

# Sensitivity and Computation of a Defective Eigenvalue

Zhonggang Zeng\*

April 9, 2015

## Abstract

A multiple and defective eigenvalue is well documented to be hypersensitive to data perturbations and round-off errors, making it a formidable challenge in numerical computation particularly when entries of the matrix are not known exactly. This paper establishes a finitely bounded sensitivity of a defective eigenvalue with respect to perturbations that preserve the geometric multiplicity and the smallest Jordan block size. Based on this perturbation theory, numerical computation of a defective eigenvalue is regularized as a well-posed least squares problem so that it can be accurately carried out using floating point arithmetic even if the matrix is perturbed.

## 1 Introduction

Computing matrix eigenvalues is one of the fundamental problems in theoretical and numerical linear algebra. The highly celebrated Francis QR algorithm effectively solves the problem for simple and non-defective eigenvalues in practical applications. However, it is well documented that multiple and defective eigenvalues are hypersensitive to both data perturbations on the matrix and the inevitable round-off. For an eigenvalue of a matrix  $A$  associated with the largest Jordan block size  $l \times l$  while  $A$  is perturbed by  $\Delta A$ , the error bound [2, p. 58][3, 13] on the eigenvalue deviation is a constant multiple of  $\|\Delta A\|_2^{\frac{1}{l}}$ , implying that the accuracy of the computed eigenvalue in the number of correct digits is only a fraction  $\frac{1}{l}$  of the accuracy of the matrix data. As a result, numerical computation of defective eigenvalues from empirical matrix data remains a formidable challenge.

On the other hand, it has long been known that a defective eigenvalue disperses into a cluster of simple eigenvalues when the matrix is under arbitrary perturbations but the mean of the cluster is not hypersensitive [10, 17]. In his seminal technical report[9], Kahan proved that the sensitivity of an  $m$ -fold eigenvalue is actually bounded by  $\frac{1}{m}\|P\|_2$  where  $P$  is the spectral projector associated with the eigenvalue as long as the perturbation is

---

\*Department of Mathematics, Northeastern Illinois University, Chicago, IL 60625. *Email:* zzeng@neiu.edu

29 constrained so that the multiplicity is preserved. The same proof and the same sensitivity  
 30 also apply to the mean of the eigenvalue cluster emanating from the  $m$ -fold eigenvalue  
 31 with respect to arbitrary perturbations. Indeed, using cluster means as approximations  
 32 to defective eigenvalues has been extensively applied to numerical computation of Jordan  
 33 Canonical Forms and staircase forms, provided that the clusters can be sorted out from  
 34 the spectrum. This approach includes works of Ruhe [16], Sdridhar and Jordan [20], and  
 35 culminated in Golub and Wilkinson's review [6] as well as Kågström and Ruhe's JNF [7, 8].  
 36 Theoretical issues have been analyzed in, e.g. works of Demmel [4, 5] and Wilkinson [22, 23].  
 37 Perturbations on eigenvalue clusters are also studied as pseudospectra of matrices in works  
 38 of Trefethon and Ebreë [21] as well as Rump [18, 19].

39 In this paper we elaborate a different measurement of the sensitivity of a defective eigenvalue  
 40 with respect to perturbations constrained to preserve the geometric multiplicity and the  
 41 smallest Jordan block size. We prove that such a sensitivity is also finitely bounded even  
 42 if the multiplicity is not preserved, and it is large only if either the geometric multiplicity  
 43 or the smallest Jordan block size can be increased by a small perturbation on the matrix.  
 44 This sensitivity can be small even if the spectral projector norm is large, or vice versa.

45 In general, perturbations are expected to be arbitrary without preserving either the mul-  
 46 tiplicity or what we refer to as the multiplicity support. We prove that a certain type  
 47 of pseudo-eigenvalue uniquely exists, is Lipschitz continuous, is backward accurate and ap-  
 48 proximates the defective eigenvalue with a forward accuracy in the same order of the data  
 49 accuracy. Namely, we prove that finding a defective eigenvalue via computing a pseudo-  
 50 eigenvalue is a well-posed problem. Based on this analysis, we develop an iterative algorithm  
 51 that is capable of accurate computation of defective eigenvalues using floating point arith-  
 52 metic from empirical matrix data even if the spectral projector norm is large and thus the  
 53 cluster mean is inaccurate.

## 54 2 Notation

55 The space of dimension  $n$  column vectors is  $\mathbb{C}^n$  and the vector space of  $m \times n$  matrices  
 56 is  $\mathbb{C}^{m \times n}$  as usual. Matrices are denoted by upper case letters  $A$ ,  $X$ , and  $G$ , etc, with  
 57  $O$  representing a zero matrix whose dimensions can be derived from the context. Boldface  
 58 lower case letters such as  $\mathbf{x}$  and  $\mathbf{y}$  represent column vectors. Particularly, the zero vector  
 59 in  $\mathbb{C}^n$  is denoted by  $\mathbf{0}_n$  or simply  $\mathbf{0}$  if the dimension is clear. The conjugate transpose of  
 60 a matrix or vector  $(\cdot)$  is denoted by  $(\cdot)^H$ , and the pseudo-inverse of a matrix  $(\cdot)$  is  $(\cdot)^\dagger$ .  
 61 The submatrix formed by entries in rows  $i_1, \dots, i_2$  and columns  $j_1, \dots, j_2$  of a matrix  $A$   
 62 is denoted by  $A_{i_1:i_2, j_1:j_2}$ . The kernel and range of a matrix  $(\cdot)$  are denoted by  $\mathcal{Ker}(\cdot)$  and  
 63  $\mathcal{Ran}(\cdot)$  respectively. The notation  $\mathit{eig}(\cdot)$  represents the spectrum of a matrix  $(\cdot)$ .

64 We also consider vectors in vector spaces such as  $\mathbb{C} \times \mathbb{C}^{m \times k}$  throughout this paper. In  
 65 such cases, the vector 2-norm is the square root of the sum of squares of all components.  
 66 For instance, a vector  $(\lambda, X) \in \mathbb{C} \times \mathbb{C}^{m \times k}$  can be arranged as a column vector  $\mathbf{u}$  in  $\mathbb{C}^{n \times k+1}$   
 67 and  $\|(\lambda, X)\|_2 = \|\mathbf{u}\|_2$  regardless of the ordering. A zero vector in such a vector space is  
 68 also denoted by  $\mathbf{0}$ .

69 Let  $\lambda_*$  be an eigenvalue of a matrix  $A \in \mathbb{C}^{n \times n}$ . Its algebraic multiplicity can be parti-  
 70 tioned into a non-increasing sequence of nonnegative integers called the *Segre characteristic*  
 71  $\{l_1 \geq l_2 \geq \dots\}$  that are the sizes of elementary Jordan blocks. In other words, there is a  
 72 matrix  $X_* \in \mathbb{C}^{n \times m}$  of full column rank such that

$$A X_* = X_* \begin{bmatrix} J_{l_1}(\lambda_*) & & \\ & J_{l_2}(\lambda_*) & \\ & & \ddots \end{bmatrix}$$

73 where

$$J_k(\lambda_*) = \begin{bmatrix} \lambda_* & 1 & & \\ & \lambda_* & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_* \end{bmatrix} \in \mathbb{C}^{k \times k}$$

74 denotes the  $k \times k$  elementary Jordan block associated with an eigenvalue  $\lambda_*$ . For conve-  
 75 nience, a Segre characteristic is infinite in formality and the number of nonzero entries is the  
 76 geometric multiplicity. The last nonzero component of a Segre characteristic, namely the  
 77 size of the smallest elementary Jordan block associated with  $\lambda_*$ , is of particular importance  
 78 in our perturbation analysis and we shall call it the *Segre characteristic anchor*.

79 For instance, if  $\lambda_*$  is an eigenvalue of  $A$  associated with elementary Jordan blocks  
 80  $J_5(\lambda_*), J_5(\lambda_*), J_4(\lambda_*), J_4(\lambda_*)$  and  $J_3(\lambda_*)$ , its Segre characteristic is  $\{5, 5, 4, 4, 3, 0, \dots\}$   
 81 with an anchor 3. The geometric multiplicity is 5 since the Segre characteristic anchor  
 82 is the fifth component. A Segre characteristic along with its conjugate that is called the  
 83 Weyr characteristic can be illustrated by a Ferrer's diagram in Fig. 1, where the geometric  
 84 multiplicity and the Segre characteristic anchor represent the dimensions of the base rectangle  
 85 occupied by the equal leading entries of the Weyr characteristic.

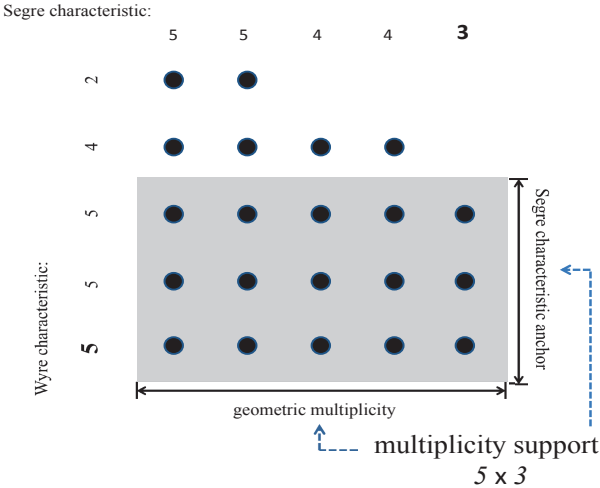


Figure 1: Illustration of the multiplicity support for a defective eigenvalue

86 For a matrix  $A$ , we shall say the *multiplicity support* of its eigenvalue  $\lambda_*$  is  $m \times k$  if the  
 87 geometric multiplicity of  $\lambda_*$  is  $m$  and the Segre characteristic anchor is  $k$ . In this case,

88 there is a unique  $X_* \in \mathbb{C}^{n \times k}$  satisfying the equation

$$\begin{aligned} (A - \lambda_* I) X_* &= X_* J_k(0) \\ C^H X_* &= T \end{aligned}$$

89 with proper choices of  $C \in \mathbb{C}^{n \times m}$  and  $T \in \mathbb{C}^{m \times k}$ , as we shall prove in Lemma 3.1. Here  
90  $J_k(0)$  is a nilpotent upper-triangular matrix of rank  $k - 1$  and can be replaced with any  
91 matrix of such kind. For integers  $m, k \leq n$ , we define a holomorphic mapping

$$\begin{aligned} \mathbf{g} : \mathbb{C}^{n \times n} \times \mathbb{C} \times \mathbb{C}^{n \times k} &\longrightarrow \mathbb{C}^{n \times k} \times \mathbb{C}^{m \times k} \\ (G, \lambda, X) &\longmapsto \begin{pmatrix} (G - \lambda I) X - X S \\ C^H X - T \end{pmatrix} \end{aligned} \quad (1)$$

92 that depends on parameters  $C \in \mathbb{C}^{n \times m}$ ,  $T \in \mathbb{C}^{m \times k}$  and an upper-triangular nilpotent  
93 matrix

$$S = \begin{bmatrix} 0 & s_{12} & \cdots & s_{1k} \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & s_{k-1,k} \\ 0 & \cdots & \cdots & 0 \end{bmatrix} \quad \text{with} \quad s_{12}s_{23} \cdots s_{k-1,k} \neq 0. \quad (2)$$

94 We shall also denote the Jacobian and partial Jacobian

$$\begin{aligned} \mathbf{g}_{G\lambda X}(G_0, \lambda_0, X_0) &= \left. \frac{\partial \mathbf{g}(G, \lambda, X)}{\partial (G, \lambda, X)} \right|_{(G, \lambda, X) = (G_0, \lambda_0, X_0)} \\ \mathbf{g}_{\lambda X}(G_0, \lambda_0, X_0) &= \left. \frac{\partial \mathbf{g}(G, \lambda, X)}{\partial (\lambda, X)} \right|_{(G, \lambda, X) = (G_0, \lambda_0, X_0)} \end{aligned}$$

