

## STOCHASTIC PERTURBATION THEORY\*

G. W. STEWART†

**Abstract.** In this paper classical matrix perturbation theory is approached from a probabilistic point of view. The perturbed quantity is approximated by a first-order perturbation expansion, in which the perturbation is assumed to be random. This permits the computation of statistics estimating the variation in the perturbed quantity. Up to the higher-order terms that are ignored in the expansion, these statistics tend to be more realistic than perturbation bounds obtained in terms of norms. The technique is applied to a number of problems in matrix perturbation theory, including least squares and the eigenvalue problem.

**Key words.** perturbation theory, random matrix, linear system, least squares, eigenvalue, eigenvector, invariant subspace, singular value

**AMS(MOS) subject classifications.** 15A06, 15A12, 15A18, 15A52, 15A60

**1. Introduction.** Let  $A$  be a matrix and let  $F$  be a matrix valued function of  $A$ . Two principal problems of matrix perturbation theory are the following. Given a matrix  $E$ , presumed small,

1. Approximate  $F(A + E)$ ,
2. Bound  $\|F(A + E) - F(A)\|$  in terms of  $\|E\|$ .

Here  $\|\cdot\|$  is some norm of interest.

The first problem is usually, but not always, solved by assuming that  $F$  is differentiable at  $A$  with derivative  $F'_A$ . Then

$$F(A + E) = F(A) + F'_A(E) + o(\|E\|),$$

so that for  $E$  sufficiently small  $F'_A(E)$  is the required approximation. The problem then reduces to finding tractable expressions for  $F'_A(E)$ , which in itself is often a nontrivial task. The second problem may be treated in a variety of ways; but if the results are to be sharp, for small  $E$  they have to approach a bound that could be obtained by manipulating  $F'_A(E)$ .

For example, it is well known that if  $A$  is nonsingular, then

$$(1.1) \quad (A + E)^{-1} = A^{-1} - A^{-1}EA^{-1} + O(\|E\|^2).$$

Moreover, if in some norm  $\|A^{-1}\|\|E\| < 1$  then

$$(1.2) \quad \|(A + E)^{-1} - A^{-1}\| \leq \frac{\|A^{-1}\|^2\|E\|}{1 - \|A^{-1}\|\|E\|}.$$

Except for the denominator, which approaches 1 as  $E \rightarrow 0$ , the inequality (1.2) could be derived from (1.1) by ignoring the quadratic term and taking norms.

The formulas (1.1) and (1.2) represent two extremes. If the higher-order term can be ignored, equation (1.1) tells the entire story, but in overabundant detail: it is not easy to interpret. On the other hand, the bound (1.2) makes a clear statement

\* Received by the editors February 8, 1989; accepted for publication (in revised form) January 19, 1990.

† Department of Computer Science and Institute for Advanced Computer Studies, University of Maryland, College Park, Maryland 20742. This work was supported in part by the Air Force Office of Scientific Research under contract AFOSR-87-0188.

about the size of the perturbation, but it is likely to be an overestimate, since the submultiplicative inequality for norms was used in its derivation.

In this paper we will consider a third approach that is in some sense intermediate to the other two. We will take  $E$  to be a stochastic matrix and compute expectations of quantities derived from the perturbation expansion (1.1). This represents a compression of the information, but up to higher-order terms it gives nothing away.

For example, let

$$\check{A}^{-1} = A^{-1} - A^{-1}EA^{-1}$$

be the first-order approximation to  $(A + E)^{-1}$ . Suppose that the elements of  $E$  are uncorrelated with mean zero and standard deviation  $\sigma$ . Let us agree to measure the size of a random matrix by the function  $\|\cdot\|_S$  defined by

$$(1.3) \quad \|E\|_S^2 = \mathbf{E}(\|E\|_F^2),$$

where  $\mathbf{E}$  is the expectation operator and  $\|\cdot\|_F$  is the Frobenius norm. Then from the results of §3.3, it can be shown that

$$(1.4) \quad \|\check{A}^{-1} - A^{-1}\|_S = \sigma \|A^{-1}\|_F^2.$$

The equality (1.4) has much the same form as (1.2), when the latter is stripped of its denominator. The left-hand side of both is a measure of the size of the perturbation. The right-hand side of both consists of a measure of the size of the error times the square of a norm of  $A^{-1}$ . However, there are two important differences. First, (1.4) is an equality—there is no question of sharpness here. Second, if  $\|\cdot\|$  in (1.2) is the Frobenius norm, then the right-hand side of (1.4) will generally be smaller than (1.2), since  $\|E\|_S = n\sigma$ .<sup>1</sup>

A person accustomed to using norms to bound errors may feel uncomfortable with a probabilistic statement like (1.4). A statistician would have no such qualms, and in fact might feel uncomfortable with an inequality like (1.2). Even outside statistics, rigorous bounds are often supplemented by informal probabilistic statements, as when we say that rounding error in the sum in  $n$  numbers grows as the square root of  $n$ , although the best upper bound grows as  $n$ . To be realistic, we must prune away the unlikely. What is left is necessarily a probabilistic statement.

Stochastic perturbation theory, as we shall understand it, consists of two steps. First, the perturbation in  $F(A)$  is estimated by the first-order expansion  $F(A) + F'_A(E)$ , a strictly conventional procedure. However, instead of going on to bound  $F'_A(E)$ , we assume that  $E$  is random and compute  $\|F'_A(E)\|_S$ .

To realize this program fully, we must address three questions.

1. How do we compute the stochastic norm of  $\|F'_A(E)\|_S$ ?
2. What does a knowledge of the stochastic norm tell us about the actual error?
3. What is the justification for ignoring higher-order terms?

These questions will be answered in the next section, which is the technical heart of the paper; however, it is appropriate to sketch the answers here.

<sup>1</sup> Actually, this exaggerates the difference in our favor, since  $\|E\|_F$  in (1.2) could be replaced by the spectral norm defined below by (1.5). However, a result on the limiting behavior of the spectral norm of stochastic matrices [15] shows that  $\sqrt{2n}\sigma$  is a reasonable estimate of  $\|E\|$ , so that (1.2) will still be an overestimate.

In principle, the answer to the first question is that given the first and second moments of  $E$  the calculation of  $\|F'_A(E)\|_S$  is a straightforward, if tedious procedure (see Theorem 2.5). However, this answer ignores the fact that the object of any perturbation theory is usually insight rather than a specific numerical bound. In order to obtain interpretable formulas we must put restrictions on  $E$ . In the next section we will introduce the class of *cross-correlated* matrices, whose structure is at the same time sufficiently rich to be useful and sufficiently simple to be tractable. The approach through cross-correlated matrices has the added advantage that it incorporates the scaling of the error into the final results.

The second question is answered by an appeal to the Chebyshev inequality, which asserts that it is improbable that a random matrix be much larger than its stochastic norm. It should be stressed that the bounds given by the Chebyshev inequality are very weak; for a given distribution the situation may be much better than they indicate.

The third question involves subtle issues in probability theory. The crux of the matter is that  $F(A + E) - F(A)$  can fail to have even a mean, much less a stochastic norm. Nonetheless, we will show that provided the second moments of  $E$  are small enough the distributions of  $F(A + E) - F(A)$  and  $F'_A(E)$  are close, so that any statement about the size of the latter can be transferred to the former. Moreover, this result is independent of the distribution of  $E$ .

This paper is organized as follows. In the next section we give the necessary probabilistic background and address the three questions raised above. The next two sections are devoted to the application of these results, first to the pseudo-inverse and least squares problems, then to the eigenvalue problem and the singular value decomposition. These sections are of independent interest, since they collect a number of perturbation expansions that have lain scattered about in the literature. The last section is devoted to a brief summary.

Throughout this paper  $\|\cdot\|_F$  will denote the Frobenius norm defined by

$$\|A\|_F^2 = \text{trace}(A^T A),$$

and  $\|\cdot\|_S$  will denote the stochastic norm defined by (1.3). The norm  $\|\cdot\|$  denotes the Euclidean vector norm and the spectral matrix norm defined by

$$(1.5) \quad \|A\| = \max_{\|x\|=1} \|Ax\|.$$

In dealing with perturbations of a matrix function  $F(A)$ , we will write  $\tilde{A}$  for  $A + E$  and  $\tilde{F}$  for  $F(\tilde{A})$ . If  $F$  is differentiable at  $A$ , we will write  $\check{F}$  for  $F(A) + F'_A(E)$ . Note that  $\check{F}$  is not just any approximation of  $\tilde{F}$  that is accurate up to terms of the first-order; it is the unique first order approximation that is linear in  $E$ .

**Notes and references.** For general surveys of perturbation theory for matrices and linear operators, see [25], [42]. The idea of using first-order expansions of nonlinear functions of random variables is by no means new. Gauss [12]–[14], used the technique to approximate the variances of parameters from nonlinear least squares fits. Hotelling, writing in 1940 [22], refers to a “method of differentials,” with the implication that the practice was widespread. Recently Chatelin [3]–[5] has used first-order expansions and random matrices to analyze the effects of rounding error on numerical calculations. As we have pointed out, the chief difficulty with this approach is that

the quantities being approximated may not have means or variances, so that the interpretation of the means and variances of the approximations becomes problematical. Theorem 2.8 provides a resolution of this difficulty.

The perturbation theory developed in this paper should not be confused with results on the properties of random matrices. For example, Demmel [9] considers the distance of a random matrix from a manifold of degenerate problems. Here the random matrices are not small, and the concern is not with perturbations of a matrix function. The work of Weiss et al. [45] is closer to ours in that they assume their random errors are small enough to ignore higher-order terms; but their concern is with evaluating average condition numbers, not with perturbation theory as such.

**2. The probabilistic background.** In this section we will introduce the ideas and techniques from probability theory that will be used throughout the rest of the paper. We will assume that the reader is familiar with the basic concepts of multivariate probability theory—distributions, expectations, independence, etc.

The expectation operator will be denoted by  $\mathbf{E}$ . The covariance of random vectors  $x$  and  $y$  will be written

$$\mathbf{C}(x, y) \stackrel{\text{def}}{=} \mathbf{E} [(x - \mathbf{E}(x))(y - \mathbf{E}(y))^T],$$

and the variance of a random vector  $x$  will be written

$$\mathbf{V}(x) \stackrel{\text{def}}{=} \mathbf{C}(x, x).$$

If  $\mathbf{C}(x, y) = 0$ , the random vectors  $x$  and  $y$  are said to be *uncorrelated*.

We will denote by  $\mathcal{G}^n$  the space of all random  $n$ -vectors whose components have finite second moments. Note that  $\mathcal{G}^n$  is a vector space under addition and multiplication by a scalar. The zero element is the vector with mean and variance zero. We write

$$x \sim \mathcal{G}^n(u, \Sigma)$$

to say that  $x \in \mathcal{R}^n$  has mean  $u$  and variance  $\Sigma$ . If  $x \sim \mathcal{G}^n(u, \Sigma)$ , then  $x$  can be written in the form

$$x = u + \Sigma^{1/2}e,$$

where  $e \sim \mathcal{G}^n(0, I)$ .

**2.1. Random matrices.** We will denote by  $\mathcal{G}^{m \times n}$  the space of all random  $m \times n$  matrices whose elements have finite second moments. As we pointed out in the introduction, random matrices are difficult to manipulate in this generality. Hence we introduce a more tractable class—the cross-correlated matrices—which cover many actual applications.

**DEFINITION 2.1.** A random matrix  $A \in \mathcal{R}^{m \times n}$  is *cross-correlated* with mean  $U$ , row scale  $S_r$ , and column scale  $S_c$  if it can be written in the form

$$(2.1) \quad A = U + S_c H S_r^T,$$

where  $H$  is a random matrix whose elements are uncorrelated with mean zero and variance 1. We write

$$(2.2) \quad A \sim \mathcal{T}^{m \times n}(U; S_r, S_c).$$

The matrices  $S_r$  and  $S_c$  are called scales, because they represent row and column scalings of the matrix  $H$ . Their relation to the variance of  $A$  is the following. Let  $S_r^{(2)} = S_r^T S_r$  and let  $S_c^{(2)} = S_c^T S_c$ . It can be shown (see Theorem 2.3 below) that if  $A \sim \mathcal{T}^{m \times n}(U; S_r, S_c)$ , then the covariance of the  $i$ th and  $j$ th columns of  $A$  is  $(S_r^{(2)})_{ij} S_c^{(2)}$ . Consequently, if we let  $\text{vec}(A)$  denote the vector formed by stacking the columns of  $A$  in their natural order, then  $\mathbf{C}[\text{vec}(A)] = S_r^{(2)} \otimes S_c^{(2)}$ , i.e., the tensor or Kronecker product of  $S_r^{(2)}$  and  $S_c^{(2)}$ . Hence the symbol  $\mathcal{T}$  in (2.2).

Definition 2.1 has been phrased with an eye to applications in which the row and columns scales are known. For theoretical work, the following characterization leads to a more compact notation.

