

ON THE SENSITIVITY OF SINGULAR AND ILL-CONDITIONED
LINEAR SYSTEMS*

ZHONGGANG ZENG†

Abstract. Solving a singular linear system for an individual vector solution is an ill-posed problem with a condition number infinity. From an alternative perspective, however, the general solution of a singular system is of a bounded sensitivity as a unique element in an affine Grassmannian. If a singular linear system is given through empirical data that are sufficiently accurate with a tight error bound, a properly formulated general numerical solution uniquely exists in the same affine Grassmannian, enjoys Lipschitz continuity, and approximates the underlying exact solution with an accuracy in the same order as the data. Furthermore, any backward accurate numerical solution vector is an accurate approximation to one of the solutions of the underlying singular system.

Key words. condition number, linear system, Grassmannian

AMS subject classifications. 65F22, 65F35, 15A12, 15A06

DOI. 10.1137/18M1197990

1. Introduction. Solving linear systems in the matrix-vector form $A\mathbf{x} = \mathbf{b}$ is one of the most fundamental problems in scientific computing. In the literature of numerical analysis, linear systems are always assumed to be nonsingular with few exceptions. Numerical solutions of singular systems are almost never mentioned directly in textbooks. The following rare remark in Meyer's textbook [25, pp. 217–218] accurately reflects the current state of knowledge:

If \mathbf{A} is singular, . . . even a stable algorithm can result in a significant loss of information. . . . the small perturbation \mathbf{E} due to roundoff makes the possibility that $\text{rank}(\mathbf{A} + \mathbf{E}) > \text{rank}(\mathbf{A})$ very likely. *The moral is to avoid floating-point solutions of singular systems*" [emphasis added].

In applications such as deblurring images and discrete inverse problems, rank-deficient and highly ill-conditioned linear systems are approached using the Tikhonov regularization [11, 12, 13, 27]. As Neumaier states [27, p. 637]:

Though frequently needed in applications, the adequate handling of such ill-posed linear problems is hardly ever touched upon in numerical analysis text books.

Singular linear systems are unavoidable in scientific computing and often need to be solved without knowledge of the exact matrices and vectors, as shown in our case studies in section 3. The obvious difficulty in solving a singular linear system from empirical data is the condition number infinity, and thus the error is unbounded when solving for an individual vector solution. While this error analysis is in itself

*Received by the editors July 2, 2018; accepted for publication (in revised form) by J. L. Barlow April 24, 2019; published electronically August 1, 2019. This work was performed by an employee of the U.S. Government or under U.S. Government contract. The U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or allow others to do so, for U.S. Government purposes. Copyright is owned by SIAM to the extent not limited by these rights.

<https://doi.org/10.1137/18M1197990>

Funding: This work was partially supported by National Science Foundation grant DMS-1620337.

†Department of Mathematics, Northeastern Illinois University, Chicago, IL 60625 (zzeng@neiu.edu).

impeccable, the solution of a singular system is more than an individual vector. The very notion of the numerical solution to a given system $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ needs clarification when entries of $(\tilde{A}, \tilde{\mathbf{b}})$ serve as empirical data for an underlying singular linear system $A\mathbf{x} = \mathbf{b}$.

This paper attempts to analyze the accuracy and sensitivity of solving singular linear systems from the following different perspective: The solution of a singular linear system is either an empty set or an affine subspace as a unique element in an affine Grassmannian rather than a vector. Using this point of view, the condition number becomes bounded. A properly formulated general numerical solution in a certain affine Grassmannian is of a sensitivity proportional to $\|A\|_2 \|A^\dagger\|_2$, never infinity, with respect to either constrained or arbitrary perturbations, where A^\dagger is the Moore–Penrose inverse of A . Such a numerical solution of a perturbed system $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ within a viable error tolerance accurately solves the underlying singular system $A\mathbf{x} = \mathbf{b}$, and the ratio of solution accuracy to the data error is bounded by a factor of $\|A\|_2 \|A^\dagger\|_2$, not $\|\tilde{A}\|_2 \|\tilde{A}^{-1}\|_2$, assuming the data error is small with an attainable tight bound.

We shall further demonstrate that the sensitivity of a singular linear system $A\mathbf{x} = \mathbf{b}$ is measured by $\|A\|_2 \|A^\dagger\|_2$ rather than infinity from multiple perspectives, including homogeneous cases, under constrained perturbations preserving the singularity and consistency, solving for the general numerical solutions in an affine Grassmannian, and solving for a single particular solution. Furthermore, every backward accurate numerical (vector) solution of a singular consistent linear system accurately approximates a particular exact solution regardless of the algorithm used. The “error” largely falls harmlessly in the kernel of A . This result extends what Peters and Wilkinson discovered in [29] beyond inverse power iterations. While any numerical (single-vector) solution may be inaccurate to a linear system that is genuinely nonsingular and highly ill-conditioned, we shall prove that a stable numerical (affine subspace) solution may exist and contain an accurate approximation to the exact solution. For practical computation, efficient and robust algorithms already exist for general numerical solutions in affine Grassmannians. Regularization algorithms, such as the Tikhonov method and truncated singular value decomposition (TSVD) [9, sect. 5.5.4], [10], produce the accurate vector component, and numerical rank-revealing algorithms [2, 7], [9, sect. 5.4.6], [18, 19, 20, 31] provide the numerical kernel as the remaining component.

For continuity of presentation, lemmas and long proofs are listed in the appendices. Additional computational results and software demonstration are given in the supplementary material.

2. Preliminaries. Column vectors are denoted by boldface lowercase letters such as \mathbf{b} , \mathbf{x} , \mathbf{y} , etc., with $\mathbf{0}$ being a zero vector whose dimension can be derived from the context. The vector space of n -dimensional complex column vectors is denoted by \mathbb{C}^n . The vector space of $m \times n$ matrices with complex entries is denoted by $\mathbb{C}^{m \times n}$. Matrices are denoted by uppercase letters such as A , B , X , etc., while O and I denote a zero matrix and an identity matrix, respectively. The range, kernel, rank, and Hermitian transpose of a matrix A are denoted by $\text{Range}(A)$, $\text{Kernel}(A)$, $\text{rank}(A)$, and A^H , respectively. In this paper, we consider general $m \times n$ linear systems in the form of $A\mathbf{x} = \mathbf{b}$, and we say the system is *singular* when $\text{rank}(A) < n$ so that $\text{Kernel}(A) \neq \{\mathbf{0}\}$, including nonsquare cases where $m < n$ or $m > n$. The system is *consistent* if $\mathbf{b} \in \text{Range}(A)$.

For any matrix $A \in \mathbb{C}^{m \times n}$, the j th largest singular value of a matrix A is denoted

by $\sigma_j(A)$. The *numerical rank* of a matrix A within an error tolerance $\theta > 0$ is defined as

$$\text{rank}_\theta(A) := \min_{\|B-A\|_2 < \theta} \text{rank}(B) \equiv \max_{\sigma_j(A) > \theta} j,$$

assuming θ is not equal to any singular value of A . Let $U \Sigma V^H$ be the SVD of A , where $U = [\mathbf{u}_1, \dots, \mathbf{u}_m]$ and $V = [\mathbf{v}_1, \dots, \mathbf{v}_n]$. If $\text{rank}_\theta(A) = r$ within θ , then the θ -projection A_θ of A is defined as

$$A_\theta := \sigma_1(A) \mathbf{u}_1 \mathbf{v}_1^H + \dots + \sigma_r(A) \mathbf{u}_r \mathbf{v}_r^H = \sum_{\sigma_j(A) > \theta} \sigma_j(A) \mathbf{u}_j \mathbf{v}_j^H.$$

In this case, the *numerical kernel* of A within θ is $\text{Kernel}(A_\theta) = \text{span}\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$, where $\text{span}\{\dots\}$ denotes the vector space spanned by vectors in the list. The entities $\text{rank}_\theta(A)$, A_θ , and $\text{Kernel}(A_\theta)$ are undefined if θ is a singular value of A . The *Moore–Penrose inverse* of A , denoted by A^\dagger , is the unique matrix satisfying the Moore–Penrose conditions $A A^\dagger A = A$, $A^\dagger A A^\dagger = A^\dagger$, $(A A^\dagger)^H = A A^\dagger$, and $(A^\dagger A)^H = A^\dagger A$. Using the singular value decomposition as above and assuming $\text{rank}(A) = r$, we see that the identity [9, sect. 5.5.2]

$$A^\dagger \equiv \frac{1}{\sigma_1(A)} \mathbf{v}_1 \mathbf{u}_1^H + \dots + \frac{1}{\sigma_r(A)} \mathbf{v}_r \mathbf{u}_r^H = \sum_{\sigma_j(A) > 0} \frac{1}{\sigma_j(A)} \mathbf{v}_j \mathbf{u}_j^H$$

holds, and $X = A^\dagger$ is the minimum Frobenius norm matrix such that $A X$ and $X A$ are orthogonal projections from \mathbb{C}^m onto $\text{Range}(A)$ and from \mathbb{C}^n onto $\text{Range}(A^H)$, respectively. We shall frequently use $\|A^\dagger\|_2^{-1}$ as an alternative notation for the smallest positive singular value $\sigma_r(A)$ of A with rank r .

The set of k -dimensional subspaces of \mathbb{C}^n is called the *Grassmannian* [6], [17, p. 52] of index k of \mathbb{C}^n denoted by $\mathcal{G}_k(\mathbb{C}^n)$. For any $\mathcal{P}, \mathcal{Q} \in \mathcal{G}_k(\mathbb{C}^n)$, let $P, Q \in \mathbb{C}^{n \times k}$ be matrices whose columns form orthonormal bases for \mathcal{P} , \mathcal{Q} , respectively, while $\hat{P}, \hat{Q} \in \mathbb{C}^{n \times (n-k)}$ such that $[P, \hat{P}]^H [P, \hat{P}] = [Q, \hat{Q}]^H [Q, \hat{Q}] = I$. The Grassmannian $\mathcal{G}_k(\mathbb{C}^n)$ is a metric space with the distance [9, sect. 2.5.3]

$$\text{dist}(\mathcal{P}, \mathcal{Q}) := \|P P^H - Q Q^H\|_2 \equiv \|P^H \hat{Q}\|_2 \equiv \|Q^H \hat{P}\|_2.$$

The set of k -dimensional affine subspaces of \mathbb{C}^n is called the *affine Grassmannian* [16, sect. 7.1], [21, 22] of index k of \mathbb{C}^n and is denoted by

$$\mathcal{A}_k(\mathbb{C}^n) := \{\mathbf{u} + \mathcal{V} \subset \mathbb{C}^n \mid \mathbf{u} \in \mathbb{C}^n, \mathcal{V} \in \mathcal{G}_k(\mathbb{C}^n)\}.$$

Here, for any vector $\mathbf{u} \in \mathbb{C}^n$ and subspace $\mathcal{V} \in \mathcal{G}_k(\mathbb{C}^n)$, the *affine subspace*

$$\mathbf{u} + \mathcal{V} := \{\mathbf{u} + \mathbf{v} \in \mathbb{C}^n \mid \mathbf{v} \in \mathcal{V}\}$$

can be written as $\hat{\mathbf{u}} + \mathcal{V}$ with a unique $\hat{\mathbf{u}} \in \mathcal{V}^\perp \cap (\mathbf{u} + \mathcal{V})$ of the minimum norm, where $(\cdot)^\perp$ denotes the unitary complement of any subspace (\cdot) . The metric

$$(1) \quad \begin{aligned} \text{dist}(\mathbf{u}_1 + \mathcal{V}_1, \mathbf{u}_2 + \mathcal{V}_2) \\ := \max_{\hat{\mathbf{u}}_j \in \mathcal{V}_j^\perp \cap (\mathbf{u}_j + \mathcal{V}_j), j=1,2} \{\|\hat{\mathbf{u}}_1 - \hat{\mathbf{u}}_2\|_2, \text{dist}(\mathcal{V}_1, \mathcal{V}_2)\} \end{aligned}$$

for every $\mathbf{u}_1 + \mathcal{V}_1, \mathbf{u}_2 + \mathcal{V}_2 \in \mathcal{A}_k(\mathbb{C}^n)$ is a distance in $\mathcal{A}_k(\mathbb{C}^n)$. For every $(A, \mathbf{b}) \in \mathbb{C}^{m \times n} \times \mathbb{C}^m$, denote the set of vector solutions to the system $A \mathbf{x} = \mathbf{b}$ by

$$\text{sol}(A, \mathbf{b}) := \{\mathbf{u} \in \mathbb{C}^n \mid A \mathbf{u} = \mathbf{b}\}.$$

For $r = \text{rank}(A)$, the set $\text{sol}(A, \mathbf{b})$ as the solution of $A\mathbf{x} = \mathbf{b}$ uniquely exists as either \emptyset or an element in the affine Grassmannian $\mathcal{A}_{n-r}(\mathbb{C}^n)$. The dimension of $\text{sol}(A, \mathbf{b})$ is either $n-r$ if it is in $\mathcal{A}_{n-r}(\mathbb{C}^n)$ or -1 if it is empty [3, p. 6]. We define $\text{dist}(\emptyset, \emptyset) = 0$ so that the deviation of solutions can be measured if and only if they are of the same dimension.

The *condition number* of a square matrix A in the context of solving a linear system $A\mathbf{x} = \mathbf{b}$ is well known to be $\kappa(A) = \|A\|_2 \|A^{-1}\|_2$, with a convention $\kappa(A) = \infty$ when A is singular [9, p. 87]. This condition number is based on the attainable error of the solution as an individual vector. The infinity convention can be justified by $\limsup_{G \rightarrow A} \|G\|_2 \|G^\dagger\|_2 = \infty$ when A is singular, and by the interpretation of the condition number as the reciprocal of the distance to the singularity [14, Theorem 6.5]. For a rectangular matrix A , it is natural to generalize the condition number as $\kappa(A) = \|A\|_2 \|A^\dagger\|_2$ (see, e.g., [14, p. 382]). We present arguments from multiple perspectives that the infinity convention may be unnecessary even if A is square and singular.

It is easy to see that $\kappa(A) = \|A\|_2 \|A^\dagger\|_2$ is discontinuous at any rank deficient matrix A and cannot be approximated from empirical data \tilde{A} since $\kappa(\tilde{A}) = \|\tilde{A}\|_2 \|\tilde{A}^\dagger\|_2$ can be arbitrarily large when $\|\Delta A\|_2 = \|\tilde{A} - A\|_2$ is small. For any error tolerance θ with $0 < \theta < \|A^\dagger\|_2^{-1}$, however, the asymptotic bound

$$\kappa(A) - 2\|\Delta A\|_2 + O(\|\Delta A\|_2^2) \leq \kappa(\tilde{A}_\theta) \leq \kappa(A) + 2\|\Delta A\|_2 + O(\|\Delta A\|_2^2)$$

follows from [9, Corollary 8.6.2] when the data matrix \tilde{A} is sufficiently accurate so that $\|\Delta A\|_2 < \|A^\dagger\|_2^{-1} - \theta$. Assuming an error bound $\beta > \|\Delta A\|_2$ is attainable and is sufficiently tight so that $\beta < \|A^\dagger\|_2^{-1} - \|\Delta A\|_2$, the condition number $\kappa(\tilde{A}_\theta)$ of the θ -projection \tilde{A}_θ of the data matrix \tilde{A} , not $\kappa(\tilde{A})$, is an approximation to the underlying condition number $\kappa(A) = \|A\|_2 \|A^\dagger\|_2 < \infty$.

3. Models of singular linear systems. We shall elaborate some case studies to show that solving singular linear systems is not only unavoidable in scientific computing but also crucial in many applications. It may even be beneficial for the systems to be singular. Moreover, singular linear systems are often not known with exact matrices and right-hand side vectors in practical computation, and need to be solved from empirical data.