95 that can be considered linear transformations

$$\begin{aligned} \mathbf{g}_{G\lambda X}(G_0, \lambda_0, X_0) : \mathbb{C}^{n \times n} \times \mathbb{C} \times \mathbb{C}^{n \times k} &\longrightarrow \mathbb{C}^{n \times k} \times \mathbb{C}^{m \times k} \\ (G, \lambda, X) &\longmapsto \begin{pmatrix} (G - \lambda I) X_0 + (G_0 - \lambda_0 I) X - X S \\ C^H X \end{pmatrix} \end{aligned}$$

96 and

$$\begin{aligned} \mathbf{g}_{\lambda X}(G_0, \lambda_0, X_0) : \mathbb{C} \times \mathbb{C}^{n \times k} &\longrightarrow \mathbb{C}^{n \times k} \times \mathbb{C}^{m \times k} \\ (\lambda, X) &\longmapsto \begin{pmatrix} -\lambda X_0 + (G_0 - \lambda_0 I) X - X S \\ C^H X \end{pmatrix} \end{aligned}$$

97 respectively. The actual matrices representing the Jacobians depend on the ordering of the  
98 bases for the domains and codomains of those linear transformations. The pseudo-inverse  
99 of a linear transformation such as  $\mathbf{g}_{\lambda X}(G_0, \lambda_0, X_0)^\dagger$  is the linear transformation whose  
100 matrix representation is the pseudo-inverse of the matrix representation of  $\mathbf{g}_{\lambda X}(G_0, \lambda_0, X_0)$   
101 corresponding to the same bases.

### 102 3 Properties of the multiplicity support

103 The following lemma sets the foundation for the sensitivity analysis and algorithm develop-  
104 ment for a defective eigenvalue with certain multiplicity support by laying out the critical  
105 properties of the holomorphic mapping in (1).

106 **Lemma 3.1** Let  $A \in \mathbb{C}^{n \times n}$  with an eigenvalue  $\lambda_*$  of multiplicity support  $m_* \times k_*$  and  
 107  $\mathbf{g}$  be the mapping in (1). The following assertions hold.

108 (i) Let  $S$  be as in (2) and

$$T = \begin{bmatrix} 1 & \mathbf{0}_{k-1}^\top \\ \mathbf{0}_{m-1} & O \end{bmatrix}. \quad (3)$$

109 For almost all  $C \in \mathbb{C}^{n \times m}$ , there is an  $X_* \in \mathbb{C}^{n \times k}$  such that  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$  if  
 110 and only if  $m \leq m_*$  and  $k \leq k_*$ . Such an  $X_*$  is unique if and only if  $m = m_*$ .

111 (ii) Let  $m \leq m_*$ ,  $k \leq k_*$  and  $S$  be as in (2). For any  $T \in \mathbb{C}^{m \times k}$  with  $T_{1:m,1} \neq \mathbf{0}$   
 112 and almost all  $C \in \mathbb{C}^{n \times m}$ , assume  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$ . Then  $\mathbf{g}_{G\lambda X}(A, \lambda_*, X_*)$  is  
 113 surjective, and  $\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)$  is injective if and only if  $m = m_*$  and  $k = k_*$ .

114 (iii) Let  $m = m_*$ ,  $k \leq k_*$ ,  $T$  be in (3) and columns of  $C \in \mathbb{C}^{n \times m}$  form an orthonormal  
 115 basis for the kernel of  $A - \lambda_* I$ . Then the matrix parameter  $S$  can be chosen so that  
 116 the unique  $X_*$  satisfying  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$  has orthonormal columns.

117 **PROOF.** Let  $N \in \mathbb{C}^{n \times m_*}$  be a matrix whose columns span the kernel  $\mathcal{Ker}(A - \lambda_* I)$ .  
 118 Then for almost all  $C \in \mathbb{C}^{n \times m}$ , the matrix  $C^H N$  is of full row rank if  $m \leq m_*$  so that  
 119 there is a  $\mathbf{u} \in \mathbb{C}^m$  such that

$$(C^H N) \mathbf{u} = \begin{bmatrix} 1 \\ \mathbf{0}_{m-1} \end{bmatrix}$$

120 and  $\mathbf{u}$  is unique if and only if  $m = m_*$ . For  $m \leq m_*$ , let  $\mathbf{x}_1 = N \mathbf{u}$  and assume vectors  
 121  $\mathbf{x}_1, \dots, \mathbf{x}_j \in \mathbb{C}^n$  are obtained such that  $1 \leq j < k$  and

$$(A - \lambda_* I) [\mathbf{x}_1, \dots, \mathbf{x}_j] = [\mathbf{x}_1, \dots, \mathbf{x}_j] \begin{bmatrix} 0 & s_{12} & \cdots & s_{1j} \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & s_{j-1,j} \\ 0 & \cdots & \cdots & 0 \end{bmatrix}$$

$$C^H [\mathbf{x}_1, \dots, \mathbf{x}_j] = \begin{bmatrix} 1 & \mathbf{0}_{j-1}^\top \\ \mathbf{0}_{m-1} & O \end{bmatrix}.$$

122 Then the system

$$(A - \lambda_* I) \mathbf{x} = s_{1,j+1} \mathbf{x}_1 + \cdots + s_{j,j+1} \mathbf{x}_j$$

123 has an affine space solution  $\{\mathbf{u} + N \mathbf{v} \mid \mathbf{v} \in \mathbb{C}^m\}$  and a unique solution

$$\mathbf{x}_{j+1} = \mathbf{u} - N (C^H N)^{-1} C^H \mathbf{u} \in \mathbb{C}^n$$

124 such that  $C^H \mathbf{x}_{j+1} = \mathbf{0}$  when  $m = m_*$ . Consequently, there is an  $X_* = [\mathbf{x}_1, \dots, \mathbf{x}_k] \in \mathbb{C}^{n \times k}$   
 125 such that  $(\lambda_*, X_*)$  is a solution to the system  $\mathbf{g}(A, \lambda, X) = \mathbf{0}$  and  $X_*$  is unique if and  
 126 only if  $m = m_*$ . The assertion (i) is proved.

127 For any parameter  $T \in \mathbb{C}^{m \times k}$  of  $\mathbf{g}$  with  $T_{1:m,1} \neq \mathbf{0}$  and almost all  $C \in \mathbb{C}^{n \times m}$ , assume  
 128  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$ . Write  $X_* = [\mathbf{x}_1, \dots, \mathbf{x}_k]$ . Then  $\mathbf{x}_1 \neq \mathbf{0}$  and

$$\mathbf{x}_j \in \mathcal{Ker}((A - \lambda_* I)^j) \setminus \mathcal{Ker}((A - \lambda_* I)^{j-1})$$

129 for  $j = 1, \dots, k$ . Thus  $X_*$  is of full column rank. Furthermore, by proper column  
 130 permutations the matrix representation of the Jacobian  $\mathbf{g}_{G\lambda X}(A, \lambda_*, X_*)$  is in the form of

$$\begin{bmatrix} X_*^\top & & & * & \cdots & * & * \\ & \ddots & & \vdots & \ddots & \vdots & \vdots \\ & & X_*^\top & * & \cdots & * & * \\ & & & C^H & & & * \\ & & & & \ddots & & \vdots \\ & & & & & C^H & * \end{bmatrix}$$

131 Therefore  $\mathbf{g}_{G\lambda X}(A, \lambda_*, X_*)$  is surjective.

132 Let  $(A, \lambda_*, \hat{X})$  be a zero of  $\mathbf{g}$  and assume  $m = m_*$  and  $k = k_*$ . Then, for almost all  
 133  $C \in \mathbb{C}^{m \times m}$ , the solution  $\mathbf{u} = \hat{\mathbf{u}}$  of the equation  $(C^H N) \mathbf{u} = T_{1:m,1}$  is unique and the first  
 134 column of  $\hat{X}$  is

$$\hat{\mathbf{x}}_1 = N \hat{\mathbf{u}} \in \left( \bigcap_{j=1}^{k_*-1} \mathcal{Ran}((A - \lambda_* I)^j) \right) \setminus \mathcal{Ran}((A - \lambda_* I)^{k_*}) \quad (4)$$

135 since the smallest elementary Jordan block is  $J_{k_*}(\lambda_*)$  of size  $k_* \times k_*$ . Assume, for certain  
 136  $(\sigma, Y) \in \mathbb{C} \times \mathbb{C}^{n \times k}$ , its image of the linear transformation  $\mathbf{g}_{\lambda X}(A, \lambda_*, \hat{X})(\sigma, Y) = \mathbf{0}$ . Namely

$$-\sigma \hat{X} + (A - \lambda_* I)Y - YS = O \quad (5)$$

$$C^H Y = O. \quad (6)$$

137 Right-multiplying both sides of the equation (5) by  $S$  yields

$$\begin{aligned} YS^2 + \sigma \hat{X}S &= (A - \lambda_* I)YS \\ &= (A - \lambda_* I)^2 Y - \sigma (A - \lambda_* I) \hat{X} \\ &= (A - \lambda_* I)^2 Y - \sigma \hat{X}S \end{aligned}$$

138 Continuing the process of recursive right-multiplying the above equation by  $S$  leads to

$$\begin{aligned} (A - \lambda_* I)^k Y &= YS^k + k\sigma \hat{X}S^{k-1} \\ &= k\sigma s_{12} s_{23} \cdots s_{k-1,k} [O_{n \times (k-1)}, \hat{\mathbf{x}}_1] \end{aligned}$$

139 with  $s_{12} s_{23} \cdots s_{k-1,k} \neq 0$ . Hence  $\sigma$  must be zero due to (4) and  $k = k_*$ . Denote columns  
 140 of  $Y$  as  $\mathbf{y}_1, \dots, \mathbf{y}_k \in \mathbb{C}^n$ . Then the first column of the equations (5) and (6) are

$$(A - \lambda_* I) \mathbf{y}_1 = \mathbf{0}, \quad C^H \mathbf{y}_1 = 0$$

141 that implies  $\mathbf{y}_1 = \mathbf{0}$ . For  $1 \leq j < k$ , using  $\sigma = 0$  and  $\mathbf{y}_1 = \dots = \mathbf{y}_j = \mathbf{0}$  on the  
 142  $(j+1)$ -th column of the equations (5) and (6) we have  $\mathbf{y}_{j+1} = \mathbf{0}$ . Thus  $Y = O$ . As a  
 143 result,  $(A, \lambda_*, \hat{X})$  is a zero of  $\mathbf{g}$  with injective partial Jacobian  $\mathbf{g}_{\lambda X}(A, \lambda_*, \hat{X})$ .

