

# 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

**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. A rare remark in Meyer's textbook [25, page 218] accurately reflects the state of knowledge: "If  $A$  is singular, ... even a stable algorithm can result in a significant loss of information. ... [T]he small perturbation  $E$  due to roundoff makes the possibility that  $\text{rank}\kappa(A + E) > \text{rank}\kappa(A)$  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]: "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 knowing the exact matrices and vectors, as shown in case studies in §3. The obvious difficulty in solving a singular linear system from empirical data is the condition number infinity so that the error is unbounded when solving for an individual vector solution. While this error analysis in itself is 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 a 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

---

\*Department of Mathematics, Northeastern Illinois University, Chicago, Illinois 60625, USA. email: [zzeng@neiu.edu](mailto:zzeng@neiu.edu). Research is supported in part by NSF under grant DMS-1620337.

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 SVD [9, §5.5.4][10] produce the accurate vector component and numerical rank-revealing algorithms [2, 7][9, §5.4.6][18, 19, 20, 31] provide the numerical kernel as the remaining component.

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

**2. Preliminaries.** Column vectors are denoted by boldface lower case 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 upper case letters such as  $A$ ,  $B$ ,  $X$ , etc with  $O$  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  $\mathcal{R}ange(A)$ ,  $\mathcal{K}ernel(A)$ ,  $rank(A)$  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  $rank(A) < n$  so that  $\mathcal{K}ernel(A) \neq \{\mathbf{0}\}$ , including non-square cases where  $m < n$  or  $m > n$ . The system is *consistent* if  $\mathbf{b} \in \mathcal{R}ange(A)$ .

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

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

assuming  $\theta$  does not equal to any singular value of  $A$ . Let  $U\Sigma V^H$  be the singular value decomposition of  $A$  where  $U = [\mathbf{u}_1, \dots, \mathbf{u}_m]$  and  $V = [\mathbf{v}_1, \dots, \mathbf{v}_n]$ . If  $rank_\theta(A) = r$  within  $\theta$ , then the  $\theta$ -*projection*  $A_\theta$  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  $\theta$  is  $\mathcal{K}ernel(A_\theta) = span\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$  where  $span\{\dots\}$  denotes the vector space spanned by vectors in the list. The entities  $rank_\theta(A)$ ,  $A_\theta$  and  $\mathcal{K}ernel(A_\theta)$  are undefined if  $\theta$  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  $AA^\dagger A = A$ ,  $A^\dagger AA^\dagger = A^\dagger$ ,  $(AA^\dagger)^\mathfrak{H} = AA^\dagger$  and  $(A^\dagger A)^\mathfrak{H} = A^\dagger A$ . Using the singular value decomposition as above and assuming  $\text{rank}(A) = r$ , the identity [9, §5.5.2]

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

holds and  $X = A^\dagger$  is the minimum Frobenius norm matrix such that  $AX$  and  $XA$  are orthogonal projections from  $\mathbb{C}^m$  onto  $\mathcal{R}\text{ange}(A)$  and from  $\mathbb{C}^n$  onto  $\mathcal{R}\text{ange}(A^\mathfrak{H})$  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, page 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  $[P, \hat{P}]^\mathfrak{H} [P, \hat{P}] = [Q, \hat{Q}]^\mathfrak{H} [Q, \hat{Q}] = I$ . The Grassmannian  $\mathcal{G}_k(\mathbb{C}^n)$  is a metric space with the distance [9, §2.5.3]

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

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

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

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, page 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 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 shall make 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 can not 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  $\theta$  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  $\theta$ -projection  $\tilde{A}_\theta$  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  $\alpha$  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  $\alpha$  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  $A(t)X + X B(t) = C(t)$  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] (c.f. supplementary material). Let

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

When  $t$  varies continuously, the system becomes singular but still 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  $\tilde{t} \approx 0.6666$  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, §1.3][36]), such that the Bézout identity

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

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  $\mathbb{P}_k$  as 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, the 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 in turning 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 the equation (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 the equation (5) in the given data is

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

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 (c.f. supplementary material):

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

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  $\lambda_*$  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  $\lambda_*$  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 presense of annihilators [12, §2.4 and page 83] (c.f. an example in 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 nonsingular 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  $\mathcal{K}_{\text{kernel}}(\tilde{A}_\theta)$  within a proper error tolerance  $\theta > 0$  is an approximation to  $\mathcal{K}_{\text{kernel}}(A) = \text{sol}(A, \mathbf{0})$  in  $\mathcal{G}_{n-r}(\mathbb{C}^n)$  (c.f. Lemma 11 in §A in appendix). For every  $G \in \mathbb{C}^{m \times n}$ , we define

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

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  $\theta$  where  $r = \text{rank}_\theta(G)$ . Numerical methods for computing  $\mathcal{K}_{\text{kernel}}(G_\theta)$  as  $\text{sol}_\theta(G, \mathbf{0})$  are well-established, including the singular value decomposition 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  $\Delta A$  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.* A straightforward verification 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  $r = \text{rank}(A)$ ,  $\mu = \sigma_{r+1}(\hat{A}) \leq \|\Delta A\|_2$  and  $\eta = \sigma_r(\hat{A}) \geq \sigma_r(A) - \|\Delta A\|_2$ .  $\square$

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  $\text{sol}_\theta(\tilde{A}, \mathbf{0}) = \mathcal{Ker}(\tilde{A}_\theta)$  within  $\theta$ .*

*Proof.* A straightforward verification using Theorem 1.  $\square$

The error tolerance  $\theta$  in Theorem 1 is an operational parameter that needs to be set up for solving Problem 1. If we assume the underlying application allows the data error to a certain extent, say  $\|A - \tilde{A}\|_2 < \hat{\theta}$ , the data error bound  $\beta$  in Corollary 2 is expected to be below  $\hat{\theta}$ . The inequality (13) ensures there is a window  $(\mu, \eta)$  for setting the operational error tolerance  $\theta$  at  $\beta$  or slightly larger. Using the notation of Problem 1, it is reasonable to assume the data error bound  $\beta$  on  $\|A - \tilde{A}\|_2$  is known or can be estimated. The crucial criterion for operational purpose is to set  $\theta$  at or slightly above  $\|A - \tilde{A}\|_2$  according to Theorem 1, part (iii). The error tolerance  $\theta$  should not exceed  $\|A^\dagger\|_2^{-1} - \|A - \tilde{A}\|_2$  whose exact value or estimation is not needed if the data error bound  $\beta$  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 cases of  $m < n$  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  $\Delta A$  and the condition  $\|A\|_2 \|A^\dagger\|_2$  is discontinuous. The assertions of Theorem 1 remain the same either  $m < n$  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 is originated in [15] by Kahan who suggests 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  $\text{sol}(A, \mathbf{b}) = \mathbf{x}_0 + \mathcal{Ker}(A)$  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

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

Notice that the component  $\text{dist}(\mathcal{K}\text{ernel}(A), \mathcal{K}\text{ernel}(B)) \leq 1$  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  $\omega$  on the component  $\|A^\dagger \mathbf{b} - B^\dagger \mathbf{d}\|_2$ . We choose not to use weights for the sake of simplicity of elaborations and for the reason that the weight  $\omega$  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

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

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 notations 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  $\mathbf{b} \in \mathcal{R}\text{ange}(A)$ . Assume the perturbation  $(\Delta A, \Delta \mathbf{b})$  is constrained such that:  $\text{rank}(\tilde{A}) = \text{rank}(A)$  and  $\tilde{\mathbf{b}} \in \mathcal{R}\text{ange}(\tilde{A})$  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  $\mathcal{K}\text{ernel}(A)$ . 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. Detailed proof is 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 (44) 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  $\mathcal{M}_r^{m \times n} := \{(A, \mathbf{b}) \in \mathbb{C}^{m \times n} \times \mathbb{C}^m \mid \text{rank}(A) = r, \mathbf{b} \in \mathcal{R}\text{ange}(A)\}$  is a complex analytic manifold in  $\mathbb{C}^{m \times n} \times \mathbb{C}^m$ . Although the problem of solving a singular linear system in general is ill-posed, Theorem 3 implies 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 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

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

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

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 in the same affine Grassmannian or both empty. It is natural to choose  $\hat{A} = \tilde{A}_\theta$  within a proper  $\theta$  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  $\mathcal{R}\text{ange}(\tilde{A}_\theta)$ . 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  $\theta$  is defined as

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

where  $G_\theta$  is the  $\theta$ -projection of  $G$  and  $\mathbf{d}_\theta = G_\theta G_\theta^\dagger \mathbf{d}$  is the orthogonal projection of  $\mathbf{d}$  onto the range  $\mathcal{R}\text{ange}(G_\theta)$  of  $G_\theta$ .

