

# TOWARDS A BACKWARD PERTURBATION ANALYSIS FOR DATA LEAST SQUARES PROBLEMS

X.-W. CHANG\*, G. H. GOLUB†, AND C.C. PAIGE‡

*Xiao-Wen Chang and Chris Paige dedicate this  
in memory of their warm, generous and inspirational friend Gene Golub.*

**Abstract.** Given an approximate solution to a data least squares (DLS) problem, we would like to know its minimal backward error. Here we derive formulas for what we call an “extended” minimal backward error, which is at worst a lower bound on the minimal backward error. When the given approximate solution is a good enough approximation to the exact solution of the DLS problem (which is the aim in practice), the extended minimal backward error is the actual minimal backward error, and this is also true in other easily assessed and common cases. Since it is computationally expensive to compute the extended minimal backward error directly, we derive a lower bound on it and an asymptotic estimate for it, both of which can be evaluated less expensively. Simulation results show that for reasonable approximate solutions the lower bound has the same order as the extended minimal backward error, and the asymptotic estimate is an excellent approximation to the extended minimal backward error.

**Key words.** data least squares, backward errors, numerical stability, perturbation analysis, asymptotic estimate, iterative methods, stopping criteria.

**AMS subject classifications.** 15A06, 65F20, 65G50.

**1. Introduction.** Given an approximate solution to a problem, the aim of backward perturbation analysis is to find a minimum size perturbation in the data such that the approximate solution is an exact solution of the perturbed problem. In the analysis one tries to find a formula for, or good bounds on, the size of the minimal perturbation (to be referred to as the minimal backward error) and design an efficient algorithm to evaluate or estimate the formula or the bounds. If the relative minimal backward error (i.e., the size of the minimal perturbation divided by an acceptable measure of the size of the data) is of the order of the unit round-off then we say that the approximate solution is a (normwise) backward stable solution. Backward perturbation analyses are useful in practice. Sometimes we may not know if an algorithm for solving a problem is numerically stable, e.g., the backward numerical stability of some fast algorithms for structured matrix problems is unknown. But if we know that a computed solution of a specific problem is a backward stable solution, we are satisfied with this computed solution. Also when we solve a large scale problem by an iterative algorithm, the results of a backward perturbation analysis can often be used to design effective stopping criteria, see, for example, [1], [20] and [25].

There has been a lot of work on the backward perturbation analysis of linear systems, especially in recent years. For example, for consistent linear systems, see [14], [25], [30], [31], [32], [34], [37]; for unconstrained least squares problems, see [9], [12], [17]–[19], [26]–[30], [35]; and for constrained least squares problems, see [4], [18] and [19].

---

\*School of Computer Science, McGill University, Montreal, Quebec, Canada, H3A 2A7. Email: chang@cs.mcgill.ca. The work of this author was supported by NSERC of Canada Grant RGPIN217191-03.

†Born 29 February 1932, died 16 November 2007. Previously: Department of Computer Science, Stanford University, CA 94305-9025. His work was in part supported by an NSF grant QAACT.

‡School of Computer Science, McGill University, Montreal, Quebec, Canada, H3A 2A7. Email: paige@cs.mcgill.ca. The work of this author was supported by NSERC of Canada Grant RGPIN9236.

The main purpose of this paper is to give a norm-wise backward perturbation analysis for the general linear data least squares (DLS) problem. As a result the structure of the matrix and magnitudes of individual elements of the matrix in the DLS problem will not be considered. We derive formulas for an “extended” minimal backward error in section 2. This extended minimal backward error is at worst a lower bound on the minimal backward error. But we show that when the given approximate solution is a good enough approximation to the exact solution of the DLS problem (which is the aim in practice), the extended minimal backward error is the actual minimal backward error. Section 2.1 deals with perturbations in both  $A$  and  $b$ , while section 2.2 considers perturbations in  $A$  alone, and shows how these are limiting cases of those in section 2.1. Since computing the extended minimal backward error directly is time consuming, in section 3 we derive a lower bound on, and in section 4 an asymptotic estimate for, this extended minimal backward error. We give numerical examples in section 5. Finally a summary is given in section 6.

We use  $I = [e_1, \dots, e_n]$  to denote the unit matrix. For any matrix  $B \in \mathbb{R}^{m \times n}$ , its column range is denoted by  $\mathcal{R}(B)$ , its Moore-Penrose generalized inverse is denoted by  $B^\dagger$ , its smallest singular value (the  $p$ -th largest singular value with  $p = \min\{m, n\}$ ) by  $\sigma_{\min}(B)$ , and its condition number in the 2-norm is denoted by  $\kappa_2(A)$ . For any symmetric  $B \in \mathbb{R}^{n \times n}$ , its eigenvalues are labeled in non-decreasing order:  $\lambda_{\min} = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ , but when only  $\lambda_{\min}$  is of interest we will write  $\lambda = \lambda_{\min}$ . For any vector  $v \in \mathbb{R}^n$ , its Moore-Penrose generalized inverse is

$$v^\dagger \equiv \begin{cases} 0 & \text{if } v = 0, \\ v^T / \|v\|_2^2 & \text{if } v \neq 0; \end{cases} \quad \|v\|_2 \equiv (v^T v)^{\frac{1}{2}}.$$

Note that  $vv^\dagger$  is the orthogonal projector onto  $\mathcal{R}(v)$  and  $I - vv^\dagger$  is the orthogonal projector onto the orthogonal complement of  $\mathcal{R}(v)$ .

**2. Backward perturbation analysis.** Given  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ , the data least squares (DLS) problem defined by DeGroat and Dowling [5] is:

$$(2.1) \quad \sigma_D \equiv \min_{E, x} \|E\|_F \quad \text{subject to} \quad (A + E)x = b, \quad \|E\|_F \equiv [\text{trace}(E^T E)]^{\frac{1}{2}}.$$

See also, for example, [21, 22]. The purpose of the DLS problem is to find the optimal  $x$ . For applications of the DLS method to some signal processing problems, see [5]. Let  $\mathcal{U}_{\min}(A)$  be the left singular vector subspace of  $A$  corresponding to its minimum singular value  $\sigma_{\min}(A)$ . In [21] it was explained that a satisfactory condition for building the theory for the DLS problem (2.1) is the condition that we will now assume holds:

$$(2.2) \quad A \text{ has full column rank, and } b \notin \mathcal{U}_{\min}(A).$$

With this condition the solution to (2.1) must exist and be unique. From  $(A + E)x = b$  we have  $Ex = b - Ax$ . Thus the minimal  $E$  must satisfy

$$(2.3) \quad E = (b - Ax)x^\dagger.$$

But (2.2) implies  $b \neq 0$ , so the solution  $x$  must be nonzero and this allows us to eliminate  $E$  and reformulate the DLS problem (2.1) as

$$(2.4) \quad \sigma_D \equiv \min_x \|(b - Ax)x^\dagger\|_F = \min_x \frac{\|b - Ax\|_2}{\|x\|_2}.$$

From [21, (5.14)–(5.17)],  $\hat{x}$  solves the DLS problem (2.1) if and only if

$$(2.5) \quad A^T(b - A\hat{x}) = -\hat{x} \frac{\|b - A\hat{x}\|_2^2}{\|\hat{x}\|_2^2},$$

$$(2.6) \quad \frac{\|b - A\hat{x}\|_2}{\|\hat{x}\|_2} < \sigma_{\min}(A).$$

Differentiating the objective function in (2.4) and setting the result to zero gives (2.5), corresponding to a stationary point. The global minimum also satisfies (2.6).

The DLS formulation is designed for problems where the right hand side  $b$  is accurately known, but the matrix  $A$  is only known approximately. Given a nonzero approximate solution  $y \in \mathbb{R}^n$  to  $Ax \approx b$ , two questions are of particular interest here:

Q1: Is  $y$  a feasible (not necessarily DLS) solution, given the accuracy of the data?

Q2: Is  $y$  a backward stable solution to the DLS problem for the given data  $A, b$ ?

Q1 will often be easy to check: for example if it is known that the given data matrix  $A$  approximates an unknown ideal matrix  $\hat{A}$  to within  $\|A - \hat{A}\|_{\{F \text{ or } 2\}} \leq \alpha$  while  $b$  is accurately known, then from (2.3) we need only check that  $\|b - Ay\|/\|y\| \leq \alpha$ . If the answer to Q1 is positive and we are not interested in the DLS solution, we might accept  $y$ . But then in practice there will be an infinite set of  $y$  satisfying Q1, and we will often seek some additional criterion, for example “does  $y$  make sense physically?”—a difficult question we might ask of an ill-posed problem. Here we consider the more generally approachable question Q2, since if we can answer this affirmatively we will know that  $y$  is a desirable *computational* solution to (2.1). Even if the answer to Q1 is “no” we might still check Q2, since it is possible for  $y$  to satisfy Q2 but not Q1, in that  $y$  can be a DLS solution for  $A + \Delta A, b + \Delta b$  for very small  $\Delta A$  and  $\Delta b$ , but the minimal norm  $E$  in  $(A + \Delta A + E)y = b + \Delta b$  can be too large for Q1. This would indicate that there are difficulties with the data.

To answer Q2 we would like to solve the minimal backward error problem:

$$(2.7) \quad \min_{\Delta A, \Delta b} \|\Delta A, \Delta b\|_F \quad \text{subject to} \quad y = \arg \min_x \frac{\|b + \Delta b - (A + \Delta A)x\|_2}{\|x\|_2},$$

see (2.4), where the chosen scalar  $\theta \geq 0$  allows a different emphasis on each data error.

From (2.5) and (2.6) we see that  $[\Delta A, \Delta b]$  is a backward perturbation for the DLS problem with the given solution  $y$  if and only if it is in the set  $\mathcal{C}_{A,b}$  where

$$(2.8) \quad \mathcal{C}_{A,b}^+ \equiv \left\{ [\Delta A, \Delta b] : (A + \Delta A)^T [b + \Delta b - (A + \Delta A)y] = -y \frac{\|b + \Delta b - (A + \Delta A)y\|_2^2}{\|y\|_2^2} \right\},$$

$$(2.9) \quad \mathcal{C}_{A,b} \equiv \left\{ [\Delta A, \Delta b] : [\Delta A, \Delta b] \in \mathcal{C}_{A,b}^+ \ \& \ \frac{\|b + \Delta b - (A + \Delta A)y\|_2}{\|y\|_2} < \sigma_{\min}(A + \Delta A) \right\}.$$

The inequality in (2.9) makes it difficult to derive a general expression for  $[\Delta A, \Delta b] \in \mathcal{C}_{A,b}$ , so we initially ignore it and consider the larger set  $\mathcal{C}_{A,b}^+$ , which we will show is also useful. The following result from Theorem 5.1 of [3] characterizes  $[\Delta A, \Delta b] \in \mathcal{C}_{A,b}^+$ .

LEMMA 2.1. *If  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ , and  $y \in \mathbb{R}^n$  is nonzero, then  $[\Delta A, \Delta b] \in \mathcal{C}_{A,b}^+$  in (2.8) if and only if there exist  $w \in \mathbb{R}^n$  and  $Z \in \mathbb{R}^{m \times n}$  such that*

$$(2.10) \quad A + \Delta A = (b + \Delta b - w)y^\dagger + (I - ww^\dagger)Z(I - yy^\dagger), \quad (b + \Delta b)^T w = 0.$$

**2.1. Allowing backward perturbations in  $A$  and  $b$ .** Based on Lemma 2.1 we will first find a computable expression for  $\mu_F(y, \theta)$  in the following ‘‘extended’’ minimal backward error problem:

$$(2.11) \quad \mu_F(y, \theta) \equiv \min_{[\Delta A, \Delta b] \in \mathcal{C}_{A,b}^+} \|\Delta A, \Delta b \theta\|_F, \text{ for } \mathcal{C}_{A,b}^+ \text{ in (2.8).}$$

We call  $\mu_F(y, \theta)$  the extended minimal backward error because we minimize over the extended set  $\mathcal{C}_{A,b}^+$ , giving at worst a lower bound on the minimal backward error.