**THEOREM 2.2.** *If  $A \sim \mathcal{T}^{m \times n}(U; S_r, S_c)$ , then  $A$  can be written in the form*

$$A = U + S'_c H' S'_r$$

where  $S'_c$  and  $S'_r$  are positive semidefinite and the elements of  $H'$  are uncorrelated with mean zero and variance 1.

*Proof.* We will show how to replace  $S_c$  with a positive semidefinite matrix, leaving the modification of  $S_r$  as an exercise. Without loss of generality we may assume that  $S_c$  has at least as many columns as rows (if not augment  $S_c$  with zero columns while augmenting  $H$  with rows of uncorrelated elements). Then  $S_c$  has a polar factorization  $S_c = S'_c Q^T$ , where  $S'_c$  is positive semidefinite and  $Q$  has orthonormal columns.<sup>2</sup> The result now follows on setting  $H' = Q^T H$ .  $\square$

For the rest of this paper we will assume that the matrices  $S_r$  and  $S_c$  are positive semidefinite. In particular, this permits us to write  $S_r^2$  for  $S_r^T S_r$ , and similarly for  $S_c$ .

If a matrix is cross-correlated, certain quadratic forms involving it may be easily computed, as the following theorem shows.

**THEOREM 2.3.** *Let  $E \sim \mathcal{T}^{m \times n}(0; S_r, S_c)$ ,  $B \in \mathcal{R}^{m \times m}$ , and  $C \in \mathcal{R}^{n \times n}$ . Then*

$$(2.3) \quad \mathbf{E}(E^T B E) = \text{trace}(S_c B S_c) S_r^2 = \text{trace}(S_c^2 B) S_r^2 = \text{trace}(B S_c^2) S_r^2,$$

$$(2.4) \quad \mathbf{E}(E C E) = S_c^2 C^T S_r^2.$$

*Proof.* The results will first be established for the case  $E \sim \mathcal{T}^{m \times n}(0; I, I)$ . For (2.3), let  $S = E^T B E$ . Then

$$(2.5) \quad s_{ij} = \sum_{k,l} e_{ki} e_{lj} b_{kl}.$$

Since the elements of  $E$  are uncorrelated,  $\mathbf{E}(s_{ij}) = 0$  unless  $i = j$ . Moreover, if  $i = j$ , the expectations of all terms in the sum (2.5) are zero, except those for which  $k = l$ . Thus

$$\mathbf{E}(s_{ii}) = \mathbf{E} \left( \sum_k x_{ki} x_{ki} b_{kk} \right) = \text{trace}(B),$$

and  $\mathbf{E}(S) = \text{trace}(B)I$ , which is just (2.3) when  $E \sim \mathcal{T}^{m \times n}(0; I, I)$ .

<sup>2</sup> Namely, let  $S_c = U \Psi V^T$  be the singular value factorization of  $S_c$  [17]. Then the factorization  $S_c = (U \Psi U^T)(U V^T)$  is the required polar decomposition.

Similarly, the  $(i, j)$ -element of  $T = XCX$  has the form

$$t_{ij} = \sum_{k,l} x_{ik} x_{lj} c_{kl}.$$

The only term in this sum having nonzero expectation occurs when  $j = k$  and  $i = l$ . Thus  $\mathbf{E}(t_{ij}) = c_{ji}$  or  $\mathbf{E}(T) = C^T$ . This is just the form (2.4) assumes when  $E \sim \mathcal{T}^{m \times n}(0; I, I)$ .

Turning now to the general case, write  $E = S_c H S_r$ , where the elements of  $H$  are uncorrelated with mean zero and variance 1. Then

$$\begin{aligned} \mathbf{E}(E^T B E) &= \mathbf{E}(S_r H^T S_c B S_c H S_r) \\ &= S_r \mathbf{E}(H^T S_c B S_c H) S_r \\ &= S_r [\text{trace}(S_c B S_c) I] S_r \\ &= \text{trace}(S_c B S_c) S_r^2. \end{aligned}$$

The other inequalities in (2.3) follow from the fact that the trace of a product of two matrices is independent of the order of multiplication.

The derivation of (2.4) goes as follows:

$$\begin{aligned} \mathbf{E}(E C E) &= \mathbf{E}(S_c H S_r C S_c H S_r) \\ &= S_c \mathbf{E}(H S_r C S_c H) S_r \\ &= S_c (S_r C S_r)^T S_r \\ &= S_c^2 C^T S_r^2. \end{aligned}$$

□

For later reference we note that when  $B = A^T A$ , equation (2.3) reduces to

$$(2.6) \quad \mathbf{E}(E^T A^T A E) = \|S_c A\|_F^2 S_r^2.$$

**2.2. Properties of the stochastic norm.** The purpose of this subsection is to establish the basic properties of the stochastic norm defined by (1.3). The first step is show that it is indeed a norm.

**THEOREM 2.4.** *The function  $\|\cdot\|_S$  defined by (1.3) is a norm on  $\mathcal{G}^n$  or  $\mathcal{G}^{m \times n}$ . If  $\mathbf{E}(A^T B) = 0$  then*

$$(2.7) \quad \|A + B\|_S^2 = \|A\|_S^2 + \|B\|_S^2.$$

*If  $A$  and  $B$  are independent, then*

$$(2.8) \quad \|AB\|_S \leq \|A\|_S \|B\|_S.$$

*Proof.* We will show first that  $\|\cdot\|_S$  is a norm on  $\mathcal{G}^n$ . For any  $x, y \in \mathcal{G}^n$  define

$$\langle x, y \rangle = \mathbf{E}(x^T y).$$

The function  $\langle \cdot, \cdot \rangle$  is bilinear, symmetric, and definite in the sense that

$$x \neq 0 \iff \langle x, x \rangle > 0.$$

Hence  $\langle \cdot, \cdot \rangle$  is an inner product on  $\mathcal{G}^n$ , and the function  $\langle x, x \rangle^{1/2}$  is a norm. It is easily verified that  $\langle x, x \rangle = \|x\|_S^2$ .

To establish the result for  $\mathcal{G}^{m \times n}$ , identify  $\mathcal{G}^{m \times n}$  with  $\mathcal{G}^{mn}$  and observe that the matrix and vector norms are the same.

Equation (2.7) is established as follows:

$$\begin{aligned}\|A + B\|_S^2 &= \mathbf{E}\{\text{trace}[(A + B)^T(A + B)]\} \\ &= \mathbf{E}[\text{trace}(A^T A)] + 2\mathbf{E}[\text{trace}(A^T B)] + \mathbf{E}[\text{trace}(B^T B)] \\ &= \mathbf{E}[\text{trace}(A^T A)] + \mathbf{E}[\text{trace}(B^T B)] \\ &= \|A\|_S^2 + \|B\|_S^2.\end{aligned}$$

Finally to establish (2.8),

$$\|AB\|_S^2 = \mathbf{E}(\|AB\|_F^2) \leq \mathbf{E}(\|A\|_F^2 \|B\|_F^2) = \mathbf{E}(\|A\|_F^2) \mathbf{E}(\|B\|_F^2) = \|A\|_S^2 \|B\|_S^2. \quad \square$$

The next theorem shows how to calculate the stochastic norm of a single matrix.

**THEOREM 2.5.** *Let  $A = U + E$ , where  $U$  is constant and  $E \in \mathcal{G}^{m \times n}$  has mean zero. Let  $S$  be the matrix of standard deviations of the corresponding elements of  $E$ . Then*

$$(2.9) \quad \|A\|_S^2 = \|U\|_F^2 + \|E\|_S^2 = \|U\|_F^2 + \|S\|_F^2.$$

*In particular, if  $A \in \mathcal{T}^{m \times n}(U; S_r, S_c)$  then*

$$(2.10) \quad \|A\|_S^2 = \|U\|_F^2 + \|S_r\|_F^2 \|S_c\|_F^2.$$

*Proof.* The proof of (2.9) is purely computational and will be left to the reader. For (2.10), we need to show that  $\|E\|_S^2 = \|S_r\|_F^2 \|S_c\|_F^2$ . By (1.3) and (2.3),

$$\begin{aligned}\|E\|_S^2 &= \mathbf{E}[\text{trace}(E^T E)] \\ &= \text{trace}[\mathbf{E}(E^T E)] \\ &= \text{trace}[\text{trace}(S_c^2) S_r^2] \\ &= \text{trace}(S_c^2) \text{trace}(S_r^2) \\ &= \|S_r\|_F^2 \|S_c\|_F^2.\end{aligned} \quad \square$$

There are some observations to be made about this theorem. In the first place, a stochastic perturbation theory can, in principle, be based on (2.9) alone. However, in our applications we will be concerned with sums and products of matrices. Here any attempt to use (2.9) will result in a welter of incomprehensible formulas. However, if we restrict ourselves to cross-correlated errors, then Theorem 2.3 provides the wherewithal to produce simple expressions for the stochastic norm. Fortunately, the class cross-correlated matrices is extensive enough to be suitable for a wide variety of applications.

In the sequel we will take  $U = 0$ . Since this seems to be a restriction on our theory, an explanation is in order. Returning to the notation of the introduction, we note that  $F'_A(U + E) = F'_A(U) + F'_A(E)$ . Hence by Theorem 2.5,

$$\|F'_A(U + E)\|_S^2 = \|F'_A(U)\|_F^2 + \|F'_A(E)\|_S^2.$$

Thus the stochastic norm of the error in the first-order approximation decomposes into the Frobenius norm of a constant part and the stochastic norm of a random part. The constant part is just what would be obtained by applying first-order perturbation theory to  $U$ . Thus we take  $U = 0$  to focus attention on the random part, which is what is new in this paper. However, there is nothing to keep one from adding in a constant part if the application demands.

**2.3. Interpretation of the stochastic norm.** We now turn to the interpretation of the stochastic norm; i.e., to the second question in the introduction. It is not enough to know the size of  $\|A\|_S$ . We also need to know how much larger  $A$  can be than  $\|A\|_S$ . One answer is provided by the Chebyshev inequality, which says that for any random variable  $e$  with finite second moment,

$$\mathbf{P}\{|e| \geq \lambda \mathbf{E}(e^2)^{1/2}\} \leq \frac{1}{\lambda^2}.$$

Since,  $\mathbf{E}(\|A\|_F^2)^{1/2} = \|A\|_S$ , we have the following bound.

**THEOREM 2.6.** *Let  $A \sim \mathcal{G}^{m \times n}$ . Then*

$$(2.11) \quad \mathbf{P}\{\|A\|_F \geq \lambda \|A\|_S\} \leq \frac{1}{\lambda^2}.$$

Although this result holds for general matrices, its natural application is to matrices  $E$  with mean zero. It says that the probability of observing  $\|E\|_F$  to be larger than  $10\|E\|_S$  is less than one in one hundred. It should be appreciated that (2.11) is very conservative, since it takes into account the worst possible distributions. For most distributions, the probability is much less. For example, if the elements of  $E$  are independently, normally distributed random variables with mean zero and equal variance and  $mn > 10$ , then the probability of  $\|E\|_F$  being greater than  $2.5\|E\|_S$  is less than 0.005.

**2.4. Convergence of linear approximations.** As we indicated in the introduction to this paper, we will estimate perturbations  $\tilde{F}$  in a function  $F$  by computing the perturbation in a linearization  $\check{F}$ . In such an approach, there is always the problem of determining when the linearization is a good approximation to the actual value. In considering stochastic perturbations, we have the additional problem that the distribution of  $\tilde{F}$  may not have a mean or variance. What then does a value of  $\|\tilde{F} - F\|_S$  mean?

To illustrate the problem, let  $e$  be normally distributed with mean zero and standard deviation  $\sigma$ , which is presumed small. Let  $\tilde{\varphi} = 1/(1 - e)$  be a random perturbation of the function  $\varphi(x) = 1/(1 - x)$  at  $x = 0$ . We have  $\check{\varphi} = 1 + e$ , from which it follows that  $\|\check{\varphi} - \varphi\|_S = \sigma$ . On the other hand the density function of  $e$  is nonzero and continuous at the singularity  $e = 1$  of  $\tilde{\varphi}$ ; hence  $\tilde{\varphi}$  has neither mean nor variance. Yet one feels that the number  $\sigma$  should give us some information about the distribution of  $\tilde{\varphi}$ , since when  $\sigma$  is small it is exceedingly improbable that  $e$  will be anywhere near 1.

We will solve this problem by showing that  $\check{\varphi} - \varphi$  and  $\tilde{\varphi} - \varphi$ , suitably scaled, approach each other in probability.

