other column in group C_r with a nonzero in row 7 (else node 7 would have another incident r -node). Therefore, $H_{7,3}$ (hence $H_{3,7}$) can be determined directly. Once $(3, 7)$ and $(3, 9)$ are determined, edge $(2, 3)$ can be computed: column 2 has a nonzero in row 3 and resides in group C_s . Columns 7 and 9 are the other columns in group C_s with nonzeros in row 3. However, $H_{3,7}$ and $H_{3,9}$ are now known quantities; hence, $H_{3,2}$ can be computed with 2 substitutions.

Clearly the process can be carried to completion until every edge in $F_{r,s}$ is determined. (It is easy to see that the diagonal elements can be directly determined: this follows from the fact that ϕ is a coloring.) But every pair of colors induces a forest, otherwise ϕ would not be a cyclic coloring, and therefore every nonzero can be determined by considering each pair of colors in turn.

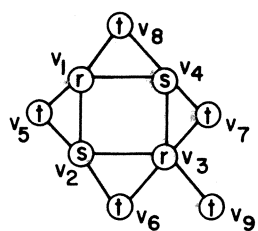


FIG. 1

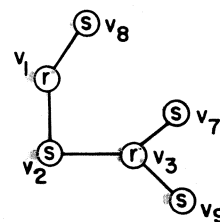
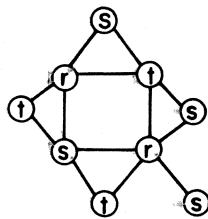


FIG. 2

The "if" part of Theorem 2.1 is proved along the lines of the example given above. The "only if" part is perhaps a bit surprising but not difficult to prove.

Proof of Theorem 2.1. First we prove that every cyclic coloring of $G(H)$ induces a substitution method. Since ϕ is a coloring, every diagonal element of H can be computed. In particular, if $\phi(v_j) = r$, then j is the only column in group C_r with a nonzero in row j (otherwise ϕ is not a coloring). But then (1.1) and (1.2) yield

$$H_{jj} = \frac{[\nabla f(x + d_r) - \nabla f(x)]_j}{h_j}.$$

Consider next $(i, j) \in U$, $i \neq j$. Suppose that column i is in group C_r and column j is in group C_s . Clearly, since ϕ is a coloring, $r \neq s$. Consider the subgraph induced by the nodes colored r and the nodes colored s , say $F_{r,s}$. Since ϕ is a cyclic coloring, $F_{r,s}$ contains no cycle and therefore is a forest. The edges in $F_{r,s}$ correspond to off-diagonal unknowns of H . They can be solved, or ordered, independent of the rest of the unknowns of H . In particular, each leaf-incident edge can be solved directly since there is no conflict. We can now "delete" all such edges and consider each new leaf-incident edge. Each such edge is now incident to known edges and can therefore be solved. Clearly the process can be repeated until an edge-less graph remains. The entire procedure can now be repeated for each pair of colors until every unknown is determined.

We now show that if ϕ induces a substitution method, then ϕ is a cyclic coloring. First it is clear that ϕ must be a valid coloring, otherwise the diagonal elements would not be determined. To see this, suppose that $(i, j) \in U$, and v_i and v_j are assigned the same color r . Hence both column j and column i are in the same group, C_r . Since column j belongs only to group C_r , it follows from (*) that H_{jj} can be determined only after either H_{ij} or H_{ji} is determined. Similarly, H_{ii} can be determined only after either H_{ij} or H_{ji} is determined. But the determination of one of H_{ij} , H_{ji} must be preceded by the determination of one of H_{ii} , H_{jj} , by (*), which is a contradiction.

Suppose then that ϕ is a coloring but is not a cyclic coloring. Hence there must be a cycle, with at least 4 edges, colored with just 2 colors, say r, s . Let (i, j) be the first edge in this cycle to be solved (ordered) and let us assume, without loss of generality, that v_i is colored r and v_j is colored s . Let node v_i be incident also to node v_h (on the cycle) and let v_j be incident also to node v_k (on the cycle), as illustrated in Fig. 3. But (i, j) cannot be determined from group C_s , because columns j and h both reside in this group with nonzeros in row i (and (i, h) is not yet known (ordered)). Similarly, (j, i) is not determined from group C_r , because columns i and k both reside in this group with nonzeros in row j (and (j, k) is not yet known (ordered)). Therefore no edge in this cycle can be solved first and ϕ cannot induce a substitution method. \square

Hence the partition problem is equivalent to the

Cyclic coloring problem: Obtain a minimum cyclic coloring of $G(H)$.