If  $b = 0$  the DLS problem (2.1) has the solution  $\hat{x} = 0$ ; if  $y = 0$  then it cannot be the DLS solution of any problem with  $b \neq 0$ . We do not need to consider these cases further. For the remainder of this paper we will assume the conditions and notation of the following theorem.

The following will simplify the presentation:

$$(2.12) \quad \rho \equiv 1/(1 + \theta^2 \|y\|_2^2), \quad \text{so that } \rho \theta^2 \|y\|_2^2 = 1 - \rho \quad \text{and} \quad 0 \leq \rho \leq 1.$$

**THEOREM 2.2.** *Suppose that we are given  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ , nonzero  $y \in \mathbb{R}^n$  and  $\theta \geq 0$ ; and suppose that (2.2) holds. Let  $r \equiv b - Ay$ ,  $\rho \equiv 1/(1 + \theta^2 \|y\|_2^2)$ , and*

$$(2.13) \quad N \equiv N(\theta) \equiv [A(I - yy^\dagger), \rho^{\frac{1}{2}} \theta \|r\|_2 (I - rr^\dagger), b\theta],$$

$$(2.14) \quad M \equiv M(\theta) \equiv A(I - yy^\dagger)A^T - r\rho\theta^2 r^T + b\theta^2 b^T = NN^T - \rho\theta^2 \|r\|_2^2 I.$$

Then  $M(\theta)$  has at most one negative eigenvalue, and the DLS extended minimal backward error  $\mu_F(y, \theta)$  in (2.11) satisfies

$$(2.15) \quad \mu_F^2(y, \theta) = \begin{cases} \rho\theta^2 \|r\|_2^2 & \text{if } \lambda_{\min}(M(\theta)) \geq 0, \\ \rho\theta^2 \|r\|_2^2 + \lambda_{\min}(M(\theta)) = \sigma_{\min}^2(N(\theta)) & \text{if } \lambda_{\min}(M(\theta)) < 0. \end{cases}$$

Furthermore  $\mu_F(y, \theta)$  is given by the backward perturbations  $\widehat{\Delta A}$  and  $\widehat{\Delta b}$  in

$$(2.16) \quad A + \widehat{\Delta A} = \begin{cases} A + r(1 - \rho)y^\dagger & \text{if } \lambda_{\min}(M(\theta)) \geq 0, \\ (I - w_\theta w_\theta^\dagger)[A + r(1 - \rho)y^\dagger] + w_\theta w_\theta^\dagger A y y^\dagger & \text{if } \lambda_{\min}(M(\theta)) < 0, \end{cases}$$

$$(2.17) \quad b + \widehat{\Delta b} = \begin{cases} b - r\rho & \text{if } \lambda_{\min}(M(\theta)) \geq 0, \\ (I - w_\theta w_\theta^\dagger)(b - r\rho) & \text{if } \lambda_{\min}(M(\theta)) < 0, \end{cases}$$

where  $w_\theta$  is the (‘unique’ when  $\lambda_{\min}(M(\theta)) < 0$ ) eigenvector of  $M(\theta)$  corresponding to  $\lambda_{\min}(M(\theta))$ , also the left singular vector of  $N(\theta)$  corresponding to  $\sigma_{\min}(N(\theta))$ .

*Proof.*  $A(I - yy^\dagger)A^T + b\theta^2 b^T$  is nonnegative definite, so from [15, Theorem 4.3.4(b)] with  $k=1$ ,  $M(\theta)$  has at most one negative eigenvalue.

Now we want to determine the optimal  $w$ ,  $Z$ ,  $\Delta A$  and  $\Delta b$  in (2.10) to minimize  $\|\Delta A, \Delta b \theta\|_F$ . In the following we discuss two cases separately.

Case 1: The optimal  $w = 0$ . Let  $Y = [y/\|y\|_2, Y_2] \in \mathbb{R}^{n \times n}$  be an orthogonal matrix. From (2.10) we have  $(b + \Delta b)^T w = 0$  automatically, and

$$\begin{aligned} \Delta AY &= (b + \Delta b)y^\dagger [y/\|y\|_2, Y_2] + Z(I - yy^\dagger)[y/\|y\|_2, Y_2] - A[y/\|y\|_2, Y_2] \\ &= [(b + \Delta b)/\|y\|_2, 0] + [0, ZY_2] - [Ay/\|y\|_2, AY_2] \\ &= [(r + \Delta b)/\|y\|_2, (Z - A)Y_2]. \end{aligned}$$

It follows that

$$\|\Delta A, \Delta b \theta\|_F^2 = \|\Delta AY\|_F^2 + \theta^2 \|\Delta b\|_2^2 = \frac{1}{\|y\|_2^2} \left\| \begin{bmatrix} I \\ \theta \|y\|_2 I \end{bmatrix} \Delta b + \begin{bmatrix} r \\ 0 \end{bmatrix} \right\|_2^2 + \|(Z - A)Y_2\|_F^2.$$

Thus  $\|[\Delta A, \Delta b \theta]\|_F$  is minimized when

$$(2.18) \quad \Delta b = \widehat{\Delta b} \equiv - \begin{bmatrix} I \\ \theta \|y\| I \end{bmatrix}^\dagger \begin{bmatrix} r \\ 0 \end{bmatrix} = -r\rho, \quad Z = \widehat{Z} \equiv A,$$

and from (2.10) we see that the optimal  $\Delta A$  must satisfy

$$(2.19) \quad \Delta A = \widehat{\Delta A} \equiv (b - r\rho)y^\dagger + A(I - yy^\dagger) - A = r(1 - \rho)y^\dagger,$$

$$(2.20) \quad \|[\widehat{\Delta A}, \widehat{\Delta b} \theta]\|_F^2 = (1 - \rho)^2 \|ry^\dagger\|_2^2 + \rho^2 \theta^2 \|r\|_2^2 = \rho \theta^2 \|r\|_2^2.$$

In Case 2 we will show that if  $\lambda_{\min}(M) \geq 0$  then  $w = 0$  is optimal.

Case 2: The optimal  $w \neq 0$ . Let  $Y$  be as in Case 1, and  $W = [w/\|w\|_2, W_2] \in \mathbb{R}^{m \times m}$  be an orthogonal matrix. Since  $w^T(b + \Delta b) = 0$  we can write

$$(2.21) \quad b + \Delta b = W_2 s \quad \text{for some } s \in \mathbb{R}^{m-1}.$$

From (2.10) we have

$$\begin{aligned} W^T \Delta A Y &= \begin{bmatrix} w^T/\|w\|_2 \\ W_2^T \end{bmatrix} [(b + \Delta b - w)y^\dagger + (I - ww^\dagger)Z(I - yy^\dagger) - A] [y/\|y\|_2, Y_2] \\ &= \left[ \begin{array}{c|c} -\|w\|_2/\|y\|_2 - w^T Ay/(\|w\|_2\|y\|_2) & -w^T AY_2/\|w\|_2 \\ \hline s/\|y\|_2 - W_2^T Ay/\|y\|_2 & W_2^T ZY_2 - W_2^T AY_2 \end{array} \right]. \end{aligned}$$

Thus the objective function can be written as five additive nonnegative terms:

$$(2.22) \quad \|[\Delta A, \Delta b \theta]\|_F^2 = [\|w\|_2/\|y\|_2 + w^T Ay/(\|w\|_2\|y\|_2)]^2 + \|w^T AY_2\|_2^2/\|w\|_2^2 \\ + \|s - W_2^T Ay\|_2^2/\|y\|_2^2 + \|W_2^T(Z - A)Y_2\|_F^2 + \theta^2 \|W_2 s - b\|_2^2.$$

To minimize this we take  $Z = \widehat{Z} \equiv A$  and note the sum of terms involving  $s$  is

$$(2.23) \quad \begin{aligned} \phi(s) &\equiv \|s - W_2^T Ay\|_2^2/\|y\|_2^2 + \theta^2 \|W_2 s - b\|_2^2 \\ &= [\|W_2^T(W_2 s - Ay)\|_2^2 - (w^T Ay)^2/\|w\|_2^2 + \theta^2 \|y\|_2^2 \|W_2 s - b\|_2^2]/\|y\|_2^2 \\ &= \frac{1}{\|y\|_2^2} \left\| \begin{bmatrix} I \\ \theta \|y\|_2 I \end{bmatrix} W_2 s - \begin{bmatrix} Ay \\ b \theta \|y\|_2 \end{bmatrix} \right\|_2^2 - \left( \frac{w^T Ay}{\|w\|_2\|y\|_2} \right)^2. \end{aligned}$$

The normal equations for  $\hat{s}$ , the optimal  $s$ , give  $(1 + \theta^2 \|y\|_2^2) \hat{s} = W_2^T (Ay + b \theta^2 \|y\|_2^2)$ . Therefore

$$(2.24) \quad \hat{s} = W_2^T [Ay\rho + b(1 - \rho)] = W_2^T (b - r\rho).$$

Substituting this in the first line of (2.23) gives with  $W_2 W_2^T = I - ww^\dagger$

$$\begin{aligned} \phi(\hat{s}) &= \|W_2^T r(1 - \rho)\|_2^2/\|y\|_2^2 + \theta^2 \|ww^\dagger b + W_2 W_2^T r\rho\|_2^2 \\ &= \rho \theta^2 \|(I - ww^\dagger)r\|_2^2 + \theta^2 \|ww^\dagger b\|_2^2. \end{aligned}$$

Then from (2.22), we obtain

$$(2.25) \quad \min_{[\Delta A, \Delta b] \in \mathcal{C}_{A,b}^+} \|[\Delta A, \Delta b \theta]\|_F^2 = \min_w [\psi_1(w) + \psi_2(w)],$$

$$(2.26) \quad \psi_1(w) \equiv [\|w\|_2/\|y\|_2 + w^T Ay/(\|w\|_2\|y\|_2)]^2,$$

$$(2.27) \quad \psi_2(w) \equiv \|w^T AY_2\|_2^2/\|w\|_2^2 + \rho \theta^2 \|(I - ww^\dagger)r\|_2^2 + \theta^2 \|ww^\dagger b\|_2^2.$$

We will minimize  $\psi_2(w)$ , which is a function of  $w/\|w\|_2$  alone, and then show that we can set  $\psi_1(w)$  to zero by scaling  $w$ , leading to the optimal  $w$ . Since  $Y_2 Y_2^T = I - yy^\dagger$ ,

$$(2.28) \quad \begin{aligned} \psi_2(w) &= \frac{w^T A(I - yy^\dagger)A^T w}{w^T w} + \rho\theta^2 \|r\|_2^2 \frac{w^T (I - rr^\dagger)w}{w^T w} + \theta^2 \frac{w^T bb^T w}{w^T w} \\ &= \frac{w^T NN^T w}{w^T w} = \rho\theta^2 \|r\|_2^2 + \frac{w^T Mw}{w^T w}, \end{aligned}$$

whose minimum is  $\rho\theta^2 \|r\|_2^2 + \lambda_{\min}(M)$  given by  $w = w_\theta \alpha$  for any nonzero  $\alpha \in \mathbb{R}$  and  $w_\theta$  satisfying  $Mw_\theta = w_\theta \lambda_{\min}(M)$ ,  $\|w_\theta\|_2 = 1$ , since we assumed  $w \neq 0$ .

If  $\lambda_{\min}(M) \geq 0$  the above with (2.20) in Case 1 show that  $w = 0$  is optimal for minimizing  $\|[\Delta A, \Delta b \theta]\|_F^2$ , giving the minimum value  $\rho\theta^2 \|r\|_2^2$ . So from (2.18), (2.19) and (2.20) we see that the top equalities in each of (2.15), (2.16) and (2.17) hold. Only when  $\lambda_{\min}(M) < 0$  do we need to consider the possibility that  $w \neq 0$ .