**DEFINITION 2.7.** For each  $\lambda$  in an index set with limit point  $\mu$ , let  $e_\lambda$  be a random vector. Then  $e$  converges in probability to a random vector  $e$  if for every  $\epsilon > 0$

$$\lim_{\lambda \rightarrow \mu} \mathbf{P}\{\|e_\lambda - e\| \geq \epsilon\} = 0.$$

We write

$$\text{plim}_{\lambda \rightarrow \mu} e_\lambda = e.$$

It is easily verified that the definition is independent of the norm; in fact  $e_\lambda$  converges in probability to  $e$  if and only if the individual components of  $e_\lambda$  converge in probability



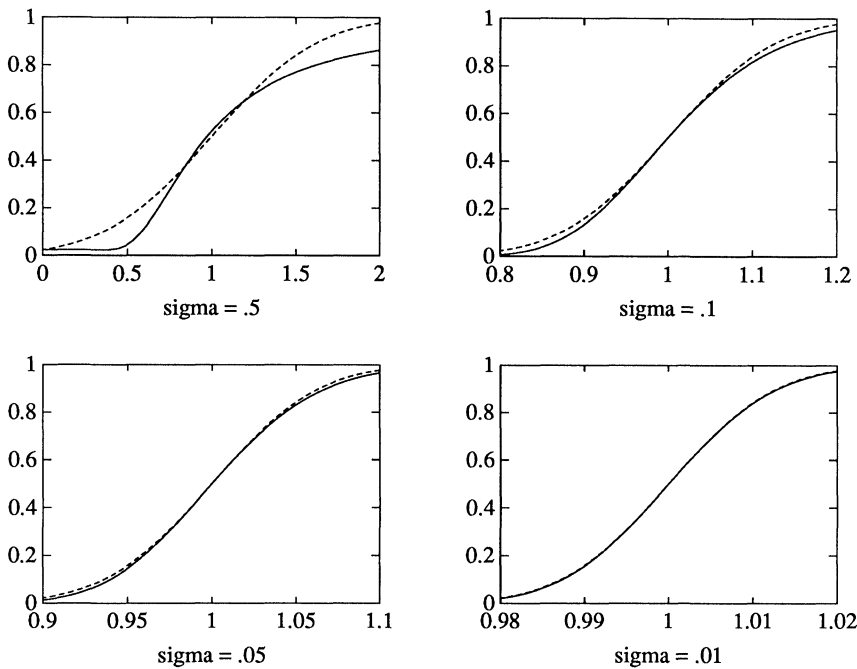


FIG. 2.1. Distributions of  $\check{\phi}$  (solid line) and  $\tilde{\phi}$  (dashed line) for normal  $e$ .

to the corresponding coefficients of  $e$ . Moreover, if  $f$  is continuous at the point  $e$ , then  $\text{plim } e_\lambda = e$  implies that  $\text{plim } f(e_\lambda) = f(e)$ .

For our problem, the critical fact is that if  $\text{plim}_{\lambda \rightarrow \mu} e_\lambda = e$ , then the distribution function of  $e_\lambda$  converges to that of  $e$  at all points of continuity. This has the following consequence for our example. We will show (Theorem 2.13 below) that

$$(2.12) \quad \text{plim}_{\sigma \rightarrow 0} \frac{\tilde{\phi} - 1}{\sigma} = \frac{\check{\phi} - 1}{\sigma}.$$

Suppose we use the fact that  $(\check{\phi} - 1)/\sigma$  is normally distributed with mean zero and variance 1 to predict that  $\check{\phi}$  lies in the interval  $(1 - 3.3\sigma, 1 + 3.3\sigma)$  with probability greater than 0.999. Then ultimately the same holds true for  $\tilde{\phi}$ . Figure 2.1 illustrates convergence of the distribution of  $\tilde{\phi}$  to that of  $\check{\phi}$  for the case where  $e$  is normal.

A formal justification of the above claims is provided by the following theorem.

**THEOREM 2.8.** *Let  $f : \mathcal{R}^n \rightarrow \mathcal{R}$  have a Frechet derivative  $f'_x$  at the point  $x$  and let  $e \sim \mathcal{G}^n(0, \Sigma)$ . Then*

$$(2.13) \quad \text{plim}_{\Sigma \rightarrow 0} \frac{f(x + e) - f(x) - f'_x{}^T e}{\|\Sigma^{1/2}\|_F} = 0.$$

Moreover, if there is an  $\alpha > 0$  such that as  $\Sigma \rightarrow 0$

$$(2.14) \quad \|\Sigma^{1/2} f'_x\| \geq \alpha \|\Sigma^{1/2}\| \|f'_x\| > 0,$$

then

$$(2.15) \quad \text{plim}_{\Sigma \rightarrow 0} \frac{f(x + e) - f(x) - f'_x{}^T e}{\|\Sigma^{1/2} f'_x\|} = 0.$$

*Proof.* We will prove (2.15), the proof of (2.13) being similar. For  $e \neq 0$  define

$$\gamma(e) = \frac{|f(x+e) - f(x) - f'_x{}^T e|}{\|e\|}$$

and set  $\gamma(0) = 0$ . Then  $\gamma$  is continuous at zero and

$$(2.16) \quad \lim_{\Sigma \rightarrow 0} \gamma(e) = 0.$$

Now let  $\epsilon, \delta > 0$  be given. It is sufficient to show that for  $\Sigma$  sufficiently small

$$\mathbf{P}\{\|e\|\gamma(e) \geq \|\Sigma^{1/2} f'_x\| \epsilon\} < \delta.$$

In view of (2.14) this will be true if

$$(2.17) \quad \mathbf{P}\{\|e\|\gamma(e) \geq \alpha \|\Sigma^{1/2}\| \|f'_x\| \epsilon\} < \delta.$$

By the Chebyshev inequality (cf. Theorem 2.6) there is a  $\beta \geq 1$  so that

$$\mathbf{P}\{\|e\| \geq \beta \|\Sigma^{1/2}\|\} < \frac{\delta}{2},$$

independently of  $\Sigma$ . From (2.16) it follows that for all  $\Sigma$  sufficiently small

$$\mathbf{P}\{\gamma(e) \geq \alpha \beta^{-1} \|f'_x\| \epsilon\} < \frac{\delta}{2}.$$

It follows that

$$\begin{aligned} \mathbf{P}\{\|e\|\gamma(e) \geq \alpha \|\Sigma^{1/2}\| \|f'_x\| \epsilon\} &\leq \mathbf{P}\{\|e\| \geq \beta \|\Sigma^{1/2}\|\} + \mathbf{P}\{\gamma(e) \geq \alpha \beta^{-1} \|f'_x\| \epsilon\} \\ &< \frac{\delta}{2} + \frac{\delta}{2} = \delta, \end{aligned}$$

which establishes the theorem.  $\square$

There are some technical comments and some general observations to be made about this theorem. We will begin with the technical comments.

The denominator  $\|\Sigma^{1/2} f'_x\|$  is the standard deviation of  $\check{f}$  and serves the same role as the denominator  $\sigma$  in (2.12). Condition (2.14) says that this standard deviation must not decrease more rapidly than  $\Sigma^{1/2}$ , as could happen when  $\Sigma$  is degenerate or when  $f'_x = 0$ . Equation (2.17) shows clearly that convergence will be delayed when either  $\alpha$  or  $f'_x$  is small.

Another way of looking at this is to realize that if the standard deviation of  $\check{f}$  is zero, it is impossible to scale the distribution  $\check{f} - f$ . However, in this case (2.13) says that the distribution of  $\check{f}$  degenerates superlinearly—which is almost as good as having zero variance.

Condition (2.14) can be replaced by

$$\alpha_{\text{opt}} \stackrel{\text{def}}{=} \liminf_{\Sigma \rightarrow 0} \frac{\|\Sigma^{1/2} f'_x\|}{\|\Sigma^{1/2}\| \|f'_x\|} > 0.$$

In general  $\alpha_{\text{opt}}$  will depend on  $\Sigma$  and its relation to  $f'_x$ . However if the condition number  $\kappa(\Sigma^{1/2}) = \|\sigma^{1/2}\| \|\Sigma^{-\frac{1}{2}}\|$  is uniformly bounded, we may take

$$\alpha = \liminf \kappa^{-1}(\Sigma^{1/2})$$

as a lower bound on  $\alpha_{\text{opt}}$ .

The first general observation to be made is that this is a distribution-free result. Not only does it not assume that  $e$  has a particular distribution, but it does not assume that  $e$  belong to a particular class of distributions (e.g., normal) as  $\Sigma \rightarrow 0$ .

The price to be paid for the generality of the theorem is that it does not give explicit error bounds, something it shares with many asymptotic results from probability theory.<sup>3</sup> In the sequel we will use results from perturbation theory to evaluate the domain of applicability of the theorem.

One of the referees has suggested that sharper results may be obtained by assuming that  $e$  is uniformly distributed in a sphere not containing a singularity of  $f$ , in which case  $f(e)$  has second moments. Of course, if this is the distribution appropriate to the application at hand, then one should use it. But many applications require normal distributions (see §3.4 below), or even distributions with heavier tails.

The notion that a uniform distribution will produce sharper bounds is worth a closer examination. Since the stochastic norm depends only on first and second moments, it is effectively independent of the form of distribution, which enters only via its effect on the rate of convergence of the linear approximation. Now the proof of Theorem 2.8 shows that simply excluding singularities from support of the distribution is not enough. The crux of the matter is whether the distribution is concentrated in a region where a linear approximation is valid. In this respect the uniform distribution is at a disadvantage compared to distributions, like the normal distribution, whose density drops off very rapidly away from its mean. However, we should not make too much of this. Although a comparison of Figs. 2.1 and 2.2 shows that convergence is slower for the uniform distribution, it is not very much slower.

**2.5. Complex values.** Since some of the objects we will be treating, such as eigenvalues and their eigenvectors, can have complex values, it is important to indicate how the results of this section are affected by the switch from real to complex numbers.

The calculus of expectations remains unchanged as long as we replace the product  $xy$  by  $\bar{x}y$ , so that  $x^2$  becomes  $|x|^2$ . In particular if we replace the transpose by the conjugate transpose, the results on cross-correlated matrices remain unaltered. Since the results on the stochastic norm and the convergence theorem deal with real-valued quantities, they also remain unaltered.

**Notes and references.** The background for this section will be found in almost any probability or statistics book that treats multivariate distributions. Elementary treatments may be found in [21],[30].

The notation  $\mathcal{G}(M, \Sigma)$  was suggested by the use of the letter G in queuing theory to stand for a general distribution, a practice started by Kendal [26].

§2.1. The material in this section appeared in some lecture notes by the author (c. 1982). Theorem 2.3 has been published by Neudecker and Wansbeek [28]. Although their paper treats normal matrices, their proof is quite general.

<sup>3</sup> However, its proof does provide hints about what makes for fast convergence. For example,  $\alpha$  should be large, and the distribution should have small tails so that  $\beta$  is small. Moreover,  $\gamma$  should not grow swiftly; i.e., the first-order approximation should be good.

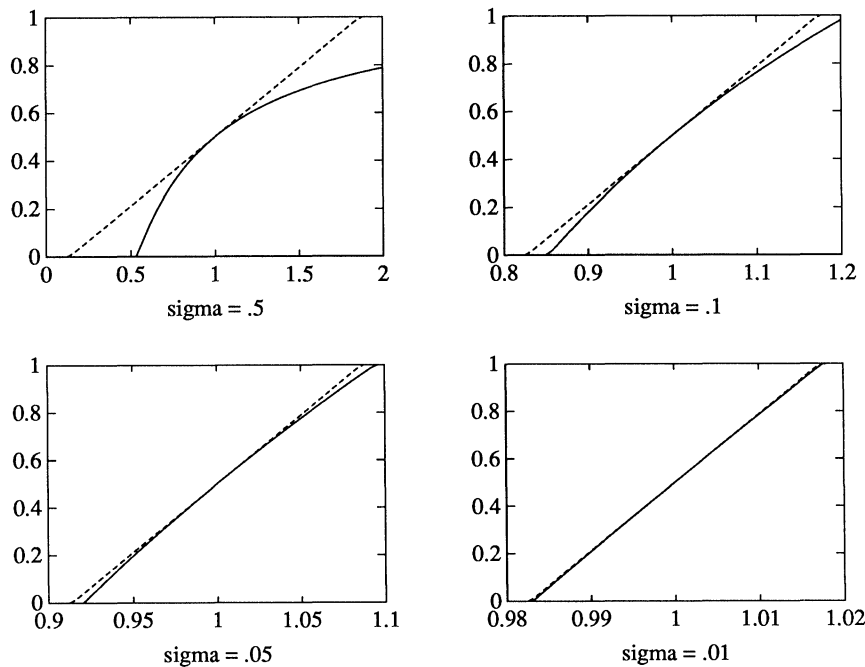


FIG. 2.2. Distributions of  $\check{\phi}$  (solid line) and  $\tilde{\phi}$  (dashed line) for uniform  $e$ .

§2.2. The formal use of the function  $\mathbf{E}[\text{trace}(X^T X)]$  as a norm on random matrices appears to be new. Its major problem is that the submultiplicative inequality (2.8) can fail. For example, if  $e$  is distributed normally with mean zero and variance 1, then

$$\|e \cdot e\|_S^2 = \mathbf{E}(e^4) = 3 > 1 = \mathbf{E}(e^2)\mathbf{E}(e^2) = \|e\|_S^2 \|e\|_S^2.$$

The inequality can even fail for uncorrelated matrices.