Example 1 (multiplicity of a singular solution to a nonlinear system). For a system of nonlinear equations in the form of $\mathbf{f}(\mathbf{x}) = \mathbf{0}$, where $\mathbf{f} = (f_1, \dots, f_m)$ and $f_j : \mathbb{C}^n \rightarrow \mathbb{C}$ is an analytic function for $j = 1, \dots, m$, a zero \mathbf{x}_* of \mathbf{f} is multiple if the Jacobian of \mathbf{f} at \mathbf{x}_* is rank-deficient. At such a multiple \mathbf{x}_* , there is a vector space called the dual space $\mathcal{D}_{\mathbf{f}, \mathbf{x}_*}$ that forms the multiplicity structure of the zero and the dimension of $\mathcal{D}_{\mathbf{f}, \mathbf{x}_*}$ is the multiplicity. The multiplicity structure can be determined by solving a sequence of homogeneous linear systems

$$(2) \quad S_\alpha(\mathbf{x}_*) \mathbf{c} = \mathbf{0} \quad \text{for } \alpha = 1, 2, \dots,$$

where $S_\alpha(\mathbf{x}_*)$ is the Macaulay matrix whose entries are derivatives of f_j 's of orders up to α evaluated at \mathbf{x}_* . The solution $\text{sol}(S_\alpha(\mathbf{x}_*), \mathbf{0})$ of (2) in a proper Grassmannian is isomorphic to the desired dual space $\mathcal{D}_{\mathbf{f}, \mathbf{x}_*}$ when α reaches the so-called depth. See, e.g., [4] for detailed elaborations and the supplementary material for a computing demo. The exact Macaulay matrix $S_\alpha(\mathbf{x}_*)$ is almost never available since \mathbf{x}_* is generally known approximately through a certain $\tilde{\mathbf{x}} \approx \mathbf{x}_*$ within an error bound. The model is to solve the singular system (2) for the solution in a Grassmannian rather than individual vectors from empirical data matrix $S_\alpha(\tilde{\mathbf{x}}) \approx S_\alpha(\mathbf{x}_*)$.

Example 2 (Sylvester equation). This is an application arising in control and system theory [1] in the form of the Sylvester matrix equation $A(t)X + XB(t) = C(t)$, where $A(t)$, $B(t)$, and $C(t)$ are matrices depending on a parameter t . The system may inevitably become singular when the parameter t varies continuously and passes through a certain t_* whose value may only be obtained approximately. The following illustrative example is slightly modified from [1] (cf. the supplementary material). Let

$$(3) \quad A(t) = \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}, \quad B(t) = \begin{bmatrix} -\frac{5}{3} + t & 1 \\ -1 & -\frac{1}{3} + 2t \end{bmatrix}, \quad \text{and} \quad C(t) = \begin{bmatrix} 1 & 0 \\ 2 & -1 \end{bmatrix}.$$

When t varies continuously, the system becomes singular but still is consistent when t hits the value $t_* = \frac{2}{3}$ with the general solution

$$(4) \quad X_* = \frac{1}{4} \begin{bmatrix} 1 & -1 \\ -3 & -1 \end{bmatrix} + \alpha_1 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \alpha_2 \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix}, \quad \alpha_1, \alpha_2 \in \mathbb{C}.$$

Suppose we know that $\tilde{t} \approx 0.6666$ with an error bound 0.0001. Can we find a numerical solution \tilde{X} of the perturbed system at the parameter value $t = \tilde{t}$ approximating X_* in (4) of the underlying system at $t = t_*$ with an accuracy $\|\tilde{X} - X_*\|_2$ roughly 0.0001?

Example 3 (Bézout coefficients). For polynomials f_1, \dots, f_n , with a greatest common divisor g , there exist polynomials u_1, \dots, u_n , known as the Bézout coefficients (see, e.g., [26, sect. 1.3], [36]) such that the Bézout identity

$$(5) \quad u_1 f_1 + \cdots + u_n f_n = g$$

holds. Solving the linear equation (5) for the Bézout coefficients appears in many applications, such as computing the Smith normal form in linear control theory [36], and the systems are often singular for $n \geq 3$. Denote by \mathbb{P}_k the vector space of polynomials with degrees up to k . For instance, let f_1, f_2, f_3 be polynomials of degrees, say 4, 7, 6, with degree of g , say 2. Equation (5) for $(u_1, u_2, u_3) \in \mathbb{P}_3 \times \mathbb{P}_1 \times \mathbb{P}_2$ is consistent and rank-deficient by 2. The rank-deficiency is, in fact, a blessing because it transforms the general solution

$$(u_1, u_2, u_3) = (u_{01}, u_{02}, u_{03}) + t_1 (u_{11}, u_{12}, u_{13}) + t_2 (u_{21}, u_{22}, u_{23})$$

into an invertible transformation

$$(6) \quad \begin{bmatrix} u_{01} & u_{02} & u_{03} \\ u_{11} & u_{12} & u_{13} \\ u_{21} & u_{22} & u_{23} \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} = \begin{bmatrix} g \\ 0 \\ 0 \end{bmatrix}.$$

The exact coefficients of polynomial parameters f_1, \dots, f_n and g may be unknown beyond their empirical data, say,

$$\begin{aligned} \tilde{f}_1 &= 2.5714 + 3.8571x - 3x^2 - 6.4286x^3 - 2.1429x^4, \\ \tilde{f}_2 &= -1.7143 - 1.7143x + 0.4286x^2 + 0.4286x^3 - 3.4286x^5 - 5.1429x^6 - 1.7143x^7, \\ \tilde{f}_3 &= 0.8571 + 1.2857x + 2.1429x^2 + 2.5714x^3 + 3.4286x^4 + 3.8571x^5 + 1.2857x^6, \\ \tilde{g} &= 4.6667 + 7x + 2.3333x^2, \end{aligned}$$

with coefficientwise error bound $\varepsilon = 0.5 \times 10^{-4}$. Can we accurately calculate the general solution for $(u_1, u_2, u_3) \in \mathbb{P}_3 \times \mathbb{P}_1 \times \mathbb{P}_2$ of (5) using the imperfect data $\tilde{f}_1, \tilde{f}_2, \tilde{f}_3$, and \tilde{g} within an error in the same order of the data? A computation/software demo

for this example is given in the supplementary material. The matrix-vector representation $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ of (5) in the given data is

$$(7) \quad \begin{aligned} & \left[\begin{array}{ccccccccc} 2.5714 & 0 & 0 & 0 & -1.7143 & 0 & 0.8571 & 0 & 0 \\ 3.8571 & 2.5714 & 0 & 0 & -1.7143 & -1.7143 & 1.2857 & 0.8571 & 0 \\ -3.0000 & 3.8571 & 2.5714 & 0 & 0.4286 & -1.7143 & 2.1429 & 1.2857 & 0.8571 \\ -6.4286 & -3.0000 & 3.8571 & 2.5714 & 0.4286 & 0.4286 & 2.5714 & 2.1429 & 1.2857 \\ -2.1429 & -6.4286 & -3.0000 & 3.8571 & 0 & 0.4286 & 3.4286 & 2.5714 & 2.1429 \\ 0 & -2.1429 & -6.4286 & -3.0000 & -3.4286 & 0 & 3.8571 & 3.4286 & 2.5714 \\ 0 & 0 & -2.1429 & -6.4286 & -5.1429 & -3.4286 & 1.2857 & 3.8571 & 3.4286 \\ 0 & 0 & 0 & -2.1429 & -1.7143 & -5.1429 & 0 & 1.2857 & 3.8571 \\ 0 & 0 & 0 & 0 & 0 & -1.7143 & 0 & 0 & 1.2857 \end{array} \right] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \end{bmatrix} \\ & = [4.6667 \quad 7.0000 \quad 2.3333 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0]^T, \end{aligned}$$

with respect to monomial bases, and the condition number $\kappa(\tilde{A}) \gtrsim 2.29 \times 10^6$. The system (7) in the conventional sense is highly ill-conditioned since $\varepsilon \kappa(\tilde{A}) > 1$.

Applications are abundant involving singular linear systems. The output regulation problem arises in the application of neural networks [23] for finding the matrix pair (X, U) satisfying the so-called regulator equations, whose solutions are not necessarily unique. An illustrative example is as follows (cf. the supplementary material):

$$(8) \quad \begin{cases} X \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 2 & -1 & 0 \end{bmatrix} X + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} U + \begin{bmatrix} -2 & 1 \\ -1 & 1 \\ 0 & 0 \end{bmatrix} \\ [0 \ 0] = [1 \ 0 \ -1] X + [-1 \ 0], \end{cases}$$

where the unknowns X and U are matrices. The system is rank-deficient by one. Furthermore, the matrix parameters are not known exactly but given by estimation. For a matrix A with a defective eigenvalue λ_* and an associated eigenvector \mathbf{z}_* , a generalized eigenvector satisfies the singular system $(A - \lambda_* I)\mathbf{x} = \mathbf{z}_*$ for $\mathbf{x} \in \mathbb{C}^n$. The value of λ_* and \mathbf{z}_* generally can only be known approximately. The problem is to solve the underlying system by solving $(\tilde{A} - \tilde{\lambda} I)\mathbf{x} = \tilde{\mathbf{z}}$ from the data $\tilde{A} \approx A$, $\tilde{\lambda} \approx \lambda_*$, and $\tilde{\mathbf{z}} \approx \mathbf{z}_*$. More applications include solving the singular homogeneous linear systems of Ruppert matrices in numerical factorization of polynomials [8, 37], numerical elimination of polynomial variables [38], etc. The generalized Lyapunov equation $E^H A X + A^H X E = -G$ with given matrices A, E , and G is singular when E is rank-deficient [33]. A singular linear system that models the atmospheric path delay and the water vapor constant estimation is given in [30]. Linear systems derived from discretizing the Fredholm and Volterra integral equations can be considered empirical data of singular systems in the presence of annihilators [12, sect. 2.4 and p. 83] (cf. an example in the supplementary material).

4. Homogeneous systems with empirical data. A problem is *well-posed* if its solution satisfies existence, uniqueness, and Lipschitz continuity with respect to the data or, otherwise, it is an *ill-posed problem*. For an $m \times n$ singular homogeneous linear system $A\mathbf{x} = \mathbf{0}$, the problem

$$(9) \quad \text{Solve } A\mathbf{x} = \mathbf{0} \text{ for a single-vector solution } \mathbf{x} \text{ in } \mathbb{C}^n$$

is obviously ill-posed, as its solutions are not unique. However, the problem (9) is not precisely the problem to be solved in standard linear algebra, where all the solutions are in question. There is a unique solution to the problem

$$(10) \quad \text{Solve } A\mathbf{x} = \mathbf{0} \text{ for the solution } \text{sol}(A, \mathbf{0}) \text{ in the Grassmannian } \mathcal{G}_{n-r}(\mathbb{C}^n),$$

where $r = \text{rank}(A)$. The problem may become somewhat confounding when the exact A is unknown but given through empirical data in \tilde{A} , as illustrated in Example 1. What really is at stake is a nontrivial solution $\text{sol}(A, \mathbf{0}) \equiv \text{Kernel}(A)$ in

the Grassmannian $\mathcal{G}_{n-r}(\mathbb{C}^n)$, but the data system $\tilde{A}\mathbf{x} = \mathbf{0}$ is almost always non-singular with $\text{sol}(\tilde{A}, \mathbf{0}) = \{\mathbf{0}\} \in \mathcal{G}_0(\mathbb{C}^n)$ when $m \geq n$. The condition number $\kappa(\tilde{A}) = O(\|A - \tilde{A}\|_2^{-1})$ can be huge as well if $r < \min\{m, n\}$. The very problem of solving a homogeneous linear system from empirical data needs clarification.

PROBLEM 1 (numerical solution of a homogeneous linear system). *Let \tilde{A} be an $m \times n$ matrix serving as empirical data for an underlying homogeneous system $A\mathbf{x} = \mathbf{0}$, where entries of A may or may not be known exactly. Identify the rank r of A using \tilde{A} , and find a numerical solution of $\tilde{A}\mathbf{x} = \mathbf{0}$ in the Grassmannian $\mathcal{G}_{n-r}(\mathbb{C}^n)$ in the form of an orthonormal basis $\{\mathbf{z}_1, \dots, \mathbf{z}_{n-r}\}$ so that*

$$(11) \quad \text{dist}(\text{span}\{\mathbf{z}_1, \dots, \mathbf{z}_{n-r}\}, \text{sol}(A, \mathbf{0})) = O\left(\frac{\|A - \tilde{A}\|_2}{\|A\|_2}\right).$$

From Wedin's perturbation analysis [34], the numerical kernel $\text{Kernel}(\tilde{A}_\theta)$ within a proper error tolerance $\theta > 0$ is an approximation to $\text{Kernel}(A) = \text{sol}(A, \mathbf{0})$ in $\mathcal{G}_{n-r}(\mathbb{C}^n)$ (cf. Lemma 11 in Appendix A). For every $G \in \mathbb{C}^{m \times n}$, we define

$$\text{sol}_\theta(G, \mathbf{0}) := \text{Kernel}(G_\theta) \equiv \text{sol}(G_\theta, \mathbf{0})$$

as the *numerical solution* of the homogeneous system $G\mathbf{x} = \mathbf{0}$ in the Grassmannian $\mathcal{G}_{n-r}(\mathbb{C}^n)$ within an error tolerance θ , where $r = \text{rank}_\theta(G)$. Numerical methods for computing $\text{Kernel}(G_\theta)$ as $\text{sol}_\theta(G, \mathbf{0})$ are well established, including the SVD and other numerical rank-revealing methods (see, e.g., [11, 20]). The following theorem summarizes the properties of the numerical solution as a generalization of the exact solution to the homogeneous system and as a well-posed computing problem that solves the underlying system in Problem 1. The essence and underlying substance of Theorem 1 are based on Wedin [34].

THEOREM 1. *Let $A \in \mathbb{C}^{m \times n}$. The following properties hold for the numerical solution of a homogeneous system:*

- (i) The exact solution is a special case of the numerical solution

$$\text{sol}(A, \mathbf{0}) \equiv \text{sol}_\theta(A, \mathbf{0}) \quad \text{for all } \theta \in (0, \|A^\dagger\|_2^{-1}).$$

- (ii) Computing the numerical solution is a well-posed problem: *If $\text{sol}_\theta(A, \mathbf{0})$ is well-defined within $\theta > 0$, then $\text{sol}_\theta(A + \Delta A, \mathbf{0})$ uniquely exists in the same Grassmannian as $\text{sol}_\theta(A, \mathbf{0})$ and enjoys Lipschitz continuity with*

$$(12) \quad \begin{aligned} & \text{dist}(\text{sol}_\theta(A + \Delta A, \mathbf{0}), \text{sol}_\theta(A, \mathbf{0})) \\ & \leq \frac{\|A_\theta\|_2 \|A_\theta^\dagger\|_2}{1 - \|A_\theta^\dagger\|_2 (\|A - A_\theta\|_2 + \|\Delta A\|_2)} \frac{\|\Delta A\|_2}{\|A\|_2} \end{aligned}$$

for all ΔA with sufficiently small $\|\Delta A\|_2$ satisfying

$$\|\Delta A\|_2 \leq \min\left\{\frac{1}{2} \left(\|A_\theta^\dagger\|_2^{-1} - \|A - A_\theta\|_2\right), \theta - \|A - A_\theta\|_2, \|A_\theta^\dagger\|_2^{-1} - \theta\right\}.$$

- (iii) A homogeneous system can be solved from empirical data with an accuracy in the same order as the data: *For any $A + \Delta A$ serving as empirical data of A with $\|\Delta A\|_2 < \frac{1}{2} \|A^\dagger\|_2^{-1}$, there exist $\mu, \eta > 0$ with*

$$(13) \quad \mu \leq \|\Delta A\|_2 < \|A^\dagger\|_2^{-1} - \|\Delta A\|_2 \leq \eta$$

such that the numerical solution $\text{sol}_\theta(A + \Delta A, \mathbf{0})$ within any error tolerance $\theta \in (\mu, \eta)$ is in the same Grassmannian as the exact solution $\text{sol}(A, \mathbf{0})$ and

$$(14) \quad \text{dist}(\text{sol}_\theta(A + \Delta A, \mathbf{0}), \text{sol}(A, \mathbf{0})) \leq \frac{\|A\|_2 \|A^\dagger\|_2}{1 - \|A^\dagger\|_2 \|\Delta A\|_2} \frac{\|\Delta A\|_2}{\|A\|_2}.$$

Proof. The proof is straightforward from Wedin's error bound [34] on singular subspaces (see Lemma 11 in Appendix A) along with the identity $A_\theta \equiv A$ for $0 < \theta < \|A^\dagger\|_2^{-1} = \sigma_r(A)$, where $r = \text{rank}(A)$, $\mu = \sigma_{r+1}(\tilde{A}) \leq \|\Delta A\|_2$, and $\eta = \sigma_r(\tilde{A}) \geq \sigma_r(A) - \|\Delta A\|_2$. \square

By Theorem 1, Problem 1 is solvable if the data are sufficiently accurate and a tight error bound on data is attainable, as asserted in the following corollary.