Assume that  $\lambda_{\min}(M) < 0$ . It is easy to verify that  $\psi_1(\hat{w}) = 0$  if

$$(2.29) \quad \hat{w} \equiv -w_\theta (w_\theta^T A y) \neq 0.$$

Suppose  $w_\theta^T A y = 0$ , then from  $w_\theta^T M w_\theta = \lambda_{\min}(M) < 0$  and (2.14)

$$0 > \lambda_{\min}(M) = w_\theta^T A A^T w_\theta - (w_\theta^T r)^2 \rho\theta^2 + (w_\theta^T b)^2 \theta^2 = w_\theta^T A A^T w_\theta + (w_\theta^T b)^2 (1 - \rho)\theta^2$$

which is impossible since the right hand side is nonnegative, see (2.12), proving that the inequality in (2.29) holds. Therefore from (2.25) we see that when  $\lambda_{\min}(M) < 0$  the extended minimal backward error  $\mu_F(y, \theta)$  satisfies the two bottom equalities in (2.15). The bottom equality in (2.17) follows immediately from (2.21) and (2.24), and substituting this with  $Z = A$  and (2.29) in (2.10) gives

$$\begin{aligned} A + \Delta A &= [(I - w_\theta w_\theta^\dagger)(b - r\rho) + w_\theta (w_\theta^T A y)] y^\dagger + (I - w_\theta w_\theta^\dagger) A (I - yy^\dagger) \\ &= (I - w_\theta w_\theta^\dagger) [A + (b - r\rho - A y) y^\dagger] + w_\theta (w_\theta^T A y) y^\dagger \\ &= (I - w_\theta w_\theta^\dagger) [A + r(1 - \rho) y^\dagger] + w_\theta (w_\theta^T A y) y^\dagger, \end{aligned}$$

to prove the bottom equation in (2.16).  $\square$

**REMARK 2.1.** *The criterion  $\lambda_{\min}(M(\theta)) \geq 0$  appears in (2.15)–(2.17). But if, as is usual,  $m > n+1$ , then  $M(\theta)$  has at least  $m-n-1$  zero eigenvalues corresponding to eigenvectors spanning  $\mathcal{R}([A, b])^\perp$ . Eigenvalues of a parameterized matrix that are zero independent of the parameter (here  $\theta$ ) will be called “trivial zero eigenvalues”. Because they remain zero, their limiting behavior is trivial.*

**2.2. Allowing a backward perturbation in  $A$  alone.** In DLS problems only the matrix  $A$  is assumed to have uncertainty, so it is also important to consider the case where there is a backward perturbation in  $A$  alone. Then the corresponding minimal backward error problem becomes

$$(2.30) \quad \min_{\Delta A \in \mathcal{C}_A} \|\Delta A\|_F \quad \text{where} \quad \mathcal{C}_A^+ \equiv \left\{ \Delta A : (A + \Delta A)^T [b - (A + \Delta A)y] = -y \frac{\|b - (A + \Delta A)y\|_2^2}{\|y\|_2^2} \right\},$$

$$(2.31) \quad \mathcal{C}_A \equiv \left\{ \Delta A : \Delta A \in \mathcal{C}_A^+ \ \& \ \frac{\|b - (A + \Delta A)y\|_2}{\|y\|_2} < \sigma_{\min}(A + \Delta A) \right\},$$

and these two sets are just  $\mathcal{C}_{A,b}^+$  and  $\mathcal{C}_{A,b}$  in (2.8) and (2.9) with  $\Delta b = 0$ , so that  $\Delta A \in \mathcal{C}_A^+ \Rightarrow [\Delta A, 0] \in \mathcal{C}_{A,b}^+$ . We can force  $\Delta b = 0$  by taking the limit as  $\theta \rightarrow \infty$  in (2.11), giving for the *extended* minimal backward error in this more limited case

$$(2.32) \quad \mu_F(y) \equiv \min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_F = \lim_{\theta \rightarrow \infty} \mu_F(y, \theta).$$

Here we have abused the notation a little by using both  $\mu_F(\cdot)$  and  $\mu_F(\cdot, \cdot)$ . The proof using  $\theta \rightarrow \infty$  (we use  $\varepsilon \equiv \theta^{-2} \searrow 0$ ) is made possible by a beautiful classical result.

LEMMA 2.3. (*Rellich [24, pp. 29–37], see also Kato [16, pp. 121–2]*). For  $\varepsilon \in \mathbb{R}$  suppose  $H(\varepsilon) = H(\varepsilon)^T \in \mathbb{R}^{n \times n}$  is analytic about  $\varepsilon = 0$ , then its eigenvalues and an orthonormal set of eigenvectors can be chosen analytic about  $\varepsilon = 0$ .

We need some more results to prepare for Theorem 2.6.

LEMMA 2.4. If  $M = SW S^T$  with  $S \in \mathbb{R}^{m \times n}$ ,  $W = W^T \in \mathbb{R}^{n \times n}$ ,  $m \geq n$ , then

(a)  $M$  has no more positive (negative) eigenvalues than  $W$ .

(b) If  $W = I - y\alpha y^\dagger$  then  $M$  has at most one negative eigenvalue.

*Proof.* Part (a) was proven in [15, §4.5.11] for the case  $m = n$ . For  $m > n$  writing  $M = [S, 0] \text{diag}(W, 0) [S, 0]^T$  with square  $[S, 0]$  proves that (a) still holds. Since  $I - y\alpha y^\dagger$  has eigenvalues 1 when  $y = 0$ , and  $1 - \alpha, 1, \dots, 1$  when  $y \neq 0$ , (b) follows from (a).  $\square$

To make later analysis easier, we use  $\varepsilon$  to replace  $\theta^{-2}$ . From (2.12)

$$(2.33) \quad \rho = \varepsilon / (\varepsilon + \|y\|_2^2), \quad \varepsilon + \rho \|y\|_2^2 = \varepsilon(2 - \rho), \quad \varepsilon \equiv \theta^{-2}.$$

In our limits we only consider  $\varepsilon \geq 0$ , so  $\varepsilon \rightarrow 0$  will always mean  $\varepsilon \searrow 0$ .

THEOREM 2.5. With (2.33) if  $A \in \mathbb{R}^{m \times n}$ ,  $0 \neq y \in \mathbb{R}^n$ ,  $0 \neq b \in \mathbb{R}^m$ ,  $r \equiv b - Ay$ ,

$$(2.34) \quad H(\varepsilon) \equiv \varepsilon A(I - yy^\dagger)A^T - r\rho r^T + bb^T, \quad \varepsilon \in \mathbb{R}, \quad H(\varepsilon)w(\varepsilon) = w(\varepsilon)\lambda(\varepsilon), \quad w(\varepsilon) \in \mathbb{R}^m,$$

$m \geq 2$  and  $\lambda(0) \equiv \lambda_{\min}(H(0))$ , then for small enough  $\varepsilon \geq 0$  the minimum eigenvalue  $\lambda(\varepsilon)$  and its normalized eigenvector  $w(\varepsilon)$  can be chosen analytic with the forms

$$(2.35) \quad \lambda(\varepsilon) = \lambda_1 \varepsilon + \lambda_2 \varepsilon^2 + \dots; \quad w(\varepsilon) = w_0 + w_1 \varepsilon + w_2 \varepsilon^2 + \dots, \quad b^T w_0 = 0, \quad \|w(\varepsilon)\|_2 = 1.$$

$$(2.36) \quad \text{If } T(\varepsilon) \equiv \varepsilon^{-1} P_b^\perp H(\varepsilon) P_b^\perp, \quad \text{then } T(0) \equiv \lim_{\varepsilon \rightarrow 0} T(\varepsilon) = P_b^\perp A(I - y2y^\dagger)A^T P_b^\perp.$$

Let  $\lambda_*(0) \equiv \lambda_{\min}(T(0))$ , then for small enough  $\varepsilon \geq 0$  the minimum eigenvalue  $\lambda_*(\varepsilon)$  and its normalized eigenvector  $w_*(\varepsilon)$  can be chosen analytic in

$$(2.37) \quad T(\varepsilon)w_*(\varepsilon) = w_*(\varepsilon)\lambda_*(\varepsilon), \quad w_*(\varepsilon) \in \mathbb{R}^m, \quad \|w_*(\varepsilon)\|_2 = 1.$$

Finally  $\lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon)$  exists (see (2.35)), and for  $\lambda_*(0)$  and  $w_*(0)$  in (2.37),

$$(2.38) \quad \lim_{\varepsilon \rightarrow 0} b^T w(\varepsilon) = 0, \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^{-\frac{1}{2}} b^T w(\varepsilon) = 0, \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} b^T w(\varepsilon) = b^T w_1;$$

$$(2.39) \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) < 0 \quad \Rightarrow \quad \lambda_*(0) < 0 \quad \Rightarrow \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) = \lambda_*(0) \quad \& \quad w(0) = \pm w_*(0);$$

$$(2.40) \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) \geq 0 \quad \Leftrightarrow \quad \lambda_*(0) = 0;$$

$$(2.41) \quad \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) = 0 \quad \Rightarrow \quad \lambda_*(0) = 0 \quad \& \quad \begin{cases} \{\exists x \text{ such that } b = Ax \text{ \& } \angle(x, y) = \pm\pi/4\}, \\ \text{or } \{\varepsilon^{-1} \lambda(\varepsilon) = 0 \text{ in a neighborhood of } \varepsilon = 0\}. \end{cases}$$

*Proof.* The expression for  $T(0)$  in (2.36) follows from (2.34) and (2.33), and then

$$(2.42) \quad T(\varepsilon) = T(0) + P_b^\perp A y \rho y^\dagger A^T P_b^\perp = P_b^\perp A [I - y(2 - \rho)y^\dagger] A^T P_b^\perp.$$

Clearly  $H(\varepsilon)$  and  $T(\varepsilon)$  are analytic about  $\varepsilon=0$ , so  $w(\varepsilon)$ ,  $\lambda(\varepsilon)$ ,  $w_*(\varepsilon)$ , and  $\lambda_*(\varepsilon)$  can be chosen to be analytic with  $\|w(\varepsilon)\|_2 = \|w_*(\varepsilon)\|_2 = 1$ , see Lemma 2.3. Also  $H(\varepsilon)$  can have at most one negative eigenvalue, see the start of the proof of Theorem 2.2, so from Lemma 2.4  $T(\varepsilon)$  can have at most one negative eigenvalue. Since  $m \geq 2$ ,  $H(0) = bb^T$  has minimum eigenvalue  $\lambda(0) = 0$ , proving the first part of (2.35). Since  $bb^T w(0) = w(0)\lambda(0) = 0$ , we must have  $b^T w(0) = b^T w_0 = 0$ , proving the rest of (2.35). Next (2.35) proves (2.38), and we have

$$(2.43) \quad \begin{aligned} \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) &= \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} w(\varepsilon)^T H(\varepsilon) w(\varepsilon) = \lim_{\varepsilon \rightarrow 0} \varepsilon^{-\frac{1}{2}} w(\varepsilon)^T P_b^\perp H(\varepsilon) P_b^\perp w(\varepsilon) \varepsilon^{-\frac{1}{2}} \\ &= \lim_{\varepsilon \rightarrow 0} w(\varepsilon)^T T(\varepsilon) w(\varepsilon) = w(0)^T T(0) w(0) \geq \lambda_*(0) = \lambda_{\min}(T(0)), \end{aligned}$$

so  $\lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) < 0 \Rightarrow \lambda_*(0) < 0$ . When  $\lambda_*(0) < 0$  it is a singleton, see Lemma 2.4, so for small enough  $\varepsilon$ ,  $\lambda_*(\varepsilon) < 0$ ; then  $b^T w_*(\varepsilon) = 0$  from (2.37) and (2.36), giving

$$(2.44) \quad \lambda_*(\varepsilon) = \varepsilon^{-1} w_*(\varepsilon)^T P_b^\perp H(\varepsilon) P_b^\perp w_*(\varepsilon) = \varepsilon^{-1} w_*(\varepsilon)^T H(\varepsilon) w_*(\varepsilon).$$

Taking the limit as  $\varepsilon \rightarrow 0$  and using (2.43) with  $\|w(\varepsilon)\|_2 = \|w_*(\varepsilon)\|_2 = 1$  gives

$$\begin{aligned} \lambda_*(0) &= \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} w_*(\varepsilon)^T H(\varepsilon) w_*(\varepsilon) \geq \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) = w(0)^T T(0) w(0) \\ &\geq \lambda_*(0) = w_*(0)^T T(0) w_*(0), \end{aligned}$$

proving equality throughout, so that  $\lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) = \lambda_*(0)$  when  $\lambda_*(0) < 0$ . Also  $w(0)^T T(0) w(0) = w_*(0)^T T(0) w_*(0)$  is a minimum of  $w^T T(0) w$  over  $w^T w = 1$  with unique minimizer (up to sign) when  $\lambda_*(0)$  is a singleton, completing the proof of (2.39). Since  $T(\varepsilon)b = 0$  we see that  $\lambda_*(0) \leq 0$ , and (2.40) follows using (2.39).