Note that once we have found a cyclic coloring of $G(H)$, then the coloring induces a substitutable partition and the corresponding ordering of U is available, as the proof of Theorem 2.1 indicates. A tridiagonal matrix provides a simple example. The graph

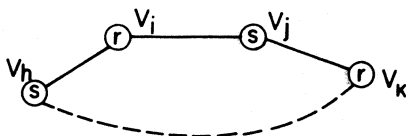


FIG. 3

is shown in Fig. 4 and a valid cyclic coloring is provided by assigning r to the even nodes and s to the odd nodes. The diagonal elements can be solved directly and the off-diagonal elements are obtained via substitution: edges $(1, 2)$ and $(n-1, n)$ are obtained first (directly), followed by $(2, 3)$ and $(n-2, n-1)$, with 1 substitution each, and so on. The middle edge will be the last determined element with approximately $\frac{1}{2}n$ dependencies or substitutions.

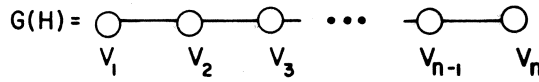


FIG. 4

Suppose we modify the above example by adding an edge from node v_1 to node v_n . A cyclic coloring would then require 3 colors; for example, we could use our previous assignment except we apply a new color, t , to node 1. Now $(1, 2)$ and $(1, n)$ can be determined directly (or, ordered first) and the remaining elements can be determined, as before, via substitution.

3. Substitution methods and roundoff error. The above two examples raise an interesting question with numerical significance: Is there a limit to the number of dependencies or substitutions? The amount of computational work as well as the potential growth of roundoff error depends, in part, on this number; therefore, a bound tighter than the total number of nonzeros in H may be consequential. Powell and Toint [1979] established that for a particular class of substitution methods, triangular substitution methods, the bound is $n-2$. The cyclic coloring characterization leads us immediately to a more general result. Every unknown can be determined by considering the forest induced by a particular pair of colors. But each forest can have at most $n-1$ edges and therefore we have the following result.

THEOREM 3.1. *Let ϕ be a substitutable partition. Then, each unknown in H is dependent on at most $n-2$ other unknowns.*

Clearly this result is the best possible worst case upper bound, if we allow any possible feasible ordering of the unknowns or edges. To see this, just consider the tridiagonal case: if the edges are solved from one end of $G(H)$ to the other, then the last edge requires $n-2$ substitutions. However, certain orderings are preferable over others. For example, in the tridiagonal case one can achieve a bound of $\lfloor \frac{1}{2}(n-2) \rfloor$ if each edge is solved by substituting from the nearest end of $G(H)$. It is not hard to see that, over different orderings, this is the best possible worst case upper bound; again, just consider the tridiagonal case.

Is it possible to order the unknowns, in general, so that the maximum number of substitutions is less than or equal to $\lfloor \frac{1}{2}(n-2) \rfloor$? In order to answer this question, consider when it is feasible, during the solution process, to solve for edge $l \triangleq (x, y)$ in $T_{r,s}$ where $T_{r,s}$ is a tree in the forest induced by the colors r and s . Note that if an edge (x, y) is removed from a tree (but the nodes x, y are not removed) then two subtrees remain: Let

$$T_{r,s}^l(x) = (V_{r,s}^l(x), E_{r,s}^l(x)), \quad T_{r,s}^l(y) = (V_{r,s}^l(y), E_{r,s}^l(y))$$

represent the two subtrees, rooted at x and y respectively, that remain when edge l is removed from $T_{r,s}$: consider Fig. 5. It is clear that (x, y) is ready to be solved *if and only if* either every edge in $T_{r,s}^l(x)$ is solved or every edge in $T_{r,s}^l(y)$ is solved. Furthermore, (x, y) requires at least

$$\text{mincost}(x, y) \triangleq \min \{ |E_{r,s}^l(x)|, |E_{r,s}^l(y)| \}$$

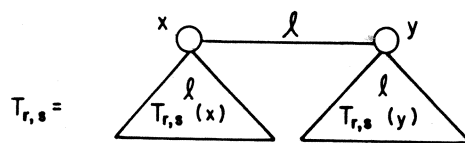


FIG. 5

substitutions. Note that an ordering that computes edge (x, y) using $\text{mincost}(x, y)$ substitutions, for each edge (x, y) , is optimal and requires less than $\lfloor \frac{1}{2}(n-2) \rfloor$ substitutions for each edge.

Such an ordering is possible and is provided by the following algorithm. Let $T \triangleq (V_T, E_T)$ be the tree under consideration, with $|V| = n_T \leq n$.

ALGORITHM *solve_tree*

```

 $T_1 \triangleq (V_1, E_1)$  where  $V_1 = V_T, E_1 = E_T$ 
for each vertex  $v \in V_1$  do  $\text{value}(v) \leftarrow 0$  endo
for  $i = 1$  to  $|E_1|$  do
  choose a leaf  $x_i$  of  $T_i$ , of smallest  $\text{value}$ 
  let  $y_i$  be the vertex such that  $(x_i, y_i) \in E_i$ 
   $\text{value}(y_i) \leftarrow \text{value}(y_i) + \text{value}(x_i) + 1$ 
  solve  $(x_i, y_i)$ 
   $E_{i+1} \leftarrow E_i - \{(x_i, y_i)\}$ 
   $V_{i+1} \leftarrow V_i - \{x_i\}$ 
   $T_{i+1} \leftarrow (V_{i+1}, E_{i+1})$ 
endo
end solve_tree

```

THEOREM 3.2. *Algorithm solve_tree solves for each edge $(x, y) \in E_T$ using the fewest possible substitutions, mincost (x, y) . Hence, solve_tree requires at most*

$$\max \{ \text{mincost}(x, y) : (x, y) \in E \} \leq \lfloor \frac{1}{2}(n_T - 2) \rfloor \leq \lfloor \frac{1}{2}(n - 2) \rfloor$$

substitutions to determine any edge of E_T .