COROLLARY 2. *Let the matrices A and \tilde{A} be as in Problem 1. Assume the data in \tilde{A} are sufficiently accurate such that $\|A - \tilde{A}\|_2 < \frac{1}{2} \|A^\dagger\|_2^{-1}$. Further assume a data error bound $\beta > \|A - \tilde{A}\|_2$ is known and is sufficiently tight so that $\beta < \|A^\dagger\|_2^{-1} - \|A - \tilde{A}\|_2$. Then Problem 1 is solvable by setting the error tolerance $\theta = \beta$ and finding an orthonormal basis for the numerical solution $\text{sol}_\theta(\tilde{A}, \mathbf{0}) = \text{Kernel}(\tilde{A}_\theta)$ within θ .*

Proof. The proof is straightforward from Theorem 1. \square

The error tolerance θ in Theorem 1 is an operational parameter that needs to be set up for solving Problem 1. If we assume that the underlying application allows the data error to a certain extent, say $\|A - \tilde{A}\|_2 < \hat{\theta}$, the data error bound β in Corollary 2 is expected to be below $\hat{\theta}$. The inequality (13) ensures there is a window (μ, η) for setting the operational error tolerance θ at β or slightly larger. Using the notation of Problem 1, it is reasonable to assume that the data error bound β on $\|A - \tilde{A}\|_2$ is known or can be estimated. The crucial criterion for operational purpose is to set θ at or slightly above $\|A - \tilde{A}\|_2$ according to Theorem 1(iii). The error tolerance θ should not exceed $\|A^\dagger\|_2^{-1} - \|A - \tilde{A}\|_2$, whose exact value or estimation is not needed if the data error bound β is sufficiently tight. See the supplementary material for examples of setting error tolerances.

For a rank- r matrix A , the sensitivity of solving $A\mathbf{x} = \mathbf{0}$ for $\text{sol}_\theta(\tilde{A}, \mathbf{0})$ in the Grassmannian $\mathcal{G}_{n-r}(\mathbb{C}^n)$ from a perturbed data matrix \tilde{A} is

$$\|A\|_2 \|A^\dagger\|_2 = \frac{\sigma_1(A)}{\sigma_r(A)} \approx \|\tilde{A}_\theta\|_2 \|\tilde{A}_\theta^\dagger\|_2$$

from (12) and (14), not infinity or $\kappa(\tilde{A})$. The convention $\kappa(A) = \infty$ for the square singular case and $\kappa(\tilde{A}) = \|\tilde{A}\|_2 \|\tilde{A}^\dagger\|_2$ may overestimate the sensitivity substantially. Problem 1 may not be solvable if the data error is large beyond, say, $\frac{1}{2} \|A^\dagger\|_2^{-1}$, or may not be solved accurately if the data error bound is unknown or the inherent sensitivity $\|A\|_2 \|A^\dagger\|_2$ is high.

For solving $A\mathbf{x} = \mathbf{0}$ with $A \in \mathbb{C}^{m \times n}$, there are differences between the cases $m < n$ and $m \geq n$. The solution is of a positive dimension when $m < n$ regardless of perturbations and, if $\text{rank}(A) = m$, the condition $\|A\|_2 \|A^\dagger\|_2$ is continuous with respect to small perturbations. When $m \geq n$ and $\text{sol}(A, \mathbf{0})$ is nontrivial, however, the dimension of $\text{sol}(A + \Delta A, \mathbf{0})$ degrades to zero for almost all perturbations ΔA , and the condition $\|A\|_2 \|A^\dagger\|_2$ is discontinuous. The assertions of Theorem 1 remain the same no matter whether $m < n$ or $m \geq n$.

5. Sensitivity of a consistent singular system. Solving a singular system for an individual vector solution is known to have an unbounded sensitivity under

arbitrary perturbations. From a different perspective, the infinity condition number is not the sensitivity of the *singular system* if the singularity is not maintained. There is an intrinsic stability in solving $A\mathbf{x} = \mathbf{b}$ when the rank and consistency are preserved. This point of view originated in [15], where Kahan suggested that the perceived hypersensitivity of multiple roots may be a “misconception” without maintaining the multiplicity.

A consistent $m \times n$ linear system $A\mathbf{x} = \mathbf{b}$ with $\text{rank}(A) = r$ has a unique solution $\text{sol}(A, \mathbf{b}) = \mathbf{x}_0 + \text{Kernel}(A)$ in the affine Grassmannian $\mathcal{A}_{n-r}(\mathbb{C}^n)$, where \mathbf{x}_0 is any particular solution. The sensitivity of the linear system $A\mathbf{x} = \mathbf{b}$ can be based on the deviation of the solution $\text{sol}(A, \mathbf{b})$ in $\mathcal{A}_{n-r}(\mathbb{C}^n)$ with respect to perturbations of $(A, \mathbf{b}) \in \mathbb{C}^{m \times n} \times \mathbb{C}^n$. From (1), the difference between solutions of two consistent systems of the same rank can be measured by the metric (1), namely,

$$(15) \quad \begin{aligned} & \text{dist}(\text{sol}(A, \mathbf{b}), \text{sol}(B, \mathbf{d})) \\ &= \max \{ \|A^\dagger \mathbf{b} - B^\dagger \mathbf{d}\|_2, \text{dist}(\text{Kernel}(A), \text{Kernel}(B)) \}. \end{aligned}$$

Notice that the component $\text{dist}(\text{Kernel}(A), \text{Kernel}(B)) \leq 1$ in (15), but the other component $\|A^\dagger \mathbf{b} - B^\dagger \mathbf{d}\|_2$ can be large or small. One way to avoid an imbalance is to put a weight factor ω on the component $\|A^\dagger \mathbf{b} - B^\dagger \mathbf{d}\|_2$. We choose not to use weights for the sake of simplicity and because the weight ω can be used to scale the linear system instead so that we can solve $A(\omega \mathbf{x}) = \omega \mathbf{b}$ equivalently. For convenience, we adopt a specific norm

$$(16) \quad \|(A, \mathbf{b})\| := \sqrt{\|A\|_2^2 + \|\mathbf{b}\|^2}$$

in the product space $\mathbb{C}^{m \times n} \times \mathbb{C}^n$. The theories in this paper can be adapted to other norms.

With these notation and metrics, the solution $\text{sol}(A, \mathbf{b})$ of a singular consistent $m \times n$ linear system $A\mathbf{x} = \mathbf{b}$ uniquely exists in the affine Grassmannian $\mathcal{A}_{n-r}(\mathbb{C}^n)$, and the sensitivity is proportional to $\|A\|_2 \|A^\dagger\|_2$ rather than infinity when the rank and consistency are preserved, as established in the following theorem.

THEOREM 3. The solution of a consistent linear system is Lipschitz continuous when the rank and consistency are preserved. Let $A \in \mathbb{C}^{m \times n}$ and $\mathbf{b} \in \text{Range}(A)$. Assume the perturbation $(\Delta A, \Delta \mathbf{b})$ is constrained such that $\text{rank}(\tilde{A}) = \text{rank}(A)$ and $\tilde{\mathbf{b}} \in \text{Range}(\tilde{A})$, where $\tilde{A} = A + \Delta A$ and $\tilde{\mathbf{b}} = \mathbf{b} + \Delta \mathbf{b}$. Then,

$$(17) \quad \begin{aligned} & \text{dist}(\text{sol}(\tilde{A}, \tilde{\mathbf{b}}), \text{sol}(A, \mathbf{b})) \\ & \leq \|A\|_2 \|A^\dagger\|_2 \cdot \frac{\sqrt{2 \|\mathbf{x}_*\|_2^2 + 1}}{\|A\|_2 - \sqrt{2 \|A\|_2 \|A^\dagger\|_2 \|\Delta A\|_2}} \|(\Delta A, \Delta \mathbf{b})\|, \end{aligned}$$

where $\mathbf{x}_* = A^\dagger \mathbf{b}$ whenever $\sqrt{2 \|A^\dagger\|_2 \|\Delta A\|_2} < 1$.

Proof sketch. The kernel component of the distance in (17) is bounded by Wedin’s error estimate [34] (see Lemma 11 in Appendix A). Let N be a matrix whose columns form an orthonormal basis for $\text{Kernel}(A)$. Then the minimum norm solution \mathbf{x}_* is the unique least squares solution of the system

$$\begin{bmatrix} \mu N^H \\ A \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} \quad \text{for any } \mu > 0,$$

and the standard error bound [24, Theorem 1.4.6] applies. A detailed proof is given in Appendix B. \square

As a result of (17), the intrinsic sensitivity of solving a singular system $A\mathbf{x} = \mathbf{b}$ for the general solution $\text{sol}(A, \mathbf{b})$ is a constant multiple of

$$\|A\|_2 \|A^\dagger\|_2 = \frac{\sigma_1(A)}{\sigma_r(A)} < \infty$$

when the rank and consistency are preserved.

As a by-product of establishing Theorem 3, the following corollary improves the standard normwise error bound [9, Theorem 5.6.1] on the minimum norm solution of a full rank underdetermined linear system by reducing a factor from 2 to $\sqrt{2}$.

COROLLARY 4. *Let $A \in \mathbb{C}^{m \times n}$ with $\text{rank}(A) = m < n$ and $\mathbf{b} \in \mathbb{C}^m$. If \mathbf{x}_* and $\tilde{\mathbf{x}}$ are minimum norm solutions of the underdetermined linear systems $A\mathbf{x} = \mathbf{b}$ and $(A + \Delta A)\mathbf{x} = \mathbf{b} + \Delta\mathbf{b}$, respectively, with $\sqrt{2}\|A^\dagger\|_2 \|\Delta A\|_2 < 1$, then*

$$(18) \quad \frac{\|\tilde{\mathbf{x}} - \mathbf{x}_*\|_2}{\|\mathbf{x}_*\|_2} \leq \|A\|_2 \|A^\dagger\|_2 \left(\sqrt{2} \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta\mathbf{b}\|_2}{\|\mathbf{b}\|_2} \right) + O(\|(\Delta A, \Delta\mathbf{b})\|^2).$$

Proof. The inequality (18) follows from (45) in Appendix B. \square

Remark 1. The subset of all rank- r matrices is a complex analytic manifold in the topological space $\mathbb{C}^{m \times n}$ [5], with the topology derived from the Frobenius norm. Similarly, the subset $\mathcal{M}_r^{m \times n} := \{(A, \mathbf{b}) \in \mathbb{C}^{m \times n} \times \mathbb{C}^m \mid \text{rank}(A) = r, \mathbf{b} \in \text{Range}(A)\}$ is a complex analytic manifold in $\mathbb{C}^{m \times n} \times \mathbb{C}^m$. Although in general the problem of solving a singular linear system is ill-posed, Theorem 3 implies that the problem of solving $A\mathbf{x} = \mathbf{b}$ for $\text{sol}(A, \mathbf{b})$ in $\mathcal{A}_{n-r}(\mathbb{C}^n)$ is well-posed on the manifold $\mathcal{M}_r^{m \times n}$.

6. The general numerical solution. When a rank-deficient $m \times n$ linear system $A\mathbf{x} = \mathbf{b}$ is given through empirical data $(\tilde{A}, \tilde{\mathbf{b}})$, the perturbed matrix \tilde{A} is almost always of full rank and highly ill-conditioned. Furthermore, the conventional single-vector solution of the data system $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ is in \mathbb{C}^n , while the general solution of the underlying system is in the completely different space $\mathcal{A}_{n-r}(\mathbb{C}^n)$. What the problem precisely is and what the numerical solution really means need to be clarified.

PROBLEM 2 (numerical solution of a linear system). *For given \tilde{A} and $\tilde{\mathbf{b}}$ serving as empirical data for an underlying linear system $A\mathbf{x} = \mathbf{b}$ to be solved, find a numerical solution of $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ that can be identified as the exact solution $\text{sol}(\hat{A}, \hat{\mathbf{b}})$ of $\hat{A}\mathbf{x} = \hat{\mathbf{b}}$ such that both the backward error and the forward error,*

$$(19) \quad \|(\tilde{A}, \tilde{\mathbf{b}}) - (\hat{A}, \hat{\mathbf{b}})\| = O(\|(\tilde{A}, \tilde{\mathbf{b}}) - (A, \mathbf{b})\|),$$

$$(20) \quad \text{dist}(\text{sol}(\hat{A}, \hat{\mathbf{b}}), \text{sol}(A, \mathbf{b})) = O(\|(\tilde{A}, \tilde{\mathbf{b}}) - (A, \mathbf{b})\|),$$

are in the same order of the data accuracy.

The accuracy requirement (20) stipulates that both $\text{sol}(A, \mathbf{b})$ and $\text{sol}(\hat{A}, \hat{\mathbf{b}})$ are either in the same affine Grassmannian or empty. It is natural to choose $\hat{A} = \tilde{A}_\theta$ within a proper θ and $\hat{\mathbf{b}} = \tilde{\mathbf{b}}_\theta := \tilde{A}_\theta \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}$ as the orthogonal projection of $\tilde{\mathbf{b}}$ onto $\text{Range}(\tilde{A}_\theta)$. The solution $\text{sol}(\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)$ is acceptable as the numerical solution of $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ if its backward error is below the error tolerance or the empty set.

DEFINITION 5 (general numerical solution). *Let $G \in \mathbb{C}^{m \times n}$, $\mathbf{d} \in \mathbb{C}^m$, and $\theta > 0$ be an error tolerance within which $\text{rank}_\theta(G)$ is well defined. With respect*

to a norm $\|\cdot\|$ on $\mathbb{C}^{m \times n} \times \mathbb{C}^m$, the general numerical solution of the linear system $G\mathbf{x} = \mathbf{d}$ within θ is defined as

$$sol_\theta(G, \mathbf{d}) := \begin{cases} sol(G_\theta, \mathbf{d}_\theta) & \text{if } \|(G, \mathbf{d}) - (G_\theta, \mathbf{d}_\theta)\| < \theta, \\ \emptyset & \text{if } \|(G, \mathbf{d}) - (G_\theta, \mathbf{d}_\theta)\| > \theta, \end{cases}$$

where G_θ is the θ -projection of G , and $\mathbf{d}_\theta = G_\theta G_\theta^\dagger \mathbf{d}$ is the orthogonal projection of \mathbf{d} onto the range $\text{Range}(G_\theta)$ of G_θ .