144 If  $m < m_*$ , the solution  $(\lambda_*, X_*)$  of  $\mathbf{g}(A, \lambda, X) = \mathbf{0}$  is not isolated and thus  $\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)$   
 145 is not injective. Let  $m = m_*$ , we now prove the partial Jacobian  $\mathbf{g}_{\lambda X}(A, \lambda_*, \hat{X})$  is injective

146 only if  $k = k_*$ . Assume  $k < k_*$  and write  $\hat{X} = [\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_k]$ . Since  $\mathbf{g}(A, \lambda_*, \hat{X}) = \mathbf{0}$   
 147 and  $S$  is upper-triangular nilpotent, hence  $\mathbf{x}_1 \in \mathcal{Ker}(A - \lambda_* I)$ . Then  $k < k_*$  implies

$$\hat{\mathbf{x}}_1 \in \bigcap_{j=1}^{k+1} \mathcal{Ran}((A - \lambda_* I)^j) \quad (7)$$

148 and there exists an  $\hat{\mathbf{x}}_{k+1}$  such that  $(A - \lambda_* I)\mathbf{x}_{k+1} = \hat{\mathbf{x}}_k$ . For almost all  $C \in \mathbb{C}^{m \times m}$ , the ma-  
 149 trix  $\begin{bmatrix} A - \lambda_* I \\ C^H \end{bmatrix}$  is of full column rank and there is a unique solution  $\mathbf{z} = \mathbf{y}_1 \in \text{span}\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2\}$   
 150 to the linear system

$$\begin{bmatrix} A - \lambda_* I \\ C^H \end{bmatrix} \mathbf{z} = \begin{bmatrix} \hat{\mathbf{x}}_1 \\ \mathbf{0} \end{bmatrix}.$$

151 Assume we have  $\mathbf{y}_1, \dots, \mathbf{y}_j \in \text{span}\{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_{j+1}\}$  for any  $j \in \{1, \dots, k-1\}$  such that

$$\begin{aligned} -[\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_j] + (A - \lambda_* I)[\mathbf{y}_1, \dots, \mathbf{y}_j] - [\mathbf{y}_1, \dots, \mathbf{y}_j] S &= \mathbf{0} \\ C^H [\mathbf{y}_1, \dots, \mathbf{y}_j] &= \mathbf{0}. \end{aligned}$$

152 By (7), there is a unique solution  $\mathbf{z} = \mathbf{y}_{j+1} \in \text{span}\{\mathbf{x}_1, \dots, \hat{\mathbf{x}}_{j+1}\}$  to the linear system

$$\begin{bmatrix} A - \lambda_* I \\ C^H \end{bmatrix} \mathbf{z} = \begin{bmatrix} \hat{\mathbf{x}}_{j+1} + s_{1,j+1}\mathbf{y}_1 + \dots + s_{j,j+1}\mathbf{y}_j \\ \mathbf{0} \end{bmatrix}.$$

153 Write  $Y = [\mathbf{y}_1, \dots, \mathbf{y}_k]$ . We have  $\mathbf{g}_{\lambda X}(A, \lambda_*, \hat{X})(1, Y) = \mathbf{0}$  and thus the partial Jacobian  
 154  $\mathbf{g}_{\lambda X}(A, \lambda_*, \hat{X})$  is not injective. As a result, the assertion (ii) is proved.

155 We now prove (iii). Let columns of  $C \in \mathbb{C}^{n \times m}$  form an orthonormal basis for  $\mathcal{Ker}(A - \lambda_* I)$   
 156 and  $X_*$  be the unique matrix in assertion (i). Then  $C$  and  $X_*$  have the identical first  
 157 columns. Let  $X_* = QR$  be a thin QR decomposition where  $R = [r_{ij}]$  with  $r_{11} = 1$ .  
 158 Since  $T$  is as in (3), it is easy to see that

$$R = \begin{bmatrix} 1 & \mathbf{0}_{k-1}^\top \\ \mathbf{0}_{k-1} & \hat{R} \end{bmatrix} \in \mathbb{C}^{k \times k}.$$

159 Thus

$$\begin{aligned} (A - \lambda_* I)Q &= Q(RSR^{-1}) \\ C^H Q &= TR^{-1}. \end{aligned}$$

160 It is straightforward to verify that  $TR^{-1} = T$  and  $RSR^{-1}$  is an upper-triangular nilpotent  
 161 matrix of rank  $k-1$ . Consequently, the assertion (iii) is proved by setting  $\check{X}$  as  $Q$  and  
 162 replacing  $S$  with  $RSR^{-1}$ .  $\square$

## 163 4 Sensitivity of a defective eigenvalue

164 Lemma 3.1 reveals the critical importance of the geometric multiplicity combined with the  
 165 Segre characteristic anchor of an eigenvalue and its behavior under data perturbations. In  
 166 fact, the matrix and the defective eigenvalue are holomorphic functions of certain entries of  
 167 the matrix as asserted in the following corollary.

168 **Corollary 4.1** Assume  $A \in \mathbb{C}^{n \times n}$  and  $\lambda_* \in \text{eig}(A)$  with a multiplicity support  $m \times k$ .  
169 Let the mapping  $\mathbf{g}$  be defined in (1) where  $S \in \mathbb{C}^{k \times k}$  is as in (2),  $C \in \mathbb{C}^{n \times m}$  and  
170  $T \in \mathbb{C}^{m \times k}$  with  $T_{1:m,1} \neq \mathbf{0}$  so that  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$  with an injective  $\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)$ .  
171 There exist  $n^2 - mk + 1$  entries of  $G$  forming a variable  $\mathbf{z}$  in a neighborhood  $\Omega \ni \mathbf{z}_*$   
172 in  $\mathbb{C}^{n^2 - mk + 1}$ , a neighborhood  $\Sigma$  of  $(A, \lambda_*, X_*)$  in  $\mathbb{C}^{n \times n} \times \mathbb{C} \times \mathbb{C}^{n \times k}$  and holo-  
173 morphic mappings  $G : \Omega \rightarrow \mathbb{C}^{n \times n}$ ,  $\lambda : \Omega \rightarrow \mathbb{C}$  and  $X : \Omega \rightarrow \mathbb{C}^{n \times k}$  with  
174  $(G(\mathbf{z}_*), \lambda(\mathbf{z}_*), X(\mathbf{z}_*)) = (A, \lambda_*, X_*)$  such that  $\mathbf{g}(G, \lambda, X) = \mathbf{0}$  for  $(G, \lambda, X) \in \Sigma$  if and  
175 only if  $(G, \lambda, X) = (G(\mathbf{z}), \lambda(\mathbf{z}), X(\mathbf{z}))$  for  $\mathbf{z} \in \Omega$ .

176 **PROOF.** By Lemma 3.1, part (ii), the Jacobian  $\mathbf{g}_{G\lambda X}(A, \lambda_*, X_*)$  is surjective to  
177  $\mathbb{C}^{n \times k} \times \mathbb{C}^{m \times k}$  and the partial Jacobian  $\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)$  is injective from  $\mathbb{C} \times \mathbb{C}^{n \times k}$ . Conse-  
178 quently, there are  $mk - 1$  entries of the variable  $G \in \mathbb{C}^{n \times n}$  forming a variable  $\mathbf{y}$  such that  
179 the partial Jacobian  $\mathbf{g}_{\mathbf{y}\lambda X}(A, \lambda_*, X_*)$  is invertible. By the Implicit Function Theorem, the  
180 remaining entries of  $G$  excluding those forming  $\mathbf{y}$  form a variable vector  $\mathbf{z} \in \mathbb{C}^{n^2 - mk + 1}$   
181 so that the assertion holds.  $\square$

182 Denote the collection of  $n \times n$  complex matrices having an eigenvalue that shares the same  
183 multiplicity support  $m \times k$  as

$$\mathcal{E}_{m \times k}^n := \{ A \in \mathbb{C}^{n \times n} \mid \text{A multiplicity support of } A \text{ is } m \times k \}. \quad (8)$$

184 Corollary 4.1 implies that every matrix  $A$  in  $\mathcal{E}_{m \times k}^n$  has an eigenvalue  $\lambda_*$  along with a  
185 unique  $X_*$  such that  $(A, \lambda_*, X_*)$  belongs to an algebraic variety defined by the solution  
186 set of the polynomial system  $\mathbf{g}(G, \lambda, X) = \mathbf{0}$ .

187 We can now establish one of the main theorems of this paper.

188 **Theorem 4.2 (Defective Eigenvalue Sensitivity Theorem)** *The sensitivity of a de-*  
189 *fective eigenvalue is finitely bounded as long as its multiplicity support is preserved. More*  
190 *precisely, let  $A \in \mathbb{C}^{n \times n}$  and  $\lambda_* \in \text{eig}(A)$  with a multiplicity support  $m \times k$ . There*  
191 *is a neighborhood  $\Sigma$  of  $(\lambda_*, A)$  in  $\mathbb{C} \times \mathbb{C}^{n \times n}$  and a neighborhood  $\Omega \in \mathbb{C}^{n^2 - mk + 1}$  of*  
192  *$\mathbf{z}_*$  formed by  $n^2 - mk + 1$  entries  $A$  along with holomorphic mappings  $\lambda : \Omega \rightarrow \mathbb{C}$*   
193 *and  $G : \Omega \rightarrow \mathbb{C}^{n \times n}$  with  $(\lambda_*, A) = (\lambda(\mathbf{z}_*), G(\mathbf{z}_*))$  such that every  $(\tilde{\lambda}, \tilde{A}) \in \Sigma$  with*  
194  *$\tilde{\lambda} \in \text{eig}(\tilde{A})$  of multiplicity support  $m \times k$  is equal to  $(\lambda(\mathbf{z}), G(\mathbf{z}))$  for certain  $\mathbf{z} \in \Omega$ .*  
195 *Furthermore,*

$$\limsup_{\mathbf{z} \rightarrow \mathbf{z}_*} \frac{|\lambda(\mathbf{z}) - \lambda_*|}{\|G(\mathbf{z}) - A\|_F} \leq \|\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)^\dagger\|_2 \quad (9)$$

$$< \infty$$

196 where  $X_* \in \mathbb{C}^{n \times k}$  satisfies  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$  for the mapping  $\mathbf{g}$  defined in (1) that  
197 renders columns of  $X_*$  orthonormal.

198 **PROOF.** The existence of neighborhoods  $\Sigma$  and  $\Omega$  along with holomorphic mappings  $\lambda$   
199 and  $G$  is a direct consequence of Corollary 4.1, which also implies

$$\mathbf{g}(G(\mathbf{z}), \lambda(\mathbf{z}), X(\mathbf{z})) \equiv O \quad \text{for all } \mathbf{z} \in \Omega.$$



200 As a result,

$$\begin{aligned}
& \left( \frac{\partial \mathbf{g}(G(\mathbf{z}), \lambda(\mathbf{z}), X(\mathbf{z}))}{\partial \mathbf{z}} \Big|_{\mathbf{z}=\mathbf{z}_*} \right) (\hat{\mathbf{z}} - \mathbf{z}_*) \\
&= \mathbf{g}_G(A, \lambda_*, X_*) G_{\mathbf{z}}(\mathbf{z}_*) (\hat{\mathbf{z}} - \mathbf{z}_*) + \mathbf{g}_{\lambda X}(A, \lambda_*, X_*) \left( \frac{\partial(\lambda(\mathbf{z}), X(\mathbf{z}))}{\partial \mathbf{z}} \Big|_{\mathbf{z}=\mathbf{z}_*} \right) (\hat{\mathbf{z}} - \mathbf{z}_*) \\
&= \mathbf{0},
\end{aligned}$$

201 implying

$$\begin{aligned}
|\lambda(\hat{\mathbf{z}}) - \lambda_*| &\leq \|(\lambda(\hat{\mathbf{z}}), X(\hat{\mathbf{z}})) - (\lambda_*, X_*)\|_2 \\
&= \left\| \frac{\partial(\lambda(\mathbf{z}), X(\mathbf{z}))}{\partial \mathbf{z}} \Big|_{\mathbf{z}=\mathbf{z}_*} (\hat{\mathbf{z}} - \mathbf{z}_*) \right\|_2 + O\left(\|\hat{\mathbf{z}} - \mathbf{z}_*\|_2^2\right) \\
&= \|\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)^\dagger \mathbf{g}_G(A, \lambda_*, X_*) G_{\mathbf{z}}(\mathbf{z}_*) (\hat{\mathbf{z}} - \mathbf{z}_*)\|_2 + O\left(\|\hat{\mathbf{z}} - \mathbf{z}_*\|_2^2\right) \\
&\leq \|\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)^\dagger\|_2 \|G_{\mathbf{z}}(\mathbf{z}_*)\|_F + O\left(\|\hat{\mathbf{z}} - \mathbf{z}_*\|_2^2\right)
\end{aligned}$$

