



On the Successive Projections Approach to Least-Squares Problems

J. E. Dennis, Jr., Trond Steihaug

SIAM Journal on Numerical Analysis, Volume 23, Issue 4 (Aug., 1986), 717-733.

Stable URL:

<http://links.jstor.org/sici?sici=0036-1429%28198608%2923%3A4%3C717%3AOTSPAT%3E2.0.CO%3B2-B>

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

SIAM Journal on Numerical Analysis is published by Society for Industrial and Applied Mathematics. Please contact the publisher for further permissions regarding the use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/siam.html>.

SIAM Journal on Numerical Analysis
©1986 Society for Industrial and Applied Mathematics

JSTOR and the JSTOR logo are trademarks of JSTOR, and are Registered in the U.S. Patent and Trademark Office. For more information on JSTOR contact jstor-info@umich.edu.

©2002 JSTOR

ON THE SUCCESSIVE PROJECTIONS APPROACH TO LEAST-SQUARES PROBLEMS*

J. E. DENNIS, JR.† AND TROND STEIHAUG‡

§

Abstract. In this paper, we suggest a generalized Gauss–Seidel approach to sparse linear and nonlinear least-squares problems. The algorithm, closely related to one given by Elfving (1980), uses the work of Curtis, Powell, and Reid (1974) as extended by Coleman and Moré (1983) to divide the variables into nondisjoint groups of structurally orthogonal columns and then projects the updated residual into each column subspace of the Jacobian in turn. In the linear case, this procedure can be viewed as an alternate ordering of the variables in the Gauss–Seidel method. Preliminary tests indicate that this leads quickly to cheap solutions of limited accuracy for linear problems, and that this approach is promising for an inexact Gauss–Newton analog of the inexact Newton approach of Dembo, Eisenstat, and Steihaug (1982).

Key words. sparse nonlinear least squares, inexact Gauss–Newton, finite-difference Jacobians, SOR iteration

AMS(MOS) subject classifications. 65F10, 32-04

1. Introduction. The purpose of this paper is to suggest a generalized group Gauss–Seidel approach to sparse least-squares problems. This approach appears in preliminary tests to be useful in obtaining cheap solutions of limited accuracy. For nonlinear problems, the approach combines the work on linearization of Coleman and Moré (1982), (1983), developing ideas of Curtis, Powell, and Reid (1974), with the work on inexact Newton methods of Dembo, Eisenstat, and Steihaug (1982) and Steihaug (1980). This amalgam leads here to an especially convenient implementation of the inexact Gauss–Newton method suggested and analyzed in § 3.

The basic idea of this paper is simple. If we group the columns of a coefficient matrix by the methods developed for use in estimating sparse Jacobian matrices, then each group of columns is mutually orthogonal in their zero/nonzero structure regardless of the particular values of the nonzero entries. This has the advantage of leading for each group of columns to smaller cheaper least-squares problems unaffected by conditioning. The idea is to cycle through these groups solving each time with the right-hand side updated as in iterative refinement.

The graph-theoretic algorithms that produce these groupings actually partition the columns of the coefficient matrix A into disjoint groups. Notice that there is no reason to want nonintersecting groups of columns, and it will be interesting in our context to consider overlap between groups. We report some experiments in § 4 with a simple heuristic scheme for overlapping orthogonal sets of columns and for groups whose normal matrices are tridiagonal. It appears in our tests as though overlapping is sometimes worthwhile in the nonlinear case, but not worth the extra arithmetic to

* Received by the editors October 11, 1983, and in revised form June 12, 1985.

† Mathematical Sciences Department, Rice University, Houston, Texas 77251. The research of this author was sponsored by DOE DE-AS05-82ER13016, ARO DAAG-29-83-K-0035, and NSF MCS81-16779. This work was supported in part by the International Business Machine Corporation, Palo Alto Scientific Center, Palo Alto, California.

‡ Statoil, Stavanger, Norway and Mathematical Sciences Department, Rice University, Houston, Texas 77251. The research of this author was sponsored by DOE DE-AS05-82ER13016 and ARO DAAG-29-83-K-0035. This work was supported in part by the International Business Machine Corporation, Palo Alto Scientific Center, Palo Alto, California.

solve the tridiagonal systems for linear problems. On the other hand, our tests suggest that using the column grouping schemes to provide orderings for point SOR does seem promising in the linear case.

In the nonlinear case, maximal overlap gives diagnostic advantages in the computation of finite-difference Jacobian approximations. This is a potentially useful point that has not to our knowledge been pointed out and applies equally in the solution of square systems of nonlinear equations. It would require no additional vector function evaluations to use a different perturbation for a given variable for each group to which its column belongs. These values could then be used heuristically to refine the partial derivative approximation and estimate its accuracy. Such an investigation is outside the scope of this paper and so we do not consider it further here.

We do not claim to bypass the effects of ill-conditioning by using any of the grouping strategies suggested here, although we may make less use of extended precision than factorization techniques do in obtaining a given accuracy. It is true that some of the subproblems we solve are unaffected by conditioning because they are diagonal, but a linear problem with two nearly-dependent columns shows what happens in the ill-conditioned case and illustrates the algorithm. We first subtract from the right-hand side, say r^0 , its projection onto the subspace spanned by the first column a_1 . Now we subtract from the new residual r^1 its projection onto the subspace spanned by a_2 . The algorithm repeats the process and the reader will see that in order to reduce to an acceptable level the part of the residual that can be accounted for by a linear combination of a_1 and a_2 , we will need more iterations if the angle between a_1 and a_2 is small. If a_1 were orthogonal to a_2 , then one projection onto each would suffice to reduce the residual as much as possible. This is exactly the feature that we are exploiting within each column group.

Since our algorithm is based on column groupings, the reader might think that for large least-squares problems it would be at a disadvantage compared to row-oriented schemes. This is not the case either here or in the nonlinear equations problem for that matter. See Coleman and Moré (1983, p. 208), for an argument that just the opposite is likely to be the case. There are many schemes that can be used to adapt the ideas here to any segmentation of the problem into sets of rows with the elements in each set stored columnwise. Of course, the corresponding sets of rows of the residual must be calculable separately.

In this paper, we will find it useful to follow Young (1971) in distinguishing between group iterative methods and block iterative methods. He regards a group method as one based on grouping the row indices, or equation numbers, into disjoint subsets, or groups, whose set union is the complete set of row indices. If the groupings are done in such a way that the first group is $\{1, 2, \dots, n_1\}$, the second is $\{n_1 + 1, n_1 + 2, \dots, n_2\}$, etc., then he calls the corresponding method a block method. In Young's terminology, a group method can be made into a block method by simply renumbering the equations. This is not true for the methods we will consider, because we will consider intersecting groups in which some indices may belong to every group. Thus, we will call the methods here generalized group methods rather than generalized block methods.