The solution $sol_\theta(G, \mathbf{d})$ is undefined if θ is equal to a singular value of G or $\theta = \|(G, \mathbf{d}) - (G_\theta, \mathbf{d}_\theta)\|$. We can now establish the following theorem on the general numerical solution.

THEOREM 6. At any $(A, \mathbf{b}) \in \mathbb{C}^{m \times n} \times \mathbb{C}^m$, the following properties of the general numerical solution hold with respect to the norm (16):

- (i) An exact general solution is a special case of general numerical solution: The identity $sol(A, \mathbf{b}) \equiv sol_\theta(A, \mathbf{b})$ holds for all $\theta < \|A^\dagger\|_2^{-1}$ if $\mathbf{b} \in \text{Range}(A)$, or $\theta < \min\{\|A^\dagger\|_2^{-1}, \|\mathbf{b} - A A^\dagger \mathbf{b}\|_2\}$ otherwise.
- (ii) Computing the general numerical solution is a well-posed problem: Assume $sol_\theta(A, \mathbf{b})$ is well-defined within a certain $\theta > 0$. There is a $\xi > 0$ depending on A , \mathbf{b} , and θ such that for every $(\Delta A, \Delta \mathbf{b})$ with a sufficiently small norm, there exists a unique $sol_\theta(A + \Delta A, \mathbf{b} + \Delta \mathbf{b})$ satisfying the Lipschitz continuity

$$(21) \quad dist(sol_\theta(A + \Delta A, \mathbf{b} + \Delta \mathbf{b}), sol_\theta(A, \mathbf{b})) \leq \xi \|(\Delta A, \Delta \mathbf{b})\|.$$

- (iii) A singular linear system can be solved from empirical data with an accuracy in the same order as the data:

(a) Assume $\mathbf{b} \in \text{Range}(A)$, and let $\mathbf{x}_* = A^\dagger \mathbf{b}$. For any empirical data pair $(\tilde{A}, \tilde{\mathbf{b}}) = (A + \Delta A, \mathbf{b} + \Delta \mathbf{b})$ satisfying $\|(\Delta A, \Delta \mathbf{b})\| < ((\omega + 1) \|A^\dagger\|_2)^{-1}$, where $\omega = \sqrt{4 \|A^\dagger\|_2^2 \|\mathbf{b}\|_2^2 + 2}$, and for any error tolerance θ satisfying

$$(22) \quad \omega \|(\Delta A, \Delta \mathbf{b})\| < \theta < \|A^\dagger\|_2^{-1} - \|(\Delta A, \Delta \mathbf{b})\|,$$

there exists a unique general numerical solution $sol_\theta(\tilde{A}, \tilde{\mathbf{b}})$ with a backward error bound $\omega \|(\Delta A, \Delta \mathbf{b})\|$ and a forward error bound

$$(23) \quad \begin{aligned} dist\left(sol_\theta(\tilde{A}, \tilde{\mathbf{b}}), sol(A, \mathbf{b})\right) \\ \leq \|A\|_2 \|A^\dagger\|_2 \frac{\sqrt{4 \|\mathbf{x}_*\|_2^2 + 1}}{\|A\|_2 - \|A\|_2 \|A^\dagger\|_2 \|\Delta A\|_2} \|(\Delta A, \Delta \mathbf{b})\|. \end{aligned}$$

(b) Assume $sol(A, \mathbf{b}) = \emptyset$. For any $\theta < \min\{\frac{1}{2} \|A^\dagger\|_2^{-1}, \|\mathbf{b} - A A^\dagger \mathbf{b}\|_2\}$, there is a constant $\rho \in (0, \theta)$ such that $sol_\theta(\tilde{A}, \tilde{\mathbf{b}}) = sol(A, \mathbf{b}) = \emptyset$ at any empirical data pair $(\tilde{A}, \tilde{\mathbf{b}})$ satisfying $\|(\tilde{A}, \tilde{\mathbf{b}}) - (A, \mathbf{b})\| < \rho$.

Proof sketch. Assertion (i) and the unique existence in assertion (ii) directly follow from Definition 5. The Lipschitz continuity (21) is a variation of the error estimate for the truncated SVD solution by Hansen [10, inequality (26a)] as an extension of Wedin error analysis [35]. The bound on the minimum norm solution component of the distance in the inequality (23) follows from Hansen [10, inequality (27a)], and the bound on the numerical kernel is established by Wedin [34]. A detailed proof is given in Appendix B. \square

For Problem 2, assume the underlying linear system $A\mathbf{x} = \mathbf{b}$ in Problem 2 is known to be consistent in applications such as those in Example 3; the solvability of the system from empirical data $(\tilde{A}, \tilde{\mathbf{b}})$ is given in the following corollary of Theorem 6.

COROLLARY 7. *Let (A, \mathbf{b}) and $(\tilde{A}, \tilde{\mathbf{b}})$ be as in Problem 2, where the underlying linear system $A\mathbf{x} = \mathbf{b}$ is consistent. Assume the data matrix \tilde{A} is sufficiently accurate with $\|A - \tilde{A}\|_2 < \frac{1}{2}\|A^\dagger\|_2^{-1}$. Further assume that an error bound $\beta > \|A - \tilde{A}\|_2$ is attainable and is sufficiently tight so that $\beta < \|A^\dagger\|_2^{-1} - \|A - \tilde{A}\|_2$. Then Problem 2 is solvable by calculating $\text{sol}(\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)$ with the error tolerance $\theta = \beta$, where $\tilde{\mathbf{b}}_\theta$ is the orthogonal projection of $\tilde{\mathbf{b}}$ onto $\text{Range}(\tilde{A}_\theta)$. Furthermore,*

$$\begin{aligned} & \text{dist} \left(\text{sol}(\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta), \text{sol}(A, \mathbf{b}) \right) \\ & \leq \|A\|_2 \|A^\dagger\|_2 \cdot \frac{\sqrt{4\|A^\dagger\|_2^2 + 1}}{\|A\|_2 - \|A\|_2 \|A^\dagger\|_2 \|A - \tilde{A}\|_2} \|(\tilde{A}, \tilde{\mathbf{b}}) - (A, \mathbf{b})\|. \end{aligned}$$

We reiterate that the sensitivity of solving $A\mathbf{x} = \mathbf{b}$ from empirical data $(\tilde{A}, \tilde{\mathbf{b}})$ is measured by

$$\|A\|_2 \|A^\dagger\|_2 = \frac{\sigma_1(A)}{\sigma_r(A)} \approx \|\tilde{A}_\theta\|_2 \|\tilde{A}_\theta^\dagger\|_2,$$

not infinity or $\kappa(\tilde{A})$, when the underlying matrix A is singular, where $r = \text{rank}(A)$. Problem 2 may still be difficult if data are inaccurate, if the intrinsic condition $\|A\|_2 \|A^\dagger\|_2$ is large, or if the window for setting the error tolerance is too narrow. The general numerical solution can be computed using existing rank-revealing tools, such as those in [18, 20] and UTV/ULV decomposition [9, sect. 5.4.6], in the following template:

```

set the error tolerance  $\theta$  at or slightly above the error bound  $\beta \gtrsim \|\Delta A\|_2$ 
if  $r = \text{rank}_\theta(A) \approx n$  then
  – calculate  $N \in \mathbb{C}^{n \times (n-r)}$ , whose columns form an orthonormal basis for
    the numerical kernel  $\text{Kernel}(\tilde{A}_\theta)$ 
  – solve  $A\mathbf{x} = \mathbf{b}$  for a particular solution  $\mathbf{x} = \mathbf{x}_*$  by any backward accurate method, such as  $\mathbf{x}_* = (A^\text{H} A + \mu^2 N N^\text{H})^{-1} A^\text{H} \mathbf{b}$  or Tikhonov regularization
  – output  $\text{sol}_\theta(A, \mathbf{b}) = \mathbf{x}_* + \text{Range}(N)$ .
else
  – calculate a decomposition  $U S V^\text{H} = A_\theta$  with  $S \in \mathbb{C}^{r \times r}$ ,  $U^\text{H} U = I$ , and
     $V^\text{H} V = I$ 
  – solve  $S\mathbf{y} = U^\text{H} \mathbf{b}$  for  $\mathbf{y} = \mathbf{y}_*$  and obtain the truncated SVD solution
     $\mathbf{x}_* = V\mathbf{y}_*$ 
  – output  $\text{sol}_\theta(A, \mathbf{b}) = \mathbf{x}_* + \text{Range}(V)^\perp$ 
end if

```

As we shall establish in section 7, the particular solution component \mathbf{x}_* of $\text{sol}_\theta(A, \mathbf{b})$ in the above template can be computed by any backward accurate numerical algorithm, including Tikhonov regularization and truncated SVD. Computation of the general numerical solution is implemented in the MATLAB package NACLAB [40] as the functionality `LinearSolve` (cf. [39] and the supplementary material). The general guideline for the error tolerance is to set it at or slightly larger than a known data error bound $\beta > \|A - \tilde{A}\|_2$ if the application allows such an adjustment. We conclude this section with the following example.

Example 4. Revisiting the linear system in Example 3, we see that the data error bound can be estimated as $\|\Delta A\|_2 \leq \|\Delta A\|_F \leq 4.5 \times 10^{-4}$, where A is the underlying matrix since the entrywise error bound is 5×10^{-5} . The error tolerance θ can be set at or slightly larger than the error bound, say $\theta = 0.0005$. The numerical solution of the system (5) within 0.0005 in the affine Grassmannian $\mathcal{A}_7(\mathbb{C}^9)$ is a representation of

$$(u_1, u_2, u_3) \\ = \left(.90710 + .33322x + .71029x^2 + .59968x^3, \quad -.79946 + .06694x, \quad 1.12433 - .06648x + .08926x^2 \right) \\ + t_1 \left(-.27897 - .08391x - .17878x^2 + .08424x^3, \quad -.35739 - .47261x, \quad .12212 - .33612x - .63016x^2 \right) \\ + t_2 \left(-.21387 + .29319x - .18465x^2 + .46503x^3, \quad -.55471 + .18011x, \quad -.46785 + .03542x + .24016x^2 \right)$$

(cf. the supplementary material). The general numerical solution $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}})$ is of a healthy sensitivity $\|\tilde{A}_\theta\|_2 \|\tilde{A}_\theta^\dagger\|_2 \approx 17.19$ —not the infinite $\kappa(A)$ or the large $\kappa(\tilde{A}) \approx 2.29 \times 10^6$. The three components of $\text{sol}_\theta(\tilde{A})$ form an invertible polynomial transformation matrix as shown in (6) with the numerical inverse

$$\begin{bmatrix} .55101 - .91839x^2, & -2.33985 + .70128x + .30047x^2 + .00001x^3, & .71342 + 1.83982x + 1.12986x^2 + .00002x^3 \\ -.36735 + .18367x - .73471x^5, & -.43135 - 1.09668x + .33663x^2 - 1.72768x^3 + .56105x^4 + 0.24037x^5, & -.9954 + .44186x + .8831x^2 + 1.11306x^3 + 1.47187x^4 + .90389x^5 \\ .18366 + .36734x^2 + .55103x^4, & 1.58108 - .53294x + 1.17553x^2 - .42079x^3 - .18027x^4, & -.28338 - 1.39825x - 1.28672x^2 - 1.10389x^3 - .6779x^4 \end{bmatrix}.$$

Remark 2. An $m \times n$ system $A\mathbf{x} = \mathbf{b}$ with $m > n$ is inconsistent for almost all $\mathbf{b} \in \mathbb{C}^m$, and its least squares solution is usually studied in the literature. In fact, the least squares solution can be considered whenever $\mathbf{b} \notin \text{Range}(A)$ even if $m \leq n$. There are substantial differences between the conventional solution and the least squares solution. In Theorem 6 and throughout this paper, our elaboration is restricted to the conventional solution so that $\text{sol}(A, \mathbf{b}) = \emptyset$ for inconsistent systems, and the nonempty set of least squares solutions is beyond the scope of this paper. The sensitivity of the least squares solution is well known to be $\kappa(A)^2$ (see, e.g., [14, sect. 20.1]) in contrast to $\kappa(A)$ for the (conventional) solution in Theorem 6.

7. Particular solution of a singular linear system. There are many applications where only a particular solution is needed from among the infinitely many solutions of a singular linear system $A\mathbf{x} = \mathbf{b}$, and it makes little difference which particular solution is obtained. For such applications, the problem of finding a *numerical particular solution* can be stated as follows.

PROBLEM 3 (numerical particular solution). *Assume a linear system $A\mathbf{x} = \mathbf{b}$ is consistent where the entries of A and \mathbf{b} may be known through empirical data of limited accuracy. Find a numerical particular solution $\tilde{\mathbf{x}}$ that approximates an exact solution $\mathbf{x}_* \in \text{sol}(A, \mathbf{b})$ with the error $\|\tilde{\mathbf{x}} - \mathbf{x}_*\|_2$ at an acceptable level.*

There are regularization approaches, such as the Tikhonov method [9, sect. 6.1.5], [11, 27] that can produce approximate particular solutions with high backward accuracy. For any backward accurate numerical solution $\tilde{\mathbf{x}}$ of the system $A\mathbf{x} = \mathbf{b}$, in the sense that there is a pair $(\tilde{A}, \tilde{\mathbf{b}})$ such that $\tilde{A}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$ and that $\|(A, \mathbf{b}) - (\tilde{A}, \tilde{\mathbf{b}})\|$ is at an acceptable level, we call $\tilde{\mathbf{x}}$ a *numerical particular solution* of $A\mathbf{x} = \mathbf{b}$. The following theorem asserts that every numerical particular solution approximates one of the exact solutions.

THEOREM 8. *Let $A \in \mathbb{C}^{m \times n}$ with $\text{rank}(A) < n$ and $\mathbf{b} \in \text{Range}(A)$. Assume $\tilde{\mathbf{x}} \in \mathbb{C}^n$ is a backward accurate numerical solution of $A\mathbf{x} = \mathbf{b}$ in the sense that $\tilde{\mathbf{x}}$ is*

an exact solution of $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ with $\|\tilde{A} - A\|_2 \leq .46 \|A^\dagger\|_2^{-1}$. Then $\tilde{\mathbf{x}}$ approximates an exact solution $\mathbf{x}_* \in \text{sol}(A, \mathbf{b})$ with an error bound

$$(24) \quad \frac{\|\tilde{\mathbf{x}} - \mathbf{x}_*\|_2}{\|\mathbf{x}_*\|_2} \leq \frac{\|A\|_2 \|A^\dagger\|_2}{1 - \|A^\dagger\|_2 \|\Delta A\|_2} \left(2\sqrt{2} \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \right),$$

assuming $\mathbf{b} \neq \mathbf{0}$ where $\Delta A = A - \tilde{A}$, $\Delta \mathbf{b} = \mathbf{b} - \tilde{\mathbf{b}}$, or

$$(25) \quad \|\tilde{\mathbf{x}} - \mathbf{x}_*\|_2 \leq \frac{\|A\|_2 \|A^\dagger\|_2}{1 - \|A^\dagger\|_2 \|\Delta A\|_2} \left(\|\tilde{\mathbf{x}}\|_2 \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|A\|_2} \right)$$

if $\mathbf{b} = \mathbf{0}$.