Proof. First we establish that algorithm *solve_tree* terminates: Since T_i is a tree and x_i is a leaf in T_i , it follows that T_{i+1} is a tree. Therefore T_{i+1} will have a leaf (indeed at least 2) and x_{i+1} will be found. It follows that the algorithm will determine every edge.

Assume then that *solve_tree* does not solve for each edge using the fewest possible substitutions. In particular suppose that at step i vertex x_i is the chosen leaf in T_i and edge $l \triangleq (x_i, y_i)$ will be determined nonoptimally. That is, $|E^l(x_i)| > |E^l(y_i)|$ and hence $|E^l(x_i)| = \text{value}(x_i) > \lfloor \frac{1}{2}(n_T - 2) \rfloor$, since $|E^l(x_i)| + |E^l(y_i)| = n_T - 2$.

Consider a leaf, v_i , in $T^l(y_i)$ (there must be at least 1). But

$$|E^l(x_i)| > \lfloor \frac{1}{2}(n_T - 2) \rfloor \Rightarrow |E^l(y_i)| \leq \lfloor \frac{1}{2}(n_T - 2) \rfloor$$

and therefore,

$$\text{value}(v_i) \leq \lfloor \frac{1}{2}(n_T - 2) \rfloor < \text{value}(x_i).$$

Therefore (x_i, y_i) would not be chosen at step i , a contradiction. \square

In summary, Theorem 3.2 says that for an arbitrary substitutable partition algorithm *solve_tree* will compute each edge with the fewest possible substitutions (with respect to that partition) and that number is always bounded by $\lfloor \frac{1}{2}(n-2) \rfloor$.

We conclude § 3 by considering the accuracy of the estimated Hessian matrix in more detail. We will show that an error bound, similar to that achieved by Powell and

Toint [1979] for a particular class of algorithms (lower triangular substitution methods) holds for any substitution method provided the unknowns are solved for in the manner suggested by *solve_tree*.

Every substitutable partition with p groups, or cyclic p -coloring, allows for the recovery of the matrix unknowns via a back substitution process provided the differencing vectors are consistent with the coloring ϕ . In particular, let S_k denote the set of nodes (columns) colored k (i.e. in C_k) and again define

$$d_k = \sum_{i \in S_k} h_i e_i$$

where h_i is the step-length associated with column i . Let x be a given point in R^n and define $u_k = \nabla f(x + d_k) - \nabla f(x)$, for $k = 1, \dots, p$. If H denotes the approximation to $\nabla^2 f(x)$, then, since ϕ is a coloring and by (1.2),

$$H_{jj} \cdot h_j = \nabla f(x + d_k)_j - \nabla f(x)_j \quad \text{for } j \in S_k$$

and therefore, since every column belongs to a group, every diagonal element can be determined. The diagonal approximations are defined by these equations, and will not participate in any subsequent calculations. Indeed such equations usually guide the choice of h_j : h_j is chosen to balance truncation and roundoff errors in order to approximate the diagonal elements as accurately as possible (e.g. Gill, Murray, Saunders, and Wright [1983]).

Our previous analysis has shown that it is only necessary to consider 2 colors (directions) at a time when solving for the off-diagonal elements. Let us concern ourselves then with a tree, $T_{r,s}$, induced by colors r and s . Let $u_r = \nabla f(x + d_r) - \nabla f(x)$, $u_s = \nabla f(x + d_s) - \nabla f(x)$ and let

$$\hat{u}_r = u_r + \varepsilon_r, \quad \hat{u}_s = u_s + \varepsilon_s$$

denote the computed quantities (i.e. contaminated with rounding error).

The solution process is provided by algorithm *solve_tree* with the statement “*solve* (x_i, y_i)” expanded, to read

$$(3.1) \quad H_{ij} \cdot h_j = (\hat{u}_c)_i - \sum_{k \in N(i)} H_{ik} \cdot h_k$$

where we identify vertex x_i with index i , and vertex y_i with index j . $N(i)$ is the set of neighbours of node x_i in $T_{r,s}(x_i)$, and $c = \phi(y_i)$, which is one of r, s (for brevity we will write $T_{r,s}(x_i)$ instead of $T_{r,s}^{l_i}(x_i)$ where $l_i = (x_i, y_i)$). In other words, when H_{ij} is solved for, every other element in row i of columns in group C_c has already been solved for; the right-hand side of (3.1) is adjusted accordingly.