§2.4. Theorem 2.8 is the author's, who first proved (2.13) in [37]. Serfling [32, Thm. 3.3A] gives a similar theorem, but with  $e$  normal,  $\Sigma$  of the form  $b_n \Sigma_0$  for fixed  $\Sigma_0$ , and convergence in distribution. It is worth noting that if  $e = \Sigma^{1/2} e_0$ , where  $e_0 \sim \mathcal{G}^n(0, I)$  is a fixed distribution, then the convergence is with probability 1.

**3. Pseudo-inverses, least squares, and projections.** In this section we will consider perturbation of pseudo-inverses, least squares solutions, and projections. Throughout this section  $A$  will denote an  $m \times n$  matrix of rank  $n$ . The matrix  $C = A^T A$  is the cross-product matrix. The matrix  $P = A A^\dagger$  is the orthogonal projection onto the column space of  $A$  and  $P_\perp = I - P$  is its complementary projection.

We will assume that the perturbation matrix  $E$  is distributed  $\mathcal{T}^{m \times n}(0; S_r, S_c)$ . The expressions we derive will be simplest when  $S_r = \sigma I_n$  and  $S_c = I_m$  so that the elements of  $E$  are uncorrelated with variance  $\sigma^2$ . We will refer to this as a *simple perturbation*.

**3.1. The pseudo-inverse.** In this subsection we will consider perturbations of the pseudo-inverse of  $A$  defined by

(3.1) 
$$A^\dagger = (A^T A)^{-1} A^T.$$

**3.1.1. Perturbation expansion.** The perturbation expansion for  $A^\dagger$  is easily derived by replacing  $A$  by  $\tilde{A} = A + E$  in (3.1), using the linear part of the Neumann expansion of  $(\tilde{A}^T \tilde{A})^{-1}$ , dropping higher-order terms, and simplifying. The result is

$$(3.2) \quad \check{A}^\dagger = A^\dagger - A^\dagger E A^\dagger + C^{-1} E^T P_\perp.$$

**3.1.2. Range of applicability.** It is important that we have some idea of when the linear approximation (3.2) is valid. There are two conditions that must be satisfied.

In the first place, the matrix  $A + E$  must be of full rank. Most of the perturbation bounds in the literature are derived under the supposition that

$$\|A^\dagger E\| < 1$$

or the stronger condition

$$(3.3) \quad \|A^\dagger\| \|E\| < 1,$$

both of which insure that  $A$  is of full rank. In keeping with our program, let us derive the stochastic norm of  $A^\dagger E$ . We have

$$\|A^\dagger E\|_S^2 = \text{trace}(E^T A^{\dagger T} A^\dagger E).$$

By (2.3),

$$\|A^\dagger E\|_S^2 = \|A^\dagger S_c\|_F^2 \text{trace}(S_r^2).$$

Hence

$$(3.4) \quad d_1 \stackrel{\text{def}}{=} \|A^\dagger E\|_S = \|A^\dagger S_c\|_F \|S_r\|_F,$$

or in the case of a simple perturbation

$$d_1 = \sqrt{n} \sigma \|A^\dagger\|_F.$$

If  $d_1$  is near 1, we should not trust  $\check{A}^\dagger$  to approximate  $\tilde{A}^\dagger$ .

A second source of nonlinearity is bias in the cross-product matrix  $C = A^T A$ . Specifically, we have

$$\tilde{C} = (A + PE)^T (A + PE) + E^T P_\perp E.$$

Since the diagonals  $E^T P_\perp E$  are nonnegative, its addition causes an upward bias in the diagonals of  $(A + PE)^T (A + PE)$ . More generally, any quadratic form  $v^T C v$  is biased upward by  $v^T E^T P_\perp E v$ . The expected value of  $v^T E^T P_\perp E v / v^T C v$  is

$$\|P_\perp S_c^2\|_F^2 \frac{v^T S_r^2 v}{v^T C v}.$$

Setting  $w = Av$ , so that  $v = A^\dagger w$ , we see that the maximum of this expectation is

$$(3.5) \quad d_2^2 \stackrel{\text{def}}{=} \|S_r A^\dagger\|^2 \|P_\perp S_c\|_F^2,$$

or in the simple case

$$d_2^2 = (m - n) \sigma^2 \|A^\dagger\|^2.$$

Again, unless  $d_2^2$  is substantially less than 1, we should not trust the linear approximation.

It is instructive to compare the two diagnostics  $d_1$  and  $d_2^2$  in the case of a simple perturbation. The latter varies as  $\sigma^2$  and, as the error approaches zero, is dominated by the former, which varies as  $\sigma$ . On the other hand, suppose that  $\sigma$  is fixed and we add rows to  $A$  in such a way that  $\lim_{m \rightarrow \infty} m^{-1} A^T A = C_\infty$ . Then  $A^\dagger = O(1/\sqrt{m})$ , and  $d_1 \rightarrow 0$  while  $d_2^2$  remains uniformly positive. Thus  $d_1$  measures an effect that dominates for small errors, while  $d_2$  measures an effect that dominates as we increase the size of the problem, holding the size of the errors fixed.

**3.1.3. The perturbation estimate.** We now turn to the computation of  $\|\check{A}^\dagger - A^\dagger\|_S$ , where  $\check{A}$  is given by (3.2). Since  $(A^\dagger E A^\dagger)(C^{-1} E^T P_\perp)^T = 0$ , the matrices  $A^\dagger E A^\dagger$  and  $C^{-1} E^T P_\perp$  are uncorrelated, and we may bound them separately [cf. (2.7)].

We have first

$$\|A^\dagger E A^\dagger\|_S^2 = \mathbf{E}[\text{trace}(A^{\dagger T} E^T A^{\dagger T} A^\dagger E A^\dagger)] = \text{trace}[A^{\dagger T} \mathbf{E}(E^T A^{\dagger T} A^\dagger E) A^\dagger].$$

Hence by (2.6),

$$\|A^\dagger E A^\dagger\|_S^2 = \|A^\dagger S_c^2\|_F^2 \text{trace}(A^{\dagger T} S_r^2 A^\dagger) = \|A^\dagger S_c\|_F^2 \|S_r A^\dagger\|_F^2.$$

Similarly<sup>4</sup>

$$\|C^{-1} E^T P_\perp\|_S^2 = \|P_\perp S_c\|_F^2 \|S_r C^{-1}\|_F^2.$$

Hence

$$(3.6) \quad \|\check{A}^\dagger - A^\dagger\|_S = \sqrt{\|A^\dagger S_c\|_F^2 \|S_r A^\dagger\|_F^2 + \|P_\perp S_c\|_F^2 \|S_r C^{-1}\|_F^2},$$

or in the simple case

$$\|\check{A}^\dagger - A^\dagger\|_S = \sigma \sqrt{\|A^\dagger\|_F^4 + (m - n) \|C^{-1}\|_F^2}.$$

Since  $\|C^{-1}\|_F = \|A^\dagger A^{\dagger T}\|_F \leq \|A^\dagger\|_F^2$ , we have for the simple case

$$\frac{\|\check{A}^\dagger - A^\dagger\|_S}{\|A^\dagger\|_F} \leq \sigma \sqrt{m - n + 1} \|A^\dagger\|_F.$$

The right-hand side of this inequality may be further manipulated to give

$$(3.7) \quad \frac{\|\check{A}^\dagger - A^\dagger\|_S}{\|A^\dagger\|_F} \leq \kappa(A) \frac{\sqrt{m - n + 1} \sigma}{\|A\|_F},$$

where

$$\kappa(A) = \|A\|_F \|A^\dagger\|_F$$

is the *condition number* of  $A$ . In this form the bound is similar to others appearing in the literature. However, we have obtained this pretty form at the expense of sharpness, and in the sequel we will not massage our formulas beyond simple equalities, unless we are forced to do so.

<sup>4</sup> In the sequel we will omit the routine computation of stochastic norms.

**3.2. Least squares and projections.** One seldom has cause to bound perturbations of the pseudo-inverse alone, since in most applications the pseudo-inverse is invoked only to be applied to a vector or matrix. In particular, it is well known that the vector  $x = A^\dagger b$  solves the least squares problem of minimizing  $\|b - Ax\|^2$ . In this case the residual vector  $r = b - Ax$  is the projection onto the orthogonal complement of the column space of  $A$ ; that is,  $r = P_\perp b$ . We now turn to estimating the sizes of the perturbations in  $x$ ,  $P$ , and  $r$ .

**3.2.1. Least squares solutions.** A perturbation expansion for  $\check{x}$  can be easily found from the expression (3.2) for  $\check{A}$ ; namely,

$$(3.8) \quad \check{x} = x - A^\dagger E x + C^{-1} E^T r.$$

From this a perturbation estimate is easily calculated in the form

$$\|\check{x} - x\|_S = \sqrt{\|A^\dagger S_c\|_F^2 \|S_r x\|^2 + \|S_c r\|^2 \|S_r C\|_F^2}.$$

In the simple case this becomes

$$\|\check{x} - x\|_S = \sigma \sqrt{\|A^\dagger\|_F^2 \|x\|^2 + \|r\|^2 \|C\|_F^2}.$$

**3.2.2. Individual components.** It is useful to have estimates for the errors in the individual components of  $x$ . Multiplying by the transpose of the  $i$ th unit vector  $\mathbf{1}_i$ , we get

$$\check{x}_i = x_i - a_i^{(\dagger)T} E x + c_i^{(-1)T} E^T r,$$

where  $a_i^{(\dagger)T}$  is the  $i$ th row of  $A^\dagger$  and  $c_i^{(-1)T}$  is the  $i$ th row of  $B^{-1}$ . From this we get

$$(3.9) \quad \|\check{x}_i - x_i\|_S = \sqrt{\|S_c a_i^{(\dagger)}\|^2 \|S_r x\|^2 + \|S_r c_i^{(-1)}\|^2 \|S_c r\|^2}.$$

A particularly interesting special case occurs when only one column of  $A$ , say the  $j$ th, is permitted to vary; that is, when  $S_c = I$  and  $S_r = \sigma_j \mathbf{1} \mathbf{1}_j^T$ . In this case (3.9) becomes

$$(3.10) \quad \|\check{x}_i - x_i\|_S = \sigma_j \sqrt{\|a_i^{(\dagger)}\|^2 |x_j|^2 + |c_{ij}^{(-1)}|^2 \|r\|^2}.$$

Thus the quantity  $\sqrt{\|a_i^{(\dagger)}\|^2 |x_j|^2 + |c_{ij}^{(-1)}|^2 \|r\|^2}$  is a condition number for  $x_i$  with respect to perturbations in the  $j$ th column of  $A$ .

**3.2.3. Projections and the residual vector.** The perturbation expansion for the projection operator  $\check{P} = \check{A} \check{A}^\dagger$  may be found by replacing  $\check{A}^\dagger$  by  $\check{A}^\dagger$ , dropping second-order terms, and simplifying. The result is

$$\check{P} = P - P_\perp E A^\dagger - (P_\perp E A^\dagger)^T.$$

From this a perturbation estimate is easily computed:

$$\|\check{P} - P\|_S = \sqrt{2} \|S_c P_\perp\|_F \|S_r A^\dagger\|_F,$$

or for the simple case

$$\|\check{P} - P\|_S = \sqrt{2(m-n)} \|A^\dagger\|_F \sigma.$$

Since  $P_{\perp} = I - P$ , the same perturbation estimates hold also for  $P_{\perp}$ .

For the residual vector  $r = P_{\perp}b$  we have

$$\check{r} = r + P_{\perp}Ex + A^{\dagger T}E^T r.$$

Hence

$$\|\check{r} - r\|_S = \sqrt{\|S_c P_{\perp}\|_F^2 \|S_r x\|^2 + \|S_c r\|^2 \|S_r A^{\dagger}\|_F^2}.$$

In the simple case this becomes

$$\|\check{r} - r\|_S = \sigma \sqrt{(m - n)\|x\|^2 + \|r\|^2 \|A^{\dagger}\|_F^2}.$$

**3.3. The inverse matrix and linear equations.** When  $A$  is square, so that  $m = n$ , we have  $P_{\perp} = 0$ . Consequently, perturbation results for the inverse and for linear systems may be obtained trivially from those of the previous subsections by taking  $m = n$  and dropping all terms containing  $P_{\perp}$ .

**3.3.1. The inverse matrix.** The perturbation expansion for  $\tilde{A}$  may be obtained from (3.2):

$$(3.11) \quad \check{A}^{-1} = A^{-1} - A^{-1}EA^{-1}.$$

For this expansion to be valid, we must have  $d_1$  defined by (3.4) to be substantially less than 1. However, the quantity  $d_2^2$  defined by (3.5) is zero and need not be considered.