Proof sketch. Let $r = \text{rank}(A)$ and $\sigma_{r+1}(\tilde{A}) < \theta < \sigma_r(\tilde{A})$. Write $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}_1 + \tilde{\mathbf{x}}_2$, where $\tilde{\mathbf{x}}_1 = \tilde{A}_\theta^\dagger \tilde{A}_\theta \tilde{\mathbf{x}}$ and $\tilde{\mathbf{x}}_2 = (I - \tilde{A}_\theta^\dagger \tilde{A}_\theta) \tilde{\mathbf{x}}$. Choose a particular solution $\mathbf{x}_* = A^\dagger \mathbf{b} + (I - A^\dagger A) \tilde{\mathbf{x}}_2$ from $\text{sol}(A, \mathbf{b})$. Since $\tilde{\mathbf{x}}_1 = \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}$ approximates $A^\dagger \mathbf{b}$, $\tilde{\mathbf{x}}_2 \in \text{Kernel}(\tilde{A}_\theta)$, $(I - A^\dagger A) \tilde{\mathbf{x}}_2 \in \text{Kernel}(A)$, and $\text{Kernel}(\tilde{A}_\theta)$ approximates $\text{Kernel}(A)$, and hence $\tilde{\mathbf{x}}$ is an approximation to the particular solution \mathbf{x}_* of $A\mathbf{x} = \mathbf{b}$, so the theorem holds. Detailed proofs of (24) and (25) are given in Appendix B. \square

For the case $\mathbf{b} = \mathbf{0}$ in Theorem 8, the objective is to solve the homogeneous system $A\mathbf{x} = \mathbf{0}$. The inequality (25) includes the following three cases:

Case (i): $\tilde{\mathbf{b}} = \mathbf{0}$ and $\tilde{\mathbf{x}} = \mathbf{0}$. The inequality (25) is trivial and perhaps meaningless since $\tilde{\mathbf{x}} = \mathbf{x}_* = \mathbf{0}$.

Case (ii): $\tilde{\mathbf{b}} = \mathbf{0}$ and $\tilde{\mathbf{x}} \neq \mathbf{0}$. Then we can normalize $\tilde{\mathbf{x}}$ to be a unit vector so that (25) becomes

$$(26) \quad \min_{\mathbf{z} \in \text{Kernel}(A)} \|\tilde{\mathbf{x}} - \mathbf{z}\|_2 \leq \|\tilde{\mathbf{x}} - \mathbf{x}_*\|_2 \leq \frac{\|A\|_2 \|A^\dagger\|_2}{1 - \|A^\dagger\|_2 \|\Delta A\|_2} \frac{\|\Delta A\|_2}{\|A\|_2}.$$

Case (iii): $\tilde{\mathbf{b}} \neq \mathbf{0}$. This case is relevant in practical computation by setting the right-hand side $\tilde{\mathbf{b}}$ as a nonzero random vector of a moderate norm and obtaining a numerical particular solution $\tilde{\mathbf{x}}$ as an exact solution of $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ with a small $\|\tilde{A} - A\|_2$, leading to the inverse power iteration. The norm $\|\tilde{\mathbf{x}}\|_2$ is almost always large due to the condition number $\kappa(\tilde{A}) = O(\|A - \tilde{A}\|_2^{-1})$. As it pleasantly turns out, the large $\|\tilde{\mathbf{x}}\|_2$ is exactly what is needed as (25) becomes

$$(27) \quad \left\| \frac{\tilde{\mathbf{x}}}{\|\tilde{\mathbf{x}}\|_2} - \frac{\mathbf{x}_*}{\|\tilde{\mathbf{x}}\|_2} \right\|_2 \leq \frac{\|A\|_2 \|A^\dagger\|_2}{1 - \|A^\dagger\|_2 \|\Delta A\|_2} \frac{1}{\|A\|_2} \left(\|\Delta A\|_2 + \frac{\|\tilde{\mathbf{b}}\|_2}{\|\tilde{\mathbf{x}}\|_2} \right).$$

The larger the norm achieved by $\|\tilde{\mathbf{x}}\|_2$, the higher the accuracy of $\frac{\tilde{\mathbf{x}}}{\|\tilde{\mathbf{x}}\|_2}$ to a particular nontrivial solution of the homogeneous system $A\mathbf{x} = \mathbf{0}$. Once again, the sensitivity of solving a singular linear system $A\mathbf{x} = \mathbf{b}$ is $\|A\|_2 \|A^\dagger\|_2 = \frac{\sigma_1(A)}{\sigma_r(A)}$ —not infinity in the sense of finding a numerical particular solution.

Particular solutions of $A\mathbf{x} = \mathbf{b}$ can vary arbitrarily, but their deviations can only stretch in $\text{Kernel}(A)$. As the following corollary states, the high sensitivity is near a direction in $\text{Kernel}(A)$, and such a sensitivity may be harmless after all.

COROLLARY 9. *Let $A \in \mathbb{C}^{m \times n}$ with $\text{rank}(A) < n$ and $\mathbf{b} \in \text{Range}(A)$. Assume \mathbf{x}_1 and \mathbf{x}_2 are both backward accurate numerical particular solutions of $A\mathbf{x} = \mathbf{b}$ in*

the sense that $A_1 \mathbf{x}_1 = \mathbf{b}_1$ and $A_2 \mathbf{x}_2 = \mathbf{b}_2$ with sufficiently small $\|(A_1, \mathbf{b}_1) - (A, \mathbf{b})\|$ and $\|(A_2, \mathbf{b}_2) - (A, \mathbf{b})\|$. Then there is an $\mathbf{x}_* \in \text{Kernel}(A)$ such that

$$(28) \quad \begin{aligned} \|(\mathbf{x}_1 - \mathbf{x}_2) - \mathbf{x}_*\|_2 &\leq \|A\|_2 \|A^\dagger\|_2 \\ &\times \left(\frac{\|\mathbf{b} - \mathbf{b}_1\|_2}{\|A\|_2} + \frac{\|\mathbf{b} - \mathbf{b}_2\|_2}{\|A\|_2} + \frac{\|A - A_1\|_2}{\|A\|_2} \|\mathbf{x}_1\|_2 + \frac{\|A - A_2\|_2}{\|A\|_2} \|\mathbf{x}_2\|_2 \right). \end{aligned}$$

Proof. Apply the inequality (25) on $\tilde{A} = A$ and $\tilde{\mathbf{x}} = \mathbf{x}_1 - \mathbf{x}_2$ that satisfy $A(\mathbf{x}_1 - \mathbf{x}_2) = (\mathbf{b}_1 - \mathbf{b}) + (\mathbf{b} - \mathbf{b}_2) + (A - A_1)\mathbf{x}_1 + (A_2 - A)\mathbf{x}_2$. \square

Theorem 8 extends the accuracy result for the inverse iteration in spite of the large condition number. In [29, pp. 339–340], Peters and Wilkinson described what they called “exaggerated fears” in the early days of the computer age when the inverse iteration

$$(29) \quad (A - \lambda I) \mathbf{x}_{k+1} = \mathbf{x}_k \quad \text{for } k = 0, 1, \dots$$

at an approximation λ to an eigenvalue λ_* of A was proposed for calculating an eigenvector \mathbf{x}_* as a nontrivial solution to the homogeneous system $(A - \lambda_* I) \mathbf{x} = \mathbf{0}$:

Although [inverse iteration is] basically a simple concept its numerical properties have not been widely understood. If λ really is very close to an eigenvalue, the matrix $(A - \lambda I)$ is almost singular and hence a typical step in the iteration involves the solution of a very ill-conditioned set of equations. . . . The period when inverse iteration was first considered was notable for exaggerated fears concerning the instability of direct methods for solving linear systems and *ill-conditioned* systems were a source of particular anxiety. . . . few numerical analysts discuss inverse iteration with any confidence.

It is counterintuitive, but pleasantly surprising nonetheless, that ill-conditioning is not harmful in computing the eigenvector. As pointed out in [29] and by Parlett [28, sect. 4.3] as follows, errors mainly lie in $\text{Kernel}(A - \lambda_* I)$ and are not really errors at all:

... roundoff errors can give rise to completely erroneous “solutions” to very ill-conditioned systems of equations. . . . Indeed some textbooks have cautioned users not to take $[\lambda]$ too close to any eigenvalue. . . . Fortunately these fears are groundless and furnish a nice example of confusing ends with means. . . . the error $\mathbf{e} [= \mathbf{x}_{k+1} - \mathbf{x}_*]$, which may be almost as large as the exact solution of $[(A - \lambda I)^{-1} \mathbf{x}_k]$, is almost entirely in the direction of [the eigenvector]. . . . This result is alarming if we had hoped for an accurate solution of [(29)] (the means) but is a delight in the search for [the eigenvector] (the end).

Theorem 8 concludes, in fact, that the fear of solving a highly ill-conditioned linear system may also be exaggerated for nonhomogeneous systems when the underlying system $A\mathbf{x} = \mathbf{b}$ is consistent and singular, as long as the numerical solution is backward accurate and the intrinsic sensitivity measure $\|A\|_2 \|A^\dagger\|_2$ is moderate. The variation between any two numerical particular solutions can be large, but the difference falls harmlessly in the kernel of A . In other words, the “error” is actually a part of the solution.

Example 5. The system $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$ in (7) is a representation of the underlying system $A\mathbf{x} = \mathbf{b}$ with $\|\Delta A\|_2 \leq \|\Delta A\|_F \leq 4.5 \times 10^{-4} = \theta$ from the entrywise

error bound 0.5×10^{-5} . Rounded to five digits after the decimal point, two numerical particular solutions $\tilde{\mathbf{x}}_0$ and \mathbf{x}_1 by truncated SVD $\tilde{A}_\theta^\dagger \tilde{\mathbf{b}}$ and MATLAB backslash “\” operator, respectively, are

$$\begin{aligned}\tilde{\mathbf{x}}_0 &= [0.90711, 0.33322, 0.71029, 0.59968, -0.79946, 0.06694, 1.12433, -0.06648, 0.08926]^H, \\ \tilde{\mathbf{x}}_1 &= [-0.78366, 0.47296, -0.45954, 1.83637, -3.47453, -1.81379, 0.84635, -1.57209, -2.41843]^H,\end{aligned}$$

with both residuals $\|\tilde{A} \tilde{\mathbf{x}}_0 - \tilde{\mathbf{b}}\| \approx 8.1 \times 10^{-5}$ and $\|\tilde{A} \tilde{\mathbf{x}}_1 - \tilde{\mathbf{b}}\| \approx 5.3 \times 10^{-5}$ roughly within the data error bound. The two numerical particular solutions are far apart with $\|\tilde{\mathbf{x}}_0 - \tilde{\mathbf{x}}_1\| \approx 5.01$ as predicted by the large condition number $\kappa(\tilde{A}) \approx 2.29 \times 10^6$. However, the underlying system is consistent and singular with a healthy sensitivity $\|A\|_2 \|A^\dagger\|_2 \approx \|\tilde{A}_\theta\|_2 \|A_\theta^\dagger\|_2 \lesssim 17.19$. Both $\tilde{\mathbf{x}}_0$ and $\tilde{\mathbf{x}}_1$ are accurate approximations to different exact solutions with estimate error bounds 0.00186 and 0.00177, respectively, and actual relative errors are 4.49×10^{-5} and 0.93×10^{-5} in the same level of the data error.

8. Bona fide ill-conditioned linear systems. A linear system $A \mathbf{x} = \mathbf{b}$ is truly ill-conditioned when $\|A\|_2 \|A^\dagger\|_2$ is large regardless of its rank. When A is of full column rank, $\mathbf{b} \in \text{Range}(A)$, and the condition number $\kappa(A)$ is huge, the solution uniquely exists but in general cannot be computed accurately from perturbed data using any algorithm. The system is de facto rank-deficient in a practical sense. Even in such cases, a stable general numerical solution may still be attainable in an affine Grassmannian from empirical data, and the underlying solution can be accurately approximated by a vector in the affine subspace as the general numerical solution.

THEOREM 10. Assume $A \in \mathbb{C}^{m \times n}$, $\mathbf{x}_* \in \mathbb{C}^n$, and $\mathbf{b} = A \mathbf{x}_*$. Let r be any integer with $\sigma_r(A) > \sigma_{r+1}(A)$. For any $(\tilde{A}, \tilde{\mathbf{b}}) = (A + \Delta A, \mathbf{b} + \Delta \mathbf{b})$ serving as empirical data of (A, \mathbf{b}) with

$$(30) \quad \|\Delta A\|_2 < \min\{\sigma_r(A) - \sigma_{r+1}(A), (2\sqrt{3} - 3)\sigma_r(A)\},$$

there is an $\tilde{\mathbf{x}} \in \text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}})$ with $\sigma_{r+1}(\tilde{A}) < \theta < \sigma_r(\tilde{A})$ such that

$$(31) \quad \frac{\|\tilde{\mathbf{x}} - \mathbf{x}_*\|_2}{\|\mathbf{x}_*\|_2} \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\sigma_{r+1}(A) - \|\Delta A\|_2}{\sigma_r(A)}} \left((2 + \sqrt{2}) \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \right).$$

Proof sketch. Since \mathbf{x}_* is a backward accurate solution of the linear system $\tilde{A}_\theta \mathbf{x} = \tilde{\mathbf{b}}_\theta$, Theorem 8 applies from a reversed perspective. The detailed proof is given in Appendix B. \square

In the following example, the underlying system is ill-conditioned but truly non-singular. All known numerical algorithms, including regularization methods, produce solutions that are inaccurate as single vectors but highly accurate as the vector component of a general numerical solution that is perfectly conditioned and contains accurate approximations to the underlying exact solution.

Example 6. Consider the polynomial division problem in the form of the equation

$$(x + 10) q + \rho = \frac{1}{3} x^8 + 4 x^7 + \frac{23}{3} x^6 + \frac{34}{3} x^5 + 15 x^4 + \frac{56}{3} x^3 + \frac{67}{3} x^2 + 26 x + \frac{89}{3}$$

for the quotient q and the constant remainder ρ . There is a unique solution which consists of $q = \frac{1}{3}(x^7 + 2x^6 + \dots + 7x + 8)$ and $\rho = 3$. The corresponding linear system

is of the form $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{bmatrix} 1 & & & \\ 10 & 1 & & \\ & \ddots & \ddots & \\ & & 10 & 1 \end{bmatrix}, \quad \mathbf{b} = \frac{1}{3} \begin{bmatrix} 1 \\ 12 \\ \vdots \\ 89 \end{bmatrix},$$

with the exact solution $\mathbf{x}_* = \frac{1}{3}[1, 2, \dots, 9]^H$ that is attainable in symbolic computation using the exact data in rational number format. In MATLAB single precision arithmetic, the system is represented as perturbed data $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$, where $\tilde{A} = A$ and $\tilde{\mathbf{b}} = [.333333, 4.0, 7.6666665, 11.333333, 15.0, 18.6666666, 22.333334, 26.0, 29.666666]^H$. The singular values $10.9461079 > 10.7891169 > \dots > 9.0683689 > 9.9 \times 10^{-9}$ indicate that it is practically impossible to calculate the single-vector solution with any meaningful accuracy using such data. Table 1 shows three sample numerical solutions: \mathbf{x}_1 by a straightforward application of the MATLAB command $\mathbf{A}\backslash\mathbf{b}$, a Tikhonov regularization solution $\mathbf{x}_2 = (A^H A + \alpha^2 I)^{-1} A^H \tilde{\mathbf{b}}$ at, say, $\alpha = 0.001$, and the truncated SVD solution $\mathbf{x}_3 = A_\theta^\dagger \tilde{\mathbf{b}}$ with an error tolerance that is roughly $\theta = \|\mathbf{b}\|_2 \varepsilon \approx 3.18 \times 10^{-6}$, where ε is the unit roundoff. As expected from the condition number $\kappa(A) \approx 1.1 \times 10^9$, none of the solutions can be considered accurate *as a single vector*. On the other hand, the general numerical solution $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}})$ is almost perfectly conditioned at $\|\tilde{A}_\theta\|_2 \|\tilde{A}_\theta^\dagger\|_2 \approx 1.21$. The three solutions \mathbf{x}_1 , \mathbf{x}_2 , and \mathbf{x}_3 that are inaccurate as individual vectors are all accurate as the component $\tilde{\mathbf{u}}$ of the general numerical solution $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}}) = \tilde{\mathbf{u}} + \text{Kernel}(\tilde{A}_\theta) = \text{span}\{\tilde{\mathbf{v}}\}$, where

$$(32) \quad \tilde{\mathbf{v}} = [.0, -.0000001, .0000010, -.0000099, .0000995, -.0009950, .0099499, -.0994987, .9949875]^H.$$

All $\mathbf{x}_j + \text{Kernel}(A_\theta)$ for $j = 1, 2, 3$ are nearly identical in the affine Grassmannian $\mathcal{A}_1(\mathbb{C}^9)$, and each contains a particular vector $\hat{\mathbf{x}}_j$ that is an accurate approximation to the exact solution \mathbf{x}_* , as shown in the bottom part of Table 1. The errors $\frac{\|\hat{\mathbf{x}}_j - \mathbf{x}_*\|_2}{\|\mathbf{x}_*\|_2}$ are all within the bound 8.28×10^{-7} predicted by (31).