The solution  $\text{sol}_\theta(G, \mathbf{d})$  is undefined if  $\theta$  equals 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  $\text{sol}(A, \mathbf{b}) \equiv \text{sol}_\theta(A, \mathbf{b})$  holds for all  $\theta < \|A^\dagger\|_2^{-1}$  if  $\mathbf{b} \in \mathcal{R}\text{ange}(A)$ , or  $\theta < \min \{\|A^\dagger\|_2^{-1}, \|\mathbf{b} - AA^\dagger \mathbf{b}\|_2\}$  otherwise.*
- (ii) Computing the general numerical solution is a well-posed problem: *Assume  $\text{sol}_\theta(A, \mathbf{b})$  is well-defined within a certain  $\theta > 0$ . There is a  $\xi > 0$  depending on  $A, \mathbf{b}$  and  $\theta$  such that, for every  $(\Delta A, \Delta \mathbf{b})$  with a sufficiently small norm, there exists a unique  $\text{sol}_\theta(A + \Delta A, \mathbf{b} + \Delta \mathbf{b})$  satisfying the Lipschitz continuity*

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

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

(a) *Assume  $\mathbf{b} \in \mathcal{R}\text{ange}(A)$  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  $\theta$  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  $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}})$  with a backward error bound  $\omega \|(\Delta A, \Delta \mathbf{b})\|$  and a forward error bound*

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

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

*Proof sketch.* The assertion (i) and the unique existence in the 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 the bound on the numerical kernel is established by Wedin [34]. 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 Example 3, the solvability 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 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  $\mathcal{R}ange(\tilde{A}_\theta)$ . Furthermore

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

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 [18, 20] and UTV/ULV decomposition[9, §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  $\mathcal{K}ernel(A_\theta)$   
 – solve  $A\mathbf{x} = \mathbf{b}$  for a particular solution  $\mathbf{x} = \mathbf{x}_*$  by any backward accurate method such as  $\mathbf{x}_* = (A^H A + \mu^2 N N^H)^{-1} A^H \mathbf{b}$  or Tikhonov regularization  
 – **output**  $\text{sol}_\theta(A, \mathbf{b}) = \mathbf{x}_* + \mathcal{R}ange(N)$ .  
**else**  
 – calculate a decomposition  $U S V^H = A_\theta$  with  $S \in \mathbb{C}^{r \times r}$ ,  $U^H U = I$  and  $V^H V = I$   
 – solve  $S \mathbf{y} = U^H \mathbf{b}$  for  $\mathbf{y} = \mathbf{y}_*$  and obtain the truncated SVD solution  $\mathbf{x}_* = V \mathbf{y}_*$ .  
 – **output**:  $\text{sol}_\theta(A, \mathbf{b}) = \mathbf{x}_* + \mathcal{R}ange(V)^\perp$   
**end if**

As we shall establish in §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 general numerical solution is implemented in the Matlab package NACLAB [40] as the functionality `LinearSolve` (c.f. [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, 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  $5 \times 10^{-5}$ . The error tolerance  $\theta$  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

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

(c.f. 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

$$\begin{bmatrix}
.55101 - .91839 x^2, & -2.33985 + .70128 x + .30047 x^2 + .00001 x^3, & .71342 + 1.83982 x + 1.12986 x^2 + .00002 x^3 \\
-.36735 + .18367 x - .73471 x^5, & -.43135 - 1.09668 x + .33663 x^2 - 1.72768 x^3 + .56105 x^4 + 0.24037 x^5, & \\
.18366 + .36734 x^2 + .55103 x^4, & 1.58108 - .53294 x + 1.17553 x^2 - .42079 x^3 - .18027 x^4, & .90389 x^5 \\
& & -1.28338 - 1.39825 x - 1.28672 x^2 - 1.10389 x^3 - .6779 x^4
\end{bmatrix}.$$

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. The sensitivity of the least squares solution is well-known to be  $\kappa(A)^2$  (see, e.g. [14, §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 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, §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  $A\mathbf{x} = \mathbf{0}$  in the sense that there is a pair  $(\tilde{A}, \tilde{\mathbf{b}})$  such that  $\tilde{A}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$  and  $\|(A, \mathbf{b}) - (\tilde{A}, \tilde{\mathbf{b}})\|$  is at an acceptable level, we shall 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  $A^\dagger \mathbf{b}$ ,  $\tilde{\mathbf{x}}_2 \in \mathcal{Ker}(\tilde{A}_\theta)$ ,  $(I - A^\dagger A) \tilde{\mathbf{x}}_2 \in \mathcal{Ker}(A)$  and  $\mathcal{Ker}(\tilde{A}_\theta)$  approximates  $\mathcal{Ker}(A)$ , 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 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 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

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

Case (iii):  $\tilde{\mathbf{b}} \neq \mathbf{0}$ . The 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  $\|A - \tilde{A}\|_2$ , 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 turns out pleasantly, 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  $\|\tilde{\mathbf{x}}\|_2$  achieves, the more accurate  $\frac{\tilde{\mathbf{x}}}{\|\tilde{\mathbf{x}}\|_2}$  is 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  $\mathcal{Ker}(A)$ . As the following corollary states, the high sensitivity is near a direction in  $\mathcal{Ker}(A)$  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  $\mathbf{b} \in \mathcal{R}ange(A)$ . 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  $\mathbf{x}_* \in \mathcal{Ker}(A)$  such that*

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

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

Theorem 8 extends the accuracy result for the inverse iteration in spite of the large condition number. In [29], Peters and Wilkinson described what they called “exaggerated fears” in the early days of 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  $\lambda$  to an eigenvalue  $\lambda_*$  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  $\lambda$  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, and pleasantly surprising nonetheless, that ill-condition is not harmful in computing the eigenvector. As pointed out in [29] and by Parlett [28, §4.3] that errors mainly lie in  $\mathcal{X}_{\text{kernel}}(A - \lambda_* I)$  and are not really errors at all:

[R]oundoff errors can give rise to completely erroneous “solutions” to very ill-conditioned systems of equations. ... Indeed some textbooks have cautioned users not 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]. ... The 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 fears of solving a highly ill-conditioned linear system may also be exaggerated for non-homogeneous systems as well 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  $0.5 \times 10^{-5}$ . Rounded to five digits after the decimal point, two numerical particular solutions  $\tilde{\mathbf{x}}_0$  and  $\tilde{\mathbf{x}}_1$  by truncated SVD  $\tilde{A}_\theta^\dagger \tilde{\mathbf{b}}$  and Matlab “\”, 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





using such data. Table 1 shows three sample numerical solutions:  $\mathbf{x}_1$  by a straightforward application of the Matlab command  $\mathbf{A} \setminus \mathbf{b}$ , 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  $\varepsilon$  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  $\text{sol}_\theta(\tilde{A}, \tilde{\mathbf{b}}) = \tilde{\mathbf{u}} + \mathcal{Ker}(\tilde{A}_\theta)$  with  $\mathcal{Ker}(\tilde{A}_\theta) = \text{span}\{\tilde{\mathbf{v}}\}$  where

$$(32) \quad \tilde{\mathbf{v}} = [0, -0.0000001, 0.0000010, -0.0000099, 0.0000995, -0.0009950, 0.0099499, -0.0994987, 0.9949875]^H.$$

All  $\mathbf{x}_j + \mathcal{Ker}(A_\theta)$  for  $j = 1, 2, 3$  are nearly identical in the affine Grassmannian  $A_1(\mathbb{C}^9)$  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  $8.28 \times 10^{-7}$  predicted by (31).