In § 2, we will present the algorithm for the linear problem and show that it generalizes the block or group SOR method in the sense that it is that method for the normal equations of an extended linear least-squares problem. This observation will allow an easy convergence proof, even in the singular case, using results from Keller (1965). Section 3 is a discussion of the nonlinear least-squares problem that combines the algorithm of § 2 applied to the linearized problem with the inexact Gauss-Newton

approach. Section 4 presents some numerical results for several heuristic column grouping schemes as well as for the associated orderings of the columns for point SOR.

2. The algorithm for the linear least squares problem. Let A be an m by n real matrix, $m \geq n$, $b \in \mathbb{R}^m$, and consider the least-squares problem

$$(2.1) \quad \min_{x \in \mathbb{R}^n} \|Ax - b\|_2.$$

In order to illustrate our algorithm, assume that the columns of A are divided into g groups A_1, A_2, \dots, A_g , where A_i is a m by n_i submatrix and A_i may share columns with A_j . Let $x_1 \in \mathbb{R}^{n_1}, x_2 \in \mathbb{R}^{n_2}, \dots, x_g \in \mathbb{R}^{n_g}$. The least-squares problem (2.1) can now be written as:

$$(2.2a) \quad \min \{ \|A_1x_1 + A_2x_2 + \dots + A_gx_g - b\|_2 : x_i \in \mathbb{R}^{n_i}, i = 1, 2, \dots, g \}.$$

Note that (2.2a) is really an m by $n = n_1 + n_2 + \dots + n_g$ least-squares problem which we will denote using boldface type as:

$$(2.2b) \quad \min_{x \in \mathbb{R}^n} \|Ax - b\|_2,$$

where $A = (A_1|A_2|\dots|A_g)$ is an $m \times n$ matrix and $x = (x_1, x_2, \dots, x_g)^T$ is an n vector. Clearly, A has exactly the same set of distinct columns as A divided into the same groups, but it will be useful that we can ignore overlaps if we view the A_i as column groups of A . It will also be convenient to have the notation $\bar{x}_i \in \mathbb{R}^n$ and $\bar{x}_i \in \mathbb{R}^n$ to denote respectively the vectors obtained by starting with the n or n zero vectors and placing the nonzero entries of x_i in the positions corresponding to the column indices in A or A of A_i . Thus, given either $x \in \mathbb{R}^n$, or x_1, x_2, \dots, x_g with each $x_i \in \mathbb{R}^{n_i}$, we can define $x \in \mathbb{R}^n, x = \bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_g$ and write

$$Ax = A\bar{x}_1 + A\bar{x}_2 + \dots + A\bar{x}_g = A_1x_1 + A_2x_2 + \dots + A_gx_g$$

or

$$x = \bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_g,$$

$$Ax = A\bar{x}_1 + A\bar{x}_2 + \dots + A\bar{x}_g = A_1x_1 + A_2x_2 + \dots + A_gx_g.$$

Suppose we have an approximation x^k to a solution x^* to (2.1), and we divide x^k into $x_1^k, x_2^k, \dots, x_g^k$ as above. (This division will not be unique for components corresponding to column overlaps.) Then (2.2) suggests the following successive replacements iteration.

FOR $i = 1, 2, \dots, g$ DO

$$\text{Solve for } x_i^{k+1}: \min \left\{ \left\| \sum_{j=1}^i A_jx_j^{k+1} + \sum_{j=i+1}^g A_jx_j^k - b \right\|_2 : x_i^{k+1} \in \mathbb{R}^{n_i} \right\}.$$

This is a method of projections (Householder (1964)) and for the special case of nonoverlapping column groups, i.e., for formulation (2.2b), this particular iteration was suggested by Elfving (1980). We will see easily below that this is block Gauss-Seidel on the normal equations for (2.2b) and so it is a generalized group Gauss-Seidel for (2.1). Björk and Elfving (1979) and Elfving (1980) have pointed out that in this form of the Gauss-Seidel iteration, we do not need to explicitly form the normal equations.

Let s^k be the step or correction, let $r^k = Ax^k - b$ be the residual, and notice that

$$A_1x_1^{k+1} + \sum_{j=2}^g A_jx_j^k - b = A_1s_1^k + r^k.$$

So we rewrite the iteration:

FOR $i = 1, 2, \dots, g$ DO
 Solve for s_i^k : $\min \{\|A_i s_i^k + r^{k+(i-1)/g}\|_2 : s_i^k \in \mathbf{R}^{n_i}\}$;
 Update the residual: $r^{k+i/g} = r^{k+(i-1)/g} + A_i s_i^k$.

The new approximate solution is now

$$x_i^{k+1} = x_i^k + s_i^k, \quad i = 1, \dots, g; \quad x^{k+1} = \sum_{i=1}^g \bar{x}_i^{k+1}.$$

We complete this section by stating the general algorithm with relaxation factors and proving that it converges. In § 4, we will discuss termination criteria and storage requirements. We will want to assume that each A_i has full rank. This can be done without any loss of generality, since any linearly dependent group of nonzero columns can be split into smaller groups of linearly independent columns by making each column a group by itself, if necessary. Clearly, we can assume that there is no zero column, since such a column can be dropped from A without changing the least-squares problem.

Subdivide A into g groups. (Each A_i has full rank, $i = 1, 2, \dots, g$)

Choose x_i^0 , $i = 1, \dots, g$.

Compute $r^0 = Ax^0 - b$. (Choose $0 < \omega_i < 2$, $i = 1, 2, \dots, g$)

FOR $k = 0$ STEP 1 UNTIL Convergence DO

FOR $i = 1$ STEP 1 UNTIL g DO

$$(2.3) \quad \begin{aligned} s_i^k &= -(A_i^T A_i)^{-1} A_i^T r^{k+(i-1)/g}; \\ r^{k+i/g} &= r^{k+(i-1)/g} + \omega_i A_i s_i^k; \\ x^{k+i/g} &= x^{k+(i-1)/g} + \omega_i \bar{s}_i^k; \end{aligned}$$

Check for Convergence.

THEOREM 2.1. *Let the columns of A be divided into g groups A_1, \dots, A_g and let each A_i have full rank. Let $\{x^k\}$ be generated by algorithm (2.3) with any choices $0 < \omega_i < 2$, and any $x_i^0 \in \mathbf{R}^{n_i}$, $i = 1, 2, \dots, g$. Then $\{x^k\}$ converges to x^* , a solution to the least-squares problem (2.1).*

Proof. We will show that the algorithm is the block SOR iteration on the normal equations for (2.2b). We will then apply a result of Keller (1965) which proves convergence for the block SOR method applied to positive semidefinite systems. This will give convergence of $\{x^k\}$ and hence of $\{x_i^k\}$ and of $\{\bar{x}_i^k\}$ for every i . But then $\{x^k\}$ converges since it is just $\{\sum_{i=1}^g \bar{x}_i^k\}$.