Now assume  $\lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) = 0$ , ( $\lambda_1 = 0$  in (2.35)). Then  $\lambda_*(0) = 0$  in (2.41) follows from (2.40). If  $[A, b]$  has rank  $s$  then  $\varepsilon^{-1} H(\varepsilon)$  has  $m-s$  trivial zero eigenvalues. If in the limit as  $\varepsilon \rightarrow 0$  there are only trivial zero eigenvalues then by continuity  $\varepsilon^{-1} \lambda(\varepsilon) = 0$  in some neighborhood of  $\varepsilon = 0$ . Next assume there is a nontrivial zero eigenvalue, that is, an eigenpair of the form (2.35) with  $\lambda_1 = 0$  and  $A^T w_0 \neq 0$ . But  $\lambda_*(0) = 0$  shows that  $T(0)$  is positive semi-definite, and  $0 = \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \lambda(\varepsilon) = \lim_{\varepsilon \rightarrow 0} w_0^T \varepsilon^{-1} H(\varepsilon) w_0 = \lim_{\varepsilon \rightarrow 0} w_0^T T(\varepsilon) w_0$ , see (2.36), so  $0 = w_0^T T(0) w_0 = w_0^T A(I - y2y^\dagger) A^T w_0$ . Here  $I - y2y^\dagger$  is a Householder reflection. Thus

$$(2.45) \quad \|A^T w_0\|_2^2 = 2(w_0^T A y)^2 / \|y\|_2^2, \quad 0 = T(0) w_0 = P_b^\perp A(I - y2y^\dagger) A^T w_0,$$

and  $A(I - y2y^\dagger) A^T w_0 = bb^\dagger A(I - y2y^\dagger) A^T w_0 \neq 0$ , see (2.2). Then  $Ax\nu = b\nu$  where

$$x\nu = (I - y2y^\dagger) A^T w_0, \quad \nu \equiv b^T A(I - y2y^\dagger) A^T w_0 / b^T b \neq 0.$$

From this  $y^T x\nu = -y^T A^T w_0$ , and with (2.45)  $x^T x\nu^2 = \|A^T w_0\|_2^2 = 2(w_0^T A y)^2 / \|y\|_2^2 = 2(y^T x)^2 \nu^2 / \|y\|_2^2$ . This gives  $2(y^T x)^2 = y^T y \cdot x^T x$ , so  $\angle(x, y) = \pm\pi/4$ , proving (2.41).  $\square$

**REMARK 2.2.** *In (2.41) the case  $b = Ax$  &  $\angle(x, y) = \pm\pi/4$  is extremely unlikely (it has ‘‘probability zero’’), requiring  $b \in \mathcal{R}(A)$  (a highly unlikely situation when we are solving DLS problems), and  $y$  to be a terrible approximation to the unique  $x$  for which  $b = Ax$ , giving exactly  $\angle(x, y) = \pm\pi/4$ .*

We can now obtain the equivalent of Theorem 2.2 for a backward perturbation restricted to  $A$  alone.



THEOREM 2.6. *Suppose that we are given  $A \in \mathbb{R}^{m \times n}$ ,  $m \geq 2$ , with nonzero  $b \in \mathbb{R}^m$  and  $y \in \mathbb{R}^n$ ; suppose also that (2.2) holds. Let  $r \equiv b - Ay$  and*

$$(2.46) \quad N_\infty \equiv \left[ (I - bb^\dagger)A(I - yy^\dagger), \frac{\|r\|_2}{\|y\|_2}(I - bb^\dagger)(I - rr^\dagger), \frac{b}{\|b\|_2} \frac{\|r\|_2}{\|y\|_2} \right],$$

$$(2.47) \quad M_\infty = M_\infty(y) \equiv (I - bb^\dagger)A(I - yy^\dagger)A^T(I - bb^\dagger) = N_\infty N_\infty^T - \|r\|_2^2 / \|y\|_2^2 I.$$

Then  $\lambda_{\min}(M_\infty) \leq 0$  and  $M_\infty$  has at most one negative eigenvalue. Also the DLS extended minimal backward error  $\mu_F(y)$  in (2.32) satisfies

$$(2.48) \quad \mu_F^2(y) = \frac{\|r\|_2^2}{\|y\|_2^2} + \lambda_{\min}(M_\infty) = \sigma_{\min}^2(N_\infty).$$

Furthermore  $\mu_F(y)$  is given by the backward perturbation:

$$(2.49) \quad \widehat{\Delta A} = \begin{cases} ry^\dagger & \text{if } \lambda_{\min}(M_\infty) = 0, \\ ry^\dagger - w_* w_*^\dagger A(I - yy^\dagger) & \text{if } \lambda_{\min}(M_\infty) < 0, \end{cases}$$

where  $w_*$  is the eigenvector of  $M_\infty$  corresponding to  $\lambda_{\min}(M_\infty) < 0$ , or equivalently the left singular vector of  $N_\infty$  corresponding to  $\sigma_{\min}(N_\infty)$ . If  $\lambda_{\min}(M_\infty) < 0$  in (2.47), then  $\lambda_{\min}(M_\infty) = \lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta))$  in (2.14). In general (2.48)–(2.49) are the corresponding limiting values of (2.15)–(2.16) as  $\theta \rightarrow \infty$ , except possibly in the probability zero case mentioned in Remark 2.2 that could only occur if  $\lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) = 0$  is a nontrivial zero eigenvalue, see Remark 2.1.

*Proof.* Since (2.47) follows from (2.46), the results on  $N_\infty$  follow trivially from those on  $M_\infty$ . Since  $M_\infty b = 0$ ,  $\lambda_{\min}(M_\infty) \leq 0$ . From Lemma 2.4  $M_\infty$  has at most one negative eigenvalue. In fact  $M_\infty$  in (2.47) is identical to  $T(0)$  in (2.36), so with  $w_* \equiv w_*(0)$ ,  $\lambda_* \equiv \lambda_{\min}(M_\infty) \equiv \lambda_*(0)$  in (2.37),  $M_\infty w_* = w_* \lambda_*$ . From (2.14) and (2.34)

$$(2.50) \quad M(\theta) \equiv A(I - yy^\dagger)A^T - r\rho\theta^2 r^T + b\theta^2 b^T = \varepsilon^{-1}H(\varepsilon) \quad \text{with } \varepsilon \equiv \theta^{-2}.$$

Since  $M(\theta)$  can have at most one negative eigenvalue, when  $\lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) < 0$ ,  $\lambda_{\min}(M(\theta))$  is the unique minimum eigenvalue for large enough  $\theta$ , and  $w_\theta$  in (2.16) can be taken as  $w(\varepsilon)$  in Theorem 2.5. This, and noting that  $\lambda(\theta^2 H(\theta^{-2}))$  is equal to  $\lambda_{\min}(M(\theta))$  for large enough  $\theta$ , gives from (2.38)–(2.41)

$$\begin{aligned} \lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) < 0 &\Rightarrow \begin{cases} \lambda_{\min}(M_\infty) = \lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) \\ \& w_* = \pm \lim_{\theta \rightarrow \infty} w_\theta \quad \& \lim_{\theta \rightarrow \infty} \theta w_\theta^T b = 0; \end{cases} \\ \lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) \geq 0 &\Leftrightarrow \lambda_* \equiv \lambda_{\min}(M_\infty) = 0; \\ \lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) = 0 &\Rightarrow \lambda_{\min}(M_\infty) = 0 \quad \& \begin{cases} \{\exists x : b = Ax \quad \& \angle(x, y) = \pm\pi/4\}, \\ \text{or } \{\exists \theta_1 : \lambda_{\min}(M(\theta)) = 0 \quad \forall \theta > \theta_1\}. \end{cases} \end{aligned}$$

But from (2.12)  $\lim_{\theta \rightarrow \infty} \rho = 0$  and  $\lim_{\theta \rightarrow \infty} \rho\theta^2 = \|y\|_2^{-2}$ , so that (2.48) is the limiting value of both cases of (2.15) as  $\theta \rightarrow \infty$ . If  $\lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) \neq 0$  it can also be seen that in the limit the two criteria in (2.16) become the respective criteria in (2.49), where if  $\lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) < 0$ ,  $b^T w_* = \lim_{\theta \rightarrow \infty} b^T w_\theta = 0$ , so that in the limit the two expressions for  $\widehat{\Delta A}$  in (2.16) become the respective expressions in (2.49). If  $\lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) = 0$  and this corresponds to trivial zero eigenvalues only, the top row of (2.49) is clearly once again the correct limit. Only the probability zero case of  $\lim_{\theta \rightarrow \infty} \lambda_{\min}(M(\theta)) = 0$  with  $b = Ax$ ,  $\angle(x, y) = \pm\pi/4$ , allows the possibility

that  $\lambda_{\min}(M(\theta)) < 0$  for arbitrarily large  $\theta$ , and since in this case  $\lambda_{\min}(M_\infty) = 0$ , this suggests that (2.16) could fail to give the correct limiting perturbation in (2.49). That is, although (2.49) is correct, until proven otherwise there remains the possibility that taking the limit in (2.16) could lead to  $\widehat{\Delta A} = ry^\dagger - w_* w_*^\dagger A(I - y2y^\dagger)$  rather than  $ry^\dagger$  in this one strange case, see (2.49).  $\square$

Deriving  $\mu_F(y)$  directly as we did for  $\mu_F(y, \theta)$  also leads to the results down to the sentence including (2.49). But Theorem 2.5 describes the limiting behavior as well.

To parallel the remark given in [13, Sec. 20.7] for some formulas for the minimal backward error of ordinary least squares problems, computing and adding the eigenvalue in (2.15) or (2.48) is not wise computationally. Catastrophic cancellation may occur when it is negative. Furthermore the computed value may have very poor accuracy even using well-known software such as MATLAB 7.0, e.g., in (2.48) the computed value of  $\lambda_{\min}(M_\infty)$  may be smaller than  $-\|r\|_2^2/\|y\|_2^2$ . The singular value is much more reliable for computation. If we computed that using the Golub-Reinsch singular value decomposition algorithm it would need about  $8/3m^3 + 4mn^2$  flops, but one point of this paper is that we can use cheaper lower bounds or estimates instead.