The perturbation estimate can be obtained from (3.6):

$$\|\check{A}^{-1} - A^{-1}\|_S = \|A^{-1}S_c\|_F \|S_r A^{-1}\|_F,$$

or in the simple case

$$\|\check{A}^{-1} - A^{-1}\|_S = \sigma \|A^{-1}\|_F^2.$$

Note that (3.7) now becomes

$$\frac{\|\check{A}^{-1} - A^{-1}\|_S}{\|A^{-1}\|_F} = \kappa(A) \frac{\sigma}{\|A\|_F}$$

with equality instead of inequality.

**3.3.2. Linear systems.** The perturbation expansion for  $\tilde{x} = \tilde{A}^{-1}b$  may be obtained from (3.8) or directly from (3.11):

$$\check{x} = x - A^{-1}Ex.$$

The corresponding estimate is given by

$$\|\check{x} - x\|_S = \|A^{\dagger}S_c\|_F \|S_r x\|,$$

and for the simple case by

$$\|\check{x} - x\|_S = \sigma \|A^{\dagger}\|_F \|x\|.$$



**3.4. An application.** In this subsection we will give an application of the above theory to the statistical analysis of regression problems with errors in the regression matrix. The standard model for the ordinary regression problem is written<sup>5</sup>

$$b = Ax + e,$$

where  $e$  is a vector of independent normal variates with mean zero and common variance  $\sigma^2$ . The vector of estimated regression coefficients is

$$(3.12) \quad \hat{x} = A^\dagger b = x + A^\dagger e.$$

Since  $A^\dagger e$  is linear in  $e$  we can approximate the distribution of  $\hat{b}$  provided we have an independent estimate of  $\sigma$ . It turns out that

$$\hat{\sigma} = \frac{\|P_\perp b\|}{\sqrt{m-n}}$$

is just such an estimate.

It sometimes happens that  $A$  cannot be observed directly but is measured or otherwise contaminated with errors. Thus the regression matrix we have at hand is  $\tilde{A} = A + E$ . In one widely used model it is assumed that

$$E \sim \mathcal{T}^{m \times n}(0; \Sigma^{1/2}, I)$$

and is normally distributed. If  $\Sigma$  is unknown, we will be forced to work with the estimate

$$\tilde{x} = \tilde{A}^\dagger b$$

instead of  $\hat{x}$ . Obviously,  $\tilde{x}$  is a nonlinear function of  $E$ . Nonetheless, if  $E$  is reasonably small it is well behaved.

To see why this should be true, rewrite the perturbed model in the form

$$b = \tilde{A}x + (e - Ex).$$

It then follows that

$$\tilde{x} = x + \tilde{A}^\dagger(e - Ex) = x + A^\dagger(e - Ex) + F^H(e - Ex),$$

where we have written  $F^H = \tilde{A}^\dagger - A^\dagger$ . Now as  $E$  becomes small  $F$  becomes small. Consequently,

$$\tilde{x} \cong \check{x} \equiv x + A^\dagger(e - Eb).$$

Comparing this equation with (3.12), we see that

*If  $F$  is small compared with  $\tilde{A}$ , then  $\tilde{x}$  behaves as if it came from the linear model  $b = Ax + (e - Ex)$ .*

Since the variance of the components of  $e + Ex$  is  $\sigma^2 + x^H \Sigma x$ , the perturbed model behaves as if the variance of the components of  $E$  had been inflated by  $x^H \Sigma x$ . In much the same way, we can show that

$$\tilde{\sigma} = \frac{\|\tilde{P}_\perp b\|}{\sqrt{m-n}}$$

<sup>5</sup> This is the model as a numerical analyst would write it. A statistician would write something like  $y = X\beta + e$ , where  $X$  is an  $n \times p$  matrix.

is asymptotically an independent estimate of  $\sqrt{\sigma^2 + x^H \Sigma x}$ , so that the usual least squares procedures work without further alteration.

It should be noted that this is not merely a continuity result;  $E$  does not have to be so small that its errors are negligible. On the contrary, it is possible for  $e$  to be zero, so that all the variability in the problem comes from  $E$ .

Nor is  $\check{x}$  a first-order expansion. Since we do not assume that  $\sigma \rightarrow 0$ , the term  $F^H e$ , which we have thrown away, is of the first order. However, it is dominated by the term  $A^\dagger e$  and must ultimately become negligible.

To apply these observations we must determine when  $F$  is negligible compared to  $A$ . Here we may apply the theory developed previously. Specifically, if  $\check{F} = \check{A}^\dagger - A^\dagger$ , then it follows from (3.6) that

$$\frac{\|\check{F}\|_s}{\|A^\dagger\|_F} \leq \sqrt{m} \|\Sigma^{1/2} A^\dagger\| \equiv \tau.$$

Thus if  $\tau$  is reasonably less than 1, then we are secure in applying least squares analysis to the perturbed model. Note that we do not require a detailed knowledge of  $\Sigma$ ; crude information sufficient to bound  $\tau$  is sufficient. Simulations suggest that least squares analysis can be trusted when  $\tau$  is less than 0.3.

**Notes and references.** Perturbation theory for the least squares problem begins with Golub and Wilkinson [18], who produced first-order expansions for least squares solutions. For surveys of the theory, see [44], [35], [42]. For extensions to more general problems, see [10], [29].

In the statistics literature, Hodges and Moore [19] and Davies and Hutton [7] have used first order approximations to least squares solutions to assess the effects of errors in the regression matrix—the problem of errors in the variables as it is known. See also [2] for applications to econometrics.

§3.1.2. The presence of bias in regression coefficients with errors in the variables has long been known to statisticians under the name asymptotic inconsistency; e.g., see [23, §9.4], and [1], [40].

§3.2.1. As mentioned above, the perturbation expansion is due to Golub and Wilkinson [18].

§3.2.2. The quantities in (3.10) were derived in [36] as sensitivity coefficients for linear regression.

§3.2.3. The quantity  $\|\tilde{P} - P\|_F^2$  is twice the sum of squares of the sines of the canonical angles between  $\mathcal{R}(P)$  and  $\mathcal{R}(\tilde{P})$  (see [35] for definitions). Thus the bounds estimate how far the column space of  $\tilde{A}$  deviates from that of  $A$ .

§3.4. The basic idea underlying this section is due to David and Stewart [6]. The present treatment is taken from [41]. For an introduction to regression analysis with a survey of the errors-in-the-variables problem, see [31].

When  $\tau$  is too large, one must have recourse to other techniques—techniques that require a precise knowledge of  $\Sigma$ . As usual, statisticians and numerical analysts have worked on the problem without consulting one another. Fuller's book on measurement error models [11] is the definitive source for the statistical approach (it contains much more than the model treated here). On the numerical side, Golub and Van Loan [16] (also see [43]) have created a technique based on the singular value decomposition known as total least squares, which is closely related to the statisticians technique. When  $\tau$  is small, total least squares and least squares give essentially the same results [38].

**4. Eigenvalue problems.** In this section we will be concerned with stochastic perturbation theory for certain eigenvalue problems. In the first subsection we will treat the perturbation of invariant subspaces and the perturbation of the representation of an operator on an invariant subspace. In the next two subsections we will consider the simplifications that obtain first when the matrix is symmetric and then when the dimension of the invariant subspace is 1 (i.e., the perturbation of eigenvalues and eigenvectors). The section concludes with a treatment of the singular value decomposition. Throughout this section,  $A$  will denote a square matrix of order  $n$ .

**4.1. Invariant subspaces.** It is convenient to approach the perturbation of eigenvectors and eigenvalues from the problem of the perturbation of invariant subspaces, since the latter, more general problem is of independent interest. However, this approach exacts a toll. The perturbation expansion for an invariant subspace involves a linear operator that does not interact nicely with cross-correlated matrices. The consequence is that we can only give bounds on the stochastic norms of the perturbations, whereas previously we have been able to compute the norms exactly. Fortunately, in some important special cases we can restore the lost equality.

We will start with a review of some facts about invariant subspaces, which will also fix the notation to be used throughout this section.

**4.1.1. The representation of invariant subspaces.** We begin with the definition of an invariant subspace of  $A$ .

**DEFINITION 4.1.** Let  $\mathcal{X}_1$  be a  $k$  dimension subspace of  $\mathcal{C}^n$ . Then  $\mathcal{X}_1$  is an *invariant subspace* of  $A$  if

$$(4.1) \quad A\mathcal{X}_1 \stackrel{\text{def}}{=} \{Ax : x \in \mathcal{X}_1\} \subset \mathcal{X}_1.$$

In other words,  $\mathcal{X}_1$  is an invariant subspace of  $A$  if  $A$  maps  $\mathcal{X}_1$  into itself.

Let  $(X_1 \ Y_2)$  be a unitary matrix with  $\mathcal{R}(X_1) = \mathcal{X}_1$ . Then from (4.1) it follows that the columns of  $AX_1$  can be written as a linear combination of  $X_1$ ; that is,

$$(4.2) \quad AX_1 = X_1 A_{11},$$

where

$$A_{11} = X_1^H A X_1.$$

Moreover, from (4.2) and the fact that  $Y_2^H X_1 = 0$ , we have

$$(4.3) \quad Y_2^H A X_1 = 0.$$

These relations may be summarized by writing

$$(4.4) \quad \begin{pmatrix} X_1^H \\ Y_2^H \end{pmatrix} A (X_1 \ Y_2) = \begin{pmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{pmatrix},$$

which also defines  $A_{12}$  and  $A_{22}$ .

When  $A$  is regarded as an operator on  $\mathcal{X}_1$ , the matrix  $A_{11}$  is the representation of  $A$  with respect to the basis formed by the columns of  $X_1$ . It is easy to see that  $\mathcal{R}(Y_2)$  is a left invariant subspace of  $A$ , or equivalently

$$Y_2^H A = A_{22} Y_2^H.$$

Thus  $A_{22}$  is the representation of  $A$  of the left invariant subspace  $\mathcal{R}(Y_2)$  with respect to the basis formed by the rows of  $Y_2^H$ .

**4.1.2. Simple invariant subspaces.** An important example of an invariant subspace is the space spanned by an eigenvector  $x_1$ . Unfortunately, even when they are normalized, say by requiring  $\|x_1\| = 1$  and that some specific nonzero component of  $x_1$  be positive, eigenvectors need not be unique. For example, if  $A$  is Hermitian and it has an eigenvalue of multiplicity  $m$ , the corresponding eigenvectors span a space of dimension  $m$ . In perturbation theory, the usual way of getting around this problem, is to assume that the eigenvalue is simple, so that  $A$  has a unique (normalized) eigenvector.

The notion of simplicity can be generalized to invariant subspaces by observing that the set eigenvalues  $\lambda(A)$  of  $A$  is the union of  $\lambda(A_{11})$  and  $\lambda(A_{22})$ . When  $A_{11}$  is a scalar (i.e., an eigenvalue), it is simple if and only if it is not also an eigenvalue of  $A_{22}$ . This leads to the following definition.

DEFINITION 4.2. The invariant subspace  $\mathcal{X}_1$  is *simple* if

$$\lambda(A_{11}) \cap \lambda(A_{22}) = \emptyset.$$

In deriving a perturbation bound, it will be important to be able to solve Sylvester equations of the form

$$A_{22}P - PA_{11} = E_{21}.$$

It turns out that this is equivalent to requiring  $\mathcal{X}_1$  to be simple, as the following widely used theorem shows.

THEOREM 4.3. Let the linear operator  $\mathbf{T}$  be defined by

$$(4.5) \quad \mathbf{T} = P \mapsto A_{22}P - PA_{11}.$$

Then  $\mathbf{T}$  is nonsingular if and only if

$$\lambda(A_{11}) \cap \lambda(A_{22}) = \emptyset.$$

Moreover, if we set

$$(4.6) \quad \delta \stackrel{\text{def}}{=} \inf_{\|P\|_{\mathbb{F}}=1} \|\mathbf{T}P\|_{\mathbb{F}} = \|\mathbf{T}^{-1}\|^{-1},$$

then

$$(4.7) \quad \delta \leq \min\{|\lambda_1 - \lambda_2| : \lambda_1 \in A_{11}, \lambda_2 \in A_{22}\}.$$

**4.1.3. Representation of  $\tilde{\mathcal{X}}$ .** Let  $\tilde{A} = A + E$ , where  $E \sim \mathcal{T}(0, S_r, S_c)$ . We will be concerned with the effects of  $E$  on the simple invariant subspace  $\mathcal{X}_1$ . For now we will assume that there is an invariant subspace  $\tilde{\mathcal{X}}_1$  of  $\tilde{A}$  which approaches  $\mathcal{X}_1$  as  $E$  approaches zero. We will justify this assumption in §4.1.6.

In order to obtain a perturbation expansion for the invariant subspace  $\tilde{\mathcal{X}}_1$  we must address two problems. The first is to represent  $\tilde{\mathcal{X}}_1$  in such a way that we can measure its distance from  $\mathcal{X}_1$ . The second is to find a perturbation equation from which we may cast out higher-order terms.

Turning to the first problem, we will seek a basis for  $\tilde{\mathcal{X}}_1$  in the form

$$(4.8) \quad \tilde{X}_1 = X_1 + Y_2P,$$

where  $P$  is to be determined. There are two reasons for this choice.