We now state the result of Keller (1965) for completeness. Let G be a real symmetric matrix of order n of the form

$$(2.4) \quad G = D + E + E^T$$

and let W be any real nonsingular matrix such that

$$(2.5) \quad N = W^{-1}D + E$$

is nonsingular. Let f be any vector for which the system

$$(2.6) \quad Gx = f$$

has a solution. Consider the iterative method

$$(2.7) \quad Nx^{k+1} = (N - G)x^k + f$$

where x^0 is an initial guess of a solution of (2.6). The following lemma is a part of Keller (1965, Cor. 2.1).

LEMMA 2.2. Let G be a symmetric positive semidefinite matrix of the form (2.4) and let W be nonsingular such that the matrix N in (2.5) is nonsingular. Let (2.6) have a solution and let

$$P = W^{-1}D + (W^{-1}D)^T - D$$

be positive definite. Then for every x^0 the sequence $\{x^k\}$ of (2.7) converges to a solution of (2.6).

It will be useful to let $G_{ij} = A_i^T A_j$ be the (i, j) block of the n by n Gram matrix, $G = A^T A$. Define

$$x_j^k = x_j^0 + \sum_{p=0}^{k-1} \omega_j s_j^p, \quad j = 1, 2, \dots, g \quad \text{and} \quad x^{k+i/g} = \sum_{j=1}^g \bar{x}_j^k + \sum_{j=1}^i \omega_j \bar{s}_j^k.$$

First we need that

$$r^{k+i/g} = Ax^{k+i/g} - b.$$

We will prove this by induction on $k + i/g$ in steps of $1/g$. By definition, the statement is true for $k + i/g = 0$. Now assume that the statement is true for some $k + (i - 1)/g \geq 0$. Then

$$\begin{aligned} r^{k+i/g} &= r^{k+(i-1)/g} + \omega_i A_i s_i^k = Ax^{k+(i-1)/g} - b + \omega_i A_i \bar{s}_i^k \\ &= A[x^{k+(i-1)/g} + \omega_i \bar{s}_i^k] - b = Ax^{k+i/g} - b. \end{aligned}$$

Now,

$$\begin{aligned} G_{ii}[x_i^{k+1} - x_i^k] &= \omega_i G_{ii} s_i^k = -\omega_i A_i^T [Ax^{k+(i-1)/g} - b] \\ &= -\omega_i \left[\sum_{j=1}^{i-1} G_{ij} x_j^{k+1} + \sum_{j=i}^g G_{ij} x_j^k - A_i^T b \right], \end{aligned}$$

which becomes:

$$G_{ii} x_i^{k+1} + \omega_i \left[\sum_{j=1}^{i-1} G_{ij} x_j^{k+1} + \sum_{j=i+1}^g G_{ij} x_j^k \right] + (\omega_i - 1) G_{ii} x_i^k = \omega_i A_i^T b.$$

When $\omega_i = \omega$, $i = 1, 2, \dots, g$ this is the standard form given on Young (1971, p. 438) of the block SOR method applied to $Gx = A^T b$. To apply Keller's result, we write $G = D + E + E^T$ where $D = (G_{ii})$ is the $n \times n$ block diagonal of G , and $E = (G_{ij})$ is the $n \times n$ block strict lower triangle of G . Let W be the $n \times n$ block diagonal matrix whose i th block is ω_i times the $n_i \times n_i$ identity matrix. Now we rewrite the iteration as

$$\begin{aligned} Dx^{k+1} + W[Ex^{k+1} + E^T x^k] + (W - I)Dx^k &= WA^T b, \\ (D + WE)x^{k+1} + [W(D + E^T) - D]x^k &= WA^T b. \end{aligned}$$

Since no $\omega_i = 0$, we can multiply through by W^{-1} , set $N = W^{-1}D + E$, and obtain

$$Nx^{k+1} + (G - N)x^k = A^T b.$$

In order to complete the proof, we only need that N is nonsingular and that $(W^{-1}D + DW^{-1} - D)$ is positive definite. Observe that, since each A_i has full rank, D is positive definite. It follows that $W^{-1}D$ is nonsingular and so N is also. The inequality $0 < \omega_i < 2$ ensures that $(W^{-1}D + DW^{-1} - D) = \text{diag}[2\omega_i^{-1}D_{ii} - D_{ii}]$ is positive definite.

3. The inexact Gauss-Newton approach. We now consider the algorithm for the nonlinear least-squares problem. Let

$$F: \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad m \geq n \quad \text{and} \quad F = (F_1, \dots, F_m)^T,$$

and define

$$\phi(x) = \frac{1}{2}F(x)^T F(x).$$

Then the nonlinear least-squares problem is to find x^* so that for some norm and some $\varepsilon > 0$,

$$(3.1) \quad \phi(x^*) \leq \phi(x) \quad \text{for all } \|x - x^*\| < \varepsilon.$$

The simultaneous solution of n nonlinear equations in n unknowns may be viewed as solving (3.1) where $m = n$. For small-residual sparse problems, the Gauss-Newton method is very attractive. It starts with x^0 and generates a sequence of iterates $\{x^k\}$ as follows:

$$(3.2) \quad \begin{array}{l} \text{FOR } k=0 \text{ STEP 1 UNTIL Convergence DO} \\ \text{Solve for } s^k: \min \{\|F'(x^k)s + F(x^k)\|_2: s \in \mathbf{R}^n\}; \\ \text{Set } x^{k+1} = x^k + s^k. \end{array}$$

If $m \times n$ is large, solving the linearized problem (3.2) may require the techniques mentioned in the penultimate paragraph of § 1.

If we use an iterative method to solve the linearized problem, then it is important to know how accurately it must be solved in order not to impede convergence. We define the *inexact Gauss-Newton* algorithm for a given real nonnegative sequence $\{\theta_k\}$, any vector norm, and any starting point x^0 as follows:

$$(3.3) \quad \begin{array}{l} \text{FOR } k=0 \text{ STEP 1 UNTIL Convergence DO} \\ \text{Find some approximate minimizer } s^k \text{ of (3.2) so that} \\ \frac{\|F'(x^k)^T r^k\|}{\|F'(x^k)^T F(x^k)\|} \leq \eta_k \quad \text{where } r^k = F(x^k) + F'(x^k)s^k; \\ \text{Set } x^{k+1} = x^k + s^k. \end{array}$$

We will assume in this section that:

- A.1. F is continuously differentiable in an open neighborhood Ω of x^* .
- A.2. $F'(x^*)^T F(x^*) = 0$.
- A.3. $F'(x^*)$ has full rank.
- A.4. There exists $\gamma \geq 0$ so that for $x \in \Omega$,

$$(3.4) \quad \|[F'(x) - F'(x^*)]^T F(x^*)\| \leq \gamma \|x - x^*\|_*$$

where $\|\cdot\|$ is any vector norm and $\|y\|_* \equiv \|F'(x^*)^T F(x^*)y\|$ for every y .