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 $\mathbf{x}_*$	.3333333 .6666667 1.0000000 1.3333333 1.6666667 2.0000000 2.3333333 2.6666667 3.0000000	
Matlab “\” $\mathbf{x}_1$	.3333333 .6666666 1.0000044 1.3333487 1.6668429 1.9985374 2.3479633 2.5203667 4.4629993	0.2612916
Tikhonov $\mathbf{x}_2$	.3333333 .6666670 0.9999974 1.3333607 1.6662034 2.0027277 2.3060572 2.4294228 0.2724202	0.4871437
trunc. SVD $\mathbf{x}_3$	.3333333 0.6666666 0.9999967 1.3333603 1.6663028 2.0027276 2.3060572 2.4294228 0.2724206	0.4870902
	particular $\hat{\mathbf{x}}_j = \mathbf{x}_j + t_j \tilde{\mathbf{v}} \in \mathbf{x}_j + \mathcal{Ker}(A_\theta)$ nearest to $\mathbf{x}_*$ with $\tilde{\mathbf{v}}$ in (32)	error $\frac{\ \hat{\mathbf{x}}_j - \mathbf{x}_*\ _2}{\ \mathbf{x}_*\ _2}$
$\mathbf{x}_1 + \mathcal{Ker}(A_\theta)$	$\hat{\mathbf{x}}_1 = \mathbf{x}_1 + t_1 \tilde{\mathbf{v}} \approx \mathbf{x}_*$ with $t_1 = -1.4703701$	$8.8 \times 10^{-8}$
$\mathbf{x}_2 + \mathcal{Ker}(A_\theta)$	$\hat{\mathbf{x}}_2 = \mathbf{x}_2 + t_2 \tilde{\mathbf{v}} \approx \mathbf{x}_*$ with $t_2 = 2.7413113$	$1.5 \times 10^{-7}$
$\mathbf{x}_3 + \mathcal{Ker}(A_\theta)$	$\hat{\mathbf{x}}_3 = \mathbf{x}_3 + t_3 \tilde{\mathbf{v}} \approx \mathbf{x}_*$ with $t_3 = 2.7410104$	$1.9 \times 10^{-7}$

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  $\mathbf{x}_*$ , as well as the accuracies of  $\mathbf{x}_1$ ,  $\mathbf{x}_2$  and  $\mathbf{x}_3$  as a component of the general numerical solution.

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  $0.2727296 + 0.9949875t = 3$ , obtaining  $t = 2.7410097$  corresponding to a numerical solution with a relative error  $1.79 \times 10^{-7}$  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

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

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

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

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

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

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

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

and thus

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

LEMMA 12. Let  $\mathcal{U}$  be a subspace of  $\mathbb{C}^n$  and  $U$  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  $[G, \hat{G}]$  be a unitary matrix. Then, for any unit vector  $\mathbf{x} \in \mathbb{C}^n$ ,

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

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

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

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

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  $\mathcal{K}\text{ernel}(A_\theta)$ . For any  $\mu > 0$ , let  $\mathbf{b} \in \mathbb{C}^m$  and  $\mathbf{x}_*$  be the least squares solution of the linear system

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

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  $\mathcal{R}\text{ange}(A_\theta)$ . Furthermore,

$$(38) \quad \mathbf{x}_* - T T^H \mathbf{x}_* = A_\theta^\dagger \mathbf{b}_\theta$$

for any  $T \in \mathbb{C}^{n \times (n-r)}$  with  $\mathcal{R}\text{ange}(T) = \mathcal{K}\text{ernel}(A_\theta)$  and  $T^H T = I$ .

*Proof.* For the case of  $m \geq n$ , we can write  $A$  in its singular value expansion  $A = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^H + \cdots + \sigma_n \mathbf{u}_n \mathbf{v}_n^H$  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  $\mathbf{x} = x_1 \mathbf{v}_1 + \cdots + x_n \mathbf{v}_n$ . Then

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

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 the assertion (i). The case  $m < n$  is similar.

To prove (ii), write the singular value decomposition  $A = U_1 \Sigma_1 V_1^H + U_2 \Sigma_2 V_2^H$  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

$$(39) \quad (V_1 \Sigma_1^H \Sigma_1 V_1^H \mathbf{x}_* - V_1 \Sigma_1^H U_1^H \mathbf{b}) + (V_2 \Sigma_2^H \Sigma_2 V_2^H \mathbf{x}_* - V_2 \Sigma_2^H U_2^H \mathbf{b} + \mu^2 N N^H \mathbf{x}_*) = \mathbf{0},$$

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  $\mathbf{x}_* = V_1 V_1^H \mathbf{x}_* + V_2 V_2^H \mathbf{x}_*$ , hence  $A_\theta \mathbf{x}_* = (U_1 \Sigma_1 V_1^H) (V_1 V_1^H \mathbf{x}_*) = U_1 U_1^H \mathbf{b} = \mathbf{b}_\theta$ . Namely  $\mathbf{x}_*$  is a particular solution of the system  $A_\theta \mathbf{x} = \mathbf{b}_\theta$ . Also,

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

Finally, if  $\text{rank}(A) = r$ , then  $A_\theta = A$  so  $\Sigma_2 = O$  in (39), implying  $N^H \mathbf{x}_* = \mathbf{0}$ . Consequently  $\mathbf{x}_* = A^\dagger \mathbf{b}$  from (38).  $\square$

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, Hansen). *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*

$$(40) \quad \zeta = \|A_\theta^\dagger \mathbf{b}\|_2 + \frac{1 + \|A_\theta^\dagger \mathbf{b}\|_2}{1 - \|A_\theta^\dagger\|_2 \|A - A_\theta\|_2}$$

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

$$(41) \quad \|\Delta A\|_2 < \min \left\{ \frac{1}{2} (\sigma_r(A) - \sigma_{r+1}(A)), \sigma_r(A) - \theta, \theta - \sigma_{r+1}(A) \right\},$$

the following inequality holds:

$$(42) \quad \|A_\theta^\dagger \mathbf{b} - \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}\|_2 \leq \frac{\sigma_1(A)}{\sigma_r(A)} \left( \zeta \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|A\|_2} \right) + O(\|(\Delta A, \Delta \mathbf{b})\|^2).$$

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

$$(43) \quad \|A^\dagger \mathbf{b} - \tilde{A}_\theta^\dagger \tilde{\mathbf{b}}\|_2 \leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\|\Delta A\|_2}{\sigma_r(A)}} \left( 2 \|A^\dagger \mathbf{b}\|_2 \frac{\|\Delta A\|_2}{\|A\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|A\|_2} \right).$$

*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 (41) so that  $\text{rank}_{\theta}(\tilde{A}) = r$  as well. Then it is straightforward to verify (42) 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 (43) follows from [10, inequality (27a)] and  $\mathbf{b} \in \mathcal{R}\text{ange}(A)$ .  $\square$