In Theorem 2.2 we have either  $\lambda_{\min}(M(\theta)) < 0$  or  $\lambda_{\min}(M(\theta)) \geq 0$ , while in Theorem 2.6 we have either  $\lambda_{\min}(M_\infty) < 0$  or  $\lambda_{\min}(M_\infty) = 0$ . By substituting the resulting perturbations in the relevant inequalities it is straightforward to see that the inequality in (2.9) is satisfied when  $\lambda_{\min}(M(\theta)) \geq 0$  and  $\text{rank}(A + r(1 - \rho)y^\dagger) = n$  and the inequality in (2.31) is satisfied when  $\lambda_{\min}(M_\infty) = 0$  and  $\text{rank}(A + ry^\dagger) = n$ . It follows that in these two special cases the extended minimal backward error is actually the true minimal backward error, and that nothing was lost by using the ‘‘supersets’’  $\mathcal{C}_{A,b}^+$  and  $\mathcal{C}_A^+$ . We will supply further justification for the use of these supersets later.

The following result indicates that the extended minimal backward error  $\mu_F(y)$  is continuous at  $y = \hat{x}$ , where  $\hat{x}$  is the DLS solution in (2.5)–(2.6), where of course  $\mu_F(\hat{x}) = 0$ . In order to save space, here and in the rest of the paper we will only consider the case where  $A$  is perturbed, but all the results could be extended to the more general case where both  $A$  and  $b$  are perturbed.

**COROLLARY 2.7.** *With the notation and conditions of Theorem 2.6, and the DLS solution  $\hat{x}$  in (2.5)–(2.6), define  $\hat{M}_\infty \equiv M_\infty(\hat{x})$  (see (2.47)), and  $\hat{r} \equiv b - A\hat{x}$ . Then with (2.15),*

$$\lim_{y \rightarrow \hat{x}} \mu_F(y) = \mu_F(\hat{x}) = \left( \frac{\|\hat{r}\|_2^2}{\|\hat{x}\|_2^2} + \lambda_{\min}(\hat{M}_\infty) \right)^{1/2} = 0.$$

*Proof.* First we see that (2.5) just says

$$A^T \hat{r} = A^T (b - A\hat{x}) = -\hat{x} \|b - A\hat{x}\|_2^2 / \|\hat{x}\|_2^2 = -\hat{x} \|\hat{r}\|_2^2 / \|\hat{x}\|_2^2,$$

and multiplying this on the left by  $\hat{x}^T$  shows that

$$(2.51) \quad 0 = (b - A\hat{x})^T (b - A\hat{x}) + (A\hat{x})^T (b - A\hat{x}) = b^T (b - A\hat{x}) = b^T \hat{r},$$

so that  $(I - bb^\dagger)\hat{r} = \hat{r}$ . Since  $\hat{M}_\infty = M_\infty(\hat{x}) = (I - bb^\dagger)A(I - 2\hat{x}\hat{x}^\dagger)A^T(I - bb^\dagger)$ , the above give

$$\begin{aligned} \hat{M}_\infty \hat{r} &= (I - bb^\dagger)A(I - 2\hat{x}\hat{x}^\dagger)A^T \hat{r} = -(I - bb^\dagger)A(I - 2\hat{x}\hat{x}^\dagger)\hat{x} \|\hat{r}\|_2^2 / \|\hat{x}\|_2^2 \\ &= (I - bb^\dagger)A\hat{x} \|\hat{r}\|_2^2 / \|\hat{x}\|_2^2 = (I - bb^\dagger)(A\hat{x} - b) \|\hat{r}\|_2^2 / \|\hat{x}\|_2^2 \\ &= -\hat{r} \|\hat{r}\|_2^2 / \|\hat{x}\|_2^2. \end{aligned}$$

Thus by Lemma 2.4,  $-\|\hat{r}\|_2^2/\|\hat{x}\|_2^2$  is the only negative eigenvalue of  $\hat{M}_\infty$ , and  $\mu_F(\hat{x}) = 0$ . Clearly when  $y \rightarrow \hat{x}$  we have  $r = b - Ay \rightarrow \hat{r}$ ,  $M_\infty \rightarrow \hat{M}_\infty$ , and by the continuity of the eigenvalues of  $M_\infty$  in (2.47),  $\lambda_{\min}(M_\infty) \rightarrow \lambda_{\min}(\hat{M}_\infty)$ , completing the proof.  $\square$

Since in (2.32)  $\mathcal{C}_A \subseteq \mathcal{C}_A^+$ ,

$$\mu_F(y) \equiv \min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_F \leq \min_{\Delta A \in \mathcal{C}_A} \|\Delta A\|_F,$$

i.e.,  $\mu_F(y)$  is a lower bound on the minimal backward error. However we have found computationally, see section 5, that when  $y$  is a reasonable approximation to the exact solution of the DLS problem (2.1), the minimal perturbation  $\widehat{\Delta A}$  usually satisfies the inequality in (2.31). Therefore in such cases  $\mu(y)$  is actually the minimal backward error. The following result partially justifies this finding.

**THEOREM 2.8.** *For given  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$  suppose that (2.2) holds. Let  $\hat{x}$  be the DLS solution to (2.1). Then there exists an  $\epsilon > 0$  such that if  $\|y - \hat{x}\|_2 < \epsilon$ , then  $\mu_F(y)$  above is the true minimal backward error.*

*Proof.* For any given  $y$  Theorem 2.6 shows that  $\widehat{\Delta A}$  satisfying (2.49) is the minimizer of (2.32). Notice that when  $y \rightarrow \hat{x}$  we have from Corollary 2.7 that  $\widehat{\Delta A} \rightarrow 0$ . Thus

$$\lim_{y \rightarrow \hat{x}} \left( \frac{\|b - (A + \widehat{\Delta A})y\|_2}{\|y\|_2} - \sigma_{\min}(A + \widehat{\Delta A}) \right) = \frac{\|b - A\hat{x}\|_2}{\|\hat{x}\|_2} - \sigma_{\min}(A)$$

Since  $\frac{\|b - A\hat{x}\|_2}{\|\hat{x}\|_2} - \sigma_{\min}(A) < 0$ , there must exist  $\epsilon > 0$  such that when  $\|y - \hat{x}\|_2 < \epsilon$ ,

$$\frac{\|b - (A + \widehat{\Delta A})y\|_2}{\|y\|_2} - \sigma_{\min}(A + \widehat{\Delta A}) < 0.$$

Therefore  $\widehat{\Delta A} \in \mathcal{C}_A$  and  $\mu_F(y) = \min_{\Delta A \in \mathcal{C}_A} \|\Delta A\|_F$ , i.e., when  $\|y - \hat{x}\|_2 < \epsilon$ ,  $\mu_F(y)$  is the true minimal backward error.  $\square$

**3. A lower bound on  $\min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_2$ .** Since computing  $\mu_F(y)$  directly is expensive, in this section we suggest a good lower bound which can be estimated easily.

First we give the following result, which is analogous to Theorem 3.1 in [35] for ordinary least squares problems.

**THEOREM 3.1.** *With the notation and conditions of Theorem 2.6, and with  $\mathcal{C}_A^+$  in (2.30) let  $\mu_2(y) \equiv \min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_2$ . Then for  $\mu_F(y)$  in (2.32)*

$$(3.1) \quad \frac{1}{\sqrt{2}}\mu_F(y) \leq \mu_2(y) \leq \mu_F(y).$$

*Proof.* For any  $\Delta A \in \mathcal{C}_A^+$  we see from Lemma 2.1 that there exist  $w$  satisfying  $b^T w = 0$  and  $Z \in \mathbb{R}^{m \times n}$  such that

$$\begin{aligned} \Delta A &= (b - w)y^\dagger + (I - ww^\dagger)Z(I - yy^\dagger) - A \\ &= [(I - ww^\dagger)b - w]y^\dagger + (I - ww^\dagger)Z(I - yy^\dagger) \\ &\quad - (I - ww^\dagger)Ayy^\dagger - ww^\dagger Ayy^\dagger - (I - ww^\dagger)A(I - yy^\dagger) - ww^\dagger A(I - yy^\dagger) \\ &= [(I - ww^\dagger)(b - Ay) - w - ww^\dagger Ay]y^\dagger - ww^\dagger A(I - yy^\dagger) \\ &\quad + (I - ww^\dagger)(Z - A)(I - yy^\dagger). \end{aligned}$$

Denote  $\Delta A_1 = [(I - ww^\dagger)(b - Ay) - w - ww^\dagger Ay]y^\dagger$ ,  $\Delta A_2 = -ww^\dagger A(I - yy^\dagger)$ , and  $\Delta A_3 = (I - ww^\dagger)(Z - A)(I - yy^\dagger)$ . Since  $\Delta A_1 \Delta A_2^T = 0$ ,  $\Delta A_1 \Delta A_3^T = 0$  and  $\Delta A_2^T \Delta A_3 = 0$  and  $\Delta A_1$  and  $\Delta A_2$  are rank 1 matrices,

$$\begin{aligned}
\|\Delta A\|_2^2 &= \|\Delta A_1 + \Delta A_2 + \Delta A_3\|_2^2 = \|(\Delta A_1 + \Delta A_2 + \Delta A_3)(\Delta A_1 + \Delta A_2 + \Delta A_3)^T\|_2 \\
&= \|\Delta A_1 \Delta A_1^T + (\Delta A_2 + \Delta A_3)(\Delta A_2 + \Delta A_3)^T\|_2 \\
&\geq \max\{\|\Delta A_1 \Delta A_1^T\|_2, \|(\Delta A_2 + \Delta A_3)(\Delta A_2 + \Delta A_3)^T\|_2\} \\
&\geq \frac{1}{2}(\|\Delta A_1\|_2^2 + \|\Delta A_2 + \Delta A_3\|_2^2) \\
&= \frac{1}{2}(\|\Delta A_1\|_2^2 + \|(\Delta A_2 + \Delta A_3)^T(\Delta A_2 + \Delta A_3)\|_2) \\
&= \frac{1}{2}(\|\Delta A_1\|_2^2 + \|\Delta A_2^T \Delta A_2 + \Delta A_3^T \Delta A_3\|_2) \geq \frac{1}{2}(\|\Delta A_1\|_2^2 + \|\Delta A_2\|_2^2) \\
&= \frac{1}{2}(\|\Delta A_1\|_F^2 + \|\Delta A_2\|_F^2) = \frac{1}{2}\|\Delta A_1 + \Delta A_2\|_F^2 \geq \frac{1}{2} \min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_F^2,
\end{aligned}$$

where the last inequality is due to the fact that  $\Delta A_1 + \Delta A_2 \in \mathcal{C}_A^+$  (take  $Z = A$ ). Therefore

$$\min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_2^2 \geq \frac{1}{2} \min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_F^2,$$

leading to the first inequality in (3.1). The second inequality in (3.1) is easy to prove. In fact, if  $\widehat{\Delta A}$  is a minimal solution to (2.32), then

$$\mu_F(y) = \|\widehat{\Delta A}\|_F \geq \|\widehat{\Delta A}\|_2 \geq \mu_2(y).$$

□

Now we give a lower bound on  $\mu_2(y) \equiv \min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_2$ .