If the Jacobian matrix is Lipschitz continuous at x^* , i.e., there exists $L \geq 0$ so that for $x \in \Omega$,

$$\|F'(x^*) - F'(x)\| \leq L \|x - x^*\|,$$

then there exists $\gamma \geq 0$ so that assumption A.4 holds. Notice that $\gamma = 0$ for zero-residual problems. The following theorem relates $\{\eta_k\}$ to the speed of convergence of $\{x^k\}$.

THEOREM 3.1. *Assume that*

$$(3.5) \quad \gamma + \eta(1 + \gamma) < r < 1$$

for γ from (3.4) and $0 \leq \theta_k \leq \eta$ from (3.3). Then there exists some $\varepsilon > 0$ for which $\|x^0 - x^*\|_* \leq \varepsilon$ implies that the sequence of inexact Gauss-Newton iterates $\{x^k\}$ is defined and converges at least linearly to x^* in the sense that

$$\|x^{k+1} - x^*\|_* \leq r \|x^k - x^*\|_*.$$

Proof. Let

$$(3.6) \quad \mu = 2 \max \{\|F'(x^*)^T\|, \|(F'(x^*)^T F'(x^*))^{-1}\|\}$$

and let $\delta > 0$ be so that

$$(1 + \mu\delta)[\eta[1 + \gamma + (\mu + 1)\delta] + \gamma + \mu\delta] \leq r.$$

This can be done in view of (3.5). Define

$$G(x) \equiv F'(x)^T F'(x), \quad G^* \equiv F'(x^*)^T F'(x^*).$$

Choose $\varepsilon > 0$ so that if $\|x - x^*\|_* \leq \varepsilon$, then $\|G^{-1}(x)\| \leq \mu$, $\|G^* - G(x)\| \leq \delta$, $\|F'(x)^T\| \leq \mu$, and $\|F(x^*) - F(x) - F'(x)(x^* - x)\| \leq \delta\|x - x^*\|_*$. This can be done in view of (3.6) and the assumptions A.1 and A.3. Let x^+ be the new inexact Gauss-Newton iterate, i.e., x^+ satisfies

$$\frac{\|F'(x)^T r\|}{\|F'(x)^T F(x)\|} \leq \eta \quad \text{where } r = F(x) + F'(x)s, \quad x^+ = x + s.$$

Consider

$$\begin{aligned} G^*(x^+ - x^*) &= [(G^* - G)G^{-1} + I][F'(x)^T r - (F'(x) - F'(x^*))^T F(x^*) \\ &\quad + F'(x)^T (F(x^*) - F(x) - F'(x)(x^* - x))]. \end{aligned}$$

Taking norms yields

$$\begin{aligned} \|x^+ - x^*\|_* &\leq [1 + \|G(x)^{-1}\| \|G^* - G(x)\|] \\ &\quad \times [\|F'(x)^T r\| + \|(F'(x) - F'(x^*))^T F(x^*)\| \\ &\quad + \|F'(x)^T\| \|F(x^*) - F(x) - F'(x)(x^* - x)\|] \\ &\leq (1 + \mu\delta)[\eta\|F'(x)^T F(x)\| + \gamma\|x - x^*\|_* + \mu\delta\|x - x^*\|_*]. \end{aligned}$$

Consider

$$\begin{aligned} F'(x)^T F(x) &= (F'(x) - F'(x^*))^T F(x^*) - F'(x)^T (F(x^*) - F(x) - F'(x)(x^* - x)) \\ &\quad - (G(x) - G^*)(x^* - x) - G^*(x^* - x). \end{aligned}$$

Taking norms,

$$\|F'(x)^T F(x)\| \leq \gamma\|x^* - x\|_* + \mu\delta\|x^* - x\|_* + \delta\|x^* - x\|_* + \|x^* - x\|_*.$$

So

$$\|x^+ - x^*\|_* \leq (1 + \mu\delta)[\eta[1 + \gamma + (\mu + 1)\delta] + \gamma + \mu\delta]\|x - x^*\|_* \leq r\|x - x^*\|_*.$$

If $\gamma = 0$, then the inexact Gauss-Newton method is closely related to the inexact Newton methods of Dembo, Eisenstat, and Steihaug (1982) and the inexact quasi-Newton method of Steihaug (1980). In the case when $\theta_k \equiv 0$, this theorem relaxes the Dennis (1977) conditions for convergence of the Gauss-Newton method.

A simpler approach would have been to make the stronger assumptions that F is twice continuously differentiable and that the hypothesis of the Ostrowski Theorem holds. Then we could have applied that theorem to the Gauss-Newton method as in Ortega (1972). In order to establish that our assumption is weaker, we will show now that if these assumptions are satisfied, then (3.4) holds and x^* is a solution to (3.1) rather than just a critical point as was guaranteed by A.2.

In this paragraph, let F be twice continuously differentiable in an open neighborhood Ω containing x^* . For x sufficiently close to x^* , A.3 lets us define

$$N(x) = x - [F'(x)^T F'(x)]^{-1} F'(x)^T F(x).$$

Then the derivative of N exists, is continuous in a neighborhood of x^* and by applying the product rule to differentiate

$$G(x)N(x) = G(x)x - F'(x)^T F(x),$$

we obtain

$$N'(x^*) = -[F'(x^*)^T F'(x^*)]^{-1} S,$$

where

$$S = \sum_{i=1}^m F_i(x^*) \nabla^2 F_i(x^*).$$

(See in Ortega (1972, 8.1.8) and Dennis (1977).) Recall that the Gauss–Newton method is $x^{k+1} = N(x^k)$ and x^* is a point of attraction of the Gauss–Newton iteration if $\rho(N'(x^*)) < 1$ where $\rho(\cdot)$ denotes the spectral radius of the matrix argument (see in Ortega (1972, 8.1.7). Define the function $h: \Omega \rightarrow \mathbf{R}^n$ by

$$h(x) = [F'(x) - F'(x^*)]^T F(x^*).$$

The assumption that F is twice continuously differentiable in an open neighborhood of x^* and the assumption A.2 give $h'(x^*) = S$. Assume for some positive δ we have $\rho(N'(x^*)) + \delta \leq \gamma < 1$. Choose a norm $\|\cdot\|$ so that $\|N'(x^*)^T\| \leq \rho(N'(x^*)) + \delta/2$. We can find a vector norm that is consistent with the chosen matrix norm, and choose a neighborhood of radius ε so that for all $\|x - x^*\| \leq \varepsilon$ we have

$$\|h(x) - h(x^*) - h'(x^*)(x - x^*)\| \leq \frac{\delta}{2} \|x - x^*\|_*.$$

This can be done since h is continuously differentiable. Consider

$$h(x) = h'(x^*)(x - x^*) + h(x) - h(x^*) - h'(x^*)(x - x^*)$$

and note that $h'(x^*)(x - x^*) = -N'(x^*)^T [F'(x^*)^T F'(x^*)(x - x^*)]$. Taking norms yields

$$\begin{aligned} \|[F'(x) - F'(x^*)]^T F(x^*)\| &= \|h(x)\| \\ &\leq \|h'(x^*)(x - x^*)\| + \|h(x) - h(x^*) - h'(x^*)(x - x^*)\| \\ &\leq \left[\|N'(x^*)^T\| + \frac{\delta}{2} \right] \|x - x^*\|_* \\ &\leq (\rho(N'(x^*)) + \delta) \|x - x^*\|_* \leq \gamma \|x - x^*\|_* \end{aligned}$$

which shows (3.4).