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),

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

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  $\text{sol}_{\theta}(\tilde{A}, \tilde{\mathbf{b}}) = \tilde{A}_{\theta}^{\dagger} \tilde{\mathbf{b}}_{\theta} + \mathcal{K}\text{ernel}(\tilde{A}_{\theta})$  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.* (p.9) Let  $N \in \mathbb{C}^{n \times (n-r)}$  whose columns form an orthonormal basis for  $\mathcal{K}\text{ernel}(A)$ . By Lemma 12, there is an  $\tilde{N} \in \mathbb{C}^{n \times (n-r)}$  whose columns form a basis for  $\mathcal{K}\text{ernel}(\tilde{A})$  such that  $\|N - \tilde{N}\|_2 \leq \text{dist}\left(\mathcal{K}\text{ernel}(A), \mathcal{K}\text{ernel}(\tilde{A})\right)$ . For  $\zeta = \sigma_r(A)$ , denote  $B = \begin{bmatrix} \zeta N^{\text{H}} \\ A \end{bmatrix}$  and  $\tilde{B} = \begin{bmatrix} \zeta \tilde{N}^{\text{H}} \\ \tilde{A} \end{bmatrix}$ . 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 part (ii). By  $\|B - \tilde{B}\|_2 \leq \sqrt{2} \|\Delta A\|_2$  from Lemma 11 part (ii), [24, Theorem 1.4.6, page 30] and Lemma 13 part (i),

$$(44) \quad \begin{aligned} \left\| B^\dagger \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} - \tilde{B}^\dagger \begin{bmatrix} \mathbf{0} \\ \tilde{\mathbf{b}} \end{bmatrix} \right\|_2 &\leq \frac{\sigma_1(B)}{\sigma_n(B)} \frac{1}{1 - \frac{\|B - \tilde{B}\|_2}{\sigma_n(B)}} \left( \|\mathbf{x}_*\|_2 \frac{\|B - \tilde{B}\|_2}{\|B\|_2} + \frac{\|\Delta \mathbf{b}\|_2}{\|B\|_2} \right) \\ &\leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\sqrt{2} \|\Delta A\|_2}{\sigma_r(A)}} \frac{\sqrt{2} \|\mathbf{x}_*\|_2 \|\Delta A\|_2 + \|\Delta \mathbf{b}\|_2}{\|A\|_2} \end{aligned}$$

leading to (17).  $\square$

*Proof Theorem 6.* (p. 11) The assertion (i) is true because  $A_\theta = A$  and  $\mathbf{b}_\theta = \mathbf{b}$  for  $\theta \in (0, \sigma_r(A))$ . If  $\mathbf{b} \in \mathcal{R}ange(A)$ , then

$$sol_\theta(A, \mathbf{b}) = sol(A, \mathbf{b}) = A^\dagger \mathbf{b} + \mathcal{K}ernel(A).$$

Otherwise  $sol_\theta(A, \mathbf{b}) = sol(A, \mathbf{b}) = \emptyset$  if  $\theta < \min\{\sigma_r(A), \|A A^\dagger \mathbf{b} - \mathbf{b}\|_2\}$ . The assertion (ii) directly follows from Lemma 14 and Lemma 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 the assertion (iii), part (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  $\|\tilde{A} - \tilde{A}_\theta\|_2 = \min_{rank(B)=r} \|\tilde{A} - B\|_2 \leq \|\Delta A\|_2$ , we have

$$(44) \quad \begin{aligned} \|\tilde{\mathbf{b}} - \tilde{\mathbf{b}}_\theta\|_2 &= \|\tilde{P} \tilde{\mathbf{b}}\|_2 = \|\tilde{P} \tilde{\mathbf{b}} - P \mathbf{b}\|_2 \leq \|\tilde{P}\|_2 \|\tilde{\mathbf{b}} - \mathbf{b}\|_2 + \|\tilde{P} - P\|_2 \|\mathbf{b}\|_2 \\ &\leq \|\Delta \mathbf{b}\|_2 + \text{dist}(\mathcal{R}ange(\tilde{A}_\theta), \mathcal{R}ange(A)) \|\mathbf{b}\|_2 \\ \text{(by (33))} \quad &\leq \|\Delta \mathbf{b}\|_2 + \frac{\sigma_1(A)}{\sigma_r(A)} \frac{2 \|\mathbf{b}\|_2}{\|A\|_2} \|\Delta A\|_2 \\ &\leq \sqrt{4 \|A^\dagger\|_2^2 \|\mathbf{b}\|_2^2 + 1} \|(\Delta A, \Delta \mathbf{b})\| = \sqrt{\omega^2 - 1} \|(\Delta A, \Delta \mathbf{b})\| \end{aligned}$$

and

$$(45) \quad \begin{aligned} \|(\tilde{A}, \tilde{\mathbf{b}}) - (\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)\| &\leq \sqrt{\|\Delta A\|_2^2 + (\omega^2 - 1) (\|\Delta A\|_2^2 + \|\Delta \mathbf{b}\|_2^2)} \\ &\leq \omega \|(\Delta A, \Delta \mathbf{b})\| < \sigma_r(A) - \|(\Delta A, \Delta \mathbf{b})\| \end{aligned}$$

Then, for any  $\theta$  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  $rank_\theta(\tilde{A}) = r$ ,  $\|(\tilde{A}, \tilde{\mathbf{b}}) - (\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)\| < \theta$  and thus  $sol_\theta(\tilde{A}, \tilde{\mathbf{b}})$  is of the same dimension as  $sol(A, \mathbf{b})$ . Since  $sol_\theta(\tilde{A}, \tilde{\mathbf{b}}) = sol(\tilde{A}_\theta, \tilde{\mathbf{b}}_\theta)$ , the backward error of  $sol_\theta(\tilde{A}, \tilde{\mathbf{b}})$  is bounded above by  $\omega \|(\Delta A, \Delta \mathbf{b})\|$  from (45). Thus (23) follows from (33) in Lemma 11 and (43) in Lemma 14, leading to the assertion (iii).

We now prove the assertion (b) of part (iii). If  $sol(A, \mathbf{b})$  is empty, then  $sol_\theta(A, \mathbf{b}) = \emptyset$  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.* (p. 13) Let  $\tilde{A} = \tilde{U}_1 \tilde{\Sigma}_1 \tilde{V}_1^H + \tilde{U}_2 \tilde{\Sigma}_2 \tilde{V}_2^H$  be the singular value decomposition where  $\tilde{\Sigma}_1$  is  $r \times r$  with  $r = \text{rank}(A)$ . Then  $\tilde{\mathbf{x}}$  is a solution to  $\tilde{A} \tilde{\mathbf{x}} = \tilde{\mathbf{b}}$  implies  $\tilde{U}_1 \tilde{\Sigma}_1 \tilde{V}_1^H \tilde{\mathbf{x}}_1 = \tilde{U}_1 \tilde{U}_1^H \tilde{\mathbf{b}}$  and  $\tilde{U}_2 \tilde{\Sigma}_2 \tilde{V}_2^H \tilde{\mathbf{x}}_2 = \tilde{U}_2 \tilde{U}_2^H \tilde{\mathbf{b}}$  where  $\tilde{\mathbf{x}}_1 = \tilde{V}_1 \tilde{V}_1^H \tilde{\mathbf{x}}$  and  $\tilde{\mathbf{x}}_2 = \tilde{V}_2 \tilde{V}_2^H \tilde{\mathbf{x}}$ . Then  $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}_1 + \tilde{\mathbf{x}}_2$  with  $\tilde{\mathbf{x}}_1 = \tilde{A}_1^\dagger \tilde{\mathbf{b}}$  for any  $\theta$  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  $\mathcal{K}(\text{Kernel}(A))$ . Since  $\tilde{\mathbf{x}}_2 \in \mathcal{K}(\text{Kernel}(\tilde{A}_\theta))$ ,

$$\begin{aligned} \|N N^H \tilde{\mathbf{x}}_2 - \tilde{\mathbf{x}}_2\|_2 &= \min_{\mathbf{u} \in \mathcal{K}(\text{Kernel}(A))} \|\mathbf{u} - \tilde{\mathbf{x}}_2\|_2 \\ &\leq \text{dist}(\mathcal{K}(\text{Kernel}(A_\theta)), \mathcal{K}(\text{Kernel}(A))) \|\tilde{\mathbf{x}}_2\|_2 \\ (46) \quad &\leq \frac{\sigma_1(A)}{\sigma_r(A)} \frac{1}{1 - \frac{\|\Delta A\|_2}{\sigma_r(A)}} \frac{\|\Delta A\|_2}{\|A\|_2} \|\tilde{\mathbf{x}}_2\|_2 \end{aligned}$$

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

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

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

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

leading to (24). For the case  $\mathbf{b} = \mathbf{0}$ , the bound (25) follows from (46)

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

*Proof of Theorem 10.* (p. 16) 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 singular value decompositions 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

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

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. on Automatic Control, 47 (2002), pp. 1116–1119.
- [2] J. BARLOW, H. ERBAY, AND I. SLAPNICAR, *An alternative algorithm for the refinement of ULV decompositions*, SIAM J. Matrix Anal. Appl., 27 (2005), pp. 198–211.
- [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*, Mathematics of Computation, 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 Alg. and its 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.
- [7] R. D. FIERRO, P. C. HANSEN, AND P. S. K. HANSEN, *UTV Tools: Matlab templates for rank-revealing UTV decompositions*, Numerical 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*, The John Hopkins University Press, Baltimore and London, 4th ed., 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*, SIAM, Philadelphia, 1997.
- [12] P. C. HANSEN, *Discrete Inverse Problems. Insight and Algorithms*, SIAM, Philadelphia, 2010.
- [13] P. C. HANSEN, J. G. NAGY, AND D. P. O’LEARY, *Deblurring Images, Matrices, Spectra, and Filtering*, SIAM, Philadelphia, 2006.
- [14] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM, 2nd ed., 2002.
- [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.