202 since, by proper row and column permutations, the partial Jacobian  $\mathbf{g}_G(A, \lambda_*, X_*)$  has a  
203 matrix representation

$$\begin{bmatrix} X_*^\top & & \\ & \ddots & \\ & & X_*^\top \end{bmatrix}$$

204 with a unit 2-norm due to orthonormal columns of  $X_*$ , leading to the inequality (9). The  
205 norm  $\|\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)^\dagger\|_2$  is finite because  $\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)$  is injective by Lemma 3.1.  $\square$

206 In light of Theorem 4.2, we can introduce an  $m \times k$  *condition number* of an eigenvalue  
207  $\lambda_* \in \text{eig}(A)$  as

$$\begin{aligned}
\tau_{A, m \times k}(\lambda_*) &:= \inf_{C, T, S} \left\| (\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)^\dagger) \right\| \\
&< \infty
\end{aligned} \tag{10}$$

208 where  $\mathbf{g}$  is the mapping defined in (1) and the infimum is taken over all the proper choices of  
209 matrix parameters  $C, T, S$  that render the columns of the unique  $X_*$  orthonormal so that  
210  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$ . We shall refer to  $\tau_{A, m \times k}(\lambda_*)$  as the *multiplicity support condition number*  
211 if the specific  $m$  and  $k$  are irrelevant in the discussion. From Lemma 3.1, the  $m \times k$   
212 condition number is infinity *only if* either  $m$  is less than the actual geometric multiplicity  
213 or  $k$  is less than the Segre characteristic anchor. Consequently, the multiplicity support  
214 condition number  $\tau_{A, m \times k}(\lambda_*)$  is large only if  $A$  is close to a matrix  $\tilde{A}$  that possesses an  
215 eigenvalue  $\tilde{\lambda} \approx \lambda_*$  whose geometric multiplicity is larger than  $m$  or its Segre characteristic  
216 anchor is larger than  $k$ .

217 We can now revisit the question:

*Is a defective eigenvalue hypersensitive to perturbations?*

218 The answer is not as simple as the question may seem to be.

219 It is well documented in the literature that, under an *arbitrary* perturbation  $\Delta A$  on  
220 the matrix  $A$ , a defective eigenvalue of  $A$  generically disperses into a cluster of simple  
221 eigenvalues with an error bound proportional to  $\|\Delta A\|_2^{\frac{1}{l}}$  where  $l$  is the size of the largest  
222 Jordan block associated with the eigenvalue [2, p. 58][3, 13]. Similar and related sensitivity  
223 results can be found in the works such as [1, 14, 15]. This error bound implies that the  
224 asymptotic sensitivity of a defective eigenvalue is infinity, and only a fraction  $\frac{1}{l}$  of the data  
225 accuracy passes on to the accuracy of the eigenvalue. For instance, if the largest Jordan  
226 block is  $5 \times 5$ , only three correct digits can be expected from the computed eigenvalues  
227 regarding the defective eigenvalue since one fifth the hardware precision (about 16 digits)  
228 remains in the forward accuracy. It is also known that the mean of the cluster emanating  
229 from the defective eigenvalue under perturbations is not hypersensitive [10, 17].

230 Kahan is the first to discover the finite sensitivity  $\frac{1}{m} \|P\|_2$  of a multiple eigenvalue under  
231 constrained perturbations that preserve the algebraic multiplicity  $m$ , where  $P$  is the spec-  
232 tral projector associated with the eigenvalue. This spectral projector norm is large only if a  
233 small perturbation on the matrix can increase the multiplicity [9]. As pointed out by Kahan,  
234 the seemingly infinite sensitivity of a defective eigenvalue may not be a conceptually mean-  
235 ingful measurement for the condition of a *multiple eigenvalue* since arbitrary perturbations  
236 do not maintain the characteristics of the eigenvalue as being multiple. Theorem 4.2 sheds  
237 light on another intriguing and pleasant property of a defective eigenvalue: Its algebraic  
238 multiplicity does *not* need to be maintained under data perturbations for its sensitivity to  
239 be under control, as long as the geometric multiplicity *and* the Segre characteristic anchor  
240 are preserved. As a result, the condition number  $\tau_{A,m \times k}(\lambda_*)$  provides a new and different  
241 measurement on the sensitivity of a *defective* eigenvalue  $\lambda_*$  when its multiplicity support  
242 is preserved. More importantly, a defective eigenvalue can be accurately computed in nu-  
243 merical computation from imposing the constraints on the multiplicity support, as we shall  
244 demonstrate in later sections.

245 Interestingly, the same eigenvalue can be ill-conditioned in the spectral projector norm while  
246 being well conditioned in multiplicity support and vice versa (c.f. Example 4 in §11), no  
247 contradiction whatsoever.

248 Even if perturbations are unconstrained, the problem of computing a defective eigenvalue  
249 may not have to be hypersensitive at all if the problem is properly generalized, i.e. regular-  
250 ized. We shall prove in Theorem 6.1 that the  $m \times k$  condition number still provides the  
251 finitely bounded sensitivity of  $\lambda_*$  as what we call the  $m \times k$  pseudo-eigenvalue of  $A$ , and  
252 this condition number is large only if  $m$  or  $k$  can be increased by small perturbations.

253 There are further subtleties on the condition of a defective eigenvalue. The sensitivity is  
254 finitely bounded if the multiplicity or the multiplicity support of *the* eigenvalue is preserved.  
255 The set  $\mathcal{E}_{m \times k}^n$  in (8) is not a manifold so the Tubular Neighborhood Theorem does not  
256 apply. As a result, requiring the matrix to maintain a multiplicity support  $m \times k$  is not  
257 enough to dampen the sensitivity of a particular defective eigenvalue with that multiplicity  
258 support. The matrix staying on  $\mathcal{E}_{m \times k}^n$  does not guarantee the finite sensitivity of a defective

259 eigenvalue. If a matrix  $A \in \mathcal{E}_{m \times k}^n$  has two eigenvalues of the same multiplicity support  
 260  $m \times k$ , then  $A$  is in the intersection of images of two holomorphic mappings described  
 261 in Corollary 4.1. When  $A$  drifts on  $\mathcal{E}_{m \times k}^n$ , the multiplicity support  $m \times k$  may be  
 262 maintained for one eigenvalue but lost on the other. Consequently, the “other” defective  
 263 eigenvalue still disperses into a cluster.

## 264 5 Backward errors near a defective eigenvalue

265 Lemma 3.1 makes it possible to calculate a defective eigenvalue in numerical computation.  
 266 When an approximate value  $\tilde{\lambda}$  of the defective eigenvalue is obtained as a component of  
 267 the approximate zero  $(A, \tilde{\lambda}, \tilde{X})$  of  $\mathbf{g}$  in (1), a backward error measure of  $\tilde{\lambda}$  can be easily  
 268 obtained from the following lemma.

269 **Corollary 5.1** *Assume  $A \in \mathbb{C}^{n \times n}$  and the mapping  $\mathbf{g}$  is defined in (1) with any param-*  
 270 *eters  $C \in \mathbb{C}^{n \times m}$ ,  $S$  as in (2) and  $T \in \mathbb{C}^{m \times k}$  with  $T_{1:m,1} \neq \mathbf{0}$ . For any  $(\tilde{\lambda}, \tilde{X}) \in \mathbb{C} \times \mathbb{C}^{n \times k}$*   
 271 *with a full-ranked  $\tilde{X}$  and residual  $(A - \tilde{\lambda}I)\tilde{X} - \tilde{X}S = E$ , the backward error bound of*  
 272  *$\tilde{\lambda}$  is  $\|E\|_2 \|\tilde{X}^\dagger\|_2$  in the sense that  $\tilde{\lambda}$  is an exact eigenvalue of  $A - E\tilde{X}^\dagger$  associated with*  
 273 *an elementary Jordan block  $J_l(\tilde{\lambda})$  of size  $l \geq k$ , and the bound is  $\|E\|_2$  if columns of  $\tilde{X}$*   
 274 *are orthonormal.*

275 **PROOF.** Since  $\tilde{X}$  is of full column rank, we have  $\tilde{X}^\dagger \tilde{X} = I$  and thus  $E = E\tilde{X}^\dagger \tilde{X}$ ,  
 276 leading to  $(A - E\tilde{X}^\dagger - \tilde{\lambda}I)\tilde{X} = \tilde{X}S$  and the assertions follow from a straightforward  
 277 verification.  $\square$

278 Corollary 5.1 provides a backward error on the approximate eigenvalue  $\tilde{\lambda}$  in the sense  
 279 that it is an exact defective eigenvalue of the matrix  $A - E\tilde{X}^\dagger$ . The backward error bound  
 280  $\|E\|_2 \|\tilde{X}^\dagger\|_2$  points to the importance of orthonormalizing columns of  $\tilde{X}$  in achieving highest  
 281 possible accuracy. Lemma 3.1 ensures that we can choose the parameter  $C$  and  $S$  in (1)  
 282 so that columns of  $\tilde{X}$  are orthonormal or nearly orthonormal. In that way, the magnitude  
 283  $\|E\|_2$  of the the residual  $E = (A - \tilde{\lambda}I)\tilde{X} - \tilde{X}S$  directly measures such a backward  
 284 error. The method of orthonormalization is elaborated in §10 along with numerical results  
 285 confirming the accuracy improvement.

286 There is another way to measure the backward error so that the computed eigenvalue is an  
 287 exact eigenvalue of a nearby matrix having an identical multiplicity support, as asserted in  
 288 the following lemma.

289 **Corollary 5.2** *Let  $A \in \mathbb{C}^{n \times n}$  and  $\lambda_* \in \text{eig}(A)$  with a multiplicity support  $m \times k$ . As*  
 290 *parameters of  $\mathbf{g}$  in (1), assume  $S \in \mathbb{C}^{k \times k}$  is as in (2),  $T \in \mathbb{C}^{m \times k}$  is as in (3) and*  
 291  *$C \in \mathbb{C}^{n \times m}$  is full-ranked such that the mapping  $\mathbf{g}(A, \lambda_*, X)$  has a unique zero  $X_*$  and*  
 292  *$\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)$  is injective. Further assume  $\tilde{A} \in \mathbb{C}^{n \times n}$  represents the available data of  $A$*   
 293 *with sufficiently small error  $\|A - \tilde{A}\|_2$ . Then every sufficiently accurate approximation  $\tilde{\lambda}$*   
 294 *of  $\lambda_*$  is an exact eigenvalue of  $\tilde{A} + \Delta\tilde{A}$  with the identical multiplicity support  $m \times k$  and*

$$\|\Delta\tilde{A}\|_2 = \|[E_1 - (\tilde{A} - \tilde{\lambda}I)C^{\text{H}\dagger}E_2 + C^{\text{H}\dagger}E_2S](\tilde{X} - C^{\text{H}\dagger}E_2)^\dagger\|_2 \quad (11)$$

295 where  $\tilde{X}$  is the least squares solution of  $\mathbf{g}(A, \tilde{\lambda}, X) = \mathbf{0}$  with the residual

$$\mathbf{g}(\tilde{A}, \tilde{\lambda}, \tilde{X}) = \begin{bmatrix} (\tilde{A} - \tilde{\lambda}I)\tilde{X} - \tilde{X}S \\ C^H \tilde{X} - T \end{bmatrix} = \begin{bmatrix} E_1 \\ E_2 \end{bmatrix}$$