Since A.2 is the assumption that $\nabla \phi(x^*) = 0$, we only need to show that $\nabla^2 \phi(x^*)$ is positive definite to know that x^* is a minimizer of ϕ . This follows in a straightforward way from the assumption that $\rho(N'(x^*)) < 1$, but a proof by contradiction is shorter. Thus, let C be a symmetric positive-definite square root of $[F'(x^*)^T F'(x^*)]$ and suppose that for some $v \neq 0$,

$$0 \geq v^T \nabla^2 \phi(x^*) v = v^T [F'(x^*)^T F'(x^*)] v + v^T S v = (Cv)^T (Cv) + (Cv)^T C^{-1} S C^{-1} (Cv).$$

But this means that $C^{-1} S C^{-1}$ has an eigenvalue less than -1 . Since $C^{-1} S C^{-1}$ is similar to $N'(x^*)$, this contradicts the assumption that $\rho(N'(x^*)) < 1$.

In the inexact Gauss–Newton approach, we ignore the specific method we are using to find an approximate minimizer s^k of (3.2). If F' is sparse, then as in Curtis, Powell, and Reid (1974), and Coleman and Moré (1982), (1983), we may group the

columns of F' so that the columns in each group are mutually orthogonal vectors. We note that a column can be in several groups. The columns $F'(x^k)_i$ in group i may be approximated by finite differences $\Delta F(x^k)_i$ with only one extra value $F(x^k + v_i)$, where v_i is an appropriate linear combination of the corresponding standard unit vectors. For $s_i \in \mathbf{R}^{n_i}$, let $\bar{s}_i \in \mathbf{R}^n$ be constructed as in § 2. This suggests the following cycle in the inner loop:

$$\begin{aligned}
 &\text{For given } x^k, \text{ let } r^k = F(x^k), y^k = x^k; \\
 &\text{FOR } c = 1 \text{ STEP 1 UNTIL Termination DO} \\
 &\quad \text{FOR } i = 1 \text{ STEP 1 UNTIL } g \text{ DO} \\
 &\quad\quad \text{Compute } A_i^k = F'(x^k)_i \text{ or } \Delta F(x^k)_i; \\
 (3.7) \quad &\quad\quad \text{Solve for } s_i^k: \min \{ \|A_i^k s_i + r^{k+(i-1)/g}\|_2 : s_i \in \mathbf{R}^{n_i} \}; \\
 &\quad\quad \text{Set } r^{k+i/g} = r^{k+(i-1)/g} + \omega_i A_i^k s_i^k; \\
 &\quad\quad \text{Set } y^{k+i/g} = y^{k+(i-1)/g} + \omega_i \bar{s}_i^k. \\
 &\quad \text{Check Termination (3.3) (with } A^k \text{ in place of } F'(x^k)).
 \end{aligned}$$

The next iterate is now $x^{k+1} = y^{k+c} \in \mathbf{R}^n$ where c is the number of cycles, which corresponds to terminating the inner iteration after c sweeps through all of the column groups. The inner loop for the nonlinear problem is written differently here from the complete cycle for the linear problem because we save some work by only checking termination after each complete sweep through all the groups. The least-squares problem (3.7) is trivial to solve when the columns in this group are mutually orthogonal. This especially convenient way to group the columns has been discovered independently by Coleman (1984).

If $c > 1$, then the above inner loop cycle requires either recomputing $F'(x^k)_i$ or $\Delta F(x^k)_i$ when needed, or storing $F'(x^k)$ or $\Delta F(x^k)$. An alternative approach recomputes the Jacobian matrix of one particular group at a time and updates the nonlinear residual. This would suggest the following nonlinear substitution method:

$$\begin{aligned}
 &\text{Given } x^0, \text{ compute } F(x^0) \\
 &\text{FOR } k = 0 \text{ STEP 1 UNTIL Convergence DO} \\
 &\quad \text{FOR } i = 1 \text{ STEP 1 UNTIL } g \text{ DO} \\
 (3.8) \quad &\quad \text{Compute } A_i^{k+(i-1)/g} = F'(x^{k+(i-1)/g})_i \text{ or } \Delta F(x^{k+(i-1)/g})_i; \\
 &\quad \text{Solve for } s_i^k: \min \{ \|A_i^{k+(i-1)/g} s_i + F(x^{k+(i-1)/g})\|_2 : s_i \in \mathbf{R}^{n_i} \}; \\
 &\quad \text{Set } x^{k+i/g} = x^{k+(i-1)/g} + \omega_i \bar{s}_i^k; \\
 &\quad \text{Check Convergence.}
 \end{aligned}$$

4. Numerical results. In this section, we describe two column grouping strategies to be used with the algorithms given in §§ 2 and 3, and we present some numerical results for the Duff and Reid (1979) sparse least-squares test problems. We begin with a discussion of the problems.

4.1. The test problems. These problems are specified only in their sparsity structures which come from adjustment of survey data (Matrices 28 to 32 in the test bed).

- Problem 1: A is 219 by 85 and the survey pattern is from Holland.
- Problem 2: A is 958 by 292 and the survey pattern is from United Kingdom.
- Problem 3: A is 331 by 104 and the survey pattern is from Scotland.
- Problem 4: A is 608 by 188 and the survey pattern is from England.
- Problem 5: A is 313 by 176 and the survey pattern is from Sudan.

The specific problems used here were found by generating the nonzero matrix elements randomly in the interval $(-1, 1)$ and the components of a solution vector x randomly

in the interval $(0, 1)$. The right-hand side b was found by computing $b = Ax$. The nonlinear problems were found by replacing x_j by x_j^3 , i.e., component i in F is

$$F_i(x) = \sum_{j=1}^n A_{ij}x_j^3 - b_i.$$

Thus, our problems have zero residuals at the solution. We approximate all derivatives by finite differences in these tests.