- [19] T.-L. LEE, T.-Y. LI, AND Z. ZENG, *RankRev — A Matlab package for computing numerical ranks*, Numerical 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.
- [21] L.-H. LIM, K. S.-W. WONG, AND K. YE, *Numerical algorithms on the affine Grassmannian*, SIAM J. Matrix Anal. Appl., 40 (2019) pp. 371–393, DOI 10.1137/18M1169321.
- [22] L.-H. LIM, K. S.-W. WONG, AND K. YE, *The Grassmannian of affine subspaces*. arXiv:1807.10883.
- [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. on Neural Networks and Learning systems, 29 (2018), pp. 2271–2277.
- [24] ÅKE BJÖRCK, *Numerical Methods for Least Squares Problems*, SIAM, Philadelphia, 1996.
- [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 Univ. Press, London, 2003.
- [27] A. NEUMAIER, *Solving ill-conditioned and singular linear systems: A tutorial on regularization*, SIAM Review, 40 (1998), pp. 636–666.
- [28] B. N. PARLETT, *The Symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs, N.J., 1980.
- [29] G. PETERS AND J. H. WILKINSON, *Inverse iteration, ill-conditioned equations and Newton’s method*, SIAM Review, 21 (1979), pp. 339–360.
- [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. Geoscience and Remote Sensing, 46 (2008), pp. 3187–3195.
- [31] G. W. STEWART, *UTV decompositions*, in Numerical Analysis, 1993, D. F. Griffith and G. Watson, eds., Pitman Research Notes in Mathematical Sciences, New York, 1994, 1994.
- [32] G. W. STEWART AND J. SUN, *Matrix Perturbation Theory*, Academic Press, Inc, Boston, San Diego, New York, London, Sydney, Tokyo, Toronto, 1990.
- [33] T. STYKEL, *Numerical solution and perturbation theory for generalized Lyapunov equations*, Linear Alg. and Its 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 Computation, 46 (2011), pp. 1–12.
- [37] W. WU AND Z. ZENG, *The numerical factorization of polynomials*, J. Foundation of Computational Mathematics, 17 (2017), pp. 259–286.
- [38] Z. ZENG, *A polynomial elimination method for numerical computation*, Theoretical Computer Science, 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 Communications in Computer Algebra, 47 (2013), pp. 170–173.



# Online Supplement to “On the Sensitivity of Singular and Ill-Conditioned Linear Systems”

Zhonggang Zeng<sup>1</sup>

**Abstract.** This online supplement provides a software demo of the package NACLAB and additional computing examples for calculating numerical solutions of singular and ill-conditioned linear systems.

In this online supplementary material, we briefly introduce the software package NACLAB in the context of solving singular linear systems for the general numerical solution elaborated in the paper *On the sensitivity of singular and ill-conditioned linear systems*. All the example numbers point to the examples in the paper and all the citation numbers point to the references of the paper.

**1. NACLAB functionality LinearSolve.** NACLab<sup>2</sup> is a software package of Matlab functions for numerical algebraic computation [40]. We implemented the computation of general numerical solution  $sol_\theta(A, \mathbf{b})$  as a functionality `LinearSolve` [39] in a simple call with input  $A$ ,  $\mathbf{b}$  and  $\theta$ :

```
>> [x0, N, lcmd, res] = LinearSolve(A, b, theta)
```

The output `x0`, `N`, `lcmd` and `res` carries  $A_\theta^\dagger \mathbf{b}$ ,  $\mathcal{K}ernel(A_\theta)$  spanned by the orthonormal columns of  $N$ , the sensitivity estimate  $\|A_\theta\|_2 \|A_\theta^\dagger\|_2$  and the residual  $\max\{\|A \mathbf{x}_0 - \mathbf{b}\|_2, \|A N\|_2\}$  respectively. The functionality `LinearSolve` follows from the high-rank case of the template in §6.

Furthermore, `LinearSolve` provide a mechanism to solve a linear system in the form of

$$L(\mathbf{u}_1, \dots, \mathbf{u}_m) = (\mathbf{b}_1, \dots, \mathbf{b}_n) \quad \text{for} \quad (\mathbf{u}_1, \dots, \mathbf{u}_m)$$

where  $L : \mathcal{U}_1 \times \dots \times \mathcal{U}_m \rightarrow \mathcal{V}_1 \times \dots \times \mathcal{V}_n$  is a linear transformation and  $\mathcal{U}_1, \dots, \mathcal{U}_m, \mathcal{V}_1, \dots, \mathcal{V}_n$  are vector spaces of column vectors, matrices, or polynomials in a call syntax

```
>> [x0, N, lcmd, res] = LinearSolve({L, domain, parameter}, b, theta)
```

where `L` is a Matlab (anonymous) function for evaluating the linear transformation  $L$  along with `domain` and `parameter` in cell arrays representing the domain  $\mathcal{U}_1 \times \dots \times \mathcal{U}_m$  and parameters of  $L$ .

**2. Supplement to Example 1.** Consider a system of polynomial equations  $\mathbf{f}(x, y) = \mathbf{0}$  that is known through empirical data in the perturbed system  $\tilde{\mathbf{f}}(x, y) = \mathbf{0}$  with

$$\tilde{\mathbf{f}}(x, y) = \begin{bmatrix} x^3 + y - 0.7698 \\ x + y^3 - 0.7698 \end{bmatrix}$$

and the coefficientwise error bound  $5 \times 10^{-7}$ . A multiple zero

$$(x_*, y_*) \approx (\tilde{x}, \tilde{y}) = (0.57735, 0.57735)$$

---

<sup>1</sup>Department of Mathematics, Northeastern Illinois University, Chicago, Illinois 60625, USA. email: [zzeng@neiu.edu](mailto:zzeng@neiu.edu). Research is supported in part by NSF under grant DMS-1620337.a

<sup>2</sup><http://homepages.neiu.edu/~zzeng/naclab.html>

is computed using the NACLAB polynomial system solver `psolve` and the depth-deflation method [4] with an error bound  $\|(x_* - \tilde{x}, y_* - \tilde{y})\|_2 \leq \varepsilon = 9.46 \times 10^{-6}$

As briefly elaborated in Example 1 and in [4], the multiplicity structure can be computed via solving a sequence of homogeneous linear systems  $S_\alpha(x_*, y_*) \mathbf{c} = \mathbf{0}$  for  $\alpha = 1, 2, \dots$  from Macaulay matrices  $S_\alpha(\tilde{x}, \tilde{y})$  serving as empirical data. The NACLAB functionality `MacaulayMatrix` is built for generating the Macaulay matrices. To construct, say  $S_2(\tilde{x}, \tilde{y})$ , use the following statements:

```
>> f = {'x^3+y-0.7698', 'x+y^3-0.7698'}; % cell array of polynomials in character strings
>> var = {'x', 'y'}; % cell array of variable names in character stings
>> z = [.57735, .57735]; % approximate zero
>> M = MacaulayMatrix(z, f, var, 2); % generate the Macaulay matrix
>> single(full(M)) % display the matrix in single precision

ans =

    0    1.0000000    0.9999990         0         0    1.7320499
    0    0.9999990    1.0000000    1.7320499         0         0
    0         0         0    1.0000000    0.9999990         0
    0         0         0    0.9999990    1.0000000         0
    0         0         0         0    1.0000000    0.9999990
    0         0         0         0    0.9999990    1.0000000
```