TABLE 1

For $A\mathbf{x} = \tilde{\mathbf{b}}$ in Example 6, numerical solutions \mathbf{x}_1 , \mathbf{x}_2 , and \mathbf{x}_3 by MATLAB “\”, Tikhonov regularization, and truncated SVD, respectively, in comparison with the exact solution \mathbf{x}_* . Also shown are the accuracies of \mathbf{x}_1 , \mathbf{x}_2 and \mathbf{x}_3 as a component of the general numerical solution.

| Solution type | Numerical (single-vector) solution with incorrect digits crossed out | Error $\frac{\ \mathbf{x}_j - \mathbf{x}_*\ _2}{\ \mathbf{x}_*\ _2}$ |
|--|--|--|
| 8 digits of \mathbf{x}_* MATLAB “\” \mathbf{x}_1 Tikhonov \mathbf{x}_2 trunc. SVD \mathbf{x}_3 | .333333 .6666667 1.0000000 1.3333333 1.6666667 2.0000000 2.3333333 2.6666667 3.0000000 .333333 .6666665 1.0000044 1.333187 1.6668129 1.9985374 2.3479633 2.5203667 4.4620993 .333333 .6666674 0.9999974 1.3333607 1.6663934 2.0027277 2.3060572 2.9394238 0.2724303 .3333335 0.6666669 0.9999967 1.3333608 1.6668988 2.0027270 2.3060613 2.9398885 0.2727296 | 0.2612916 0.4871437 0.4870902 |
| $\mathbf{x}_1 + \text{Kernel}(A_\theta)$ $\mathbf{x}_2 + \text{Kernel}(A_\theta)$ $\mathbf{x}_3 + \text{Kernel}(A_\theta)$ | Particular $\hat{\mathbf{x}}_j = \mathbf{x}_j + t_j \tilde{\mathbf{v}} \in \mathbf{x}_j + \text{Kernel}(A_\theta)$ nearest to \mathbf{x}_* with $\tilde{\mathbf{v}}$ in (32) $\hat{\mathbf{x}}_1 = \mathbf{x}_1 + t_1 \tilde{\mathbf{v}} \approx \mathbf{x}_*$ with $t_1 = -1.4703701$ $\hat{\mathbf{x}}_2 = \mathbf{x}_2 + t_2 \tilde{\mathbf{v}} \approx \mathbf{x}_*$ with $t_2 = 2.7413113$ $\hat{\mathbf{x}}_3 = \mathbf{x}_3 + t_3 \tilde{\mathbf{v}} \approx \mathbf{x}_*$ with $t_3 = 2.7410104$ | Error $\frac{\ \hat{\mathbf{x}}_j - \mathbf{x}_*\ _2}{\ \mathbf{x}_*\ _2}$ 8.8 $\times 10^{-8}$ 1.5 $\times 10^{-7}$ 1.9 $\times 10^{-7}$ |

The linear system in Example 6 is nonsingular in theory but practically underdetermined in numerical computation. Suppose an additional piece of information becomes available, say the remainder $\rho = 3$. One can impose such a constraint on the general numerical solution $\{\tilde{\mathbf{u}} + t \tilde{\mathbf{v}} \mid t \in \mathbb{C}\}$ at the trailing component as $0.2727296 + .9949875t = 3$, obtaining $t = 2.7410097$ corresponding to a numerical solution with a relative error 1.79×10^{-7} in the same order of the data.

Appendix A. Lemmas.

LEMMA 11. Let $A, \tilde{A} \in \mathbb{C}^{m \times n}$ with $\Delta A = \tilde{A} - A$. Assume $\sigma_r(A) > \sigma_{r+1}(A)$.

(i) (Wedin) If $\|\Delta A\|_2 < \frac{1}{2}(\sigma_r(A) - \sigma_{r+1}(A))$, then

$$(33) \quad \begin{aligned} & \text{dist} \left(\text{Kernel}(A_\theta), \text{Kernel}(\tilde{A}_\theta) \right) \\ & \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\sigma_{r+1}(A) + \|\Delta A\|_2}{\sigma_r(A)}} \frac{\|\Delta A\|_2}{\|A\|_2} \\ & \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{2}{1 - \frac{\sigma_{r+1}(A)}{\sigma_r(A)}} \frac{\|\Delta A\|_2}{\|A\|_2} \end{aligned}$$

for any $\theta \in (\sigma_{r+1}(A), \sigma_r(A)) \cap (\sigma_{r+1}(\tilde{A}), \sigma_r(\tilde{A})) \neq \emptyset$.

(ii) If $\text{rank}(A) = \text{rank}(\tilde{A}) = r$ and $\|\Delta A\| < \sigma_r(A)$, then

$$(34) \quad \text{dist} \left(\text{Kernel}(A), \text{Kernel}(\tilde{A}) \right) \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{\|\Delta A\|_2}{\|A\|_2}.$$

Proof. Assertion (i) is established by Wedin [34] (also see [32, Theorem 4.4] and [10, Theorem 3.3]). To prove (ii), let the SVDs of A and \tilde{A} be

$$A = [U_1, U_2] \begin{bmatrix} \Sigma_1 & \\ & O \end{bmatrix} [V_1, V_2]^H \quad \text{and} \quad \tilde{A} = [\tilde{U}_1, \tilde{U}_2] \begin{bmatrix} \tilde{\Sigma}_1 & \\ & O \end{bmatrix} [\tilde{V}_1, \tilde{V}_2]^H,$$

respectively, where $\Sigma_1, \tilde{\Sigma}_1 \in \mathbb{C}^{r \times r}$. Then,

$$-\tilde{V}_2^H \Delta A^H = \tilde{V}_2^H \tilde{A}^H - \tilde{V}_2^H \Delta A^H = \tilde{V}_2^H A^H = (\tilde{V}_2^H V_1) (\Sigma_1^H U_1^H),$$

and thus,

$$\text{dist} \left(\text{Kernel}(A), \text{Kernel}(\tilde{A}) \right) = \|\tilde{V}_2^H V_1\|_2 \leq \frac{\|\Delta A\|_2}{\sigma_r(A)} \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{\|\Delta A\|_2}{\|A\|_2}. \quad \square$$

LEMMA 12. Let \mathcal{U} be a subspace of \mathbb{C}^n , and let U be a matrix whose columns form an orthonormal basis for \mathcal{U} . For every subspace \mathcal{V} of \mathbb{C}^n of the same dimension as \mathcal{U} with $\text{dist}(\mathcal{U}, \mathcal{V}) < 1$, there is a matrix V whose columns form a basis for \mathcal{V} such that $\|U - V\|_2 \leq \text{dist}(\mathcal{U}, \mathcal{V})$.

Proof. Let G be any matrix whose columns form an orthonormal basis for \mathcal{V} , and let $[G, \hat{G}]$ be a unitary matrix. Then, for any unit vector $\mathbf{x} \in \mathbb{C}^n$,

$$1 = \|([G, \hat{G}]^H U \mathbf{x})\|_2^2 = \|(G^H U) \mathbf{x}\|_2^2 + \|(\hat{G}^H U) \mathbf{x}\|_2^2 \leq \|(G^H U) \mathbf{x}\|_2^2 + \text{dist}(\mathcal{U}, \mathcal{V})^2$$

leading to $\|(G^H U) \mathbf{x}\|_2^2 \geq 1 - \text{dist}(\mathcal{U}, \mathcal{V})^2 > 0$, implying $G^H U$ is invertible so that columns of $V = G(G^H U)$ form a basis for \mathcal{V} , and $\|U - V\|_2 = \|(U U^H - G G^H) U\|_2$ that is less than or equal to $\text{dist}(\mathcal{U}, \mathcal{V})$. \square

LEMMA 13. Let $A \in \mathbb{C}^{m \times n}$ with $\sigma_r(A) > \theta > \sigma_{r+1}(A)$.

- (i) For every $\mu \in [\sigma_r(A), \sigma_1(A)]$, let $N \in \mathbb{C}^{n \times (n-r)}$ be a matrix whose columns form an orthonormal basis for $\text{Kernel}(A_\theta)$. Then

$$(35) \quad \left\| \begin{bmatrix} \mu N^H \\ A \end{bmatrix} \right\|_2 = \max \{ \sigma_1(A), \sqrt{\mu^2 + \sigma_{r+1}(A)^2} \} \\ \in \left[\|A\|_2, \sqrt{2} \|A\|_2 \right)$$

$$(36) \quad \left\| \begin{bmatrix} \mu N^H \\ A \end{bmatrix}^\dagger \right\|_2 = \max \left\{ \frac{1}{\sigma_r(A)}, \frac{1}{\sqrt{\mu^2 + \eta^2}} \right\} \leq \|A_\theta^\dagger\|_2,$$

where $\eta = \sigma_n(A)$ if $m \geq n$ or $\eta = 0$ otherwise.

- (ii) Assume columns of $N \in \mathbb{C}^{n \times (n-r)}$ span $\text{Kernel}(A_\theta)$. For any $\mu > 0$, let $\mathbf{b} \in \mathbb{C}^m$, and let \mathbf{x}_* be the least squares solution of the linear system

$$(37) \quad \begin{bmatrix} \mu N^H \\ A \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix}.$$

Then $\mathbf{x}_* = A^\dagger \mathbf{b}$ if $\text{rank}(A) = r$ or $A_\theta \mathbf{x}_* = \mathbf{b}_\theta$ if $\text{rank}(A) > r$, where $\mathbf{b}_\theta = A_\theta A_\theta^\dagger \mathbf{b}$ is the orthogonal projection of \mathbf{b} onto $\text{Range}(A_\theta)$,

$$(38) \quad \left\| \begin{bmatrix} \mu N^H \\ A \end{bmatrix} \mathbf{x}_* - \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} \right\|_2 \leq \|\mathbf{b} - \mathbf{b}_\theta\|_2,$$

$$(39) \quad \mathbf{x}_* - TT^H \mathbf{x}_* = A_\theta^\dagger \mathbf{b}_\theta,$$

for any $T \in \mathbb{C}^{n \times (n-r)}$ with $\text{Range}(T) = \text{Kernel}(A_\theta)$ and $T^H T = I$.

Proof. For the case of $m \geq n$, we can write A in its singular value expansion $A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^H + \dots + \sigma_n \mathbf{u}_n \mathbf{v}_n^H$ and $N = [\mathbf{v}_{r+1}, \dots, \mathbf{v}_n] G$ where $\mathbf{u}_1, \dots, \mathbf{u}_m$ and $\mathbf{v}_1, \dots, \mathbf{v}_n$ are left and right singular vectors, respectively, with a unitary matrix $G \in \mathbb{C}^{(n-r) \times (n-r)}$. Write $\mathbf{x} = x_1 \mathbf{v}_1 + \dots + x_n \mathbf{v}_n$. Then,

$$\begin{aligned} \left\| \begin{bmatrix} \mu N^H \\ A \end{bmatrix} \mathbf{x} \right\|_2^2 &= \left\| \begin{bmatrix} G^H & \\ & I \end{bmatrix} \begin{bmatrix} \mu [\mathbf{v}_{r+1}, \dots, \mathbf{v}_n]^H \\ A \end{bmatrix} \mathbf{x} \right\|_2^2 \\ &= \sigma_1^2 |x_1|^2 + \dots + \sigma_r |x_r|^2 + (\mu^2 + \sigma_{r+1}^2) |x_{r+1}|^2 + \dots + (\mu^2 + \sigma_n^2) |x_n|^2, \end{aligned}$$

whose extrema subject to $\|\mathbf{x}\|_2 = 1$ are $\max\{\sigma_1^2, \mu^2 + \sigma_{r+1}^2\}$ and $\min\{\sigma_r^2, \mu^2 + \sigma_n^2\}$, leading to (35) and (36) in assertion (i). The case $m < n$ is similar.

To prove (ii), write the SVD $A = U_1 \Sigma_1 V_1^H + U_2 \Sigma_2 V_2^H$, where $\Sigma_1 \in \mathbb{C}^{r \times r}$ and $\Sigma_2 \in \mathbb{C}^{(m-r) \times (n-r)}$. Then \mathbf{x}_* is the solution of the normal equation $\mu^2 N N^H \mathbf{x}_* + A^H A \mathbf{x}_* - A^H \mathbf{b} = \mathbf{0}$. Namely, we have an orthogonal decomposition

$$(40) \quad (V_1 \Sigma_1^H \Sigma_1 V_1^H \mathbf{x}_* - V_1 \Sigma_1^H U_1^H \mathbf{b}) + (V_2 \Sigma_2^H \Sigma_2 V_2^H \mathbf{x}_* - V_2 \Sigma_2^H U_2^H \mathbf{b} + \mu^2 N N^H \mathbf{x}_*) \\ = \mathbf{0},$$

implying $V_1 \Sigma_1^H \Sigma_1 V_1^H \mathbf{x}_* - V_1 \Sigma_1^H U_1^H \mathbf{b} = \mathbf{0}$ and thus $V_1^H \mathbf{x}_* = \Sigma_1^{-1} U_1^H \mathbf{b}$. Since $\mathbf{x}_* = V_1 V_1^H \mathbf{x}_* + V_2 V_2^H \mathbf{x}_*$, we have $A_\theta \mathbf{x}_* = (U_1 \Sigma_1 V_1^H)(V_1 V_1^H \mathbf{x}_*) = U_1 U_1^H \mathbf{b} = \mathbf{b}_\theta$. Namely \mathbf{x}_* is a particular solutions of the system $A_\theta \mathbf{x} = \mathbf{b}_\theta$. Also,

$$A_\theta^\dagger \mathbf{b}_\theta = (V_1 \Sigma_1^{-1} U_1^H)(U_1 U_1^H \mathbf{b}) = (V_1 \Sigma_1^{-1} U_1^H) \mathbf{b} = V_1 V_1^H \mathbf{x}_* = (I - TT^H) \mathbf{x}_*.$$

Furthermore,

$$\left\| \begin{bmatrix} \mu N^H \\ A \end{bmatrix} \mathbf{x}_* - \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} \right\|_2 \leq \left\| \begin{bmatrix} \mu N^H \\ A \end{bmatrix} V_1 V_1^H \mathbf{x}_* - \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} \right\|_2 = \|\mathbf{b} - \mathbf{b}_\theta\|_2.$$

Finally, if $\text{rank}(A) = r$, then $A_\theta = A$ so $\Sigma_2 = O$ in (40), implying $N^H \mathbf{x}_* = \mathbf{0}$. Consequently $\mathbf{x}_* = A^\dagger \mathbf{b}$ from (39). \square

The following lemma is a variation of Theorem 5.1 in [35] by Wedin and its extension in Theorem 3.4 in [10] by Hansen.