4.2. The column grouping schemes. We have already mentioned that one grouping scheme is based primarily on the ideas of Curtis, Powell, and Reid (1974) as expanded and improved by Coleman and Moré (1983). A FORTRAN code found in Coleman and Moré (1982) furnished our first pass partitioning the columns of A into disjoint groups. We will refer to this work as "CM". In § 1, we argued that there could be some advantages to allowing the groups to overlap. In our tests, we used the following heuristic to expand each group in turn. To expand a given group, we first mark all columns that have a nonzero element in the same row position as a nonzero of some column in the group. This identifies the columns that cannot be added to the group. We then add one unmarked column to the group and add to the set of marked columns all columns that have a nonzero row element in the same position as the column that was added to the group. This process is then repeated until no columns are left unmarked.

Finally, let A_i denote a resulting submatrix of columns a_j of A , $j \in I_i$, then $A_i^T A_i$ is a diagonal matrix where the diagonal elements are the squares of the l_2 -norms of the columns, so A_i has full rank, as we required in Theorem 2.1. As an example, we present in Table 1 the results of this scheme applied to Problem 3. We note that when the groups are expanded, for the last groups the increase in number of columns is larger than for the first few. This is to be expected for most sparsity structures by the way the methods of partitioning the columns work.

TABLE 1
Groups in Problem 3.

Group	Number of Columns	
	CM	Expanded
1	25	25
2	25	25
3	25	25
4	20	23
5	8	25
6	1	25

We also considered an expansion of the groups of columns beyond mutual orthogonality to the case where the normal equations for the least-squares subproblems are banded. In particular, we used the following sequential heuristic algorithm to group the columns so that $A_i^T A_i$ is tridiagonal and block diagonal. Initially, the columns are ordered according to some criterion such as increasing number of nonzeros. Choose the first column, mark it, and let all other columns be unmarked. This will be our first column in the group. As the next column in the group, choose the first of the unmarked columns that have a nonzero element in any same row position as the first column. Now mark all other columns that have an element in any same row position as the

first column. This process is now repeated until a column has no unmarked columns with an element in any same row position. At this point, either all columns are marked or there is an unmarked column. Choose the first unmarked column and continue the process until all columns are marked. We have now generated one group of columns so that the normal matrix $A_i^T A_i$ is block tridiagonal and block diagonal. Columns in different blocks in the same group are orthogonal. Unmark all columns except the columns already in a group. If there are no unmarked columns at this point, then every column is in a group and the process is over. Otherwise, choose the first unmarked column to be the first column in the next group, and repeat the process to generate that next group. We illustrate this grouping strategy on Problem 3 in Table 2.

TABLE 2
Groups in Problem 3.

Group	Number of	
	columns	blocks
1	45	5
2	41	9
3	17	11
4	1	1

4.3. Storage requirements. It is of interest to compare the storage requirements of the algorithm of § 2 applied directly to these problems to the requirements of a very good package for sparse symmetric and positive-definite systems applied to the normal equations. In Table 3, columns A and $A^T A$ give the storage required for the real nonzero elements in A and the lower triangular part of $A^T A$, as well as the associated integer pointers when we use the storage scheme of the Harwell testbed. Column L gives the storage requirements for the Yale Sparse Matrix Package (YSMP) (Eisenstat et al. (1982)) to store the lower triangular factor of $A^T A$.

For our scheme, if $A_i^T A_i$, $i = 1, 2, \dots, g$ are diagonal matrices, we do not need the vector s_i^k explicitly. Instead, when we compute the components of $A_i^T r^{k+(i-1)/g}$, we also compute the components of s_i^k and accumulate the inner product $(A_i^T r^{k+(i-1)/g})^T s_i^k$. Hence the only storage that is needed is the original data A , and b (overwritten by $r^{k+i/g}$), and the solution vector x plus some additional pointer storage for the groups. If $A_i^T A_i$, $i = 1, 2, \dots, g$ are tridiagonal, then we need the LDL^T factorizations of the tridiagonal matrices and the vector s_i^k . For the inexact Gauss-Newton method (3.7), we need to store the Jacobian matrix. For the nonlinear substitution method (3.8), we need only one extra vector of length m if the columns of the

TABLE 3
Storage.

Storage Requirement								
Problem	m	n	A		$A^T A$		L	
			Real	Int	Real	Int	Real	Int
1	219	85	438	524	304	390	520	642
2	958	292	1916	2209	1250	1543	2568	2497
3	331	104	662	767	435	540	774	812
4	608	188	1216	1405	796	985	1625	1609
5	313	176	1557	1734	1485	1662	1593	1210

Jacobian matrix or its approximant in each group have no elements in the same row positions.

4.4. Numerical experiments. Now we briefly discuss the termination criteria that we use. From the definition of $r^{k+i/g}$ and the choice of s_i^k , we have

$$\begin{aligned} \|r^{k+i/g}\|_2^2 &= \|r^{k+(i-1)/g}\|_2^2 + 2\omega_i(A_i^T r^{k+(i-1)/g})^T s_i^k + \omega_i^2(A_i^T A_i s_i^k)^T s_i^k \\ &= \|r^{k+(i-1)/g}\|_2^2 + \omega_i(2 - \omega_i)(A_i^T r^{k+(i-1)/g})^T s_i^k. \end{aligned}$$

The major work required to calculate the l_2 -norm of the residual is an extra inner product since $A_i^T r^{k+(i-1)/g}$ is already computed as in § 4.3. Since we want to compare different algorithms, we need to base our stopping criteria on a monotonically decreasing sequence. This suggests the following termination rule for the linear problem (2.1):

$$(4.1) \quad \frac{\|r^{k+i/g}\|_2}{\|r^0\|_2} \leq \varepsilon.$$

For the nonlinear problems, the outer loop is terminated when

$$(4.2) \quad \frac{\|F(x^{k+i/g})\|}{\|F(x^0)\|} \leq \varepsilon.$$

In the inexact Gauss-Newton method (3.7), we base the termination rule for the c -loop on the residual $A^T r$ in the normal equations. We note that this costs one matrix-vector product per iteration. We terminate the c -loop cycle when (3.3) holds.

As explained in § 4.2, one grouping scheme begins by using CM graph coloring to partition the columns of A , and then we use a heuristic strategy to expand the groups. Table 4 compares the CM grouping to the expanded groups that “overlap”. The entries in the tables are: in the column marked “it” for iterations are the numbers $kg + i$ in the notation of (2.3) of diagonal least-squares problems solved; we also include in the column marked “vup” the total number of variables that were updated. Since the block matrix $A_i^T A_i$ is diagonal, the CM grouping is a point SOR using the CM grouping as the ordering. The number of variables updated is therefore the number of point SOR corrections. In Table 4, we chose $\omega_i \equiv 1$.

TABLE 4
Overlap vs. CM for linear problems.