First, it is easily verified that  $\tilde{X}_1(I + P^H P)^{-1/2}$  has orthonormal columns. In other words, up to second-order terms in  $P$ , which is presumed small, the columns of  $\tilde{X}$  form an orthonormal basis for  $\tilde{\mathcal{X}}_1$ .

Second, there are many ways of choosing bases for  $\mathcal{X}_1$  and  $\tilde{\mathcal{X}}_1$ , some of which may be quite different even when  $\mathcal{X}_1$  and  $\tilde{\mathcal{X}}_1$  are near. However, if we define  $\tilde{X}_1$  by (4.8), then of all matrices whose column spaces span  $\mathcal{X}_1$ , the matrix  $X_1$  is nearest  $\tilde{X}_1$ . Specifically, we have the following theorem.

**THEOREM 4.4.** *If  $(X_1 \ Y_2)$  is unitary and  $\tilde{X}_1$  is defined by (4.8), then  $X_1$  solves the following least squares problem*

$$(4.9) \quad \begin{array}{ll} \text{minimize:} & \|\tilde{X}_1 - \bar{X}_1\|_F, \\ \text{subject to:} & \mathcal{R}(\bar{X}_1) = \mathcal{X}_1. \end{array}$$

*Proof.* Let  $\bar{X}_1 = X_1 R$ . Then the above minimization problem becomes

$$\begin{array}{ll} \text{minimize:} & \|\tilde{X}_1 - X_1 R\|_F, \\ \text{subject to:} & R \text{ is nonsingular.} \end{array}$$

If we drop the restriction that  $R$  be nonsingular, then from the theory of least squares, the minimizing value of  $R$  is given by

$$R = (X_1^H X_1)^{-1} X_1^H \tilde{X}_1.$$

But from (4.8)  $X_1^H \tilde{X}_1 = X_1^H X_1 = I$ . Hence  $R = I$  (which is nonsingular) and  $X_1$  solves (4.9).  $\square$

**4.1.4. The perturbation bound.** The second problem is to determine a perturbation equation and solve it. The key is to observe that equation (4.3) is not only a consequence of  $\mathcal{R}(X_1)$  being an invariant subspace, it is actually a sufficient condition for  $\mathcal{R}(X_1)$  to be an invariant subspace. For if  $Y_2^H (AX_1) = 0$ , then  $\mathcal{R}(AX_1)$  lies in the orthogonal complement of  $\mathcal{R}(Y_2)$ ; that is, it lies in  $\mathcal{R}(X_1)$ .

Thus our problem reduces to finding a basis for the orthogonal complement of  $\tilde{\mathcal{X}}$ . But if  $\tilde{X}_1$  is represented in the form (4.8), it is easily verified that the columns of

$$Y_2 - X_1 P^H$$

form a basis for the required space. Thus a necessary and sufficient condition for  $\mathcal{R}(\tilde{X}_1)$  to span an invariant subspace of  $\tilde{A}$  is that

$$(4.10) \quad (Y_2^H - X_1 P)(A + E)(X_1 + Y_2 P) = 0.$$

Equation (4.10) can be simplified by the introduction of the operator  $\mathbf{T}$  defined by (4.5). Specifically, let

$$(X_1 \ Y_2)^H E (X_1 \ Y_2) = \begin{pmatrix} E_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix}.$$

Then (4.10) is equivalent to

$$(4.11) \quad \mathbf{T}P + (E_{22}P - PE_{11}) = E_{21} - P(A_{12} + E_{12})P,$$

where  $A_{12} = X_1^H A Y_2$  as above.

If we drop second-order terms from (4.11), we get our first-order perturbation equation

$$(4.12) \quad \mathbf{T}\check{P} = E_{21}.$$

Since by hypothesis  $\mathcal{X}_1$  is simple,  $\mathbf{T}$  is nonsingular, and

$$(4.13) \quad \check{P} = \mathbf{T}^{-1}E_{21}.$$

This gives the perturbation bound

$$(4.14) \quad \|\check{P}\|_S \leq \frac{\|E_{21}\|_S}{\delta} = \frac{\|S_r X_1\|_F \|S_c Y_2\|_F}{\delta},$$

where  $\delta$  is defined by (4.6). In the simple case (4.14) becomes

$$(4.15) \quad \|\check{P}\|_S \leq \frac{\sqrt{k(n-k)}\sigma}{\delta}.$$

**4.1.5. Interpretation of the bound.** The bounds (4.14) and (4.15) consist of two parts: a norm of  $E_{21}$  and a divisor  $\delta$ . From (4.7) it follows that if the eigenvalues of  $A_{11}$  and  $A_{22}$  are not well separated,  $\delta$  will be small. In other words, if the eigenvalues of  $A$  corresponding to the invariant subspace  $\mathcal{X}_1$  are near the eigenvalues corresponding to the complementary invariant subspace, the bounds (4.14) and (4.15) will be large, and we can expect  $\mathcal{X}_1$  to be sensitive to perturbations in  $A$ .

The norm of  $E_{21}$  has an interesting interpretation, which is given in the following theorem, whose proof is left as an exercise.

**THEOREM 4.5.** *The Frobenius norm of  $E_{21}$  satisfies*

$$\|E_{21}\|_F = \min_{B \in \mathbb{C}^{k \times k}} \|\tilde{A}X_1 - X_1B\|_F,$$

and the minimum is attained when  $B = X_1^H \tilde{A}X_1$ .

If  $\mathcal{R}(X_1)$  were an invariant subspace of  $\tilde{A}$ , we could make the residual  $\tilde{A}X_1 - X_1B$  zero by choosing  $B$  to be equal to the representation of  $\tilde{A}$  on  $\mathcal{R}(X_1)$ ; i.e.,  $B = X_1^H \tilde{A}X_1$ . Even when  $\mathcal{R}(X_1)$  is not an invariant subspace of  $\tilde{A}$ , this choice of  $B$  minimizes the norm of the residual, whose value at the minimum is  $\|E_{21}\|_F$ .

**4.1.6. Range of applicability.** The foregoing development presupposes that there is an invariant subspace of  $\tilde{A}$  that approaches  $\mathcal{X}_1$  as  $E$  approaches zero. The following theorem shows that this is true by showing that if  $\mathcal{X}_1$  is simple and  $E$  is small enough, the perturbation equation (4.11) has a solution  $P$  that approaches zero as  $E$  approaches zero.

**THEOREM 4.6.** *If*

$$(4.16) \quad \|E_{11}\| + \|E_{22}\| < \delta$$

and

$$(4.17) \quad \frac{\|E_{21}\|_F \|A_{12} + E_{12}\|_F}{(\delta - \|E_{11}\| - \|E_{22}\|)^2} < \frac{1}{4},$$

then (4.11) has a unique solution  $P$  that satisfies

$$(4.18) \quad \|P\|_F \leq \frac{2\|E_{21}\|_F}{\delta - \|E_{11}\| - \|E_{22}\|}.$$

The condition (4.16) essentially says that the perturbations  $E_{11}$  and  $E_{22}$  do not make the operator  $\mathbf{T}$  singular [cf. (3.3)]. If we replace the norms by stochastic norms, we get

$$(4.19) \quad \|S_r X_1\|_F \|S_c X_1\|_F + \|S_r Y_2\|_F \|S_c Y_2\|_F < \delta,$$

or in the simple case

$$(4.20) \quad n\sigma < \delta.$$

If (4.19) or (4.20) show that  $E_{11}$  and  $E_{22}$  have a negligible effect in the denominator of (4.17), then we ignore it. In this case, the second condition written in terms of the stochastic norm becomes

$$(4.21) \quad \frac{\|S_c Y_2\|_F \|S_r X_1\|_F \sqrt{\|A_{12}\|_F^2 + \|S_c X_1\|_F^2 \|S_r Y_2\|_F^2}}{\delta^2} < \frac{1}{4}.$$

In the simple case this becomes

$$\frac{\sigma \sqrt{k(n-k)} \sqrt{\|A_{12}\|_F^2 + k(n-1)\sigma^2}}{\delta^2} < \frac{1}{4}.$$

If (4.20) is satisfied, this latter condition can be replaced by the simpler, but stronger condition

$$\frac{\sqrt{\|A_{12}\|_F^2 + k(n-1)\sigma^2}}{\delta} < \frac{1}{4}.$$

**4.1.7. The representation of  $\tilde{A}$  on  $\mathcal{R}(\check{X}_1)$ .** Just as  $A_{11} = X_1^H A X_1$  is the representation of  $A$  on  $\mathcal{X}_1$ , so is  $\check{X}_1^H (A + E) \check{X}_1$  the representation of  $\tilde{A}$  on  $\check{\mathcal{X}}_1$ . This matrix is readily found to be

$$\tilde{A}_{11} = A_{11} + E_{11} + A_{12} \tilde{P},$$

and its first-order approximation is

$$(4.22) \quad \check{\tilde{A}}_{11} = A_{11} + E_{11} + A_{12} \check{P},$$

where  $\check{P}$  is defined by (4.13).

From (4.22) and (4.14) we have

$$(4.23) \quad \|\check{\tilde{A}}_{11} - A_{11}\|_S \leq \|S_r X_1\|_F \|S_c X_1\|_F + \frac{\|A_{12}\| \|S_r X_1\|_F \|S_c Y_2\|_F}{\delta},$$

or in the simple case

$$\|\check{\tilde{A}}_{11} - A_{11}\|_S \leq \left( k + \frac{\sqrt{k(n-k)} \|A_{12}\|}{\delta} \right) \sigma.$$

When  $k = 1$  these inequalities provide a stochastic bound on the perturbation of an eigenvalue of  $A$ . In the general case, we need more information about  $A_{11}$  (e.g.,  $A_{11}$  is normal) to say anything about the relation of the eigenvalues of  $A_{11}$  and  $\check{\tilde{A}}_{11}$ .

**4.1.8. Other representations of  $\tilde{X}_1$ .** The representation of  $\tilde{X}_1$  in the form (4.8) amounts to imposing the normalization  $X_1^H \tilde{X}_1 = I$ . But this is not the only desirable normalization. For example, if  $A^T$  is stochastic and  $x_1$  its Peron vector, then one might wish to impose the normalization  $\mathbf{1}^T x_1 = 1$ , where  $\mathbf{1}$  is the vector of all ones. This allows the components of  $x_1$  to be regarded as probabilities, since they are nonnegative and sum to 1.

We will consider normalizations of the form

$$(4.24) \quad W_1^H \tilde{X}_1 = I,$$

where  $W_1$  satisfies

$$(4.25) \quad W_1^H X_1 = I.$$

Any matrix satisfying (4.25) can be written in the form

$$W_1 = X_1 + Y_2 Q^H$$

(to see this write  $W_1 = X_1 R + Y_2 Q^H$ , and multiply first by  $X_1^H$  and then by  $Y_2^H$ ).

The normalization (4.24) can be regarded as a block transformation of the coordinate system defined by  $(X_1 \ Y_2)$ . Let

$$(X_1 \ W_2) = (X_1 \ Y_2) \begin{pmatrix} I & Q \\ 0 & I \end{pmatrix} = (X_1 \ Y_2 + X_1 Q).$$

Then

$$(4.26) \quad (X_1 \ W_2)^{-1} = \begin{pmatrix} I & -Q \\ 0 & I \end{pmatrix} \begin{pmatrix} X_1^H \\ Y_2^H \end{pmatrix} = \begin{pmatrix} W_1^H \\ Y_2^H \end{pmatrix}.$$

Moreover,

$$\begin{pmatrix} W_1^H \\ Y_2^H \end{pmatrix} A (X_1 \ W_2) = \begin{pmatrix} A_{11} & A_{12}^{(W)} \\ 0 & A_{22} \end{pmatrix},$$

where

$$(4.27) \quad A_{12}^{(W)} = W_1 A W_2 = A_{12} + A_{11} Q - Q A_{22}.$$

Thus  $(X_1 \ W_2)$  defines a (nonunitary) similarity that reduces  $A$  to block triangular form [cf. (4.4)].

We can now repeat our entire development. The normalization (4.24) is equivalent to seeking  $\tilde{X}_1$  in the form

$$\tilde{X}_1^{(W)} = X_1 + W_2 P^{(W)}$$

[cf. (4.8)]. From (4.26) it follows that the columns of  $Y_2 - W_1 P^H$  form a basis for the orthogonal complement of  $\mathcal{R}(\tilde{X}_1)$ . Hence the perturbation equation is

$$(Y_2 - W_1 P^H)(A + E)(X_1 + W_2 P^{(W)}) = 0$$

[cf. (4.10)], from which we find that

$$(4.28) \quad \mathbf{T} \tilde{P}^{(W)} = E_{21}.$$



In this case the representation of  $A$  with respect to  $\tilde{X}_1^{(W)}$  up to first-order terms is given by

$$(4.29) \quad \check{A}_{11}^{(W)} = A_{11} + E_{11}^{(W)} + A_{12}^{(W)} \check{P}^{(W)},$$

where  $E_{11}^{(W)} = W_1^H E X_1$  and  $E_{12}^{(W)} = W_1^H E W_2$  [cf. (4.22)].