**THEOREM 3.2.** *With the notation and conditions of Theorem 2.6, and with  $\mu_2(y) \equiv \min_{\Delta A \in \mathcal{C}_A^+} \|\Delta A\|_2$  for  $\mathcal{C}_A^+$  in (2.30),*

$$(3.2) \quad \mu_2(y) \geq \mu_2^{\text{lb}}(y) \equiv \frac{2\beta_0}{\beta_1 + \sqrt{\beta_1^2 + 4\beta_0}},$$

where

$$(3.3) \quad \beta_0 \equiv \frac{\|(A^T r \|y\|_2^2 + y \|r\|_2^2)\|_2}{2\|y\|_2^3}, \quad \beta_1 \equiv \frac{\|y\|_2^3 \|A\|_2 + 3\|y\|_2^2 \|r\|_2}{2\|y\|_2^3}.$$

*Proof.* For any  $\Delta A \in \mathcal{C}_A^+$ , from (2.30) we obtain

$$(A + \Delta A)^T (r - \Delta A y) \|y\|_2^2 + y \|r - \Delta A y\|_2^2 = 0.$$

Thus we have

$$\begin{aligned}
A^T r \|y\|_2^2 + y \|r\|_2^2 &= A^T \Delta A y \|y\|_2^2 - \Delta A^T r \|y\|_2^2 + y 2(r^T \Delta A y) \\
&\quad + \Delta A^T \Delta A y \|y\|_2^2 - y \|\Delta A y\|_2^2.
\end{aligned}$$

Then taking the 2-norm on both sides of this equation, we obtain the inequality

$$\|(A^T r \|y\|_2^2 + y \|r\|_2^2)\|_2 \leq (\|y\|_2^3 \|A\|_2 + 3\|y\|_2^2 \|r\|_2) \|\Delta A\|_2 + 2\|y\|_2^3 \|\Delta A\|_2^2,$$

that is, with (3.3), the quadratic inequality in terms of  $\xi \equiv \|\Delta A\|_2$ :

$$\beta_0 \leq \beta_1 \xi + \xi^2.$$

Since  $\xi$  and  $\beta_1$  are nonnegative,  $\xi \geq \xi_+$ , where  $\xi_+$  is the positive root of  $\beta_0 = \beta_1 \xi + \xi^2$ , so

$$\xi \geq \xi_+ = (\sqrt{\beta_1^2 + 4\beta_0} - \beta_1)/2 = 2\beta_0/(\sqrt{\beta_1^2 + 4\beta_0} + \beta_1),$$

giving (3.2).  $\square$

The lower bound in (3.2) can usually be evaluated in  $\mathcal{O}(mn)$  flops, since  $\|A\|_2$  can usually be estimated by a standard norm estimator in  $O(mn)$  flops, see [13, Sec. 15.2]. In fact a good estimate of  $\|A\|_2$  might already be available from whatever method is used for obtaining  $y$ , and the cost will essentially be the  $4mn$  flops for computing  $A^T(b - Ay)$ .

Also  $\mu_2^{\text{lb}}(\hat{x}) = \mu_2(\hat{x}) = \mu_F(\hat{x}) = 0$  as desired, see (2.5), Corollary 2.7, and (3.1).

**4. An asymptotic estimate for  $\mu_F(y)$ .** Computing  $\mu_F(y)$  directly is expensive and the lower bound (3.2) may not be very tight. In this section we would like to give an asymptotic estimate by following the general approach given in [9].

Let  $f(A, y) \equiv (b - Ay)^T(b - Ay)y + (y^T y)A^T(b - Ay) = \|r\|_2^2 y + \|y\|_2^2 A^T r$ . Note that  $f(A, \hat{x}) = 0$  (see (2.5)). The extended minimal backward perturbation  $\Delta A$  is the matrix satisfying  $f(A + \Delta A, y) = 0$  and  $\mu_F(y) = \|\Delta A\|_F$ . But by Taylor's expansion, for small enough  $E \in \mathbb{R}^{m \times n}$ ,

$$f(A + E, y) \approx f(A, y) + \mathcal{J}_A f(A, y) \text{vec}(E),$$

where  $\mathcal{J}_A f(A, y) \in \mathbb{R}^{n \times mn}$  is the Jacobian matrix of  $f$  with respect to  $\text{vec}(A)$ . Thus an approximation to  $\Delta A$  is that  $E$  giving the minimum 2-norm solution to

$$(4.1) \quad f(A, y) + \mathcal{J}_A f(A, y) \text{vec}(E) = 0,$$

that is,  $E$  such that

$$(4.2) \quad \text{vec}(E) = -[\mathcal{J}_A f(A, y)]^\dagger f(A, y), \quad \tilde{\mu}_F(y) \equiv \|E\|_F = \|[\mathcal{J}_A f(A, y)]^\dagger f(A, y)\|_2.$$

**THEOREM 4.1.** *With the notation and conditions of Theorem 2.6,  $\tilde{\mu}_F(y)$  in (4.2) is an asymptotic estimate of  $\mu_F(y)$  in (2.32), i.e. for  $\hat{x}$  solving (2.1),*

$$\lim_{y \rightarrow \hat{x}} \frac{\tilde{\mu}_F(y)}{\mu_F(y)} = 1.$$

*Proof.* By Taylor's expansion,

$$(4.3) \quad 0 = f(A + \Delta A, y) = f(A, y) + \mathcal{J}_A f(A, y) \text{vec}(\Delta A) + O(\|\Delta A\|_F^2).$$

Thus from (4.2),

$$\text{vec}(E) = -[\mathcal{J}_A f(A, y)]^\dagger f(A, y) = [\mathcal{J}_A f(A, y)]^\dagger \mathcal{J}_A f(A, y) \text{vec}(\Delta A) + O(\|\Delta A\|_F^2).$$

Taking the 2-norm and noticing  $[\mathcal{J}_A f(A, y)]^\dagger \mathcal{J}_A f(A, y)$  is an orthogonal projection matrix, we obtain

$$\tilde{\mu}_F(y) \leq \mu_F(y) + O(\|\Delta A\|_F^2),$$

which, with Corollary 2.7, leads to  $\lim_{y \rightarrow \hat{x}} \tilde{\mu}_F(y)/\mu_F(y) \leq 1$ .

On the other hand, from (4.1) and (4.3) we can obtain

$$\mathcal{J}_A f(A, y) \text{vec}(\Delta A) = \mathcal{J}_A f(A, y) [\text{vec}(E) + O(\|\Delta A\|_F^2)].$$

Since  $\Delta A$  is a matrix satisfying the above equality with minimum F-norm, we must have

$$\|\text{vec}(\Delta A)\|_2 \leq \|\text{vec}(E) + O(\|\Delta A\|_F^2)\|_2 \leq \|\text{vec}(E)\|_2 + O(\|\Delta A\|_F^2)$$

which, with Corollary 2.7, leads to  $\lim_{y \rightarrow \hat{x}} \tilde{\mu}_F(y)/\mu_F(y) \geq 1$ , completing the proof.  $\square$

Theorem 4.1 is similar to [9, Cor 3.4], where a general minimal backward error problem was considered, and applying the corollary to our case will result in the asymptotic estimate  $\|[\mathcal{J}_A f(A, \hat{x})]^\dagger f(A, y)\|_2$ . We thank the referee who pointed out that a general version of Theorem 4.1 was given in [10] (with no formal proof).

In the following we will consider computing or estimating  $\tilde{\mu}_F(y)$ . First we would like to obtain an explicit expression for it. If  $f = (f_i)$  and  $g = (g_i)$  are column vectors, then we define the matrix  $\partial f / \partial g^T \equiv (\partial f_i / \partial g_j)$ . Write  $m \times n$   $A = [a_1, \dots, a_n]$ ,  $y^T = (\eta_1, \dots, \eta_m)$ , then  $\partial r / \partial a_j^T = \partial(b - Ay) / \partial a_j^T = -\eta_j I$ ,  $\partial(r^T r) / \partial a_j^T = 2r^T \partial r / \partial a_j^T = -2\eta_j r^T$ , and if  $i \neq j$ ,  $\partial(a_i^T r) / \partial a_j^T = -\eta_j a_i^T$ , while  $\partial(a_j^T r) / \partial a_j^T = r^T - \eta_j a_j^T$ , from which we see that, with  $\text{vec}(A)^T = (a_1^T, \dots, a_n^T)$  and

$$\begin{aligned} \mathcal{J}_A f(A, y) &\equiv \partial f(A, y) / \partial \text{vec}(A)^T = [\partial f(A, y) / \partial a_1^T, \dots, \partial f(A, y) / \partial a_n^T], \\ \partial f(A, y) / \partial a_j^T &= \partial(r^T r y + y^T y A^T r) / \partial a_j^T = -2\eta_j y r^T + y^T y e_j r^T - \eta_j y^T y A^T, \\ \mathcal{J}_A f(A, y) \cdot [\mathcal{J}_A f(A, y)]^T &= \sum_{j=1}^n [\partial f(A, y) / \partial a_j^T] \cdot [\partial f(A, y) / \partial a_j^T]^T \\ &= \sum_{j=1}^n (-2\eta_j y r^T + y^T y e_j r^T - \eta_j y^T y A^T) (-2\eta_j y r^T + y^T y e_j r^T - \eta_j y^T y A^T)^T \\ &= \|y\|_2^6 [A^T A + A^T r y^\dagger + (y^\dagger)^T r^T A + (\|r\|_2^2 / \|y\|_2^2) I] \\ (4.4) \quad &= \|y\|_2^6 \left[ \begin{array}{c} A + r y^\dagger \\ (\|r\|_2 / \|y\|_2)(I - y y^\dagger) \end{array} \right]^T \left[ \begin{array}{c} A + r y^\dagger \\ (\|r\|_2 / \|y\|_2)(I - y y^\dagger) \end{array} \right]. \end{aligned}$$

Here the matrix  $\left[ \begin{array}{c} A + r y^\dagger \\ (\|r\|_2 / \|y\|_2)(I - y y^\dagger) \end{array} \right]$  has full column rank. In fact if it does not then there exists nonzero  $x \in \mathbb{R}^n$  such that

$$(A + r y^\dagger)x = 0, \quad (I - y y^\dagger)x = 0,$$

and it follows that  $Ay + r = 0$ , so  $b = 0$ , contradicting our assumption (2.2). Therefore  $\mathcal{J}_A f(A, y)$  has full row rank. Then from (4.2),

$$\begin{aligned} \tilde{\mu}_F(y) &= \|[\mathcal{J}_A f(A, y)]^T \{ \mathcal{J}_A f(A, y) \cdot [\mathcal{J}_A f(A, y)]^T \}^{-1} f(A, y)\|_2 \\ (4.5) \quad &= \| \{ \mathcal{J}_A f(A, y) \cdot [\mathcal{J}_A f(A, y)]^T \}^{-1/2} f(A, y)\|_2, \end{aligned}$$

where the second equality can easily be proved by using the SVD of  $\mathcal{J}_A f(A, y)$ .

Define

$$(4.6) \quad B \equiv \left[ \begin{array}{c} A + r y^\dagger \\ (\|r\|_2 / \|y\|_2)(I - y y^\dagger) \end{array} \right], \quad c \equiv \begin{bmatrix} r \\ 0 \end{bmatrix} \in \mathbb{R}^{m+n}.$$

Note that

$$(4.7) \quad f(A, y) = \|y\|_2^2(A^T r + \|r\|_2^2(y^\dagger)^T) = \|y\|_2^2 B^T c.$$

Then from (4.5) with (4.4) and (4.7), it follows that

$$(4.8) \quad \tilde{\mu}_F(y) = \frac{\|(B^T B)^{-1/2} B^T c\|_2}{\|y\|_2} = \frac{[c^T B (B^T B)^{-1} B^T c]^{1/2}}{\|y\|_2} = \frac{\|B (B^T B)^{-1} B^T c\|_2}{\|y\|_2}.$$

Note that  $B(B^T B)^{-1} B^T$  is an orthogonal projector onto  $\mathcal{R}(B)$ .

The asymptotic estimate  $\tilde{\mu}_F(y)$  is analogous to an estimate for the minimal backward error for ordinary least squares problems whose various forms have been studied in [9], [11], [12], and [17]. One method for computing  $\tilde{\mu}_F(y)$  is to use the QR factorization. If  $B = QR$  where  $Q \in \mathbb{R}^{(m+n) \times n}$  satisfies  $Q^T Q = I_n$ , and  $R$  is upper triangular, then we see that

$$(4.9) \quad \tilde{\mu}_F(y) = \|Q^T c\|_2 / \|y\|_2.$$

If we use Householder QR factorization, this method will cost  $2(m + 2/3n^2)n^2$  flops. The other method is to use the moment method by following [27, Part I]. For brevity, we will not give details here.