Problem	g	it: number of subproblems solved -- vup: number of variables updated											
		$\epsilon = .1$				$\epsilon = .01$				$\epsilon = .001$			
		Overlap		CM		Overlap		CM		Overlap		CM	
it	vup	it	vup	it	vup	it	vup	it	vup	it	vup		
1	4	8	188	8	170	19	449	19	408	35	825	35	727
2	6	9	633	9	505	22	1524	23	1166	39	2713	41	2042
3	6	9	223	9	179	28	690	28	511	86	2122	86	1506
4	6	9	411	9	321	24	1092	31	987	52	2362	68	2160
5	10	33	703	33	602	123	2584	123	2186	242	5051	235	4160

Table 4 indicates a fast decrease in the residual in the first few iterations. This can be explained from the observation that the iterative process is somewhat related to coordinate search for the least-squares problem, where the spans of the A_i act like coordinates. Notice further, there is basically no difference in Table 4, where $\omega = 1$, between partitioning the columns and allowing overlaps if we count the number of iterations. We see a bigger difference between overlapping and partitioning for $\omega \neq 1$ than for $\omega = 1$. Perhaps this can be explained from the observation that if the column

a_t of A is in group i and $i + 1$, then for $\omega = 1$, component t of \bar{s}_{i+1}^k is 0, but for $\omega \neq 1$, this component can be nonzero. However, in terms of numbers of variables updated, we see that for the linear problem it does not pay to expand the groups. On the other hand, in the following three sets of results for the nonlinear problems and the nonlinear substitution technique (3.8), we see that overlap may be more efficient in terms of fewer iterations and function calls. Of course, this is hardly surprising, but the extra function calls used by nonlinear substitution make it less attractive than the inexact Gauss-Newton method, unless storage is the main concern.

In Tables 6 and 7 respectively, we use a large value, $\epsilon = .1$, and a small value, $\epsilon = .001$ to compare the CM grouping strategy and the grouping strategy of § 4.2 that makes every $A_i^T A_i$ tridiagonal. The strategies are compared separately to point SOR with the original ordering, which accounts for the two lines labeled "Point SOR". The entries for the grouping schemes are the numbers of variables updated to satisfy (4.1)

TABLE 5

Results for nonlinear test problems. Here $f/i/o$ denotes f = number of function calls, i = number of iterations or subproblems solved, and for the Gauss-Newton method (3.7), the third number is o = number of outer iterations.

Problem 1

		$\epsilon = .1$	$\epsilon = .01$	$\epsilon = .001$
CM ordering	Nonlinear substitution	19/9	45/22	81/40
	Inexact Gauss Newton	11/19/2	21/51/4	26/70/5
Overlap	Nonlinear substitution	19/9	45/22	81/40
	Inexact Gauss Newton	11/19/2	21/51/4	26/70/5

Problem 2

		$\epsilon = .1$	$\epsilon = .01$	$\epsilon = .001$
CM ordering	Nonlinear substitution	23/11	63/31	105/52
	Inexact Gauss Newton	15/24/2	22/44/3	36/90/5
Overlap	Nonlinear substitution	23/11	57/28	103/51
	Inexact Gauss Newton	15/22/2	29/61/4	36/84/4

Problem 3

		$\epsilon = .1$	$\epsilon = .01$	$\epsilon = .001$
CM ordering	Nonlinear substitution	27/13	57/28	115/57
	Inexact Gauss Newton	15/23/2	29/85/4	36/84/5
Overlap	Nonlinear substitution	23/11	55/27	101/50
	Inexact Gauss Newton	15/22/2	29/84/4	36/124/5

TABLE 6

Point SOR with CM ordering vs. point SOR with original column ordering and tridiagonal blocks vs. point SOR with original ordering. $\epsilon = .1$.

Problem 1

 $\epsilon = .1$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	281	217	196	170	153	170	196	238	302	387	557
Point SOR	304	238	223	210	171	210	237	275	346	434	596
Tridiagonal	255	247	210	170	162	162	210	247	295	417	587
Point SOR	296	268	238	214	191	208	269	288	346	463	625
Rel. eff.	1.07	.99	1.0	1.02	1.05	1.04	1.06	1.01	1.02	.99	1.0

Problem 2

 $\epsilon = .1$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	797	584	559	505	505	559	658	851	1089	1381	1965
Point SOR	802	634	568	552	560	664	807	985	1234	1535	2119
Tridiagonal	814	584	577	522	522	577	705	814	1106	1398	1982
Point SOR	828	664	630	569	613	747	848	975	1277	1557	2164
Rel. eff.	1.01	1.05	1.12	1.0	1.06	1.09	.98	1.03	1.02	1.00	1.01

Problem 3

 $\epsilon = .1$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	283	233	199	179	179	207	233	303	387	491	699
Point SOR	271	233	197	198	202	230	278	329	426	537	745
Tridiagonal	294	253	190	190	190	190	253	294	398	502	710
Point SOR	297	262	198	205	232	222	305	322	436	540	757
Rel. eff.	1.05	1.04	1.05	.98	1.08	1.05	1.01	1.01	1.0	.98	1.00

Problem 4

 $\epsilon = .1$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	548	423	360	321	321	360	468	548	697	924	1261
Point SOR	521	431	341	339	360	387	529	612	762	984	1322
Tridiagonal	529	453	373	341	341	373	453	529	717	905	1281
Point SOR	513	474	365	343	360	461	513	597	799	988	1360
Rel. eff.	1.02	1.03	1.03	.95	.94	1.15	1.00	1.01	1.02	1.03	1.01

Problem 5

 $\epsilon = .1$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	876	778	667	602	602	602	602	654	778	973	1306
Point SOR	802	736	614	583	589	614	633	682	816	1016	1351
Tridiagonal	824	704	621	562	562	562	562	648	738	973	1294
Point SOR	805	715	612	583	588	613	631	700	814	1028	1405
Rel. eff.	1.07	1.07	1.07	1.07	1.07	1.07	1.07	1.04	1.05	1.01	1.05

TABLE 7

Point SOR with CM ordering vs. point SOR with original ordering and point SOR with tridiagonal ordering vs. point SOR with original ordering. $\epsilon = .001$.