It is a straightforward to verify that the entrywise error on  $S_2(\tilde{x}, \tilde{y})$  is bounded by  $6\varepsilon$  on 8 entries. As a result, we have an error bound for the error tolerance

$$\|S_2(x_*, y_*) - S_2(\tilde{x}, \tilde{y})\|_2 \leq \|S_2(x_*, y_*) - S_2(\tilde{x}, \tilde{y})\|_F \leq \sqrt{8} 6 \cdot 9.46 \times 10^{-6} \approx 1.6 \times 10^{-4}.$$

Set the error tolerance slightly larger at

$$\theta = 2 \times 10^{-4} = 0.0002.$$

Then a one-line call of `LinearSolve` produces the matrix  $N$  whose columns form an orthonormal basis for the numerical solution  $\text{sol}_\theta(S_2(\tilde{x}, \tilde{y}), \mathbf{0})$  in the Grassmannian  $\mathcal{G}_3(\mathbb{C}^6)$ .

```
>> [z, N, lcond, res] = LinearSolve(M, zeros(6, 1), 2e-4); % solve M*z = 0 within 2e-4
>> single(N) % display solution basis in single precision

ans =

    1.0000000   -0.0000000    0.0000000
    0   -0.7828174    0.2320854
    0    0.5924006    0.5618970
    0    0.1099370   -0.4584060
    0   -0.1099372    0.4584060
    0    0.1099374   -0.4584059
```

The multiplicity of  $(x_*, y_*)$  is thus 3, with the dual space  $\mathcal{D}_{\mathbf{f}, (x_*, y_*)}$  accurately represented by the basis

$$\begin{aligned} &1, \quad .78282 \partial_x + .59240 \partial_y + .10994 \frac{1}{2!} \partial_{x^2} - .10994 \partial_{xy} + .10994 \frac{1}{2!} \partial_{y^2} \\ &.23209 \partial_x + .56190 \partial_y - .45840 \frac{1}{2!} \partial_{x^2} + .45840 \partial_{xy} - .45841 \frac{1}{2!} \partial_{y^2} \end{aligned}$$

so that those differential operators vanish on the entire ideal generated by the polynomial system at the zero point  $(x_*, y_*)$ .

**3. Supplement to Example 2.** The linear equation  $A(t)X + X B(t) = C(t)$

can be written as  $L(X) = C(t)$  where  $L$  is a linear transformation

$$\begin{aligned} L &: \mathbb{C}^{2 \times 2} &\longrightarrow &\mathbb{C}^{2 \times 2} \\ X &\longmapsto &A(t)X + X B(t) \end{aligned}$$

with the domain is  $\mathbb{C}^{2 \times 2}$  and parameters  $A(t)$ ,  $B(t)$  and  $C(t)$  in (3), The linear system can be solved by constructing the representation matrix and vectors for  $L$  and  $C(t)$  and solving the resulting matrix-vector equation. `LinearSolve` in `NACLAB` provides an intuitive WYSIWYG approach for solving the equation directly. The matrix-vector representation is generated internally.

At the hypothetical  $\tilde{t} \approx 0.6666$  with an error bound 0.0001 of a  $4 \times 4$  system, it clearly safe to say the 2-norm data error bound is  $\theta = 10 \cdot 0.0001 = 10^{-3}$  that can be used as the error tolerance for the general numerical solution.

```
>> L = @(X,t) [1 -1; 1 -1]*X*X*[-5/3+t 1; -1 -1/3+2*t];           % the linear transformation L
>> C = [1 0; 2 -1];                                               % the right-hand side
>> domain = {ones(2,2)};                                         % domain of 2x2 matrices
>> parameter = {0.6666};                                         % parameter t = 0.6666
>> [x0, N, lcnd, res] = LinearSolve({L,domain,parameter}, C, 1e-3); % solve L(X)=C within 1e-3
x0 =
    [2x2 double]
N =
    {1x1 cell}    {1x1 cell}
lcnd =
    1.573435327501125
res =
    2.499756923490804e-05
```

The underlying singular linear system is quite well conditioned with a sensitivity measure roughly 1.6 along with a residual  $\approx 2.5 \times 10^{-5}$ , implying the general numerical solution carried in `x0` and `N` are as accurate as the data.

```
>> x0{1}                % display the truncated SVD solution
ans =
    0.249983334213952   -0.250004166457633
   -0.750004165972904   -0.249974998284764

>> N{1}{1}             % display the 1st kernel component
ans =
   -0.662148424976858    0.483868243696442
   -0.483822831115126    0.305526519527407

>> N{2}{1}             % display the 2nd kernel component
ans =
    0.558171384891092    0.126097939805073
   -0.126073928483796    0.810339052016092
```

The result is an accurate approximation to the exact solution (4).

**4. Solving the system (8).** The linear system (8) can be considered as the equation

$$L(X, U) = (E, -F)$$

where  $L$  is the linear transformation

$$\begin{aligned} L &: \mathbb{C}^{3 \times 2} \times \mathbb{C}^{1 \times 2} &\longrightarrow &\mathbb{C}^{3 \times 2} \times \mathbb{C}^{1 \times 2} \\ (X, U) &\longmapsto & &(X A - B X - C U, D X) \end{aligned}$$

along with parameters

$$A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 2 & -1 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \quad D = [1 \ 0 \ -1],$$

$$E = \begin{bmatrix} 2 & 1 \\ -1 & 1 \\ 0 & 0 \end{bmatrix} \quad F = [-1 \ 0]$$

In preparation for calling `LinearSolve`, define the linear transformation, its domain and the parameter array:

```
>> L = @(X,U,A,B,C,D) {X*A-B*X-C*U, D*X}; % the linear transformation function
>> A = [1 1; 0 1]; B = [0 1 0; 0 0 1; 2 -1 0]; C = [0; 0; 1]; D = [1 0 -1]; % parameters A,B,C,D
>> E = [2 1; -1 1; 0 0]; F = [-1 0]; % right side E and F
>> domain = {ones(3,2), ones(1,2)}; % domain of 3x2 and 1x2 matrices
>> parameter = {A,B,C,D}; % parameter cell array
```

The data are exact but floating point arithmetic will introduce entrywise error around the unit roundoff  $\varepsilon \approx 2.2 \times 10^{-16}$  so that we can set the error tolerance slightly larger, say  $\theta = 10^{-10}$ . A simple call to find the general numerical solution within the error tolerance:

```
>> [Z, N, lcmd, res] = LinearSolve({L,domain,parameter}, {E, -F}, 1e-10) % solve L(X,U) = (E,-F)
Z =
 [3x2 double] [1x2 double]
N =
 {1x2 cell}
lcmd =
 11.987437447750866
res =
 1.332267629550188e-15
```

The sensitivity measure approximately 11.99 along with the residual  $1.33 \times 10^{-15}$  indicate that the numerical solution is accurate.

```
>> Z{1} % display X component of the truncated SVD solution
ans =
 1.9999999999999999 -0.3333333333333333
 0 0.6666666666666666
 0.9999999999999999 -0.3333333333333333

>> Z{2} % display U component of the truncated SVD solution
ans =
 -2.9999999999999999 1.9999999999999998

>> N{1}{1} % display the X component of the kernel basis|
ans =
 0 -0.577350269189626
 0 -0.577350269189626
 0 -0.577350269189626

>> N{1}{2} % display the U component of the kernel basis|
ans =
 0 0
```

Namely, the general numerical solution is an accurate approximation to

$$(X, U) = \left( \begin{bmatrix} 2 & -\frac{1}{3} \\ 0 & -\frac{1}{3} \\ 1 & -\frac{1}{3} \end{bmatrix}, [-3 \ 2] \right) + t \left( \begin{bmatrix} 0 & -\frac{1}{\sqrt{3}} \\ 0 & -\frac{1}{\sqrt{3}} \\ 0 & -\frac{1}{\sqrt{3}} \end{bmatrix}, [0 \ 0] \right)$$