Comparing (4.12) and (4.28) we obtain the following remarkable theorem.

**THEOREM 4.7.**

$$\check{P}^{(W)} = \check{P}.$$

In other words, the same  $\check{P}$  serves for all normalizations. Unfortunately, the ranges of applicability may be different, since  $P^{(W)}$  will in general be different from  $P$ . It is therefore worthwhile to give a different proof of Theorem 4.7, which shows explicitly how they differ.

*Second proof of Theorem 4.7.* Since the columns of  $(X_1 \ W_2)$  are linearly independent, we may write  $\tilde{X}_1 = X_1 U + W_2 V$  or

$$X_1 + Y_2 P = X_1 U + W_2 V.$$

Multiplying this expression by  $Y_2^H$  and using the fact that  $Y_2^H W_2 = I$ , we get  $V = P$ . Multiplying by  $X_1^H$  and using the fact that  $X_1^H W_2 = Q$ , we get  $U = I - QP$ . As  $E \rightarrow 0$ , the matrix  $U$  approaches  $I$ . Hence for  $E$  sufficiently small  $U$  is nonsingular, and we may write

$$\tilde{\mathcal{X}}_1 = \mathcal{R}[X_1 + W_2 P(I - QP)^{-1}].$$

Thus

$$P^{(W)} = P(I - QP)^{-1}.$$

The theorem follows on replacing  $P$  by  $\check{P}$  and discarding higher-order terms.  $\square$

This second proof shows that we can trust the expansion  $\check{X}_1^{(W)} = X_1 + W_2 \check{P}$  only when  $\check{P}$  is accurate and  $Q\check{P}$  is small. We will use this fact in discussing generalized Rayleigh quotients, to which we now turn.

**4.1.9. Generalized Rayleigh quotients.** One of the unsatisfactory aspects of the above development is the presence of the term  $A_{12}^{(W)} \check{P}^{(W)}$  in the expression (4.29) for  $\check{A}_{11}^{(W)}$ . Since  $W_1$ , or equivalently  $Q$ , is a free parameter, we can choose it to make  $A_{12}^{(W)} = 0$ . From (4.27) we see that this is equivalent to choosing  $Q$  to satisfy

$$(4.30) \quad \mathbf{T}^* Q = A_{12},$$

where

$$\mathbf{T}^* = Q \mapsto Q A_{22} - A_{11} Q.$$

As the notation suggests,  $\mathbf{T}^*$  is the adjoint of the operator  $\mathbf{T}$ . Hence, by Theorem 4.3 it is nonsingular and the norm of its inverse is  $\delta^{-1}$ .

We will denote by  $Y_1$  and  $X_2$  the values of  $W_1$  and  $W_2$  corresponding to the solution of (4.30). It follows that

$$\begin{pmatrix} Y_1^H \\ Y_2^H \end{pmatrix} A(X_1 \ X_2) = \begin{pmatrix} A_{11} & 0 \\ 0 & A_{22} \end{pmatrix};$$

that is, the similarity transformation reduces  $A$  to block diagonal form. In particular,  $\mathcal{R}(X_2)$  is an invariant subspace of  $A$ . Since it is complementary to  $\mathcal{X}_1$  is called the *complementary invariant subspace*. In the same way  $\mathcal{R}(Y_1)$  is a left invariant subspace complementary to  $\mathcal{R}(Y_2)$ .

We saw at the end of §4.1.8 that in addition to the conditions of Theorem 4.6 we require that  $QP$  be small in order to use the approximation  $\tilde{P}$ . Fortunately, this will be so to the extent that the conditions of Theorem 4.6 are satisfied. Specifically, from (4.30) we have  $\|Q\|_F \leq \|A_{22}\|_F/\delta$ . Hence by (4.17) and (4.18) we have  $\|Q\|\|P\| < \frac{1}{2}$ .

Since  $A_{12}^{(Y)} = 0$ , the first-order approximation of the representation of  $\tilde{A}$  on  $\tilde{\mathcal{X}}_1$  has the particularly simple form

$$\check{A}_{11}^{(Y)} = A_{11} + Y_1^H E X_1.$$

In particular

$$(4.31) \quad \|\check{A}_{11}^{(Y)} - A_{11}\|_S = \|S_r X_1\|_F \|S_c Y_1\|_F,$$

or in the simple case

$$\|\check{A}_{11}^{(Y)} - A_{11}\|_S = \sqrt{k}\sigma \|Y_1\|_F,$$

The matrix  $Y_1^H A X_1$  is sometimes called the *generalized Rayleigh quotient* of  $A$ . Moreover the matrix  $X_1 Y_1^H$  is called the *spectral projector*, since it projects a vector onto  $\mathcal{X}_1$  along the complementary invariant subspace  $\mathcal{R}(X_2)$ . With the normalization  $X_1^H X_1 = I$ , the Frobenius norm of the spectral projector is  $\|Y_1\|_F$ . Hence the above result says that the sensitivity of the generalized Rayleigh quotient is determined by the size of the spectral projector.

**4.2. Hermitian matrices.** The principle simplification that occurs when  $A$  is Hermitian is that  $A_{12} = 0$ , or equivalently  $Y_1 = X_1$ . This does not change the bound on  $\tilde{P}$ , but it does affect the range of applicability. Specifically, (4.21) becomes

$$\frac{\|S_c Y_2\|_F \|S_r X_1\|_F \|S_c X_1\|_F \|S_r Y_2\|_F}{\delta^2} < \frac{1}{4},$$

or in the simple case

$$\frac{\sqrt{k(n-1)}}{\delta} \sigma < \frac{1}{2}.$$

The interpretation of the bound becomes simpler, for when  $A_{11}$  and  $A_{22}$  are Hermitian

$$\delta = \min\{|\lambda_1 - \lambda_2| : \lambda_1 \in \lambda(A_{11}) \lambda_2 \in \lambda(A_{22})\}.$$

Thus for Hermitian matrices the separation of the eigenvalues of an invariant subspace from those of its complement determines the sensitivity of the subspace.

Finally, since  $X_1 = Y_1$ , the bound (4.31) become

$$(4.32) \quad \|\check{A}_{11}^{(Y)} - A_{11}\|_S = \|S_r X_1\|_F \|S_c X_1\|_F,$$

or in the simple case

$$(4.33) \quad \|\check{A}_{11}^{(Y)} - A_{11}\|_S = k\sigma.$$

Since  $A_{11}$  and  $\check{A}_{11}$  are both Hermitian, we can make a strong statement about the relation of eigenvalues of the two. Specifically, if  $\lambda_1(A_{11}) \leq \lambda_2(A_{11}) \leq \cdots \leq \lambda_k(A_{11})$  are the eigenvalues of  $A_{11}$  and  $\lambda_1(\check{A}_{11}) \leq \lambda_2(\check{A}_{11}) \leq \cdots \leq \lambda_k(\check{A}_{11})$  are the eigenvalues of  $\check{A}_{11}$ , then by the Hoffman–Wielandt theorem

$$\sum_{i=1}^k [\lambda_i(\check{A}_{11}) - \lambda_i(A_{11})]^2 \leq \|\check{A}_{11} - A_{11}\|_F^2,$$

and the right-hand side can be estimated by (4.32) or (4.33).

**4.3. Eigenvalues and eigenvectors.** When  $k = 1$ , the matrix  $X_1$  becomes an eigenvector of  $A$  corresponding to the eigenvalue  $A_{11}$ . To remind us of this fact, we will write  $X = (x_1 \ Y_2)$  and

$$X^H A X = \begin{pmatrix} a_{11} & a_{12}^H \\ 0 & A_{22} \end{pmatrix},$$

with a corresponding notation for  $X^H E X$ .

The principle simplification that occurs with eigenvalues and eigenvectors is that the operator  $\mathbf{T}$  now becomes a matrix

$$\mathbf{T} = a_{11}I - A_{22}.$$

This means that we can involve  $\mathbf{T}$  directly in our calculation of the stochastic norms. Specifically, we now have for the perturbation of  $x_1$ ,

$$\check{x}_1 = x_1 + Y_2 \check{p},$$

where

$$\check{p} = \mathbf{T}^{-1} Y_2^H E x_1.$$

Hence

$$\|\check{p}\|_S = \|\mathbf{T}^{-1} Y_2^H S_c\|_F \|S_c x_1\|.$$

In the simple case this becomes

$$\|\check{p}\|_S = \sigma \|\mathbf{T}^{-1}\|_F.$$

If we define

$$\delta_F = \|\mathbf{T}^{-1}\|_F^{-1},$$

then the estimate becomes

$$\|\check{p}\|_S = \frac{\sigma}{\delta_F}.$$

When this estimate is compared with (4.15) for  $k = 1$ , it is seen that, in addition to being an equality, it lacks the factor  $\sqrt{n-1}$ . However, this is compensated for by the fact that  $\delta_F$  will in general be smaller than  $\delta$ .

Turning to the Rayleigh quotient, we have

$$\check{a}_{11} = a_{11} + y_1^H E x_1.$$

Hence

$$\|\check{a}_{11} - a_{11}\|_S = \|S_c y_1\| \|S_r x_1\|,$$

or in the simple case

$$\|\check{a}_{11} - a_{11}\|_S = \sigma \|y_1\|.$$

With the normalization  $y_1^H x_1 = 1$ ,  $\|y_1\|$  is the secant of the angle between  $x_1$  and  $y_1$ . Thus as  $x_1$  and  $y_1$  become progressively more orthogonal, their eigenvalue becomes more sensitive to perturbations.

**4.4. Singular values and singular vectors.** Let  $X \in \mathcal{R}^{m \times n}$  with  $m \geq n$ . Then there are orthogonal matrices  $U$  and  $V$  of order  $m$  and  $n$  such that

$$(4.34) \quad U^T X V = \begin{pmatrix} \Psi \\ 0 \end{pmatrix},$$

where

$$\Psi = \text{diag}(\psi_1, \psi_2, \dots, \psi_n).$$

The decomposition (4.34) is called the *singular value decomposition* of  $X$ . The number  $\psi_i$  is called a singular value of  $X$  and the corresponding columns  $u_i$  and  $v_i$  of  $U$  and  $V$  are its left and right singular vectors. They are related by the formulas

$$X v_i = \psi_i u_i, \quad X^T u_i = \psi_i v_i.$$

Although the singular value decomposition is defined for complex matrices, the singular values are not differentiable functions of the elements of the matrix. This is true even for the scalar  $\xi$ , whose “singular value” is  $\psi = |\xi|$ , since in the complex plane the absolute value is not an analytic function of its argument. The implication is that if we wish to develop first-order expansions, we must restrict ourselves to real perturbations of real matrices. Even here, we must restrict ourselves to nonzero singular values, since the absolute value, regarded as a function of a real variable, is not differentiable at zero.

Therefore, we will consider the perturbation of a nonzero singular value  $\psi_1$  of a real matrix  $X$  and its corresponding right singular vector  $v_1$  under a real perturbation  $E$ . A perturbation expansion may be obtained by observing that  $\tilde{\psi}_1^2$  is an eigenvalue of  $\tilde{A} = \tilde{X}^T \tilde{X}$  with eigenvector  $\tilde{v}_1$ . Specifically, we have the following theorem, whose proof may be found in the references.

**THEOREM 4.8.** *Let*

$$U = (u_1 \ U_2) \quad V = (v_1 \ V_2)$$

and

$$\Psi = \text{diag}(\psi_1, \Psi_2).$$

Then

$$(4.35) \quad \check{\psi}_1 = \psi_1 + u^T E v.$$

Moreover, if

$$\mathbf{T} = \psi_1^2 I - \Psi_2^2,$$

then

$$\check{v}_1 = v_1 + V_2 \check{p},$$

where

$$\check{p} = \mathbf{T}^{-1}(U_2^T E v_1 + V_2^T E^T u_1).$$

From (4.35) we immediately get the following perturbation estimate for  $\|\check{\psi}_1 - \psi_1\|_S$ :

$$\|\check{\psi}_1 - \psi_1\|_S = \|S_c u_1\| \|S_r v_1\|.$$

In the simple case, this becomes

$$\|\check{\psi}_1 - \psi_1\|_S = \sigma.$$

The estimate for  $\|\check{p} - p\|_S$  is more complicated:

$$\begin{aligned} \|\check{p} - p\|_S^2 &= \|S_c U_2 \Psi_2 \mathbf{T}^{-1}\|_F^2 \|S_r v_1\| + 2\psi_1 u_1^T S_c^2 U_2 \Psi_2 \mathbf{T}^{-2} V_2^T S_r^2 v_2 \\ &\quad + \psi_1^2 \|S_r V_2 \mathbf{T}^{-1}\|_F^2 \|S_c u_1\|^2 \\ &\leq \|S_c U_2 \Psi_2 \mathbf{T}^{-1}\|_F^2 \|S_r v_1\| + \psi_1^2 \|S_r V_2 \mathbf{T}^{-1}\|_F^2 \|S_c u_1\|^2. \end{aligned}$$

However, in the simple case the cross-product term in the equality vanishes and we get

$$\begin{aligned} \|\check{p} - p\|_S &= \sigma \|(\psi_1^2 I - \Psi^2) \mathbf{T}^{-1}\|_F \\ &= \sigma \sqrt{\sum_{i=2}^n \frac{\psi_1^2 + \psi_i^2}{\psi_1^2 - \psi_i^2}}. \end{aligned}$$

Thus the distance of a singular value from its neighbors controls the sensitivity of its singular vector to perturbations.