**5. Numerical tests.** In section 2 we gave an extended minimal backward error  $\mu_F(y)$ , which is a lower bound on the minimal backward error. But if the inequality in (2.31) holds for the extended minimal backward perturbation  $\widehat{\Delta A}$  given in (2.49), then  $\mu_F(y)$  is in fact the minimal backward error. Our numerical tests indicate that if the given vector  $y$  is a reasonable approximation to the true DLS solution, the inequality in (2.31) holds, where  $\Delta A$  is the minimal  $\widehat{\Delta A}$  given in (2.16). We will give some examples in this section to illustrate this. In sections 3 we gave a lower bound  $\mu_2^{\text{lb}}(y)$  on  $\mu_2(y)$ , which is also a lower bound on  $\mu_F(y)$  (since  $\mu_2(y) \leq \mu_F(y)$ ). In section 4 we presented an asymptotic estimate  $\tilde{\mu}_F(y)$  of  $\mu_F(y)$ . We will give numerical examples to show how good  $\mu_2^{\text{lb}}(y)$  and  $\tilde{\mu}_F(y)$  are as approximations to  $\mu_F(y)$ . We carried out computations using MATLAB 7.4 on a MacBook running Mac OS X 10.4.11.

In our numerical tests the data was constructed as follows (`randn` and `rand` are two MATLAB built-in functions for generating random matrices with normal and uniform distributions, respectively):

- We use two types of test matrix for  $A$ :
  - Type 1:  $A = \tilde{A} / \|\tilde{A}\|_F$ ,  $\tilde{A} = \text{randn}(100, 40)$ . Typically  $\kappa_2(A) \leq 10$ .
  - Type 2:  $A = \tilde{A} / \|\tilde{A}\|_F$ ,  $\tilde{A} = U \Sigma V^T$ ,  $40 \times 40$   $\Sigma = \text{diag}(\sigma_i)$ ,  $\sigma_i = 10^{-4(i-1)/39}$ ,  $U \in \mathbb{R}^{100 \times 40}$  and  $V \in \mathbb{R}^{40 \times 40}$  are the Q-factors of the QR factorizations of two random matrices `randn(100, 40)` and `randn(40, 40)`, respectively. Note that  $\kappa_2(A) = 10^4$ .
- $b = (A + E)x$ ,  $x = [1, \dots, 1]^T \in \mathbb{R}^{40}$ ,  $E = \frac{\delta_A}{\sqrt{100 \times 40}} \text{rand}(100, 40)$  (note that  $\|E\|_F \leq \delta_A$ ),  $\delta_A = 10^{-7}, 10^{-6}, \dots, 10^{-1}$  for type 1 matrices  $A$ ,  $\delta_A = 10^{-7}, 10^{-6}, 10^{-5}, 10^{-4}$  for type 2 matrices  $A$ . The DLS estimate usually has no accurate digits compared with  $x$  if  $\delta_A$  is taken to be larger.
- $y = \hat{x} + \frac{1}{\sqrt{40}} \delta_{\hat{x}} \|\hat{x}\|_2 \text{rand}(40, 1)$ ,  $\hat{x}$  is the computed solution to the DLS problem (2.1) and  $\delta_{\hat{x}} = 0, 10^{-7}, 10^{-6}, \dots, 10^{-1}$ .
- For each pair of  $\delta_A$  and  $\delta_{\hat{x}}$  and each type of matrix, we generated 1000 sample problems.

The solution  $\hat{x}$  to the DLS problem satisfies (see [5]):

$$(5.1) \quad \hat{x} = \frac{b^T b}{b^T A v_D} v_D,$$

where  $v_D$  is the right singular vector corresponding to the smallest singular value of  $(I - bb^\dagger)A$ . The equality (5.1) can also be obtained from [21, Sec. 9], which suggests a way to compute  $\hat{x}$ . In our numerical tests we used the MATLAB built-in function `svd` to find  $v_D$  and then computed  $\hat{x}$ . To compute the asymptotic estimate  $\tilde{\mu}_F(y)$ , we first computed the QR factorization of  $B$  (see (4.6)) to find the Q-factor and then used (4.9).

In our numerical tests single precision was used to generate the data  $A$ ,  $b$  and  $y$ , and to compute the DLS solution  $\hat{x}$ ; and both single precision and double precision were used to compute both sides of the inequality in (2.31) (where  $\Delta A$  gives the minimum). The number of failures to satisfy the inequality for each case by single (S) and double (D) precision is reported in Table 5.1 for type 1 matrices, and in Table 5.2 for type 2 matrices. When  $\delta_A$  is small or  $\delta_{\hat{x}}$  is large, we see that the computed version of the inequality by single precision sometimes fails. In particular for ill-conditioned type 2 matrices, when  $\delta_A = 10^{-7}$  or  $\delta_{\hat{x}} = 10^{-1}$ , the failure percentage is very large. However the computed version of the inequality by double precision *always* holds for these test cases. This shows that these failures were due to rounding errors in the single precision computed version of the inequality, and for these test cases the extended minimal backward error is actually the true minimal backward error. The reason single precision rounding errors caused some tests to fail is almost certainly the following: in each failed case the gap between the smallest and the second smallest singular values of  $N_\infty$  was small, making the computation of the singular vector  $w_*$  (see (2.49)) inaccurate (see, e.g., [2, Thm 1.2.8] or [8, Thm 8.6.5] for perturbation results concerning the singular vectors). Indeed we noticed that for the failed cases  $w_*$  computed by single precision was very inaccurate compared with the one computed by double precision, leading to a large computational error in  $\widehat{\Delta A}$ , where this is needed for checking the inequality in (2.31).

In Figures 5.1 to 5.8 we give the plots corresponding to eight extreme cases in Tables 5.1 and 5.2 which exhibit  $\mu_F(y)$  (as abscissa) vs  $\mu_2^{1b}(y)$  in (3.2), and  $\mu_F(y)$  (as abscissa) vs  $\tilde{\mu}_F(y)$  in (4.2) and (4.9), represented by the points  $\cdot$  (blue) for  $\mu_2^{1b}(y)$ , and  $*$  (green) for  $\tilde{\mu}_F(y)$ . The diagonal (red) is plotted for reference. In these figures the above quantities were computed by double precision. But we can see no difference between these figures and the corresponding ones obtained by single precision.

From Figures 5.1, 5.2, 5.5 and 5.6, where each  $y$  is the computed DLS solution  $\hat{x}$ , we see that the minimal backward error  $\mu_F(\hat{x}) \approx 10^{-7}$ , which is close to the unit roundoff for single precision, so that each computed  $\hat{x}$  is a backward stable solution. It is interesting to see from Figures 5.3, 5.4, 5.7 and 5.8 that  $\mu_F(y)$  is about one or two orders of magnitude smaller than  $\delta_{\hat{x}}$ . This phenomenon also holds for other test cases.

All these figures and the figures we did not display here indicate that the lower bound  $\mu_2^{1b}(y)$  is a reasonable approximation to the minimal backward error  $\mu_F(y)$  in the sense that these two always had the same order of magnitude, although the case for type 1 matrices is worse than the case for type 2 matrices. We also see that the asymptotic estimate  $\tilde{\mu}_F(y)$  is an excellent approximation to  $\mu_F(y)$ , even when  $y$  is not close to the DLS solution  $\hat{x}$ , see Figures 5.3, 5.4, 5.7, and 5.8, where  $\delta_{\hat{x}} = 10^{-1}$ .



			$\delta_{\hat{x}}$							
			0	$10^{-7}$	$10^{-6}$	$10^{-5}$	$10^{-4}$	$10^{-3}$	$10^{-2}$	$10^{-1}$
$\delta_A$	$10^{-7}$	S	9	2	4	2	1	1	2	1
		D	0	0	0	0	0	0	0	0
	$10^{-6}$	S	4	2	1	2	4	1	4	1
		D	0	0	0	0	0	0	0	0
	$10^{-5}$	S	0	0	0	0	0	4	3	0
		D	0	0	0	0	0	0	0	0
	$10^{-4}$	S	0	0	0	0	0	0	2	1
		D	0	0	0	0	0	0	0	0
	$10^{-3}$	S	0	0	0	0	0	0	0	0
		D	0	0	0	0	0	0	0	0
	$10^{-2}$	S	0	0	0	0	0	0	0	0
		D	0	0	0	0	0	0	0	0
	$10^{-1}$	S	0	0	0	0	0	0	0	0
		D	0	0	0	0	0	0	0	0

TABLE 5.1

Number of failures to satisfy the inequality (2.31) out of 1000 samples for type 1

			$\delta_{\hat{x}}$							
			0	$10^{-7}$	$10^{-6}$	$10^{-5}$	$10^{-4}$	$10^{-3}$	$10^{-2}$	$10^{-1}$
$\delta_A$	$10^{-7}$	S	718	725	726	749	965	997	1000	997
		D	0	0	0	0	0	0	0	0
	$10^{-6}$	S	2	0	1	0	107	970	1000	997
		D	0	0	0	0	0	0	0	0
	$10^{-5}$	S	0	0	0	0	0	120	945	995
		D	0	0	0	0	0	0	0	0
	$10^{-4}$	S	0	0	0	0	0	0	123	888
		D	0	0	0	0	0	0	0	0

TABLE 5.2

Number of failures to satisfy the inequality (2.31) out of 1000 samples for type 2

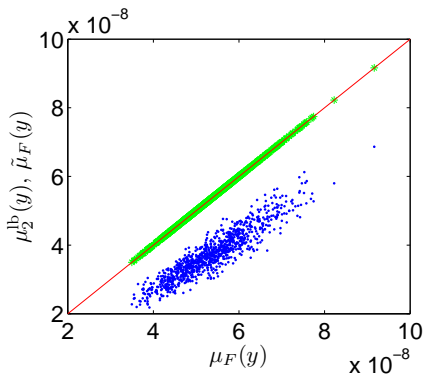


FIG. 5.1. Type 1 A,  $\delta_A = 10^{-7}$ ,  $\delta_{\hat{x}} = 0$

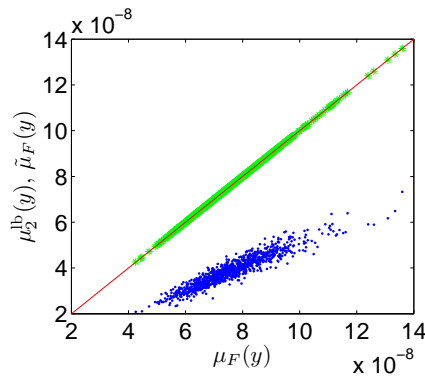
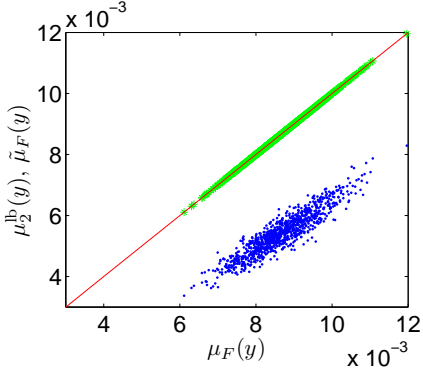
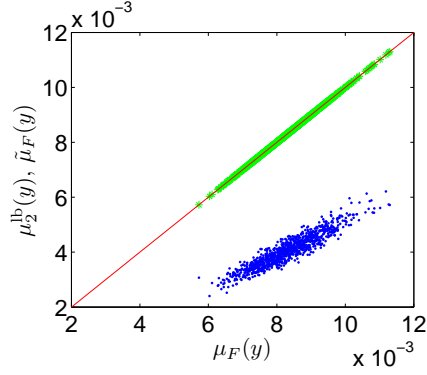
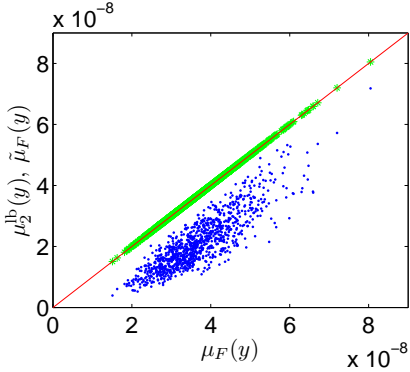
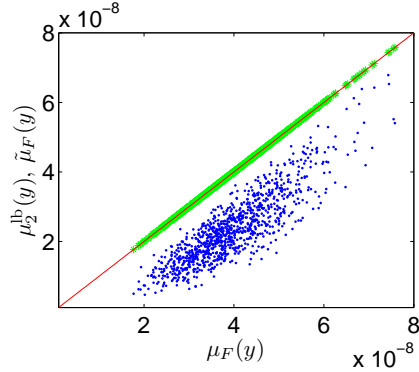
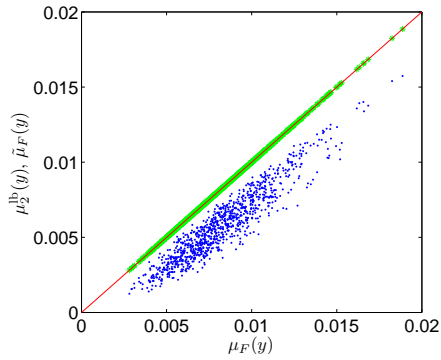
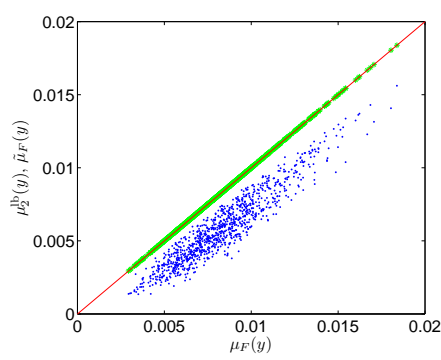