Following Powell and Toint, we define the error matrix F to be $H - \nabla^2 f$ and let

$$(3.2) \quad (\delta_c)_i = (\hat{u}_c)_i - \sum_{k \in N(i) \cup \{j\}} (\nabla^2 f(x))_{ik} \cdot h_k$$

In other words, $(\delta_c)_i$ measures the difference between the computed quantity $(\hat{u}_c)_i$ and the ideal $(\nabla^2 f(x) \cdot d_c)_i$. Hence $(\delta_c)_i$ is a composite of roundoff and truncation errors. If we assume that the second derivatives of f are Lipschitz continuous, then a standard bound is obtained:

$$\eta \triangleq \max_{c,i} \{ |(\delta_c)_i| \} \leq C \cdot \max_k \{ |h_k|^2 \} + \max_{c,i} \{ |(\varepsilon_c)_i| \}$$

where C is a positive constant.

The following result establishes a bound on the elements in the error matrix F .

THEOREM 3.3. *If H is obtained by algorithm `solve_tree` (with “`solve` (x_i, y_i)” effected by 3.1) then*

$$\begin{aligned} |F_{ij}| &\leq (|E_{r,s}(x_i)| + 1) \cdot \eta \cdot \max_{i,j,k} \left\{ \frac{|h_k|}{|h_i h_j|} \right\} \\ &\leq (\lfloor \frac{1}{2}n \rfloor) \cdot \eta \cdot \max_{i,j,k} \left\{ \frac{|h_k|}{|h_i h_j|} \right\} \end{aligned}$$

where again we identify column i with node x_i , column j with node y_i , and $\{\phi(x_i), \phi(y_i)\} = \{r, s\}$.

Proof. Combining (3.1), (3.2) and the definition of F yields

$$(\delta_c)_i = F_{ij} \cdot h_j + \sum_{k \in N(i)} F_{ik} \cdot h_k$$

which implies the bound

$$\begin{aligned} |F_{ij} h_i h_j| &\leq |(\delta_c)_i \cdot h_i| + \sum_{k \in N(i)} |h_i h_k F_{ik}| \\ &= |h_i (\delta_c)_i| + \sum_{k \in N(i)} |F_{ki} h_k h_i|. \end{aligned}$$

But this same decomposition can be applied, recursively, to each $F_{ki} h_k h_i$, for $k \in N(i)$, to yield

$$(3.3) \quad |F_{ij} h_i h_j| \leq \sum_{w \in V_{r,s}(x_i)} |(\delta_c)_w \cdot h_w|$$

where $T_{r,s}(x_i) = (V_{r,s}(x_i), E_{r,s}(x_i))$. Since the tree $T_{r,s}(x_i)$ has $(|E_{r,s}(x_i)| + 1)$ nodes, the result follows immediately from (3.3) and Theorem 3.2. \square

One can conclude from this result that the growth of roundoff error is quite limited if the steplength does not vary greatly in size. On the other hand, if there is significant variance (recall that stepsizes are chosen to accurately approximate diagonal elements) then this result may allow for unacceptable growth of error: a direct method may be preferable.

Indeed recent experiments by Coleman, Garbow, and Moré [1985] support the conclusion that unacceptable error pollution can occur when stepsizes vary noticeably in size. The resulting matrix is essentially unuseable as an approximation to the Hessian. This suggests that an automatic monitoring process that switches from an indirect method to a direct method, when necessary, might be useful. Unfortunately we have no specific suggestions at the moment as to what quantities to monitor. (Of course it is always possible to estimate the Hessian by both an indirect and a direct method, occasionally, and compare the resulting matrices.)

4. The cyclic chromatic number. How difficult is the cyclic coloring problem? We address this question in this section. The reader who is unfamiliar with the fundamentals of complexity theory and NP-completeness is urged to consult the excellent resource book *Computers and Intractability: A Guide to the Theory of NP-Completeness*, by Michael R. Garey and David S. Johnson [1979].

We will first consider the cyclic coloring decision problem (CCDP) and show that this problem is NP-complete: we do this by transforming the general graph coloring decision problem (CDP). We then conclude that the corresponding optimization problem, the cyclic coloring problem, is NP-hard. The consequence of this result is just this: if we could solve the cyclic coloring problem in polynomial time (P -time) then we could also solve the graph coloring problem in P -time (as well as a host of other

“intractable” problems). Since this is deemed highly unlikely, an expedient approach to our problem is to investigate efficient heuristic and approximation schemes (we discuss this in § 5).

It is common, when considering complexity questions related to discrete optimization problems, to consider the decision problem formulation. In this case we have the

Cyclic coloring decision problem (CCDP): Given an integer $p \geq 3$ and an arbitrary graph G , is it possible to assign a cyclic p -coloring to the nodes of G ?

We have excluded the simple cases $p = 1, 2$ since it is easy to see that polynomial algorithms exist for such cases. The following theorem shows that CCDP is not so simple for $p \geq 3$.