**Notes and references.** The general approach to invariant subspaces taken here is due to the author [33],[34]. For another view of the subject the reader is referred to Kato's work [25], which also treats the perturbation of operators in infinite-dimensional settings.

§4.1.2. The term “simple” referring to an invariant subspace does not seem to have appeared in the literature before. Theorem 4.3 is important in many of areas, and it has a number of proofs. See [17, §7.6] for a constructive approach. The quantity  $\delta$  was introduced in [33] as the function  $\text{sep}(A_{11}, A_{22})$ . This name comes from (4.7), which shows that  $\text{sep}(A_{11}, A_{22})$  is a lower bound on the separation of the spectra of  $A_{11}$  and  $A_{22}$ . Unfortunately, it can be very much smaller than the actual separation.

§4.1.3. Theorem 4.4 is not the only justification of this choice of representation. It can be shown that the singular values of  $P$  are the tangents of the canonical

angles between the subspaces  $\mathcal{X}$  and  $\tilde{\mathcal{X}}$  (see [8],[35] for definitions). Consequently,  $\|P\| = \|\tilde{X} - X\|$  is a bound on the separation of the two subspaces.

§4.1.5. Theorem 4.5 may be found in [8]. For generalization, see [24].

§4.1.6. A proof of Theorem 4.6 may be found in [34].

§4.1.8. Equation (4.24) does not represent all possible normalizations. In [27], which treats only eigenvectors, the normalizing function is allowed to be any differentiable function. Theorem 4.7 appears to be new.

§4.1.9. There is some ambiguity in the term “generalized Rayleigh quotient.” If  $\lambda$  is an eigenvalue of a Hermitian matrix then the Rayleigh quotient  $\rho(x, A) = x^H A x / x^H x$  has two properties.

- (1)  $\lambda = \rho(x, A)$ ,
- (2)  $\tilde{\lambda} = \rho(x, A + E)$ .

If for non-Hermitian matrices we require only the first property, then any quotient of the form  $w^H A x / w^H x$  generalizes the Rayleigh quotient. However, if we require both properties, then we must take  $w = y$ , the left eigenvector corresponding to  $\lambda$ .

§4.2. A completely different approach to perturbation theory for Hermitian matrices is given by Davis and Kahan [8]. The Hoffman–Wielandt theorem appears in [20], and Wilkinson [46] gives an elementary proof.

§4.4. A proof of Theorem 4.8 may be found in [39]. It is possible to develop perturbation bound for spaces of singular vectors corresponding to clusters of singular values [34]; however, the bounds are not pretty.

**5. Conclusions.** In this paper we have shown that many problems in matrix perturbation theory can be rigorously treated from a probabilistic point of view. The main advantage of this approach is that in many cases it gives estimates that are exact equalities: nothing is given away in their derivations. For example, compare the eigenvalue estimate

$$\|\tilde{a}_{11} - a_{11}\|_S = \sigma \|y_1\|$$

with the more usual bound

$$|\tilde{a}_{11} - a_{11}| \leq \|y_1\| \|E\|.$$

Not only is the first simpler, but it makes it clear that the second can be a considerable overestimate, since  $\|E\|$  will be larger than the size of a typical element. Thus stochastic perturbation theory can be used to see how well we have done with more conventional bounds.

The use of cross-correlated errors reduces the applicability of the technique, but not unduly considering the gain in simplicity. Moreover, the scale of the error appears explicitly in the bounds so they may be adjusted to the application.

The chief disadvantage of the approach is its reliance on first-order approximations. Although Theorem 2.8 provides an asymptotic justification for this, in practice we must assess when the first-order approximations are valid. In this paper we have proceeded informally, by recasting in terms of stochastic norms conditions that are necessary for the approximations to be accurate. This insures that for errors sufficiently small the probability of violating the conditions is small, and we can even use the Chebyshev inequality to bound the probability.

The theory is based on the Frobenius norm, whereas the spectral norm is more frequently used in usual approach to perturbation theory. For estimating the perturbations scalars and vectors this makes no difference, since the two norms coincide

for the estimated quantities and the exactness of the estimates assures us that any Frobenius norms in them really have to be there.

Nonetheless, one might wonder if there is a stochastic analogue of the spectral norm. Unfortunately, the natural definition

$$\max_{\|x\|=1} \|\mathbf{E}(Ex)\|$$

does not work, since it gives different results for  $E$  and  $E^T$ . If we try to restore symmetry with

$$\max_{\substack{\|x\|=1 \\ \|y\|=1}} \mathbf{E}(y^T Ex)$$

we get something that is too small. This problem can stand further investigation.

## REFERENCES

- [1] A. E. BEATON, D. B. RUBIN, AND J. L. BARONE, *The acceptability of regression solutions: Another look at computational accuracy*, J. Amer. Statist. Assoc., 71 (1976), pp. 158–168.
- [2] G. F. BROWN, J. B. KADANE, AND J. G. RAMAGE, *The asymptotic bias and mean-squared error of double K-class estimators when the disturbances are small*, Interna. Econom. Rev., 15 (1974), pp. 667–679.
- [3] F. CHATELIN, *Analyse statistique de la qualité numérique et arithmétique de la résolution approchée d'équations par calcul sur ordinateur*, Etude F.133, Centre Scientifique de Paris, 1988.
- [4] ———, *De l'utilisation en calcul matriciel de modèles probabilistes pour la simulation des erreurs de calcul*, Comptes Rendus de l'Académie des Sciences, Paris, Série I, 307 (1988), pp. 847–850.
- [5] ———, *A probabilistic round-off error propagation model. Application to the eigenvalue problem*, in Reliable Numerical Software, D. Cox and S. Hammarling, eds., Oxford University Press, Oxford, 1990, to appear.
- [6] N. DAVID AND G. W. STEWART, *Hypothesis testing with errors in the variables*, Tech. report TR-1735, Dept. of Computer Science, University of Maryland, College Park, MD, 1988.
- [7] R. B. DAVIES AND B. HUTTON, *The effects of errors in the independent variables in linear regression*, Biometrika, 62 (1975), pp. 383–391.
- [8] C. DAVIS AND W. M. KAHAN, *The rotation of eigenvectors by a perturbation. III*, SIAM J. Numer. Anal., 7 (1970), pp. 1–46.
- [9] J. W. DEMMEL, *The probability that a numerical analysis problem is difficult*, Math. Comp., 50 (1988), pp. 449–480.
- [10] L. ELDÉN, *Perturbation theory for the least squares problem with linear equality constraints*, SIAM J. Numer. Anal., 17 (1980), pp. 338–350.
- [11] W. A. FULLER, *Measurement Error Models*, John Wiley, New York, 1987.
- [12] C. F. GAUSS, *Theoria Motus Corporum Coelestium in Sectionibus Conicis Solem Ambientium*, Perthes and Besser, Hamburg, 1809. [In Latin.]
- [13] ———, *Theory of the Motion of the Heavenly Bodies Moving about the Sun in Conic Sections*, C. H. Davis, Trans. Dover, New York (1963), 1809. [In English.]
- [14] ———, *Theoria combinationum observationum erroribus minimis obnoxiae, pars prior*, in Werke, IV, Königlichen Gessellschaft der Wissenschaften zu Göttinging (1880), 1821, pp. 1–26. [In Latin.]
- [15] A. GEMAN, *A limit theorem for the norm of random matrices*, Ann. Probab., 8 (1980), pp. 252–261.
- [16] G. H. GOLUB AND C. F. VAN LOAN, *An analysis of the total least squares problem*, SIAM J. Numer. Anal., 17 (1980), pp. 883–893.
- [17] ———, *Matrix Computations*, Johns Hopkins University Press, Baltimore, MD, 1983.
- [18] G. H. GOLUB AND J. H. WILKINSON, *Note on the iterative refinement of least squares solution*, Numer. Math., 9 (1966), pp. 139–148.
- [19] S. D. HODGES AND P. G. MOORE, *Data uncertainties and least squares regression*, Appl. Statist., 21 (1972), pp. 185–195.

- [20] A. J. HOFFMAN AND H. W. WIELANDT, *The variation of the spectrum of a normal matrix*, Duke Math. J., 20 (1953), pp. 37–39.
- [21] R. V. HOGG AND A. T. CRAIG, *Introduction to Mathematical Statistics*, Fourth edition, Macmillan, New York, 1978.
- [22] H. HOTELLING, *The selection of variates for use in prediction with some comments on the general problem of nuisance parameters*, Ann. Math. Statist., 11 (1940), pp. 271–283.
- [23] J. JOHNSTON, *Econometric Methods*, Second edition, McGraw-Hill, New York, 1972.
- [24] W. KAHAN, B. N. PARLETT, AND E. JIANG, *Residual bounds on approximate eigensystems of nonnormal matrices*, SIAM J. Numer. Anal., 19 (1982), pp. 470–484.
- [25] T. KATO, *Perturbation Theory for Linear Operators*, Springer-Verlag, Berlin, New York, 1966.
- [26] D. G. KENDALL, *Stochastic processes occurring in the theory of queues and their analysis by the method of the imbedded markov chain*, Ann. Math. Statist., 24 (1953), pp. 338–354.
- [27] C. MEYER AND G. W. STEWART, *Derivatives and perturbations of eigenvectors*, SIAM J. Numer. Anal., 25 (1988), pp. 679–691.
- [28] H. NEUDECKER AND T. WANSBEEK, *Fourth-order properties of normally distributed random matrices*, Linear Algebra Appl., 97 (1987), pp. 13–22.
- [29] C. C. PAIGE, *Computer solution and perturbation analysis of generalized linear least squares problems*, Math. Comp., 33 (1979), pp. 171–184.
- [30] E. PARZEN, *Modern Probability Theory and Its Applications*, John Wiley, New York, 1960.
- [31] G. A. F. SEBER, *Linear Regression Anal.*, John Wiley, New York, 1977.
- [32] R. J. SERFLING, *Approximation Theorems of Mathematical Statistics*, John Wiley, New York, 1980.
- [33] G. W. STEWART, *Error bounds for approximate invariant subspaces of closed linear operators*, SIAM J. Numer. Anal., 8 (1971), pp. 796–808.
- [34] ———, *Error and perturbation bounds for subspaces associated with certain eigenvalue problems*, SIAM Rev., 15 (1973), pp. 727–764.
- [35] ———, *On the perturbation of pseudo-inverses, projections, and linear least squares problems*, SIAM Rev., 19 (1977), pp. 634–662.
- [36] ———, *Sensitivity coefficients for the effects of errors in the independent variables in a linear regression*, Tech. report TR-571, Dept. of Computer Science, University of Maryland, College Park, MD, 1977.
- [37] ———, *A nonlinear version of Gauss's minimum variance theorem with applications to an errors-in-the-variables model*, Tech. report TR-1263, Dept. of Computer Science, University of Maryland, College Park, MD, 1983.
- [38] ———, *On the invariance of perturbed null vectors under column scaling*, Numer. Math., 44 (1984), pp. 61–65.
- [39] ———, *A second order perturbation expansion for small singular values*, Linear Algebra Appl., 56 (1984), pp. 231–235.
- [40] ———, *Collinearity and least squares regression*, Statist. Sci., 2 (1987), pp. 68–100.
- [41] ———, *Perturbation theory and least squares with errors in the variables*, Tech. report UMIACS-TR-89-97, CS-TR 2326, Dept. of Computer Science, University of Maryland, College Park, MD, 1989. To appear in the Proceedings of the AMS Conference on Measurement Error Models, Humboldt, California.
- [42] G. W. STEWART AND G.-J. SUN, *Matrix Perturbation Theory*, Academic Press, Boston, 1990.
- [43] S. VAN HUFFEL, *Analysis of the Total Least Squares Problem and Its Use in Parameter Estimation*, Ph.D. thesis, Katholieke Universiteit Leuven, Leuven, Belgium, 1987.
- [44] P.-A. WEDIN, *Perturbation theory for pseudo-inverses*, BIT, 13 (1973), pp. 217–232.
- [45] N. WEIS, G. W. WASILKOWSKI, H. WOŹNIAKOWSKI, AND M. SHUB, *Average condition number for solving linear equations*, Linear Algebra Appl., 83 (1986), pp. 79–102.
- [46] J. H. WILKINSON, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.