296 PROOF. Since  $C \in \mathbb{C}^{n \times m}$  is of full rank and  $C^H \tilde{X} - T = E_2 = C^H C^{\dagger} E_2$ , we have

$$C^H (\tilde{X} - C^{\dagger} E_2) = T$$

297 where  $C^{\dagger}$  is the pseudo-inverse of  $C^H$ . Since  $X_*$  and  $\tilde{X}$  are least squares solutions of  
 298 linear systems  $\mathbf{g}(A, \lambda_*, X) = \mathbf{0}$  and  $\mathbf{g}(\tilde{A}, \tilde{\lambda}, X) = \mathbf{0}$  respectively, the norms  $\|X_* - \tilde{X}\|_2$ ,  
 299  $\|E_1\|_2$  and  $\|E_2\|_2$  are sufficiently small so that  $\tilde{X} - C^{\dagger} E_2$  is of full column rank. Let

$$\begin{aligned} (\tilde{A} - \tilde{\lambda}I)(\tilde{X} - C^{\dagger} E_2) - (\tilde{X} - C^{\dagger} E_2)S &= E_3 \\ &\equiv E_3 (\tilde{X} - C^{\dagger} E_2)^{\dagger} (\tilde{X} - C^{\dagger} E_2). \end{aligned}$$

300 Namely

$$[(\tilde{A} + \Delta\tilde{A}) - \tilde{\lambda}I] \hat{X} = \hat{X}S$$

301 with

$$\hat{X} = \tilde{X} - C^{\dagger} E_2 \quad \text{and} \quad \Delta\tilde{A} = -E_3 \hat{X}^{\dagger}$$

302 that satisfies (11), leading to  $\mathbf{g}(\tilde{A} + \Delta\tilde{A}, \tilde{\lambda}, \hat{X}) = \mathbf{0}$ .

303 For sufficiently small  $|\lambda_* - \tilde{\lambda}|$  and  $\|A - \tilde{A}\|_2$ , the norms  $\|E_1\|_2$ ,  $\|E_2\|_2$ ,  $\|E_3\|_2$  and  
 304  $\|\Delta\tilde{A}\|_2$  are all small. Since  $T_{2:m,1} = \mathbf{0}$ , hence

$$\begin{bmatrix} (\tilde{A} + \Delta\tilde{A}) - \tilde{\lambda}I \\ (C_{1:n,2:m})^H \end{bmatrix} \hat{X}_{1:n,1} = \mathbf{0}$$

305 with  $\hat{X}_{1:n,1} \neq \mathbf{0}$ , implying both  $\begin{bmatrix} (\tilde{A} + \Delta\tilde{A}) - \tilde{\lambda}I \\ (C_{1:n,2:m})^H \end{bmatrix}$  and  $\begin{bmatrix} A + \lambda I \\ (C_{1:n,2:m})^H \end{bmatrix}$  are rank deficient

306 with  $\|(\tilde{A} + \Delta\tilde{A} - \tilde{\lambda}I) - (A - \lambda I)\|$  sufficiently small. Thus

$$\text{nullity} \left( (\tilde{A} + \Delta\tilde{A}) - \tilde{\lambda}I \right) = \text{nullity} (A - \lambda_* I) = m$$

307 Namely the geometric multiplicity of  $\tilde{\lambda}$  as an exact eigenvalue of  $\tilde{A} + \Delta\tilde{A}$  is  $m$ .

308 Since  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$  and  $\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)$  is injective, the first column of  $X_*$  is not  
 309 in the kernel  $\mathcal{Ker}((A - \lambda_* I)^{k+1})$ . Consequently, the first column  $\hat{X}_{1:n,1}$  of  $\hat{X}$  does not  
 310 belong to the kernel  $\mathcal{Ker}((\tilde{A} + \Delta\tilde{A} - \tilde{\lambda}I)^{k+1})$  and  $\mathbf{g}_{\lambda X}(\tilde{A} + \Delta\tilde{A}, \tilde{\lambda}, \hat{X})$  is injective when  
 311  $|\lambda_* - \tilde{\lambda}|$  and  $\|A - \tilde{A}\|_2$  are sufficiently small, implying the Segre characteristic anchor of  
 312  $\tilde{\lambda} \in \text{eig}(\tilde{A} + \Delta\tilde{A})$  is exactly  $k$ .  $\square$

313 When a matrix  $A$  possesses an eigenvalue  $\lambda_*$  with a multiplicity support  $m \times k$  but is  
 314 known through empirical data in  $\tilde{A}$ , solving the system  $\mathbf{g}(\tilde{A}, \lambda, X) = \mathbf{0}$  for its nonlinear  
 315 least squares solution  $(\tilde{\lambda}, \tilde{A})$  means finding an exact eigenvalue  $\tilde{\lambda}$  of a nearby matrix

316  $\tilde{A} + \Delta\tilde{A}$  whose multiplicity support is identical to that of the underlying hidden matrix  $A$   
317 and the backward error bound on  $\Delta\tilde{A}$  is proportional to the residual  $\|\mathbf{g}(\tilde{A}, \tilde{\lambda}, \tilde{X})\|_F$  as  
318 shown in (11). Minimizing the residual reduces, if not minimizes, the distance from  $\tilde{A}$  to  
319 the set  $\mathcal{E}_{m \times k}^n$  of matrices with a multiplicity support  $m \times k$ . As a result, Theorem 4.2  
320 ensures the computed eigenvalue  $\tilde{\lambda}$  is accurate in the sense that the forward error  $|\tilde{\lambda} - \lambda_*|$   
321 is in the order of the data error  $\|A - \tilde{A}\|_2$ .

322 Requiring  $|\tilde{\lambda} - \lambda_*|$  to be sufficiently small is to ensure  $\tilde{\lambda}$  approximates the correct eigenvalue  
323  $\lambda_*$  in case there are two or more eigenvalues of  $A$  sharing the same multiplicity support.

## 324 6 A well-posed defective eigenvalue problem

325 The problem of finding an eigenvalue of a matrix in its conventional meaning is ill-posed when  
326 the eigenvalue is multiple and defective because the sensitivity of the eigenvalue is infinite  
327 with respect to arbitrary perturbations on the matrix. Lacking a Lipschitz continuity of an  
328 eigenvalue with respect to data, such a problem is not suitable for numerical computation  
329 unless the problem is properly modified, or better known as being *regularized*.

330 We can alter the problem of

*finding an eigenvalue of a matrix  $A$*

331 to

*finding a  $\tilde{\lambda}$  so that  $(\tilde{\lambda}, \tilde{X})$  is local least squares solution to  $\mathbf{g}(A, \lambda, X) = \mathbf{0}$*

332 where  $\mathbf{g}$  is the mapping defined in (1) with proper parameters. When the matrix  $A$  has  
333 an eigenvalue  $\lambda_*$  of multiplicity support  $m \times k$ , there is an  $X_*$  such that  $(\lambda_*, X_*)$  is an  
334 exact solution to  $\mathbf{g}(A, \lambda, X) = \mathbf{0}$  with zero residual. However, when  $A$  is known through  
335 its empirical data in  $\tilde{A}$ , a local least squares solution  $(\tilde{\lambda}, \tilde{X})$  to the system  $\mathbf{g}(\tilde{A}, \lambda, X) = \mathbf{0}$   
336 generally has a residual  $\|\mathbf{g}(\tilde{A}, \tilde{\lambda}, \tilde{X})\|_2 > 0$ , and  $\tilde{\lambda}$  is not an eigenvalue of either  $A$  or  $\tilde{A}$ .  
337 For the convenience of elaboration, we call such a  $\tilde{\lambda}$  an  $m \times k$  *pseudo-eigenvalue* of  $\tilde{A}$ .

338 **Theorem 6.1 (Pseudo-Eigenvalue Theorem)** *Let  $\lambda_*$  be an eigenvalue of  $A \in \mathbb{C}^{n \times n}$*   
339 *with a multiplicity support  $m \times k$  along with  $X_* \in \mathbb{C}^{n \times k}$  satisfying  $\mathbf{g}(A, \lambda_*, X_*) = \mathbf{0}$  with*  
340 *a properly defined mapping  $\mathbf{g}$  as in (1). The following assertions hold.*

- 341 (i) *The exact eigenvalue  $\lambda_*$  of  $A$  is an  $m \times k$  pseudo-eigenvalue of  $A$ .*
- 342 (ii) *There are neighborhoods  $\Phi$  of  $A$  in  $\mathbb{C}^{n \times n}$  and  $\Lambda$  of  $\lambda_*$  in  $\mathbb{C}$  respectively such*  
343 *that every matrix  $\tilde{A} \in \Phi$  has a unique  $m \times k$  pseudo-eigenvalue  $\tilde{\lambda} \in \Lambda$  that is*  
344 *Lipschitz continuous with respect to  $\tilde{A}$ .*
- 345 (iii) *For every matrix  $\tilde{A} \in \Phi$  serving as empirical data of  $A$ , there is a unique  $m \times k$*   
346 *pseudo-eigenvalue  $\tilde{\lambda} \in \Lambda$  of  $\tilde{A}$  such that*

$$|\tilde{\lambda} - \lambda_*| \leq \tau_{A, m \times k}(\lambda_*) \|\tilde{A} - A\|_2 + O\left(\|\tilde{A} - A\|_2^2\right) \quad (12)$$

347 where  $\tau_{A, m \times k}(\lambda_*)$  is the  $m \times k$  condition number of  $\lambda_* \in \text{eig}(A)$  defined in (10).

348 PROOF. The assertion (i) is a result of Lemma 3.1, part (i). Let  $\Psi$  be a small  
 349 neighborhood of  $(\lambda_*, X_*)$  in  $\mathbb{C} \times \mathbb{C}^{m \times k}$  so that  $\{A\} \times \overline{\Psi}$  is a subset of  $\Sigma$  in Corollary 4.1.  
 350 To prove the existence of a pseudo-eigenvalue of any matrix near  $A$ , assume there is a  
 351 matrix  $\tilde{A}$  with  $\|\tilde{A} - A\|_2 < \varepsilon$  for any  $\varepsilon$  such that  $\min_{(\lambda, X) \in \overline{\Psi}} \|\mathbf{g}(\tilde{A}, \lambda, X)\|_2$  are not  
 352 attainable in  $\Psi$ . Let  $\varepsilon \rightarrow 0$ . Then  $\tilde{A} \rightarrow A$  and there exists an  $(\hat{\lambda}, \hat{X}) \in \overline{\Psi} \setminus \Psi$  such that

$$\mathbf{g}(A, \hat{\lambda}, \hat{X}) = \min_{(\lambda, X) \in \overline{\Psi}} \|\mathbf{g}(A, \lambda, X)\|_2 = 0$$