FIG. 5.2. Type 1 A,  $\delta_A = 10^{-1}$ ,  $\delta_{\hat{x}} = 0$

FIG. 5.3. *Type 1 A*,  $\delta_A = 10^{-7}$ ,  $\delta_{\hat{x}} = 10^{-1}$ FIG. 5.4. *Type 1 A*,  $\delta_A = 10^{-1}$ ,  $\delta_{\hat{x}} = 10^{-1}$ FIG. 5.5. *Type 2 A*,  $\delta_A = 10^{-7}$ ,  $\delta_{\hat{x}} = 0$ FIG. 5.6. *Type 2 A*,  $\delta_A = 10^{-4}$ ,  $\delta_{\hat{x}} = 0$ FIG. 5.7. *Type 2 A*,  $\delta_A = 10^{-7}$ ,  $\delta_{\hat{x}} = 10^{-1}$ FIG. 5.8. *Type 2 A*,  $\delta_A = 10^{-4}$ ,  $\delta_{\hat{x}} = 10^{-1}$ 

**6. Summary and future work.** For a given approximate solution  $y$  to the DLS problem (2.1), we first presented formulas (2.15) in Theorem 2.2 for an extended minimal backward error  $\mu_F(y, \theta)$  for the case where backward perturbations in both  $A$  and  $b$  are allowed. Then by taking  $\theta \rightarrow \infty$  we obtained the corresponding formulas

(2.48) in Theorem 2.6 for an extended minimal backward error  $\mu_F(y)$  for the case where only backward perturbations in  $A$  are allowed—this is the case we considered later in the paper. In theory  $\mu_F(y)$  is a lower bound on the minimal backward error, but if the inequality in (2.31) is satisfied for the optimal perturbation  $\widehat{\Delta A}$  given in (2.49), it is in fact the minimal backward error. Our simulations showed that if  $y$  is a reasonable approximation to the exact DLS solution (that is, having a relative error in  $y$  of less than  $10^{-1}$  for our test cases), then, apparently the inequality in (2.31) holds in the absence of rounding errors in checking the inequality. Thus we believe in practice that  $\mu_F(y)$  can usually be used as the minimal backward error. Since the formula (2.48) for  $\mu_F(y)$  involves the minimum singular value of a matrix, it is expensive to compute directly. In order to overcome this problem we derived a lower bound  $\mu_2^{\text{lb}}(y)$  (see (3.2)) and an asymptotic estimate  $\tilde{\mu}_F(y)$  (see (4.6), (4.8) and (4.9)). These can be computed or estimated more efficiently. For our numerical test cases  $\mu_2^{\text{lb}}(y)$  always had the same order of magnitude as  $\mu_F(y)$ , and  $\tilde{\mu}_F(y)$  was an excellent approximation to  $\mu_F(y)$ . Since the computation of  $\mu_2^{\text{lb}}(y)$  is so inexpensive, it would seem to give a simple and effective indicator.

Several problems need to be investigated in the future. To check if the extended minimal backward error is the actual minimal backward error we need an efficient and reliable way to test the inequality in (2.31) (or the inequality in (2.9) when perturbations in both  $A$  and  $b$  are allowed). The relationships between  $\mu_F(y)$  and  $\tilde{\mu}_F(y)$  needs to be studied further. We would also like to incorporate the results obtained in this paper to design effective stopping criteria for iterative algorithms for solving the DLS problem, and extend the results here to total least squares problems (see [7] and [33]) and scaled total least squares problems (see [23] and [21]).

**Acknowledgments.** The first author is indebted to the second author and to Michael Saunders for their hospitality when he was on sabbatical at Stanford University where part of this research was done. He is also grateful to Zheng Su for helpful discussions. We also thank Pete Stewart and Nick Trefethen for improving our understanding of eigensystems of analytic Hermitian matrices. Two referees made many helpful suggestions, some of which were very demanding but led to necessary and substantial expansions of this paper. One referee also suggested the form of Theorem 2.8 in order to give a clearer exposition of our thinking.

#### REFERENCES

- [1] M. ARIOLI, I. DUFF AND D. RUIZ, *Stopping criteria for iterative solvers*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 1381–144.
- [2] Å. BJÖRCK, *Numerical Methods for Least Squares Problems*, SIAM, Philadelphia, PA, 1996.
- [3] X.-W. CHANG, C. C. PAIGE AND D. TITLEY-PELOQUIN, *Characterizing matrices that are consistent with given solutions*. Submitted for publication.
- [4] A.J. COX AND N.J. HIGHAM, *Backward error bounds for constrained least squares problems*, BIT, 39 (1999), pp. 210–227.
- [5] R.D. DEGROAT AND E.M. DOWLING, *The data least squares problem and channel equalization*, IEEE Trans. Signal Processing, 42 (1993), pp. 407–411.
- [6] G.H. GOLUB AND G. MEURANT, *Matrices, moment and quadrature*, in Proceedings of the 15th Dundee Conference, June–July 1993, D. F. Griffiths and G. A. Watson (eds.), Longman Scientific & Technical, 1994.
- [7] G.H. GOLUB AND C.F. VAN LOAN, *An analysis of the total least squares problem*, SIAM J. Numer. Anal., 17:883–893, 1980.
- [8] G.H. GOLUB AND C.F. VAN LOAN, *Matrix Computations*, The Johns Hopkins University Press, Baltimore MD, third ed., 1996.

- [9] J.F. GRGAR, *Optimal sensitivity analysis of linear least squares*, Technical Report LBNL-52434. Lawrence Berkeley National Laboratory, 2003.
- [10] J.F. GRGAR, *The optimum inverse problem of numerical error analysis*, Householder Symposium XVI, Seven Springs, USA, May 23-27, 2005.
- [11] J.F. GRGAR, M.A. SAUNDERS, AND Z. SU, *Estimates of optimal backward perturbations for linear least squares problems*, manuscript, 2004.
- [12] M. GU, *Backward perturbation bounds for linear least squares problems*, SIAM J. Matrix Anal. Appl., 20 (1998), pp. 363–372.
- [13] N.J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, 2nd Edn., SIAM, Philadelphia, PA, 2002.
- [14] D.J. HIGHAM AND N. J. HIGHAM, *Backward error and condition of structured linear systems*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 162–175.
- [15] R.A. HORN AND C.R. JOHNSON, *Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1985.
- [16] T. KATO, *Perturbation Theory for Linear Operators*, 2nd Edn., Springer-Verlag, New York, 1980.
- [17] R. KARLSON AND B. WALDÉN, *Estimation backward perturbation bounds for the linear least squares problem*, BIT, 37 (1997), 862-869.
- [18] A.N. MALYSHEV, *Optimal backward perturbation bounds for the LSS problems*, BIT, 41 (2001), 430–432.
- [19] A.N. MALYSHEV AND M. SADKANE, *Computation of optimal backward perturbation bounds for large sparse linear least squares problems*, BIT, 41 (2002), pp. 739–747.
- [20] C. C. PAIGE AND M. A. SAUNDERS, *LSQR: An algorithm for sparse linear equations and sparse least squares*, ACM Trans. Math. Software, 8:43–71, 1982.
- [21] C.C. PAIGE AND Z. STRAKOŠ, *Scaled total least squares fundamentals*, Numerische Mathematik, 91:117–146 (2002).
- [22] C.C. PAIGE AND Z. STRAKOŠ, *Unifying least squares, total least squares and data least squares*, in “Total Least Squares and Errors-in-Variables Modeling”, S. van Huffel and P. Lemmerling, editors, Kluwer Academic Publishers, Dordrecht, 2002, pp. 25–34.
- [23] B.D. RAO, *Unified treatment of LS, TLS and truncated SVD methods using a weighted TLS framework*. In: S. VAN HUFFEL (EDITOR), *Recent Advances in Total Least Squares Techniques and Errors-in-Variables Modelling*, pp. 11–20, SIAM Publications, Philadelphia PA, 1997.
- [24] F. RELICH, *Perturbation Theory of Eigenvalue Problems*, Gordon and Breach, New York, 1969.
- [25] J.L. RIGAL AND J. GACHES, *On the compatibility of a given solution with the data of a linear system*. JACM, 14:543–548 (1967).
- [26] G.W. STEWART, *On the perturbation of pseudo-inverses, projections, and linear least squares problems*. SIAM Rev., 19:634–662 (1977).
- [27] Z. SU, *Computational Methods for Least Squares Problems and Clinical Trials*. Ph.D Thesis, Scientific Computing & Computational Mathematics, Stanford University, 2005.
- [28] J.-G. SUN, *Optimal backward perturbation bounds for the linear least squares problem with multiple right-hand sides*, IMA J. Numer. Anal., 16 (1996), pp. 1–11.
- [29] J.-G. SUN, *On optimal backward perturbation bounds for the linear least-squares problem*, BIT, 37 (1997), pp. 179–188.
- [30] J.-G. SUN, *Bounds for the structured backward errors of Vandermonde systems*, SIAM J. Matrix Anal. Appl., 20 (1999), pp. 45–59.
- [31] J.-G. SUN, *A note on backward errors for structured linear systems*, Numerical Linear Algebra with Applications, 12(7), 2005, pp 585-603.
- [32] J.-G. SUN AND Z. SUN, *Optimal backward perturbation bounds for underdetermined systems*. SIAM J. Matrix Anal. Appl., 18 (1997), pp. 393–402.
- [33] S. VAN HUFFEL AND J. VANDEWALLE, *The Total Least Squares Problem: Computational Aspects and Analysis*, SIAM Publications, Philadelphia PA, 1991.
- [34] J.M. VARAH, *Backward error estimates for Toeplitz systems*, SIAM J. Matrix Anal. Appl., 15 (1994), pp. 408–417.
- [35] B. WALDÉN, R. KARLSON AND J.-G. SUN, *Optimal backward perturbation bounds for the linear least squares problem*, Numerical Linear Algebra with Applications, 2:271–286 (1995).
- [36] J. H. WILKINSON, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, UK (1965).
- [37] H. XIANG AND Y. WEI *On normwise structured backward errors for saddle point systems*, SIAM J. Matrix Anal. Appl., 29 (2007), pp. 838-849.