LEMMA 14 (Wedin [35], Hansen [10]). *Let $A \in \mathbb{C}^{m \times n}$ and $\mathbf{b} \in \mathbb{C}^m$. Assume, for a $\theta > 0$, we have $\text{rank}_\theta(A) = r$ and $\|\mathbf{b} - A_\theta A_\theta^\dagger \mathbf{b}\|_2 < \theta$. There is a constant*

$$(41) \quad \zeta = \|A_\theta^\dagger \mathbf{b}\|_2 + \frac{1 + \|A_\theta^\dagger \mathbf{b}\|_2}{1 - \|A_\theta^\dagger\|_2 \|A - A_\theta\|_2}$$

such that, for any $\tilde{A} = A + \Delta A \in \mathbb{C}^{m \times n}$ and $\tilde{\mathbf{b}} = \mathbf{b} + \Delta \mathbf{b} \in \mathbb{C}^m$ with

$$(42) \quad \|\Delta A\|_2 < \min \left\{ \frac{1}{2} (\sigma_r(A) - \sigma_{r+1}(A)), \sigma_r(A) - \theta, \theta - \sigma_{r+1}(A) \right\},$$

the following inequality holds:

$$(43) \quad \|A_\theta^\dagger \mathbf{b} - \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}\|_2 \leq \frac{\sigma_1(A)}{\sigma_r(A)} \left(\zeta \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|A\|_2} \right) + O(\|(\Delta A, \Delta \mathbf{b})\|^2).$$

As a special case, further assume $\text{rank}_\theta(A) = r$ and $\mathbf{b} \in \text{Range}(A)$. Then

$$(44) \quad \|A^\dagger \mathbf{b} - \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}\|_2 \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\|\Delta A\|_2}{\sigma_r(A)}} \left(2 \|A^\dagger \mathbf{b}\|_2 \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|A\|_2} \right).$$

Proof. The assumption $\text{rank}_\theta(A) = r$ implies $\sigma_{r+1}(A) < \theta < \sigma_r(A)$, and thus $\sigma_{r+1}(\tilde{A}) < \theta < \sigma_r(\tilde{A})$, following (42) so that $\text{rank}_\theta(\tilde{A}) = r$ as well. Then it is straightforward to verify (43) from the inequality (26a) in [10] using

$$\|A(A_\theta^\dagger \mathbf{b}) - \mathbf{b}\|_2 = \|A_\theta A_\theta^\dagger \mathbf{b} - \mathbf{b}\|_2 < \theta < \sigma_r(A).$$

The inequality (44) follows from [10, inequality (27a)] and $\mathbf{b} \in \text{Range}(A)$. \square

LEMMA 15. *At any $(A, \mathbf{b}) \in \mathbb{C}^{m \times n} \times \mathbb{C}^n$ and $\theta > 0$ within which $\text{sol}_\theta(A, \mathbf{b})$ is well-defined, there is a $\delta > 0$ such that $\text{sol}_\theta(A + \Delta A, \mathbf{b} + \Delta \mathbf{b})$ is well-defined with the same dimension as $\text{sol}_\theta(A, \mathbf{b})$ if $\|(\Delta A, \Delta \mathbf{b})\| < \delta$.*

Proof. Write $\tilde{A} = A + \Delta A$ and $\tilde{\mathbf{b}} = \mathbf{b} + \Delta \mathbf{b}$. Since $r = \text{rank}_\theta(A)$ is well-defined, we have $\sigma_{r+1} < \theta < \sigma_r(A)$. Thus $\|\Delta A\|_2 < \min\{\sigma_r - \theta, \theta - \sigma_{r+1}\}$ ensures $\text{rank}_\theta(\tilde{A}) = r$. Let $P = I - A_\theta A_\theta^\dagger$ and $\tilde{P} = I - \tilde{A}_\theta \tilde{A}_\theta^\dagger$. Then

$$\tilde{A} - \tilde{A}_\theta = \tilde{P} \tilde{A} = \tilde{P} \Delta A + (\tilde{P} - P) A + (A - A_\theta),$$

and by (33),

$$\|P - \tilde{P}\|_2 = \text{dist} \left(\text{Range}(A_\theta), \text{Range}(\tilde{A}_\theta) \right) \leq \eta \frac{\|\Delta A\|_2}{\|A\|_2},$$

where, assuming $\|\Delta A\|_2 \leq \frac{1}{2} (\sigma_r(A) - \sigma_{r+1}(A))$,

$$\eta = \frac{\sigma_1(A)}{\sigma_r(A)} \frac{2}{1 - \frac{\sigma_{r+1}(A)}{\sigma_r(A)}} = \frac{2 \|A_\theta\|_2 \|A_\theta^\dagger\|_2}{1 - \|A_\theta^\dagger\|_2 \|A - A_\theta\|_2},$$

implying

$$\|A - A_\theta\|_2 - (\eta + 1) \|\Delta A\|_2 \leq \|\tilde{A} - \tilde{A}_\theta\|_2 \leq \|A - A_\theta\|_2 + (\eta + 1) \|\Delta A\|_2,$$

and similarly

$$\begin{aligned} \|\mathbf{b} - \mathbf{b}_\theta\|_2 - \eta \frac{\|\Delta A\|_2}{\|A\|_2} \|\mathbf{b}\|_2 - \|\Delta \mathbf{b}\|_2 &\leq \|\tilde{\mathbf{b}} - \tilde{\mathbf{b}}_\theta\|_2 \\ &\leq \|\mathbf{b} - \mathbf{b}_\theta\|_2 + \eta \frac{\|\Delta A\|_2}{\|A\|_2} \|\mathbf{b}\|_2 + \|\Delta \mathbf{b}\|_2, \end{aligned}$$

where $\mathbf{b}_\theta = \mathbf{b} - P\mathbf{b}$ and $\tilde{\mathbf{b}}_\theta = \tilde{\mathbf{b}} - \tilde{P}\tilde{\mathbf{b}}$. If $\text{sol}_\theta(A, \mathbf{b})$ is empty, then $\|(A, \mathbf{b}) - (A_\theta, \mathbf{b}_\theta)\| > \theta$, and thus $\|(\tilde{A}, \tilde{\mathbf{b}}) - (\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)\| > \theta$ when $\|(\Delta A, \Delta \mathbf{b})\|$ is sufficiently small so that $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}}) = \emptyset$ as well. When $\|(A, \mathbf{b}) - (A_\theta, \mathbf{b}_\theta)\| < \theta$ and $\|(\Delta A, \Delta \mathbf{b})\|$ is sufficiently small, we also have $\|(\tilde{A}, \tilde{\mathbf{b}}) - (\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)\| < \theta$ and $\sigma_{r+1}(\tilde{A}) < \theta < \sigma_r(\tilde{A})$ so that $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}}) = \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}_\theta + \text{Kernel}(\tilde{A}_\theta)$ has the identical dimension $n - r$ as $\text{sol}_\theta(A, \mathbf{b})$. \square

Appendix B. Proofs of theorems and corollaries.

Proof of Theorem 3. Let $N \in \mathbb{C}^{n \times (n-r)}$, whose columns form an orthonormal basis for $\text{Kernel}(A)$. By Lemma 12, there is an $\tilde{N} \in \mathbb{C}^{n \times (n-r)}$, whose columns form a basis for $\text{Kernel}(\tilde{A})$ such that $\|N - \tilde{N}\|_2 \leq \text{dist}(\text{Kernel}(A), \text{Kernel}(\tilde{A}))$. For $\zeta = \sigma_r(A)$, denote $B = \begin{bmatrix} \zeta N^H \\ A \end{bmatrix}$ and $\tilde{B} = \begin{bmatrix} \zeta \tilde{N}^H \\ \tilde{A} \end{bmatrix}$. Then $A^\dagger \mathbf{b} = B^\dagger \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix}$ and $\tilde{A}^\dagger \tilde{\mathbf{b}} = \tilde{B}^\dagger \begin{bmatrix} \mathbf{0} \\ \tilde{\mathbf{b}} \end{bmatrix}$ by Lemma 13(ii). By $\|B - \tilde{B}\|_2 \leq \sqrt{2} \|\Delta A\|_2$ from Lemma 11(ii), [24, Theorem 1.4.6, p. 30], and Lemma 13(i),

$$\begin{aligned} \|B^\dagger \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} - \tilde{B}^\dagger \begin{bmatrix} \mathbf{0} \\ \tilde{\mathbf{b}} \end{bmatrix}\|_2 &\leq \frac{\sigma_1(B)}{\sigma_n(B)} \frac{1}{1 - \frac{\|B - \tilde{B}\|_2}{\sigma_n(B)}} \left(\|\mathbf{x}_*\|_2 \frac{\|B - \tilde{B}\|_2}{\|B\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|B\|_2} \right) \\ (45) \quad &\leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\sqrt{2} \|\Delta A\|_2}{\sigma_r(A)}} \frac{\sqrt{2} \|\mathbf{x}_*\|_2 \|\Delta A\|_2 + \|\Delta \mathbf{b}\|_2}{\|A\|_2}, \end{aligned}$$

leading to (17). \square

Proof Theorem 6. Assertion (i) is true because $A_\theta = A$ and $\mathbf{b}_\theta = \mathbf{b}$ for $\theta \in (0, \sigma_r(A))$. If $\mathbf{b} \in \text{Range}(A)$, then

$$\text{sol}_\theta(A, \mathbf{b}) = \text{sol}(A, \mathbf{b}) = A^\dagger \mathbf{b} + \text{Kernel}(A).$$

Otherwise $\text{sol}_\theta(A, \mathbf{b}) = \text{sol}(A, \mathbf{b}) = \emptyset$ if $\theta < \min\{\sigma_r(A), \|A A^\dagger \mathbf{b} - \mathbf{b}\|_2\}$. Assertion (ii) directly follows from Lemmas 14 and 15 with

$$\xi = \|A_\theta\|_2 \|A_\theta^\dagger\|_2 \frac{\sqrt{\zeta^2 + 1}}{\|A\|_2 - \|A\|_2 \|A_\theta^\dagger\|_2 \|A - A_\theta\|_2} + \varepsilon$$

for any $\varepsilon > 0$.

We now prove assertion (iii)(a). Let $\tilde{\mathbf{b}}_\theta = \tilde{A}_\theta \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}$, $P = I - A A^\dagger$, and $\tilde{P} = I - \tilde{A}_\theta \tilde{A}_\theta^\dagger$. From $\|\tilde{A} - \tilde{A}_\theta\|_2 = \min_{\text{rank}(B)=r} \|\tilde{A} - B\|_2 \leq \|\Delta A\|_2$, we have

$$\begin{aligned} \|\tilde{\mathbf{b}} - \tilde{\mathbf{b}}_\theta\|_2 &= \|\tilde{P} \tilde{\mathbf{b}}\|_2 = \|\tilde{P} \tilde{\mathbf{b}} - P \mathbf{b}\|_2 \leq \|\tilde{P}\|_2 \|\tilde{\mathbf{b}} - \mathbf{b}\|_2 + \|\tilde{P} - P\|_2 \|\mathbf{b}\|_2 \\ &\leq \|\Delta \mathbf{b}\|_2 + \text{dist}(\text{Range}(\tilde{A}_\theta), \text{Range}(A)) \|\mathbf{b}\|_2 \\ (by \ (33)) \quad &\leq \|\Delta \mathbf{b}\|_2 + \frac{\sigma_1(A)}{\sigma_r(A)} \frac{2 \|\mathbf{b}\|_2}{\|A\|_2} \|\Delta A\|_2 \\ &\leq \sqrt{4 \|A^\dagger\|_2^2 \|\mathbf{b}\|_2^2 + 1} \|(\Delta A, \Delta \mathbf{b})\| = \sqrt{\omega^2 - 1} \|(\Delta A, \Delta \mathbf{b})\| \end{aligned}$$

and

$$(46) \quad \begin{aligned} \|(\tilde{A}, \tilde{\mathbf{b}}) - (\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)\| &\leq \sqrt{\|\Delta A\|_2^2 + (\omega^2 - 1)(\|\Delta A\|_2^2 + \|\Delta \mathbf{b}\|_2^2)} \\ &\leq \omega \|\Delta A, \Delta \mathbf{b}\| < \sigma_r(A) - \|\Delta A, \Delta \mathbf{b}\|. \end{aligned}$$

Then, for any θ satisfying (22),

$$\sigma_{r+1}(\tilde{A}) \leq \|\Delta A, \Delta \mathbf{b}\| < \theta < \sigma_r(A) - \|\Delta A, \Delta \mathbf{b}\| \leq \sigma_r(\tilde{A})$$

so $\text{rank}_\theta(\tilde{A}) = r$, $\|(\tilde{A}, \tilde{\mathbf{b}}) - (\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)\| < \theta$, and thus $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}})$ is of the same dimension as $\text{sol}(A, \mathbf{b})$. Since $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}}) = \text{sol}(\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)$, the backward error of $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}})$ is bounded above by $\omega \|\Delta A, \Delta \mathbf{b}\|$ from (46). Thus (23) follows from (33) in Lemma 11 and (44) in Lemma 14, leading to assertion (iii).

We now prove assertion (iii)(b). If $\text{sol}(A, \mathbf{b})$ is empty, then $\text{sol}_\theta(A, \mathbf{b}) = \emptyset$ whenever $\theta < \min\{\sigma_r(A), \|\mathbf{b} - A A^\dagger \mathbf{b}\|_2\}$. By Lemma 15, there is a $\delta_1 > 0$ such that $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}}) = \emptyset$ for every $(\tilde{A}, \tilde{\mathbf{b}})$ with $\|(\tilde{A}, \tilde{\mathbf{b}}) - (A, \mathbf{b})\| < \delta_1$. Thus $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}}) = \text{sol}(A, \mathbf{b})$ with both backward and forward errors as zero. \square

Proof of Theorem 8. Let $\tilde{A} = \tilde{U}_1 \tilde{\Sigma}_1 \tilde{V}_1^\mathsf{H} + \tilde{U}_2 \tilde{\Sigma}_2 \tilde{V}_2^\mathsf{H}$ be the SVD where $\tilde{\Sigma}_1$ is $r \times r$ with $r = \text{rank}(A)$. Then $\tilde{\mathbf{x}}$ is a solution to $\tilde{A} \mathbf{x} = \mathbf{b}$, implying $\tilde{U}_1 \tilde{\Sigma}_1 \tilde{V}_1^\mathsf{H} \tilde{\mathbf{x}}_1 = \tilde{U}_1 \tilde{U}_1^\mathsf{H} \tilde{\mathbf{b}}$ and $\tilde{U}_2 \tilde{\Sigma}_2 \tilde{V}_2^\mathsf{H} \tilde{\mathbf{x}}_2 = \tilde{U}_2 \tilde{U}_2^\mathsf{H} \tilde{\mathbf{b}}$ where $\tilde{\mathbf{x}}_1 = \tilde{V}_1 \tilde{V}_1^\mathsf{H} \tilde{\mathbf{x}}$ and $\tilde{\mathbf{x}}_2 = \tilde{V}_2 \tilde{V}_2^\mathsf{H} \tilde{\mathbf{x}}$. Then $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}_1 + \tilde{\mathbf{x}}_2$ with $\tilde{\mathbf{x}}_1 = \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}$ for any θ between $\sigma_{r+1}(\tilde{A})$ and $\sigma_r(A) - \|\Delta A\|_2$. By Lemma 14 with $\hat{\mathbf{x}} = A^\dagger \mathbf{b}$, we have

$$\|\tilde{\mathbf{x}}_1 - \hat{\mathbf{x}}\|_2 \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{\|\hat{\mathbf{x}}\|_2}{1 - \frac{\|\Delta A\|_2}{\sigma_r(A)}} \left(2 \frac{\|\Delta A\|}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \right).$$