353 and  $\hat{\lambda} \neq \lambda(\mathbf{z}_*)$ . This is a contradiction to Corollary 4.1. As a result, there are neighborhoods  
 354  $\Phi$  of  $A$  and  $\Lambda$  of  $\lambda_*$  respectively such that every matrix in  $\Phi$  has an  $m \times k$  pseudo-  
 355 eigenvalue in  $\Lambda$ . By the local convergence of the Gauss-Newton iteration, we can assume  
 356  $\Phi$  and  $\Psi$  are sufficiently small so that the residual  $\|\mathbf{g}(G, \lambda, X)\|_2$  is small for all  $G \in \Phi$   
 357 and  $(\lambda, X) \in \Psi$ . Consequently, for every  $\tilde{A} \in \Phi$ , the Gauss-Newton iteration converges  
 358 to a local minimum point  $(\tilde{\lambda}, \tilde{X})$  of  $\|\mathbf{g}(\tilde{A}, \lambda, X)\|_2$  from any initial iterate  $(\lambda_0, X_0) \in \Psi$ .  
 359 As a result, this local minimum points is unique in  $\Psi$  and is thus the absolute minimum  
 360 point in  $\Psi$ . To prove the Lipschitz continuity of the pseudo-eigenvalue, let  $\tilde{A}, \check{A} \in \Phi$  with  
 361 corresponding minimum points  $(\tilde{\lambda}, \tilde{X})$  and  $(\check{\lambda}, \check{X})$  of  $\|\mathbf{g}\|_2$  in  $\Psi$ . The Gauss-Newton  
 362 iteration from  $(\tilde{\lambda}, \tilde{X})$  toward  $(\check{\lambda}, \check{X})$  in a single step

$$(\lambda_1, X_1) = (\tilde{\lambda}, \tilde{X}) - \mathbf{g}_{\lambda X}(\tilde{A}, \tilde{\lambda}, \tilde{X})^\dagger \mathbf{g}(\tilde{A}, \tilde{\lambda}, \tilde{X})$$

363 yields  $\|(\lambda_1, X_1) - (\check{\lambda}, \check{X})\|_2 \leq \mu \|(\tilde{\lambda}, \tilde{X}) - (\check{\lambda}, \check{X})\|_2$  with  $0 \leq \mu < 1$ . Using the identity  
 364  $(\tilde{\lambda}, \tilde{X}) = (\tilde{\lambda}, \tilde{X}) - \mathbf{g}_{\lambda X}(\tilde{A}, \tilde{\lambda}, \tilde{X})^\dagger \mathbf{g}(\tilde{A}, \tilde{\lambda}, \tilde{X})$ , there is a constant  $\gamma$  such that

$$\begin{aligned} \|(\tilde{\lambda}, \tilde{X}) - (\check{\lambda}, \check{X})\|_2 &\leq \frac{1}{1 - \mu} \|(\lambda_1, X_1) - (\tilde{\lambda}, \tilde{X})\|_2 \\ &\leq \frac{1}{1 - \mu} \left( \|\mathbf{g}_{\lambda X}(\tilde{A}, \tilde{\lambda}, \tilde{X})^\dagger\|_2 \|\mathbf{g}(\tilde{A}, \tilde{\lambda}, \tilde{X}) - \mathbf{g}(\check{A}, \tilde{\lambda}, \tilde{X})\|_2 \right. \\ &\quad \left. + \|\mathbf{g}_{\lambda X}(\tilde{A}, \tilde{\lambda}, \tilde{X})^\dagger - \mathbf{g}_{\lambda X}(\check{A}, \tilde{\lambda}, \tilde{X})^\dagger\|_2 \|\mathbf{g}(\tilde{A}, \tilde{\lambda}, \tilde{X})\|_2 \right) \\ &\leq \gamma \|\tilde{A} - \check{A}\|_2 \end{aligned}$$

365 for all  $\tilde{A}, \check{A} \in \Phi$ . Namely, the  $m \times k$  pseudo-eigenvalue is Lipschitz continuous with respect  
 366 to the matrix. Furthermore, by setting  $(\tilde{A}, \tilde{\lambda}, \tilde{X}) = (A, \lambda_*, X_*)$  in the above inequalities  
 367 we have (12) because the residual  $\|\mathbf{g}(A, \lambda_*, X_*)\|_2 = 0$  and thus  $\mu = 0$ .  $\square$

368 The above Pseudo-Eigenvalue Theorem establishes a rigorous and thorough regularization of  
 369 the ill-posed problem in computing a defective eigenvalue so that the problem of computing  
 370 a pseudo-eigenvalue enjoys unique existence and Lipschitz continuity of the solution that  
 371 approximates the underlying defective eigenvalue with an error bound proportional to the  
 372 data error. Furthermore, this theorem reaffirms the  $m \times k$  condition number as a bona  
 373 fide sensitivity measure of defective eigenvalues.

374 This regularization makes it possible to compute defective eigenvalues accurately using float-  
 375 ing point arithmetic even if the matrix data are perturbed, and we shall present such an  
 376 algorithm in next section.

## 7 An algorithm for computing a defective eigenvalue

The Pseudo-eigenvalue Theorem sets the foundation for accurate computation of a defective eigenvalue even if the the matrix data are empirical, provided that the multiplicity support can be obtained. The computation is under the assumptions that the given matrix  $A$  is the data representation of a underlying matrix possessing a defective eigenvalue and an initial estimate  $\lambda_0$  is close to that eigenvalue. Assuming the multiplicity support  $m \times k$  is known, identified or estimated, we also need to set up the matrix parameters  $C \in \mathbb{C}^{n \times m}$  and  $S \in \mathbb{C}^{k \times k}$ , while using the parameter  $T \in \mathbb{C}^{m \times k}$  as in (3). To elaborate the strategy for setting up  $C$ , we consider  $\lambda_* \in \text{eig}(A)$  with a multiplicity support  $m \times k$ . The requirements for the parameter  $C$  are two-fold:

(i) The parameter matrix  $C$  of  $\mathbf{g}$  in (1) must be generic enough so that the matrix  $\begin{bmatrix} A - \lambda_* I \\ C^H \end{bmatrix}$  is injective and  $\begin{bmatrix} A - \lambda_0 I \\ (C_{1:n,2:m})^H \end{bmatrix}$  is rank-deficient.

(ii) The parameter  $C$  must be generic enough so that the linear system

$$\begin{bmatrix} A - \lambda_* I \\ C^H \end{bmatrix} \mathbf{z} = \begin{bmatrix} \mathbf{0} \\ T_{1:m,1} \end{bmatrix} \quad (13)$$

has a unique solution that does not belong to the kernel  $\mathcal{Ker}((\tilde{A} - \tilde{\lambda} I)^{k+1})$ .

Let  $N \in \mathbb{C}^{n \times m}$  be a matrix whose columns form an orthonormal basis for the kernel  $\mathcal{Ker}(A - \lambda_* I)$ . A naive choice would be setting  $C = N$ , which certainly satisfies the requirement (i), and the equation (13) yields the solution  $\mathbf{z}$  that is identical to the first column of  $N$  and may be closed to  $\mathcal{Ker}((\tilde{A} - \tilde{\lambda} I)^{k+1})$ . Instead, let  $\mathbf{d} = Q \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix}$  be a random unit column vector and its QR decomposition. Then  $\mathbf{z} = N \mathbf{d}$  is a random vector in the kernel of  $A - \lambda_* I$  and, for almost all such random vector  $\mathbf{d}$ , does not belong to the kernel  $\mathcal{Ker}((\tilde{A} - \tilde{\lambda} I)^{k+1})$ . Thus we can set  $C = N Q$  so that the solution  $\mathbf{z}$  of the equation (13) equals to  $N \mathbf{d}$ .

Of course, we do not have the exact value of  $\lambda_*$  and we likely do not have the exact data for  $A$  either. We can assume the data are reasonably accurate and the initial estimate  $\lambda_0$  is close to  $\lambda_*$ . We can practically ensure these requirements of  $C$  to be satisfied by the following simple process: Let  $N \in \mathbb{C}^{n \times m}$  be the matrix whose columns form an orthonormal basis for the numerical kernel  $\mathcal{Ker}_\theta(A - \lambda_0 I)$  of dimension  $m$ . Namely, columns of  $N$  span the vector space spanned by the last  $m$  right singular vectors of  $A - \lambda_0 I$ . Such an  $N$  can be computed by a rank-revealing method such as [12] or a singular value decomposition. Set  $C = N Q$  where  $Q R = \mathbf{d}$  is the QR decomposition of a random vector  $\mathbf{d}$ .

407 With  $C$  and  $T$  already available, we can set up

$$\mathbf{x}_1^{(0)} = C \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix} \quad (14)$$

$$\mathbf{x}_{j+1}^{(0)} = \alpha_j \begin{bmatrix} A - \lambda_0 I \\ C^H \end{bmatrix}^\dagger \begin{bmatrix} \mathbf{x}_j^{(0)} \\ \mathbf{0} \end{bmatrix} \quad \text{for } j = 1, \dots, k-1 \quad (15)$$

$$S = \left[ \begin{array}{c|ccc} 0 & \alpha_1 & & \\ \vdots & & \ddots & \\ 0 & & & \alpha_{k-1} \\ \hline 0 & 0 & \dots & 0 \end{array} \right] \quad (16)$$

408 where, for  $j = 1, \dots, k-1$ , the scalar  $\alpha_j$  scales  $\mathbf{x}_{j+1}^{(0)}$  into a unit vector. Denote  
 409  $X_0 = [\mathbf{x}_1^{(0)}, \dots, \mathbf{x}_k^{(0)}]$ . Then  $\mathbf{g}(A, \lambda_0, X_0) \approx \mathbf{0}$  and we can apply the Gauss-Newton  
 410 iteration

$$(\lambda_{j+1}, X_{j+1}) = (\lambda_j, X_j) - \mathbf{g}_{\lambda X}(A, \lambda_j, X_j)^\dagger \mathbf{g}(A, \lambda_j, X_j) \quad (17)$$

for  $j = 0, 1, \dots$

411 that converges to  $(\lambda_*, X_*)$  assuming the initial estimate  $\lambda_0$  is sufficiently close to  $\lambda_*$ .  
 412 When the iteration stops at the  $j$ -th step, a QR decomposition of the matrix representing  
 413  $\mathbf{g}_{\lambda X}(A, \lambda_j, X_j)$  is available and thus an estimate  $\|\mathbf{g}_{\lambda X}(A, \lambda_j, X_j)^\dagger\|_2$  of the  $m \times k$  condition  
 414 number can be computed by a couple of steps of inverse iteration [12] with a negligible cost.

415 A pseudo-code of Algorithm EIGENITERATION is given in Fig. 2.

## 416 8 Taking advantage of the Jacobian structure

417 The main computing cost of Algorithm EIGENITERATION occurs at solving the linear system

$$\mathbf{g}_{\lambda X}(A, \lambda_j, X_j)(\sigma, Y) = \mathbf{g}(A, \lambda_j, X_j)$$

418 where the partial Jacobian  $\mathbf{g}_{\lambda X}(A, \lambda_j, X_j)$  is of dimensions  $(nk + mk) \times (nk + 1)$  and its  
 419 QR decomposition is needed. The partial Jacobian is pleasantly structured with a proper  
 420 arrangement so that the cost of QR decomposition can be reduced substantially.

421 Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_k]$ , the image  $\mathbf{g}(G, \lambda, X) \in \mathbb{C}^{n \times k} \times \mathbb{C}^{m \times k}$  of the mapping  $\mathbf{g}$  can be  
 422 arranged as

$$\mathbf{g}(G, \lambda, [\mathbf{x}_1, \dots, \mathbf{x}_k]) = \begin{bmatrix} C^H \mathbf{x}_k & & & & & -T_{1:m,k} \\ (G - \lambda I) \mathbf{x}_k & -s_{k-1,k} \mathbf{x}_{k-1} & -s_{k-2,k} \mathbf{x}_{k-2} & \cdots & -s_{1k} \mathbf{x}_1 & \\ & C^H \mathbf{x}_{k-1} & & & & -T_{1:m,k-1} \\ & (G - \lambda I) \mathbf{x}_{k-1} & -s_{k-2,k-1} \mathbf{x}_{k-2} & \cdots & -s_{1,k-1} \mathbf{x}_1 & \\ & & \ddots & & \vdots & \vdots \\ & & & \ddots & \vdots & \vdots \\ & & & & \ddots & \vdots \\ & & & & & -s_{12} \mathbf{x}_1 & \vdots \\ & & & & & & C^H \mathbf{x}_1 & -T_{1:m,1} \\ & & & & & & (G - \lambda I) \mathbf{x}_1 & \end{bmatrix}.$$



**Algorithm EIGENITERATION**

INPUT: matrix  $A$ , initial eigenvalue estimate  $\lambda_0$ , geometric multiplicity  $m$ , Segre characteristic anchor  $k$ .