**5. Supplement to Example 3 and Example 4.** The problem of calculating

the Bézout coefficients in Example 3 and Example 4 can be written as

$$L(u_1, u_2, u_3) = g \quad \text{for} \quad (u_1, u_2, u_3) \in \mathbb{P}_3 \times \mathbb{P}_1 \times \mathbb{P}_2$$

through the linear transformation

$$\begin{aligned} L : \mathbb{P}_3 \times \mathbb{P}_1 \times \mathbb{P}_2 &\longrightarrow \mathbb{P}_8 \\ (u_1, u_2, u_3) &\longmapsto u_1 f_1 + u_2 f_2 + u_3 f_3. \end{aligned}$$

NACLAB provides an interface for handling polynomials as character strings in WYSIWYG manner. Polynomial parameters are entered as character strings:

```
>> f1 = '2.5714 + 3.8571*x - 3*x^2 - 6.4286*x^3 - 2.1429*x^4'; % polynomials as character strings
>> f2 = '-1.7143 - 1.7143*x + 0.4286*x^2 + 0.4286*x^3 - 3.4286*x^5 - 5.1429*x^6 - 1.7143*x^7';
>> f3 = '0.8571 + 1.2857*x + 2.1429*x^2 + 2.5714*x^3 + 3.4286*x^4 + 3.8571*x^5 + 1.2857*x^6';
>> g = '4.6667 + 7*x + 2.3333*x^2'; % the known numerical gcd
```

The package NACLAB provides a library of polynomial operation functionalities, such as `pplus(...)` for adding polynomials and `ptimes(...)` for multiplying polynomials, so that the linear transformation can be defined as a Matlab (anonymous) function:

```
>> L = @(u1,u2,u3,f1,f2,f3) ... % linear transformation (u1,u2,u3) -> u1*f1 + u2*f2 + u3*f3
    pplus(ptimes(u1,f1),ptimes(u2,f2),ptimes(u3,f3));
```

To set the error tolerance, consider the entrywise error bound  $0.5 \times 10^{-4}$  on 55 nonzero entries and

$$\|A - \tilde{A}\|_2 \leq \|A - \tilde{A}\|_F \leq \sqrt{55} \cdot 0.5 \times 10^{-4} \approx 3.8 \times 10^{-4}.$$

Thus the error tolerance can be set slightly larger at, say  $5 \times 10^{-4}$ . We can then define the domain and parameter cell arrays and execute `LinearSolve` to calculate the Bézout coefficients:

```
>> domain = {'1+x+x^2+x^3','1+x','1+x+x^2'}; % domain of polynomials of degrees 3, 1, 2
>> parameter = {f1, f2, f3} % parameter cell array
>> [z0,N,lcnd,res] = LinearSolve([L,domain,parameter], g, 5e-4) % solve L(u1,u2,u3) = g
z0 =
'0.907108855304999 + 0.333222892924586*x + 0.710289197713311*x^2 + 0.599677838683852*x^3'
'-0.799463013829436 + 0.0669420537219249*x' '1.12432524246405 - 0.0664832652437786*x'
+ 0.0892574807423333*x^2'
N =
{1x3 cell} {1x3 cell}
lcnd =
20.302846223563613
res =
1.832045500993470e-05
```

The sensitivity is healthy at 20.3 with a residual  $1.8 \times 10^{-5}$  so that the error on the computed general numerical solution is at the same order of the data error. The output `z0` carries the numerical truncated SVD solution  $(u_{01}, u_{02}, u_{03})$  as shown above. The components  $(u_{11}, u_{12}, u_{13})$  and  $(u_{21}, u_{22}, u_{23})$  of the general numerical solution are in the output `N` consists of an orthonormal basis of the numerical kernel of the linear transformation  $L$  such as the result shown in Example 4. Notice that the output in `z0` and `N` carries polynomial in WYSIWYG style as character strings. There is no need for a reverse representation and interpretation of a solution in column vectors. Computation of the numerical inverse of the polynomial transformation matrix shown at the end of Example 4 is out the scope of this paper.

### 6. An application in solving an integral equation with an annihilator.

Consider a Volterra integral equation of the first kind in the form of

$$(47) \quad \int_0^s k(s-t)x(t) dt = g(s), \quad 0 \leq s \leq 1$$

for finding  $x(t)$  on the interval  $[0, 1]$  from the given kernel function  $k$  and the right-hand side function  $g$  defined on the same interval. The equation is singular if there exists an *annihilator*  $\phi(t)$  such that  $\int_0^s k(s-t)\phi(t) dt \equiv 0$  for  $0 \leq s \leq 1$ . As described in [12, pp. 82-83], the kernel<sup>3</sup>

$$(48) \quad k(\tau) = \frac{\tau^{-\frac{3}{2}} e^{-\frac{1}{4\kappa^2\tau}}}{2\kappa\sqrt{\pi}}$$

corresponds to an annihilator  $\delta(t-1)$ , the delta function at  $t=1$ . For integer  $n > 0$ , stepsize  $h = \frac{1}{n}$  and nodes  $t_i = \frac{i}{n}$ ,  $i = 0, 1, \dots, n$ , the equation (47) can be discretized by the linear spline approximation

$$x(t) \approx \begin{cases} \left(1 - \frac{t-t_{i-1}}{h}\right) z_{i-1} + \frac{t-t_{i-1}}{h} z_i \\ \text{for } t_{i-1} \leq t \leq t_i, \quad i = 1, 2, \dots, n \end{cases}$$

and represented by a linear system  $A\mathbf{z} = \mathbf{b}$  where the variable  $\mathbf{z} = [z_0, \dots, z_n]^\top$  in  $\mathbb{C}^{n+1}$ , the right-hand side vector  $\mathbf{b} = [b_1, \dots, b_n]^\top \in \mathbb{C}^n$  and the coefficient matrix  $A = [a_{ij}] \in \mathbb{C}^{n \times (n+1)}$  with

$$\begin{aligned} a_{i1} &= \int_{t_0}^{t_1} k(t_i-t) \left(1 - \frac{t-t_0}{h}\right) dt \\ a_{ij} &= \int_{t_{j-1}}^{t_j} k(t_i-t) \frac{t-t_{j-1}}{h} dt + \int_{t_j}^{t_{j+1}} k(t_i-t) \left(1 - \frac{t-t_j}{h}\right) dt \\ &\quad j = 1, 2, \dots, i+1 \\ a_{i,i+1} &= \int_{t_{i-1}}^{t_i} k(t_i-t) \frac{t-t_{i-1}}{h} dt \\ b_i &= g(t_i) \\ &\quad \text{for } i = 1, 2, \dots, n. \end{aligned}$$

As an experiment with the kernel (48) where  $\kappa = 4$ , the right-hand side function

$$g(s) = \int_0^s k(s-t) dt$$

of the equation (47) corresponds to a known general solution

$$x(t) = 1 + c\delta(t-1)$$

where  $c$  is an arbitrary constant. By any standard numerical integration method such as the composite Simpson's rule, the matrix  $A$  and the right-hand side vector

---

<sup>3</sup>There is apparently a typo in [12, p. 83] about the kernel (48).

$\mathbf{b}$  can be generated easily.

Notice that the system  $A\mathbf{z} = \mathbf{b}$  is underdetermined with the size  $n \times (n+1)$  in addition to being empirical data of a singular equation (47). We choose not to add an extra equation to square the system for the consideration of the inherent singularity in the underlying problem.

For  $n = 1024$ , a simple call of `LinearSolve` produces the truncated SVD solution  $\mathbf{z}$ , the matrix  $\mathbf{K}$  whose columns form an orthonormal basis for  $\mathcal{K}ernel(A_\theta)$ , the sensitivity measure `lcnd` defined as  $\|A_\theta\|_2 \|A_\theta^\dagger\|_2$ , and the residual `res`:

```
>> [z,K,lcnd,res] = LinearSolve(A, b, 1e-6); % solve A*z = b
>> lcnd % the sensitivity
lcnd =
    2.469428269074639e+04

>> res % the residual
res =
    8.371530784514738e-10
```

The numerical kernel  $\mathcal{K}ernel(A_\theta)$  is of dimension three. Figure 1 shows the plot of the truncated SVD solution  $\mathbf{z}$  and the numerical kernel vectors  $\mathbf{u}$ ,  $\mathbf{v}$  and  $\mathbf{w}$ .