THEOREM 4.1. *CCDP is NP-complete.*

Proof. The first step is to show that CCDP is in the class NP. In particular, we must show that we can validate, in P -time, whether or not a particular assignment of p colors is indeed a cyclic coloring. To do this, one must merely consider each pair of colors, in turn, and decide whether or not the induced graph is a forest. Clearly this is a polynomial time operation.

We now proceed to transform the general coloring problem (CDP), which is known to be NP-complete, to CCDP. Consider an arbitrary graph $G = (V, E)$ and integer $p \geq 3$. Let $|V| = n$ and $|E| = m$. We construct a new graph, $G' = (V', E')$ as follows. For each edge $e_l = (v_i, v_j) \in E$, define a bipartite graph G'_l with vertices

$$\{v_i, v_j, w_1^{(l)}, \dots, w_p^{(l)}\}$$

and edges

$$(v_i, w_k^{(l)}), \quad (v_j, w_k^{(l)}), \quad k = 1, \dots, p.$$

Graphically, this transformation is shown in Fig. 6.

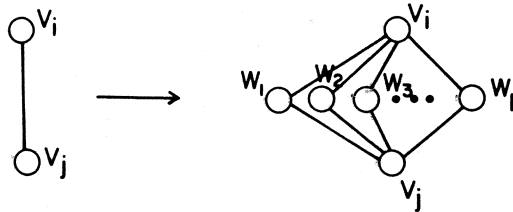


FIG. 6

Now define a bipartite graph G' by setting

$$V' = V(G) \cup \{w_k^{(l)} : 1 \leq k \leq p, 1 \leq l \leq m\}$$

and

$$E' = \bigcup_{l=1}^m E(G'_l).$$

We now show that if G can be p -colored using an assignment ϕ , then G' can be assigned a cyclic p -coloring, ϕ' . In particular, for each $v \in V$ let $\phi'(v) = \phi(v)$. Hence if we consider any G'_l , induced by $e_l = (v_i, v_j) \in E$, then $\phi'(v_i) \neq \phi'(v_j)$. Let ϕ' assign vertices $w_j^l, j = 1, \dots, p$ any color different from $\phi'(v_i)$ and $\phi'(v_j)$.

We claim that ϕ' is a cyclic p -coloring of G' : Clearly any cycle in G' must contain a path $(v_i, w_k^{(l)}, v_j)$ for some $1 \leq l \leq m$ and $1 \leq k \leq p$ where $(v_i, v_j) \in E$. But ϕ' assigns 3 colors to each such path and hence every cycle uses at least 3 colors. Moreover, the transformation from G to G' can obviously be done in P -time.

Finally we show that if G' can be assigned a cyclic p -coloring, then G can be p -colored. Assume that ϕ' is a cyclic coloring of G' . Define

$$\phi: \phi(v_i) = \phi'(v_i), \quad 1 \leq i \leq n.$$

We claim that ϕ is a p -coloring of G . Suppose instead that $\phi(v_i) = \phi(v_j)$ where $e_i = (v_i, v_j) \in E$. Then ϕ' must assign a different color to each w_k^i , $1 \leq k \leq p$; otherwise, there is a bi-colored cycle in G' . But it follows that ϕ' uses at least $p+1$ colors, a contradiction. \square

The proof above has actually established a stronger result than indicated by the statement of Theorem 4.1, since the constructed graph G' is bipartite.

COROLLARY 4.2. *The cyclic coloring decision problem on bipartite graphs is NP-complete.*

Since the *cyclic coloring problem* is the optimization version of CCDP, it follows that it cannot be an easier problem. Hence the cyclic coloring problem is NP-hard (even if we restrict our attention to bipartite graphs).

There is a marked similarity between the proof given above and the NP-completeness proof provided by Coleman and Moré [1984] with respect to the path coloring problem (symmetric *direct* problem). Indeed it turns out that the transformation given above will also establish that the path coloring decision problem is NP-complete. (Recall: $\phi: V \rightarrow \{1, 2, \dots, p\}$ is a *path p -coloring* of a graph G if ϕ is a p -coloring and if ϕ is not a 2-coloring for any path in G of length 3 edges).

We conclude this section with a short discussion on the relationship between path colorings and cyclic colorings. Let $\chi(G)$, $\chi_\pi(G)$, $\chi_0(G)$ denote the chromatic number, the path chromatic number, and the cyclic chromatic number of graph G , respectively. That is, $\chi(G)$ is the smallest integer p such that G has a p -coloring. Similarly, $\chi_\pi(G)$ [$\chi_0(G)$] is the smallest integer p such that G has a path p -coloring [cyclic p -coloring].

The first observation is that a path coloring is a cyclic coloring. To see this, consider any cycle O in G and suppose that ϕ is a path coloring. Clearly if O has only three edges then, since ϕ is a coloring, O must be assigned 3 colors. If O has more than 3 edges, then O contains a path connecting 4 distinct vertices and hence at least 3 colors are assigned by ϕ . Therefore,

$$(4.1) \quad \chi_0(G) \leq \chi_\pi(G)$$

for any graph G .