- obtain the last  $m$  right singular vectors of  $A - \lambda_0 I$  to form  $N \in \mathbb{C}^{m \times m}$
- obtain a random vector  $\mathbf{d} \in \mathbb{C}^m$  and its QR decomposition  $\mathbf{d} = QR$ .
- set  $C = NQ$  and  $\mathbf{x}_1^{(0)}$  as its first column.
- for  $j = 1, 2, \dots, k-1$  do

solve  $\begin{bmatrix} A - \lambda_0 I \\ C^H \end{bmatrix} \mathbf{u} = \begin{bmatrix} \mathbf{x}_j^{(0)} \\ \mathbf{0} \end{bmatrix}$  for the least squares solution  $\mathbf{u}$ .

set  $\mathbf{x}_{j+1}^{(0)} = \alpha_j \mathbf{u}$  so that  $\|\mathbf{x}_{j+1}^{(0)}\|_2 = 1$ .

end do

- set  $X_0 = [\mathbf{x}_1^{(0)}, \dots, \mathbf{x}_k^{(0)}]$  and set  $S$  as in (16)
- set the mapping  $\mathbf{g}$  as in (1) using  $C$  and  $S$
- for  $j = 0, 1, \dots$  do

\* solve the linear system  $\mathbf{g}_{\lambda X}(A, \lambda_j, X_j)(\sigma, Y) = \mathbf{g}(A, \lambda_j, X_j)$  for the least squares solution  $(\sigma, Y)$

\* set  $\lambda_{j+1} = \lambda_j - \sigma$ ,  $X_{j+1} = X_j - Y$ .

\* if  $\|\mathbf{g}(A, \lambda_j, X_j)\|_2 < \|\mathbf{g}(A, \lambda_{j+1}, X_{j+1})\|_2$  then

set  $(\hat{\lambda}, \hat{X}) = (\lambda_j, X_j)$ , break the loop. end if

end do

OUTPUT: eigenvalue  $\hat{\lambda}$ , backward error bound  $\|\mathbf{g}(A, \hat{\lambda}, \hat{X})\|_2 \|X^\dagger\|_2$ ,  $m \times k$

condition number  $\|\mathbf{g}_{\lambda X}(A, \hat{\lambda}, \hat{X})^\dagger\|_2$

Figure 2: Algorithm EIGENITERATION

423 As a result, the partial Jacobian matrix in a blockwise upper-triangular form

$$\frac{\partial \mathbf{g}(G, \lambda, X)}{\partial (\mathbf{x}_k, \dots, \mathbf{x}_1, \lambda)} = \begin{bmatrix} C^H & O & O & \dots & O & \mathbf{0} \\ G - \lambda I & -s_{k-1,k} I & -s_{k-2,k} I & \dots & -s_{1k} I & -\mathbf{x}_k \\ & C^H & O & \dots & O & \mathbf{0} \\ & G - \lambda I & -s_{k-2,k-1} I & \dots & -s_{1,k-1} I & -\mathbf{x}_{k-1} \\ & & \ddots & \ddots & \vdots & \vdots \\ & & & \ddots & \vdots & \vdots \\ & & & & C^H & O \\ & & & & G - \lambda I & -s_{12} I \\ & & & & & C^H \\ & & & & & G - \lambda I & -\mathbf{x}_1 \end{bmatrix}.$$

424 We can further assume the matrix  $G$  is already reduced to a Hessenberg form. Then the  
 425 matrix block

$$\begin{bmatrix} C^H \\ G - \lambda I \end{bmatrix} = \begin{bmatrix} * & * & \cdots & * \\ \vdots & \vdots & \ddots & \vdots \\ * & * & \cdots & * \\ & * & \cdots & * \\ & & \ddots & \vdots \\ & & & * \end{bmatrix}$$

426 is nearly upper-triangular with  $m + 1$  subdiagonal lines of nonzero entries. The QR  
 427 decomposition of the partial Jacobian  $\mathbf{g}_{\mathbf{x}_k \dots \mathbf{x}_1 \lambda}(G, \lambda, X)$  can then be carried out by a  
 428 sequence of standard textbook Householder transformations.

## 429 9 Identifying the multiplicity support

430 Identification of the geometric multiplicity can be carried out with numerical rank-revealing.  
 431 Let  $\lambda_0$  be an initial estimate of  $\lambda_* \in \text{eig}(A)$  in Lemma 3.1 and assume

$$|\lambda_0 - \lambda_*| < \theta < \min_{\lambda \in \text{eig}(A) \setminus \{\lambda_*\}} |\lambda - \lambda_0|.$$

432 The geometric multiplicity of  $\lambda_*$  can be computed as the numerical nullity of  $A - \lambda_0 I$   
 433 within the error tolerance  $\theta$  defined as

$$m = \max \{j \mid \sigma_{n-j+1}(A - \lambda_0 I) < \theta\} \quad (18)$$

434 where  $\sigma_i(\cdot)$  is the  $i$ -th largest singular value of  $(\cdot)$ . A misidentification of the geometric  
 435 multiplicity can be detected during the computation. Underestimating  $m$  results in an  
 436 undersized  $C$  in (1) so that both  $\begin{bmatrix} A - \lambda_* I \\ C^H \end{bmatrix}$  and the partial Jacobian  $\mathbf{g}_{\lambda X}(A, \lambda_*, X_*)$   
 437 are rank-deficient. Overestimating  $m$  renders the system

$$\begin{bmatrix} A - \lambda_* I \\ C^H \end{bmatrix} \mathbf{u} = \begin{bmatrix} 0 \\ T_{1:m,1} \end{bmatrix}$$

438 inconsistent with a large residual norm. During an iteration in which  $(\lambda_j, X_j)$  approaches  
 439  $(\lambda_*, X_*)$ , a large condition number of the partial Jacobian  $\mathbf{g}_{\lambda, X}(A, \lambda_j, X_j)$  indicates a  
 440 likely underestimated geometric multiplicity and a large residual  $\|\mathbf{g}(A, \lambda_j, X_j)\|_2$  suggests  
 441 a possible overestimation.

442 If the geometric multiplicity is identified, it is possible to find the Segre characteristic anchor  
 443 by a searching scheme based on the condition number of the Jacobian  $\mathbf{g}_{\lambda X}$ .

444 **Example 1** Let

$$A = \begin{bmatrix} 0 & 4 & 0 & -4 & 0 & -2 & 1 & 0 & 0 & -1 & -1 & -1 & 2 & 1 & 0 & 0 & -1 & 0 & 0 \\ 0 & 3 & 3 & -4 & 1 & 0 & 4 & 1 & -1 & -1 & 0 & 0 & 0 & 0 & -2 & 0 & 0 & -1 & -1 \\ 1 & -4 & 2 & 10 & -3 & 1 & -7 & -2 & -1 & 3 & 2 & 0 & 0 & -1 & 0 & 3 & -1 & 1 & 1 & 2 \\ -1 & -1 & 2 & 5 & -2 & -1 & -5 & -1 & -1 & 2 & 0 & -1 & 0 & 1 & 0 & 2 & -1 & -1 & -1 & 1 \\ -1 & -2 & 2 & 1 & 1 & 1 & -1 & 0 & -2 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & -1 & -1 & 0 & 0 \\ 1 & 4 & 1 & -12 & 4 & 2 & 13 & 3 & 0 & -4 & 0 & 0 & -2 & -1 & 1 & -6 & 1 & 1 & 0 & -3 \\ -1 & -1 & 1 & 5 & -2 & 0 & -4 & -2 & 0 & 1 & -1 & 0 & 0 & 1 & 0 & 4 & 0 & -1 & -1 & 2 \\ 1 & 2 & -4 & 0 & 1 & 0 & 1 & 1 & 4 & -2 & -1 & 1 & 0 & -1 & 0 & 1 & 2 & 1 & 0 & 1 \\ 0 & -5 & 2 & 10 & -5 & -1 & -10 & -2 & -1 & 6 & 3 & -2 & 0 & 0 & 0 & 3 & -3 & 0 & 1 & 2 \\ 1 & 1 & 1 & -1 & 2 & 2 & 4 & 0 & 1 & -1 & -1 & 2 & 0 & -1 & 0 & 0 & 2 & 1 & -1 & 0 \\ 1 & -1 & 0 & 2 & 1 & 2 & 1 & 0 & 1 & -1 & 3 & 2 & 0 & -1 & 0 & 0 & 1 & 1 & -1 & 0 \\ -1 & -3 & 0 & 5 & -1 & 2 & -4 & -1 & 0 & 1 & -1 & 4 & 4 & 1 & -2 & 2 & 0 & -1 & 0 & 1 \\ -2 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 3 & 2 & 0 & 0 & 0 & -1 & 0 & 0 & 0 \\ -3 & 4 & -1 & -4 & 0 & -2 & 1 & 0 & 0 & -1 & -1 & -2 & 5 & 2 & 0 & 0 & -1 & 1 & 0 & 0 \\ -2 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 0 & 2 & 3 & 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & -2 & 3 & 1 & 2 & -1 & -1 & 2 & -2 & 2 & 0 & 0 & 0 & 5 & 2 & 0 & 0 & 0 & 1 \\ 6 & 3 & -6 & 3 & 6 & 4 & 7 & 0 & 7 & -7 & 1 & 5 & -2 & -6 & 1 & 0 & 8 & 6 & -1 & 0 & 0 \\ 0 & 2 & -4 & -4 & 1 & -1 & 4 & 1 & 0 & -1 & 0 & -1 & -1 & 0 & 1 & -2 & 0 & 3 & 4 & -1 & 0 \\ 1 & -4 & -1 & 11 & -4 & 1 & -8 & -3 & -1 & 3 & 2 & 0 & 0 & -1 & 0 & 4 & -1 & 1 & 4 & 3 & 0 \\ 0 & 0 & -1 & 1 & -2 & 0 & -1 & -2 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 4 & 0 \end{bmatrix}$$

445 with  $\text{eig}(A) = \{2, 3\}$  of nonzero Segre characteristics  $\{4, 3, 3\}$  and  $\{5, 5\}$  respectively.

446 Applying the Francis QR algorithm implemented in Matlab yields computed eigenvalues  
447 scattered around  $\lambda_1 = 2.0$  and  $\lambda_2 = 3.0$ :

```

2.000118556521482 + 0.000118397929590i      3.000398490901253 + 0.001224915665189i
2.000118556521482 - 0.000118397929590i      3.000398490901253 - 0.001224915665189i
1.999881443477439 + 0.000118714860725i      3.000646066870935 + 0.000469627646058i
1.999881443477439 - 0.000118714860725i      3.000646066870935 - 0.000469627646058i
2.000013778528383 + 0.000018105742295i      3.001287762162967 + 0.000000000000000i
2.000013778528383 - 0.000018105742295i      2.999753002133234 + 0.000759191914332i
2.000008786021464 + 0.000020979720849i      2.999753002133234 - 0.000759191914332i
2.000008786021464 - 0.000020979720849i      2.998957628017279 + 0.000757681758834i
1.999977435451235 + 0.000002873978888i      2.998957628017279 - 0.000757681758834i
1.999977435451235 - 0.000002873978888i      2.999201861991639 + 0.000000000000000i

```

448 Using two computed eigenvalues above, say

$$\tilde{\lambda}_0 = 1.999881443477439 - 0.000118714860725i \quad \text{and} \quad \hat{\lambda}_0 = 3.001287762162967 + 0.000000000000000i$$

449 as initial estimates of the defective eigenvalues, smallest singular values of  $A - \tilde{\lambda}_0 I$  and  
450  $A - \hat{\lambda}_0 I$  can be computed using a rank-revealing method as

$\sigma_j(A - \tilde{\lambda}_0 I) :$	...	$\sigma_j(A - \hat{\lambda}_0 I) :$	...
0.084065699924186		0.070046163725993	
0.049368630759014		0.054661269198836	
<b>0.0000000000003635</b>		0.036234932328447	
<b>0.0000000000000280</b>		<b>0.0000000000000001</b>	
<b>0.0000000000000001</b>		<b>0.0000000000000000</b>	

451 indicating the geometric multiplicities 3 and 2 respectively.