$i$	$t_i$	trunc. SVD solution $z_i$	basis for $\mathcal{K}ernel(A_\theta)$		
			$u_i$	$v_i$	$w_i$
0	0	1.0000025	0.0000000	-0.0000000	-0.0000000
1	0.0009766	0.9999977	-0.0000000	0.0000000	0.0000000
2	0.0019531	1.0000017	-0.0000000	-0.0000000	-0.0000000
3	0.0029297	0.9999994	0.0000000	0.0000000	-0.0000000
4	0.0039063	1.0000004	-0.0000000	-0.0000000	0.0000000
...	...	...	...	...	...
1007	0.9833984	1.0000002	-0.0000001	-0.0000000	0.0000000
1008	0.9843750	0.9999998	0.0000004	0.0000001	-0.0000000
1009	0.9853516	1.0000010	-0.0000011	-0.0000002	0.0000000
1010	0.9863281	0.9999978	0.0000032	0.0000004	-0.0000000
1011	0.9873047	1.0000067	-0.0000093	-0.0000013	0.0000000
1012	0.9882813	0.9999813	0.0000268	0.0000037	-0.0000001
1013	0.9892578	1.0000541	-0.0000770	-0.0000106	0.0000003
1014	0.9902344	0.9998449	0.0002211	0.0000303	-0.0000009
1015	0.9912109	1.0004460	-0.0006352	-0.0000870	0.0000025
1016	0.9921875	0.9987195	0.0018246	0.0002500	-0.0000072
1017	0.9931641	1.0036784	-0.0052409	-0.0007182	0.0000206
1018	0.9941406	0.9894347	0.0150535	0.0020628	-0.0000591
1019	0.9951172	1.0303428	-0.0432346	-0.0059230	0.0001698
1020	0.9960938	0.9129785	0.1240610	0.0169516	-0.0004872
1021	0.9970703	1.2462977	-0.3529420	-0.0470151	0.0013861
1022	0.9980469	0.3926255	0.9206328	0.0892117	-0.0036175
1023	0.9990234	0.0127561	-0.1016541	0.9947379	0.0003994
1024	1.0000000	0.0000008	0.0039290	0	0.9999923

TABLE 2

The general numerical solution  $\mathbf{z} + \alpha \mathbf{u} + \beta \mathbf{v} + \gamma \mathbf{w}$  that approximates the exact underlying solution  $x(t) = 1 + c \cdot \delta(t-1)$  for the equation (47)

Table 2 shows the actual digits in single precision of the general numerical solution.

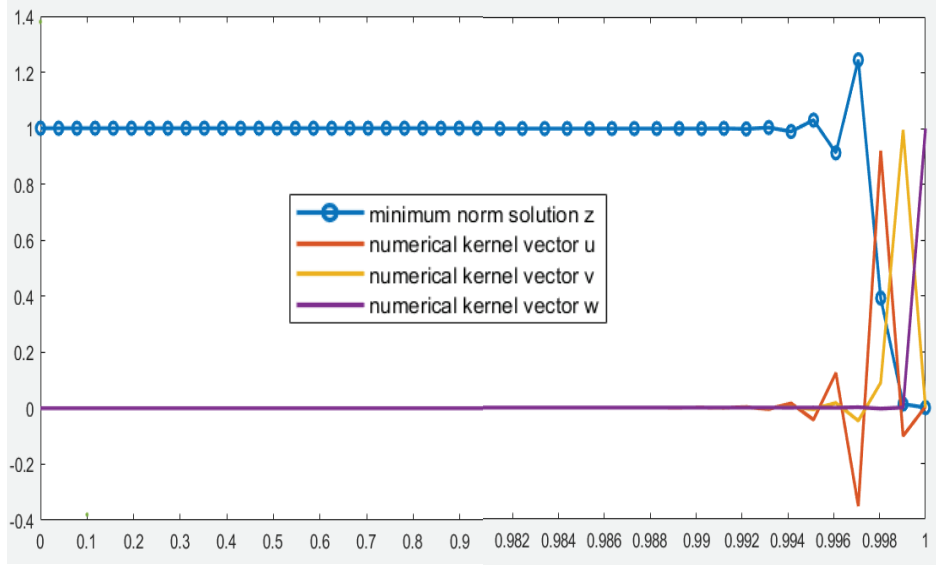


FIG. 1. The general numerical solution  $\mathbf{z} + \alpha \mathbf{u} + \beta \mathbf{v} + \gamma \mathbf{w}$  that approximates the exact underlying solution  $x(t) = 1 + c \cdot \delta(t - 1)$  for the equation (47)

The condition number

$$\kappa(A) = \frac{\sigma_1(A)}{\sigma_n(A)} \approx 3.3 \times 10^{25}$$

is huge whereas the sensitivity of the general numerical solution  $\mathbf{z} + \alpha \mathbf{u} + \beta \mathbf{v} + \gamma \mathbf{w}$  is manageable and much lower at  $2.5 \times 10^4$ . Given the residual  $8.4 \times 10^{-10}$ , we can make a rough error estimate of the general numerical solution as

$$(\|A_\theta\|_2 \|A_\theta^\dagger\|_2) (\|A\mathbf{z} - \mathbf{b}\|_2 + \|A[\mathbf{u}, \mathbf{v}, \mathbf{w}]\|_2) \approx 2.1 \times 10^{-5}$$

that can be considered accurate for such an application.

The accuracy estimate can also be justified as follows. The obvious particular solution  $x_0(t) = 1$  of the equation (47) can be approximated by a numerical particular solution:

```
>> y = K\ (1-z); % solve K*y+z = 1
y =
    0.561958104442570
    1.049493278421724
    0.997798939791854

>> norm(Z+K*u-1,1)*h % error of the numerical particular solution
ans =
    1.090344305094233e-07
```

Namely, the particular solution  $x_0(t) = 1$  can be accurately approximated by a numerical particular solution with error bound  $1.1 \times 10^{-7}$ . On the other hand, the



solution to the homogeneous equation corresponding to (47) is

$$\delta(t-1) = \lim_{\varepsilon \rightarrow 0^+} \delta_\varepsilon(t-1) \quad \text{where} \quad \delta_\varepsilon(t-1) = \frac{e^{-\left(\frac{t-1}{\varepsilon}\right)^2}}{\varepsilon \sqrt{\pi}}$$

The following Matlab statement sequence shows that a numerical kernel vector approximates  $\delta_{0.00001}(t-1)$  with an error measure  $7.76 \times 10^{-5}$ , as the error estimate suggests.

```
>> y = K\1-z; % solve K*y+z = 1
y =
>> epsilon = 1e-5; % a tiny epsilon
>> h = 1/n; % stepsize
>> t = 0:h:1; % the nodes in t
>> v = K\1./(epsilon*sqrt(pi)*exp(((t-1)/epsilon).^2)) % solve K*v = delta_epsilon
v =
 1.0e+04 *
 0.022167287609456
 0.000000000000000
 5.641852287123326

>> norm(K*v -1./(epsilon*sqrt(pi)*exp(((t-1)/epsilon).^2)),1)*h
ans =
 7.764167138103457e-05
```

Although the underlying system is underdetermined in addition to being singular, the general numerical solution  $\mathbf{z} + \alpha \mathbf{u} + \beta \mathbf{v} + \gamma \mathbf{w}$  accurately reveals that the equation (47) with the specific kernel (48) can be accurately solved in an interval  $[0, 1 - \varepsilon)$  for a small  $\varepsilon > 0$  since the annihilator represented by  $\alpha \mathbf{u} + \beta \mathbf{v} + \gamma \mathbf{w}$  is identically zero in that interval. The singularity of the underlying equation compounded by the representation linear system being underdetermined is not detrimental at all if we compute the general numerical solution and, in particular, consider the numerical kernel as an integral part of the solution.