Of course a cyclic coloring is not necessarily a path coloring: a cycle O of arbitrary large circumference needs only 1 vertex to be assigned a third color and effect a valid cyclic 3-coloring however this assignment is not a valid path coloring in general. This raises an interesting question: how large can $\chi_\pi(G)/\chi_0(G)$ be? This ratio can be arbitrarily close to 2 for band graphs; however we have been unable to prove (or disprove) that this is an upper bound. It seems reasonable to hypothesize that 2 is an upper bound because a band graph G is, in a certain sense, the worst possible graph for path coloring and the best possible graph for cyclic coloring. Specifically, the first and the last inequalities become equalities in (4.2) below for all band graphs sufficiently large. (It is easy to verify the first equality, and Coleman and Moré [1984] proved the latter.)

Finally, since every cyclic coloring of G is a coloring of G , and every coloring of G^2 is a path coloring of G , we can stretch both ends of (4.1) to get

$$(4.2) \quad \chi(G) \leq \chi_0(G) \leq \chi_\pi(G) \leq \chi(G^2).$$

Note that a partition that induces a direct method that ignores symmetry is equivalent

to a coloring of G^2 and has at least $\chi(G^2)$ groups (Coleman and Moré [1983]); a partition that induces a direct method that uses symmetry is equivalent to a path coloring of G and has at least $\chi_\pi(G)$ groups; a partition that induces a substitution method is equivalent to a cyclic coloring and has at least $\chi_0(G)$ groups. One final comment on (4.2): each inequality can be made strict by choosing appropriate graphs.

5. Algorithms. The NP-completeness result of the previous section indicates that an efficient heuristic, or approximation scheme, is required. In particular, since it is not crucial that the absolute *fewest* groups be found (though it is desirable), we are willing to settle for an efficient procedure that produces near optimal results in practise. Indeed, such procedures have been suggested by Powell and Toint [1979] and Coleman and Moré [1984]. Furthermore, Coleman and Moré report extensive experimental results. In this section we will interpret such procedures in the light of the new characterization described in this paper. In addition, we will discuss an important computational concern: Given a substitutable partition, is it possible to recover the matrix unknowns (i.e. solve for the edges) in an amount of space proportional to the number of matrix unknowns (i.e. the number of edges)?

We wish to obtain a cyclic coloring of $G(A)$ using few colors. Since efficient heuristic approaches to the ordinary graph coloring problem (i.e. no cyclic restriction) are available, a natural approach is to transform our problem to a general graph coloring problem. In particular, consider adding edges to the given graph $G = (V, E)$ to obtain a completed graph $\bar{G} = (V, \bar{E})$ such that a coloring of \bar{G} is a cyclic coloring of G . Consider the following

```

ALGORITHM add_edge
  let  $\pi: V \rightarrow \{1, \dots, n\}$  be an invertible map, initialize  $\bar{E}$  to be the set  $E$ 
  for  $i = n, \dots, 2$  do
    if  $v_j, v_k$  are neighbours of  $\pi^{-1}(i)$  in  $G$  and  $\pi(v_j), \pi(v_k) < i$  then
       $\bar{E} \leftarrow E \cup \{(v_j, v_k)\}$ 
    endif
  endo
end add_edge

```

To see that *add_edge* does the job, consider any cycle O in G . Let v_i be the vertex of largest value π on O and let v_j, v_k denote the neighbours of v_i on O . Clearly $(v_j, v_k) \in \bar{E}$ and hence O will need at least 3 colors when \bar{G} is colored.

It is clear that the initial ordering π will affect the resulting graph \bar{G} and consequently the number of colors used. For example, if G is the wheel graph on 9 vertices shown in Fig. 7, and if the center vertex is ordered last, then \bar{G} is a complete graph and requires 9 colors. On the other hand, if the center vertex is ordered first, and the outer vertices are ordered sequentially, then \bar{G} is constructed from G by adding an edge between v_2 and v_9 ; \bar{G} requires just 4 colors in this case.

A successful heuristic labelling rule, suggested by Powell and Toint, is the following. Assume that the vertices $\pi^{-1}(n), \dots, \pi^{-1}(n-k)$ have been found. Choose as the vertex to be ordered $n-k-1$, the vertex of smallest degree in $G - \{\pi^{-1}(n), \dots, \pi^{-1}(n-k)\}$. This algorithm is known as the smallest last ordering (*slo*)

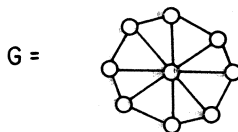


FIG. 7

and has a number of interesting properties. For further information consult Coleman and Moré [1984] and Matula and Beck [1983].