Let columns of N form an orthonormal basis for $\text{Kernel}(A)$. Since $\tilde{\mathbf{x}}_2 \in \text{Kernel}(\tilde{A}_\theta)$,

$$(47) \quad \begin{aligned} \|N N^\mathsf{H} \tilde{\mathbf{x}}_2 - \tilde{\mathbf{x}}_2\|_2 &= \min_{\mathbf{u} \in \text{Kernel}(A)} \|\mathbf{u} - \tilde{\mathbf{x}}_2\|_2 \\ &\leq \text{dist}(\text{Kernel}(A_\theta), \text{Kernel}(A)) \|\tilde{\mathbf{x}}_2\|_2 \\ &\leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\|\Delta A\|_2}{\sigma_r(A)}} \frac{\|\Delta A\|_2}{\|A\|_2} \|\tilde{\mathbf{x}}_2\|_2 \end{aligned}$$

by Lemma 11, which, combined with $\|\Delta A\|_2 \leq .46 \sigma_r(A) < (2\sqrt{3}-3) \sigma_r(A)$, implies $\text{dist}(\text{Kernel}(A_\theta), \text{Kernel}(A)) < \frac{\sqrt{3}}{2}$, and thus

$$\begin{aligned} \|N N^\mathsf{H} \tilde{\mathbf{x}}_2\|_2 &= \|N^\mathsf{H} \tilde{\mathbf{x}}_2\|_2 = \|(N^\mathsf{H} \tilde{V}_2) \tilde{V}_2^\mathsf{H} \tilde{\mathbf{x}}_2\|_2 \\ &\geq \sqrt{1 - \text{dist}(\text{Kernel}(A_\theta), \text{Kernel}(A))^2} \|\tilde{V}_2^\mathsf{H} \tilde{\mathbf{x}}_2\|_2 \geq \frac{1}{2} \|\tilde{\mathbf{x}}_2\|_2. \end{aligned}$$

Let $\mathbf{x}_* = \hat{\mathbf{x}} + N N^\mathsf{H} \tilde{\mathbf{x}}_2$. Then \mathbf{x}_* is a particular solution to $A \mathbf{x} = \mathbf{b}$ and $\|\mathbf{x}_*\|_2^2 = \|\hat{\mathbf{x}}\|_2^2 + \|N N^\mathsf{H} \tilde{\mathbf{x}}_2\|_2^2$. We have

$$\begin{aligned} \|\tilde{\mathbf{x}} - \mathbf{x}_*\| &\leq \|\tilde{\mathbf{x}}_1 - \hat{\mathbf{x}}\|_2 + \|\tilde{\mathbf{x}}_2 - N N^\mathsf{H} \tilde{\mathbf{x}}_2\|_2 \\ &\leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\|\Delta A\|_2}{\sigma_r(A)}} \left(\frac{\|\Delta A\|}{\|A\|_2} (2 \|\hat{\mathbf{x}}\|_2 + 2 \|N N^\mathsf{H} \tilde{\mathbf{x}}_2\|_2) + \|\mathbf{x}_*\| \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \right) \\ &\leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{\|\mathbf{x}_*\|_2}{1 - \frac{\|\Delta A\|_2}{\sigma_r(A)}} \left(2\sqrt{2} \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \right), \end{aligned}$$

leading to (24). For the case $\mathbf{b} = \mathbf{0}$, the bound (25) follows from (47)

$$\|\tilde{\mathbf{x}}_1\|_2 \leq \frac{\|U_1^H \tilde{\mathbf{b}}\|_2}{\sigma_r(\tilde{A})} \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\|\Delta A\|_2}{\sigma_r(A)}} \frac{\|\Delta \tilde{\mathbf{b}}\|_2}{\|A\|_2}. \quad \square$$

Proof of Theorem 10. Let $U_1 \Sigma_1 V_1^H + U_2 \Sigma_2 V_2^H$ and $\tilde{U}_1 \tilde{\Sigma}_1 \tilde{V}_1^H + \tilde{U}_2 \tilde{\Sigma}_2 \tilde{V}_2^H$ be SVDs of A and \tilde{A} , respectively, where $\Sigma_1, \tilde{\Sigma}_1 \in \mathbb{C}^{r \times r}$. Denote $A_1 = U_1 \Sigma_1 V_1^H$, $\tilde{A}_1 = \tilde{U}_1 \tilde{\Sigma}_1 \tilde{V}_1^H$, $\mathbf{x}_1 = A_1^\dagger \mathbf{b}$, $\mathbf{x}_2 = \mathbf{x}_* - \mathbf{x}_1$, $\tilde{\mathbf{x}}_1 = \tilde{A}_1^\dagger \tilde{\mathbf{b}}$, and $\mathbf{r} = A \mathbf{x}_1 - \mathbf{b}$. Then, with $\mathbf{r} = U_2 U_2^H \mathbf{b} = U_2 \Sigma_2 V_2^H \mathbf{x}_2$,

$$\begin{aligned} \tilde{\mathbf{x}}_1 - \mathbf{x}_1 &= \tilde{A}_1^\dagger (\mathbf{b} + \Delta \mathbf{b}) - \mathbf{x}_1 = \tilde{A}_1^\dagger (A \mathbf{x}_1 - \mathbf{r} + \Delta \mathbf{b}) - \mathbf{x}_1 \\ &= \tilde{A}_1^\dagger ((\tilde{A} - \Delta A) \mathbf{x}_1 - \mathbf{r} + \Delta \mathbf{b}) - \mathbf{x}_1 \\ &= \tilde{A}_1^\dagger (-\Delta A \mathbf{x}_1 - \mathbf{r} + \Delta \mathbf{b}) - (I - \tilde{A}_1^\dagger \tilde{A}_1) \mathbf{x}_1 \\ &= \tilde{A}_1^\dagger (-\Delta A \mathbf{x}_1 + \Delta \mathbf{b}) - \tilde{V}_1 \tilde{\Sigma}_1^{-1} \tilde{U}_1^H U_2 \Sigma_2 V_2^H \mathbf{x}_2 - \tilde{V}_2 \tilde{V}_2^H \mathbf{x}_1, \end{aligned}$$

leading to

$$\begin{aligned} &\|(\tilde{\mathbf{x}}_1 + \tilde{V}_2 \tilde{V}_2^H \mathbf{x}_1) - \mathbf{x}_1\|_2 \\ &\leq \|\tilde{A}_1^\dagger\|_2 (\|\Delta A\|_2 \|\mathbf{x}_1\|_2 + \|\Delta \mathbf{b}\|_2) + \|\tilde{\Sigma}_1^{-1}\|_2 \|\Sigma_2\|_2 \|\tilde{U}_1^H U_2\|_2 \|\mathbf{x}_2\|_2 \\ &\leq \|\tilde{A}_1^\dagger\|_2 (\|\Delta A\|_2 \|\mathbf{x}_1\|_2 + \|\Delta \mathbf{b}\|_2) + \text{dist}(\text{Range}(U_1), \text{Range}(\tilde{U}_1)) \|\mathbf{x}_2\|_2 \\ &\leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{\|\mathbf{x}_*\|_2}{1 - \frac{\sigma_{r+1} + \|\Delta A\|_2}{\sigma_r(A)}} \left(\sqrt{2} \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \right). \end{aligned}$$

Let $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}_1 + \tilde{\mathbf{x}}_2 \in \text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}})$ with

$$\tilde{\mathbf{x}}_2 = \tilde{V}_2 \tilde{V}_2^H \mathbf{x}_1 + \tilde{V}_2 \tilde{V}_2^H \mathbf{x}_2 \in \text{Kernel}(\tilde{A}_\theta).$$

Then

$$\|\tilde{\mathbf{x}} - \mathbf{x}_*\|_2 \leq \|(\tilde{\mathbf{x}}_1 + \tilde{V}_2 \tilde{V}_2^H \mathbf{x}_1) - \mathbf{x}_1\|_2 + \|\tilde{V}_2 \tilde{V}_2^H \mathbf{x}_2 - \mathbf{x}_2\|_2$$

while, similar to the proof of Theorem 8 from (30),

$$\|\tilde{V}_2 \tilde{V}_2^H \mathbf{x}_2 - \mathbf{x}_2\|_2 \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\sigma_{r+1} + \|\Delta A\|_2}{\sigma_r(A)}} \frac{\|\Delta A\|_2}{\|A\|_2} 2 \|\tilde{V}_2 \tilde{V}_2^H \mathbf{x}_2\|_2$$

leading to (31). \square

REFERENCES

- [1] K. E. AVRACHENKOV AND J. B. LASSERRE, *Analytic perturbation of Sylvester matrix equations*, IEEE Trans. Automat. Control, 47 (2002), pp. 1116–1119.
- [2] J. L. BARLOW, H. ERBAY, AND I. SLAPNIČAR, *An alternative algorithm for the refinement of ULV decompositions*, SIAM J. Matrix Anal. Appl., 27 (2005), pp. 198–211, <https://doi.org/10.1137/S0895479801372621>.
- [3] M. COORNAERT, *Topological Dimension and Dynamical Systems*, Springer, Switzerland, 2015.
- [4] B. DAYTON, T.-Y. LI, AND Z. ZENG, *Multiple zeros of nonlinear systems*, Math. Comp., 80 (2011), pp. 2143–2168.
- [5] J. W. DEMMEL AND A. EDELMAN, *The dimension of matrices (matrix pencils) with given Jordan (Kronecker) canonical forms*, Linear Algebra Appl., 230 (1995), pp. 61–87.

- [6] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality constraints*, SIAM J. Matrix Anal. Appl., 20 (1998), pp. 303–353, <https://doi.org/10.1137/S0895479895290954>.
- [7] R. D. FIERRO, P. C. HANSEN, AND P. S. K. HANSEN, *UTV Tools: Matlab templates for rank-revealing UTV decompositions*, Numer. Algorithms, 20 (1999), pp. 165–194.
- [8] S. GAO, *Factoring multivariate polynomials via partial differential equations*, Math. Comp., 72 (2003), pp. 801–822.
- [9] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 4th ed., The Johns Hopkins University Press, Baltimore, 2013.
- [10] P. C. HANSEN, *The truncated SVD as a method for regularization*, BIT, 27 (1987), pp. 534–553.
- [11] P. C. HANSEN, *Rank-Deficient and Discrete Ill-Posed Problems: Numerical Aspects of Linear Inversion*, Math. Model. Comput. 4, SIAM, Philadelphia, 1998, <https://doi.org/10.1137/1.9780898719697>.
- [12] P. C. HANSEN, *Discrete Inverse Problems. Insight and Algorithms*, Fund. Alg. 7, SIAM, Philadelphia, 2010, <https://doi.org/10.1137/1.9780898718836>.
- [13] P. C. HANSEN, J. G. NAGY, AND D. P. O’LEARY, *Deblurring Images: Matrices, Spectra, and Filtering*, Fund. Alg. 3, SIAM, Philadelphia, 2006, <https://doi.org/10.1137/1.9780898718874>.
- [14] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, 2nd ed., SIAM, Philadelphia, 2002, <https://doi.org/10.1137/1.9780898718027>.
- [15] W. KAHAN, *Conserving Confluence Curbs Ill-Condition*, Technical Report 6, Computer Science, University of California, Berkeley, 1972.
- [16] D. A. KLAIN AND G.-C. ROTA, *Introduction to Geometric Probability*, Cambridge University Press, Cambridge, 1997.
- [17] V. LAKSHMIBAI AND J. BROWN, *The Grassmannian Variety*, Springer, New York, 2015.
- [18] T.-L. LEE, T.-Y. LI, AND Z. ZENG, *A rank-revealing method with updating, downdating, and applications, Part II*, SIAM J. Matrix Anal. Appl., 31 (2009), pp. 503–525, <https://doi.org/10.1137/07068179X>.
- [19] T.-L. LEE, T.-Y. LI, AND Z. ZENG, *RankRev—A Matlab package for computing the numerical rank and updating/downdating*, Numer. Algorithms, 77 (2018), pp. 559–576.
- [20] T.-Y. LI AND Z. ZENG, *A rank-revealing method with updating, downdating and applications*, SIAM J. Matrix Anal. Appl., 26 (2005), pp. 918–946, <https://doi.org/10.1137/S0895479803435282>.
- [21] L.-H. LIM, K. S.-W. WONG, AND K. YE, *Numerical Algorithms on the Affine Grassmannian*, preprint, <https://arxiv.org/abs/1607.01833>, 2018.
- [22] L.-H. LIM, K. S.-W. WONG, AND K. YE, *The Grassmannian of Affine Subspaces*, preprint, <https://arxiv.org/abs/1807.10883>, 2018.
- [23] T. LIU AND J. HUANG, *A discrete-time recurrent neural network for solving rank-deficient matrix equations with an application to output regulation of linear systems*, IEEE Trans. Neural Netw. Learn. Syst., 29 (2018), pp. 2271–2277.
- [24] Å. BJÖRCK, *Numerical Methods for Least Squares Problems*, SIAM, Philadelphia, 1996, <https://doi.org/10.1137/1.9781611971484>.
- [25] C. D. MEYER, *Matrix Analysis and Applied Linear Algebra*, SIAM, Philadelphia, 2000.
- [26] T. MORA, *Solving Polynomial Equation Systems I: The Kronecker-Duval Philosophy*, Cambridge University Press, London, 2003.
- [27] A. NEUMAIER, *Solving ill-conditioned and singular linear systems: A tutorial on regularization*, SIAM Rev., 40 (1998), pp. 636–666, <https://doi.org/10.1137/S0036144597321909>.
- [28] B. N. PARLETT, *The Symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs, NJ, 1980; republished as Classics Appl. Math. 20, SIAM, 1998, <https://doi.org/10.1137/1.9781611971163>.
- [29] G. PETERS AND J. H. WILKINSON, *Inverse iteration, ill-conditioned equations and Newton’s method*, SIAM Rev., 21 (1979), pp. 339–360, <https://doi.org/10.1137/1021052>.
- [30] A. SAQELLARI-LIKOKA AND V. KARATHANASSI, *An approach for solving rank-deficient systems that enable atmospheric path delay and water vapor content estimation*, IEEE Trans. Geosci. Remote Sensing, 46 (2008), pp. 3187–3195.
- [31] G. W. STEWART, *UTV decompositions*, in Numerical Analysis 1993 (Dundee, 1993), D. F. Griffiths and G. A. Watson, eds., Pitman Res. Notes Math. Ser., Longman Sci. Tech., Harlow, 1994, pp. 225–236.
- [32] G. W. STEWART AND J. SUN, *Matrix Perturbation Theory*, Academic Press, New York, 1990.
- [33] T. STYKEL, *Numerical solution and perturbation theory for generalized Lyapunov equations*, Linear Algebra Appl., 349 (2002), pp. 155–185.

- [34] P.-Å. WEDIN, *Perturbation bounds in connection with singular value decomposition*, BIT, 12 (1972), pp. 99–111.
- [35] P.-Å. WEDIN, *Perturbation theory for pseudo-inverses*, BIT, 13 (1973), pp. 217–232.
- [36] J. WILKENING AND J. YU, *A local construction of the Smith normal form of a matrix polynomial*, J. Symbolic Comput., 46 (2011), pp. 1–12.
- [37] W. WU AND Z. ZENG, *The numerical factorization of polynomials*, J. Found. Comput. Math., 17 (2017), pp. 259–286.
- [38] Z. ZENG, *A polynomial elimination method for numerical computation*, Theoret. Comput. Sci., 409 (2008), pp. 318–331.
- [39] Z. ZENG, *Intuitive interface for solving linear and nonlinear system of equations*, in Mathematical Software — ICMS 2018, J. H. Davenport, M. Kauers, G. Labahn, and J. Urban, eds., LNCS 10931, Springer International AG, 2018, pp. 495–506.
- [40] Z. ZENG AND T.-Y. LI, *NAClab: A Matlab toolbox for numerical algebraic computation*, ACM Commun. Comput. Algebra, 47 (2013), pp. 170–173.