Problem 1

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	1598	1258	961	727	557	493	536	680	897	1173	1683
Point SOR	1870	1530	1273	1020	850	680	610	762	977	1277	1747
Tridiag. Ord.	1400	1145	935	765	635	465	550	672	890	1182	1692
Point SOR	1906	1530	1273	1020	903	735	677	777	982	1281	1756

Problem 2

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	4162	3286	2603	2042	1606	1534	1898	2410	3066	4162	5815
Point SOR	4381	3505	2886	2311	1850	1875	2268	2751	3377	4389	6009
Tridiag. Ord.	3734	3041	2457	1982	1581	1398	1873	2329	3041	4026	5778
Point SOR	4360	3563	2899	2317	2020	1783	2299	2743	3448	4344	6068

Problem 3

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	2175	1947	1714	1506	1298	1115	927	832	1090	1447	2071
Point SOR	1762	1559	1450	1247	1141	1034	926	939	1224	1560	2184
Tridiag. Ord.	1438	1230	1039	918	773	669	669	831	1085	1438	2062
Point SOR	1763	1560	1352	1349	1142	1038	933	972	1224	1560	2186

Problem 4

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	3476	2992	2536	2160	1864	1488	1220	1551	1972	2631	3744
Point SOR	2821	2316	1881	1676	1496	1308	1354	1725	2138	2797	3911
Tridiag. Ord.	2521	2145	1769	1469	1205	1017	1205	1501	1957	2629	3725
Point SOR	2823	2388	1960	1685	1488	1309	1448	1709	2159	2851	3891

Problem 5

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6	1.7
CM	7602	6156	5040	4160	3418	2834	2306	2268	2674	3418	4688
Point SOR	7126	5744	4811	4133	3632	3456	3280	3095	2959	3663	4861
Tridiag. Ord.	8306	7384	6574	5842	5138	4434	3758	3112	2702	3464	4962
Point SOR	7125	5744	4812	4133	3684	3456	3280	3095	2959	3663	4855

TABLE 8
Point SOR with CM ordering, tridiagonal ordering and original ordering. $\epsilon = .001$.

Problem 1

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6
CM	1598	1252	961	726	548	490	536	679	878	1171
Tridiagonal	1662	1332	1066	842	642	470	546	678	895	1191
Original	1869	1529	1273	1020	849	679	610	756	947	1277

Problem 2

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6
CM	4106	3242	2576	2027	1587	1469	1853	2349	3063	4109
Tridiagonal	4411	3551	2848	2271	1802	1560	1901	2444	3118	4169
Original	4306	3479	2813	2309	1832	1765	2189	2669	3361	4343

Problem 3

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6
CM	2175	1947	1714	1505	1298	1115	927	832	1069	1443
Tridiagonal	1957	1749	1541	1334	1189	1021	877	856	1104	1471
Original	1762	1557	1352	1246	1141	1032	842	939	1193	1560

Problem 4

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6
CM	3479	2989	2536	2160	1861	1486	1220	1516	1949	2625
Tridiagonal	3139	2587	2139	1769	1468	1205	1191	1545	1954	2613
Original	2821	2259	1878	1621	1317	1300	1354	1673	2119	2797

Problem 5

$\epsilon = .001$

ω	.7	.8	.9	1.0	1.1	1.2	1.3	1.4	1.5	1.6
CM	7596	6144	5030	4150	3418	2824	2296	2261	2670	3403
Tridiagonal	8322	7421	6651	5900	5163	4461	3807	3131	3011	3749
Original	7125	5744	4801	4104	3632	3455	3277	3095	2949	3628

with the specified ϵ . For point SOR with the original ordering, the entries are the numbers of variable updates needed to achieve the same accuracy as the corresponding grouping scheme. The reason for separately comparing each one to point SOR with the original ordering is to minimize the dependence of the comparisons on the specific value of ϵ . Specifically, the tridiagonal method sometimes greatly exceeds the required accuracy when it first reaches that accuracy. Thus, we thought it worthwhile in Table 6 to add the final line of the subtables, which gives the relative efficiency in point SOR corrections of the two methods, e.g. for Problem 1 with $\omega = 1.1$ and $\epsilon = .1$, it is $(191/162)/(171/153) = 1.05$.

The arithmetic needed by CM and point SOR with original ordering has the same cost per variable update. Naturally, the tridiagonal case costs more per variable update.

However, the dominating cost for all the methods is the two matrix-vector products for each sweep through all the columns.

In Table 8, we tried point SOR in every case, but the orderings used were the orderings suggested by the column grouping schemes discussed previously. Except for Problem 5, it appears as though the alternate orderings are very good ones. The rows labeled "Original" correspond to the original ordering.

Acknowledgments. We sincerely thank Kathryn Turner and both referees for carefully reading the manuscript and suggesting improvements in the presentation.

REFERENCES

- Å. BJÖRK AND T. ELFVING (1979), *Accelerated projection methods for computing pseudoinverse solutions of systems of linear equations*, BIT, 19, pp. 145-163.
- T. F. COLEMAN (1984), *Large Sparse Numerical Optimization*, Lecture Notes in Computer Science 165, Springer-Verlag, Berlin.
- T. F. COLEMAN AND J. J. MORÉ (1982), *Software for estimating sparse Jacobian matrices*, Cornell Computer Science TR 82-502.
- , (1983), *Estimation of sparse Jacobian matrices and graph coloring problems*, this Journal, 20, pp. 187-209.
- A. R. CURTIS, M. J. D. POWELL AND J. K. REID (1974), *On the estimation of sparse Jacobian matrices*, J. Inst. Math. Appl., 13, pp. 117-119.
- R. S. DEMBO, S. C. EISENSTAT AND T. STEIHAUG (1982), *Inexact Newton methods*, this Journal, 19, pp. 400-408.
- J. E. DENNIS JR. (1977), *Nonlinear least squares and equations* in The State of the Art in Numerical Analysis, D. Jacobs, ed., Academic Press, London, pp. 269-312.
- I. S. DUFF AND J. K. REID (1979), *Performance evaluation of codes for sparse matrices*, in Performance Evaluation of Numerical Software, L. D. Fosdick, ed. North-Holland, Amsterdam, pp. 121-135.
- S. C. EISENSTAT, M. C. GURSKY, M. H. SCHULTZ AND A. H. SHERMAN (1982), *Yale Sparse Matrix Package I: The symmetric codes*, Internat. J. Numer. Meth. Engrg., 18, pp. 1141-1151.
- T. ELFVING (1980), *Block iterative methods for consistent and inconsistent linear systems*, Numer. Math., 35, pp. 1-12.
- A. S. HOUSEHOLDER (1964), *The Theory of Matrices in Numerical Analysis*, Blaisdell, New York.
- A. S. HOUSEHOLDER AND F. L. BAUER (1960), *On certain iterative methods for solving linear systems*, Numer. Math., 2, pp. 55-59.
- H. B. KELLER (1965), *On the solution of singular and semidefinite linear systems by iteration*, this Journal, 2, pp. 281-290.
- J. M. ORTEGA (1972), *Numerical Analysis: A Second Course*, Academic Press, New York.
- J. K. REID (1973), *Least squares solution of sparse systems of non-linear equations by a modified Marquardt algorithm*, in Proc. NATO Conference at Cambridge, July 1972, North-Holland, Amsterdam, pp. 437-445.
- T. STEIHAUG (1980), *Quasi-Newton methods for large scale nonlinear problems*, Ph.D dissertation, SOM Technical Report #49, Yale University, New Haven, CT.
- G. W. STEWART (1973), *Introduction to Matrix Computations*, Academic Press, New York.
- D. M. YOUNG (1971), *Iterative Solution of Large Linear Systems*, Academic Press, New York.