The algorithms of Coleman and Moré [1984] and Powell and Toint [1979] both implicitly perform *add_edge/slo* followed by a \bar{G} -coloring step. Here they differ: the latter authors apply colors in a greedy fashion by considering the nodes in the given order (i.e. $\pi^{-1}(1), \dots, \pi^{-1}(n)$), Coleman and Moré apply a greedy algorithm over several different (cleverly chosen!) orderings. It has been proven that the coloring problem restricted to the class of graphs derived from the *add_edge/slo* completion process is NP-complete. However, an important question remains: Does an *optimal* coloring of such a completed graph always solve the cyclic coloring problem? If the answer is yes then one may conclude that it is not necessary to consider algorithms outside this framework. The answer is no.

To see that the cyclic coloring problem may not be solved by an *optimal* coloring of a completion produced by *add_edge*, consider Fig. 8.

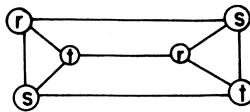


FIG. 8

The assignment shown is a valid cyclic 3-coloring; however, since every vertex is of degree 3 it follows that a coloring of a graph completed by algorithm *add_edge* will use at least 4 colors regardless of the ordering, π , of the nodes.

This example suggests that it may be worthwhile investigating heuristic algorithms for the partition/substitution problem, based on the cyclic coloring characterization perhaps, but not of the *add_edge* variety. At this point we do not know whether there is a practical gain to be made; the answer lies with further experimentation.

One final observation before discussing the solution process: A slight modification of the algorithm *add_edge* yields a procedure for the path coloring problem (symmetric direct method). In particular, change the conditional to read

“if v_j, v_k are neighbours of $\pi^{-1}(i)$ in G and $\pi(v_k) < i$ then”

and it follows that a color assignment of \bar{G} is a path coloring of G . (To our knowledge this heuristic has not been suggested or experimented with previously.) If the conditional is further changed to read

“if v_j, v_k are neighbours of v_i in G then”

it follows that a coloring of \bar{G} is a coloring of G^2 (and hence is also a path coloring of G).

An important computational concern is this: Given that $\phi: V \rightarrow \{1, \dots, p\}$ is a cyclic coloring, is it possible to compute the actual matrix elements in space proportional to $|E|$? It turns out that we can answer this question in the affirmative without imposing any additional structure on ϕ . In particular we do not assume that ϕ is necessarily consistent with the algorithm *add_edge/slo* (Coleman, Garbow, and Moré [1985] discuss, in detail, a FORTRAN 77 implementation of *add_edge/slo*, followed by a graph coloring step. Their substitution process operates in space $O(|E|)$; however it relies heavily on the regular matrix structure produced by *add_edge/slo*).

We will assume that H is stored as a sparse matrix and hence the space required is $O(|E|)$ where it is assumed that $|E| \cong n$. We will not discuss time complexity here since it is very difficult to present a convincing argument without discussing detailed

data structure and implementation requirements: we prefer to provide this analysis in a subsequent paper describing a specific implementation along with numerical results. Our purpose here is to support the claim that excessive space is not required. This discussion can be given at a fairly abstract level.

The first job is to determine $Hd_j = \nabla f(x + d_j) - \nabla f(x) \triangleq u_j$ and to save the significant information (nonzeros), for $j = 1, \dots, p$, where p is the number of colors used by ϕ . The vector d_j must be consistent with the color j : if S_j is the set of columns (nodes) in the j th group (color) then $d_j = \sum_{i \in S_j} h_i e_i$, where h_i is the steplength associated with column i . Assume that we have the vector u_j on tap. If $(u_j)_i$ is a nonzero then this quantity is stored as follows:

```

for each  $k$  such that  $H_{ik}$  is a nonzero do
  if  $\phi(v_k) = j$  then  $H_{ik} \leftarrow (u_j)_i$  endif
endo

```

We note that it is not really necessary to replicate the information in H as we have done here; however, not doing so requires a more complicated indexing scheme than we wish to describe here. The key point here is that the vector pairs (d_j, u_j) are processed sequentially and so only $2n$ space is required.

When the process is complete, H is fully assigned but the numbers do not correspond to the actual Hessian quantities: we must now effect a substitution process.

For any pair of colors r, s , we can extract a bi-colored tree, $T_{r,s}$ from the representation of G and store $T_{r,s}$ as a tree structure in space $O(n)$. Algorithm *solve_tree* can now be used: we need only be more specific about step *solve* (x_i, y_i) . The idea is simply to effect (3.1) with the knowledge that the difference results u_j are stored in H (as indicated above). In particular, *solve* (x_i, y_i) should read

$$H_{ij} \leftarrow \frac{H_{ij} - \sum_{k \in N(i)} H_{ik} \cdot h_k}{h_j},$$

$$H_{ji} \leftarrow H_{ij}$$

where $N(i)$ is the index set of neighbours of vertex x_i in $T_{r,s}(x_i)$ (i.e. all neighbours of x_i in $T_{r,s}$ except y_i). The reason this works is that when H_{ij} is solved, all other elements in row i (of columns in the same group as column j) have already been resolved.

It follows that the space required to resolve all unknowns is $O(|E|)$.

6. Concluding remarks. We have analyzed a class of methods for estimating sparse Hessian matrices, namely, substitution methods. In particular we have shown that there is an easy and elegant graph theoretic characterization of all substitution procedures based on a partition of columns of the symmetric matrix H . This characterization has allowed for a rich understanding of the combinatorial nature of the problem: we have analyzed the complexity of the partition problem, as well as suggested efficient procedures to effect the substitution process.

We have restricted our attention, in this paper, to substitution procedures based on a *partition* of columns. Indeed this is more restrictive than need be: Powell and Toint [1979], in their example (5.3), demonstrated that allowing the assignment of a column to several groups can reduce the number of required gradient evaluations. This example is particularly interesting because the solution procedure remains a "substitution process" requiring no matrix factorization. However, since the procedure allows a column to belong to several groups, it is not a method based on a partition of columns and does not belong in the class of substitution methods considered in this paper.

Indeed, a more general scheme than even this is possible provided matrix factorizations are acceptable. Newsam and Ramsdell [1983] have explored this general "elimination" option (Coleman [1984] summarizes this idea on page 49). While such methods may occasionally yield a reduction in the number of gradient evaluations, it is not clear that they provide a net benefit, in general, since they require the solution of n square dense (but relatively small) systems of equations to recover the true information.

Two other works should be mentioned. Thapa [1984] has also suggested a direct/partition method for estimating sparse Hessian matrices. Goldfarb and Toint [1984] have proposed specific (optimal) substitution procedures for specific common "mesh structures": such procedures are, of course, efficient algorithms for obtaining and using optimal cyclic colorings for particular regular structures.

We end with a comment on parallelism. There is a high degree of parallelism in the Hessian estimation problem. Specifically, each estimation Hd_j , $j = 1, \dots, p$ can be done independently, and thus in parallel. Since this work is sometimes the dominant expense in a numerical problem, exploiting this concurrency may be quite profitable. Note that the number of processors would usually be quite modest, even for large problems, since a cyclic coloring typically uses much fewer than n colors. Moreover, the substitution process also allows for parallel computation: each bi-colored tree can be processed entirely independently of the others.

Acknowledgment. We are grateful to Jorge Moré whose many suggestions have improved this paper. Specifically, he deserves credit for the elegant 6-vertex graph in § 5 that replaces our original graph of 25 vertices.

REFERENCES

- T. F. COLEMAN [1984], *Large Sparse Numerical Optimization*, Lecture Notes in Computer Sciences 165, Springer-Verlag, New York.
- T. F. COLEMAN AND J. J. MORÉ [1983], *Estimation of sparse Jacobian matrices and graph coloring problems*, SIAM J. Numer. Anal., 20, pp. 187-209.
- [1984], *Estimation of sparse Hessian matrices and graph coloring problems*, Math. Programming, 28, pp. 243-270.
- T. F. COLEMAN, B. GARBOW AND J. J. MORÉ [1984], *Software for estimating sparse Jacobian matrices*, ACM Trans. Mathematical Software, 10, pp. 329-347.
- [1985], *Software for estimating sparse Hessian matrices*, Technical Report 43, Argonne National Laboratory, Argonne, IL. Also available as TR 85-660, Dept. of Computer Science, Cornell Univ., Ithaca, NY.
- A. R. CURTIS, M. J. D. POWELL AND J. K. REID [1974], *On the estimation of sparse Jacobian matrices*, J. Inst. of Math. Appl., 13, pp. 117-119.
- M. R. GAREY AND D. S. JOHNSON [1979], *Computers and Intractability, A Guide to the Theory of NP-Completeness*, W. H. Freeman, San Francisco.
- P. E. GILL, W. MURRAY, M. A. SAUNDERS AND M. WRIGHT [1983], *Computing forward-difference intervals for numerical optimization*, SIAM J. Sci. Statist. Comp., 4, pp. 310-321.
- D. GOLDFARB AND PH. L. TOINT [1984], *Optimal estimation of Jacobian and Hessian matrices that arise in finite difference calculations*, Math. Comput., 43, pp. 69-88.
- D. W. MATULA AND L. L. BECK [1981], *Smallest-last ordering and clustering and graph coloring algorithms*, J. Assoc. Comput. Mach., 30, pp. 417-427.
- S. T. MCCORMICK [1983], *Optimal approximation of sparse Hessians and its equivalence to a graph coloring problem*, Math. Programming, 26, pp. 153-171.
- G. N. NEWSAM AND J. D. RAMSDELL [1983], *Estimation of sparse Jacobian matrices*, this Journal, 4, pp. 404-418.
- M. J. D. POWELL AND PH. L. TOINT [1979], *On the estimation of sparse Hessian matrices*, SIAM J. Numer. Anal., 16, pp. 1060-1074.
- M. N. THAPA [1984], *Optimization of unconstrained functions with sparse Hessian matrices—Newton-type methods*, Math. Programming, 29, pp. 156-186.