452 Set the geometric multiplicities for the initial eigenvalue estimate  $\tilde{\lambda}_0$  and  $\hat{\lambda}_0$  as 3 and 2  
453 respectively. Apply Algorithm EIGENITERATION with increasing input  $k = 1, 2, \dots, 6$  as  
454 estimated Segre characteristic anchors, we list the computed eigenvalues, Jacobian condition  
455 numbers and residual norms in Table 1.

456 At  $\lambda_1$ , for instance, underestimated values  $k = 1, 2$  render the condition numbers of the  
457 Jacobians as large as  $10^8$  and the residuals to be tiny, while the overestimated value  $k = 4$   
458 leads to a drastic increase of residual from  $10^{-16}$  to  $10^{-3}$  but maintains the moderate  
459 Jacobian condition number, as shown in Table 1. Similar effect of increasing estimated  
460 values of the Segre characteristic anchor at  $\lambda_2$  can be observed consistently. □

test		at $\lambda_1 = 2$ , Segre characteristic anchor $k = 3$		
$k$	value	computed eigenvalue	condition number	residual norm
$k = 1$		1.999881443477439 - 0.000118714860725i	560995239.6	0.000000000000001
$k = 2$		1.999999993438010 - 0.00000011324234i	147603979.2	0.000000000000001
$\rightarrow k = 3$		2.000000000000000 - 0.000000000000000i	58.7	0.000000000000006 ←
$k = 4$		2.109885640097783 - 0.004348977611146i	24.1	0.007
test		at $\lambda_2 = 3$ , Segre characteristic anchor $k = 5$		
$k$	value	computed eigenvalue	condition number	residual norm
$k = 1$		3.001287762162967	2161090332264.6	0.000000000000003
$k = 2$		3.001287762162967	796260062.8	0.000000000000005
$k = 3$		3.001287762162967	4556940.4	0.000000003
$k = 4$		3.000000013572103	687859583.9	0.000000000000007
$\rightarrow k = 5$		3.000000000000000	33.9	0.000000000000007 ←
$k = 6$		3.002451613695432	34.1	0.007

Table 1: Effect of increasing estimated Segre characteristic anchors: Underestimated values yield large condition numbers of the Jacobian and overestimated values lead to large residual norms. The results using the correct anchors are pointed out with arrows.

## 10 Improving accuracy with orthonormalization

Algorithm EIGENITERATION uses a simple nilpotent matrix  $S$  with only one superdiagonal line of nonzero entries. By Lemma 3.1 (iii), we can modify  $C$  and  $S$  as parameters of  $\mathbf{g}$  so that the matrix component  $\tilde{X}$  of the solution to  $\mathbf{g}(A, \lambda_*, \tilde{X}) = \mathbf{0}$  has orthonormal columns. Assuming the parameter matrix  $T$  is set as in (3), the orthonormalization process can be carried out simply by the following process:

- Execute Algorithm EIGENITERATION and obtain output  $\hat{\lambda}, \hat{X}, C, S$ .
- Obtain the thin QR decomposition  $\hat{X} = QR$ .
- Reset  $C$  and  $S$  with  $Q^H C$  and  $RSR^{-1}$  respectively in the mapping  $\mathbf{g}$ .
- Starting from  $(\lambda_0, X_0) = (\hat{\lambda}, Q)$  and execute the Gauss-Newton iteration (17) and obtain a refined eigenvalue.

The reason for such an orthonormalization is intuitively clear. When we solve for the least squares solution  $(A, \tilde{\lambda}, \tilde{X})$  of the equation  $\mathbf{g}(A, \lambda, X) = \mathbf{0}$  minimizing the magnitude of the residual  $(A - \tilde{\lambda}I)\tilde{X} - \tilde{X}S = E$ , the backward error given in Corollary 5.1 is  $\|E\|_2 \|\tilde{X}^\dagger\|_2$ . When the norm  $\|\tilde{X}^\dagger\|_2$  is large, minimizing the residual norm  $\|E\|_2$  may not achieve the highest backward accuracy. If the columns of  $\tilde{X}$  are orthonormal, however, the norm  $\|\tilde{X}^\dagger\|_2 = 1$  and the least squares solution that minimizing the residual norm  $\|E\|_2$  directly minimizes the backward error bound.

**Example 2** Consider the matrix

$$A = \begin{bmatrix} 2 & & & & \\ & -8 & & & \\ & & 1 & & \\ & & & 2 & \\ & & & & 1 \\ & -10000 & 1000 & -100 & 12 \end{bmatrix} \quad (19)$$

480 with an exact eigenvalue  $\lambda_* = 2$  and the multiplicity support  $1 \times 5$ . A straightforward  
 481 application of Algorithm EIGENITERATION in Matlab yields

$$\begin{aligned} \tilde{\lambda} &= 1.999999999999748 \\ S &= \begin{bmatrix} 0 & 0.100686223197184 & 0 & 0 & 0 \\ 0 & 0 & 0.680272615629152 & 0 & 0 \\ 0 & 0 & 0 & 0.786924421181882 & 0 \\ 0 & 0 & 0 & 0 & 0.922632632948520 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\ \tilde{X} &= \begin{bmatrix} 1.00502786434024 & 0.10210319200724 & .07627239342106 & .06542640851275 & .06584192219606 \\ -0.00000000000025 & 0.10119245986833 & .06945800549083 & .06002060904501 & .06036453955047 \\ -0.000000000000253 & 1.01192459868319 & .76341851426484 & .65486429121738 & .65902236805904 \\ -0.000000000000000 & -0.00000000000051 & .68838459356550 & .60075267245722 & .60419916522969 \\ -0.000000000000000 & -0.00000000000000 & -.00000000000052 & .54170664784191 & .55427401993992 \end{bmatrix} \end{aligned}$$

482 The residual norm

$$\|(A - \tilde{\lambda}I) \tilde{X} - \tilde{X} S\|_F \approx 4.5 \times 10^{-14}$$

483 can not be minimized further with the unit round-off about  $10^{-16}$  considering  $\|A\|_2 \approx 10^4$ .

484 The backward error

$$\|(A - \tilde{\lambda}I) \tilde{X} - \tilde{X} S\|_F \|\tilde{X}^\dagger\|_2 \approx 1.3 \times 10^{-9}$$

485 is not small enough. After orthonormalization and resetting the resulting parameter  $C$   
 486 and  $S$  in  $\mathbf{g}$  in (1), we apply the Gauss-Newton iteration again and obtain

$$\begin{aligned} \hat{\lambda} &= 2.000000000000000 \\ S &= \begin{bmatrix} 0 & 0.09950371902 & -0.00990049999 & 0.00099000050 & -0.99498744208 \\ 0 & 0 & 1.00493781395 & -0.00098508732 & 0.99004950866 \\ 0 & 0 & 0 & 1.00004900870 & -0.09850873917 \\ 0 & 0 & 0 & 0 & 10050.38307728113 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\ \hat{X} &= \begin{bmatrix} -1.0 & -0.00000000002531 & -0.000000000000000 & -0.000000000000000 & -0.000000000000000 \\ 0.0 & -0.099503719021067 & 0.009900499987341 & -0.000990000499944 & 0.994987442082474 \\ 0.0 & -0.995037190210673 & -0.000990049999192 & 0.000099000049994 & -0.099498744209908 \\ 0.0 & -0.000000000000000 & -0.999950498725976 & -0.000009900005712 & 0.009949874421338 \\ 0.0 & -0.000000000000000 & -0.000000000000000 & -0.999999505000536 & -0.000994987010945 \end{bmatrix} \end{aligned}$$

487 The residual practically stays about the same magnitude

$$\|(A - \hat{\lambda}I) \hat{X} - \hat{X} S\|_F \approx 1.25 \times 10^{-14}$$

488 but the backward error improves substantially to

$$\|(A - \hat{\lambda}I) \hat{X} - \hat{X} S\|_F \|\hat{X}^\dagger\|_2 \approx 1.25 \times 10^{-14}$$

489 as  $\|\hat{X}^\dagger\|_2 \approx 1$ . More importantly, the forward accuracy of the computed eigenvalue improves  
 490 by 3 additional accurate digits.  $\square$

491 When the given matrix represents perturbed data, the orthonormalization seems to be more  
 492 significant in improving the accuracy, as shown in the example below.

493 **Example 3** Using a random perturbation of magnitude about  $10^{-5}$ , let

$$\tilde{A} = A + 10^{-5} \begin{bmatrix} -0.092 & -0.653 & -0.201 & -0.416 & -0.787 \\ -0.135 & -0.218 & 0.054 & -0.136 & -0.255 \\ 0.651 & 0.663 & -0.166 & -0.969 & -0.603 \\ -0.833 & 0.607 & 0.314 & 0.969 & -0.020 \\ -0.733 & -0.879 & 0.256 & -0.665 & -0.321 \end{bmatrix} \quad (20)$$

494 be the data representation of the matrix  $A$  in (19). Table 2 lists the computed eigenvalues,  
 495 residual norms, backward errors and forward errors before and after orthonormalization.  
 496 The results show a substantial improvement on the both forward and backward errors even  
 497 though the residual magnitudes roughly stay the same.  $\square$

	before orthonormalization	after orthonormalization
computed eigenvalue	<b>2.004413315474177</b>	<b>2.000000343999377</b>
residual norm	$2.3 \times 10^{-6}$	$2.9 \times 10^{-6}$
backward error	$6.7 \times 10^{-2}$	$2.9 \times 10^{-6}$
forward error	$4.4 \times 10^{-3}$	$3.4 \times 10^{-7}$

Table 2: Comparison between computing results with or without orthonormalization of the  $X$  component of the least squares solution to  $\mathbf{g}(G, \lambda, X) = \mathbf{0}$  for the matrix  $\tilde{A}$  in (20) at the eigenvalue  $\lambda = 2$ . Correct digits of computed eigenvalues are highlighted in boldface.

## 11 What kind of eigenvalues are ill-conditioned, and in what sense?

The well documented claim that a defective eigenvalue is infinitely sensitive to perturbations requires an oft-missing clarification: It's unbounded sensitivity is with respect to *arbitrary* perturbations. The sensitivity of a defective eigenvalue is finitely bounded by the spectral projector norm divided by the multiplicity if the perturbation is constrained to maintain the multiplicity, or by the multiplicity support condition number if the multiplicity support remains unchanged.

Furthermore, the above sensitivity assertions and clarifications are applicable on the problem of finding eigenvalues in its strictly narrow sense. In the sense of computing a multiple eigenvalue via a cluster mean provided that the cluster can be grouped correctly, the sensitivity is still bounded by spectral projector norm divided by the multiplicity. The problem of finding a defective eigenvalue in the sense of computing a pseudo-eigenvalue elaborated in this paper also enjoys a finitely bounded sensitivity in terms of the multiplicity support condition number.

Of course, the problem can still be ill-conditioned even if the sensitivity is bounded. In the following example, the matrix  $A$  has an eigenvalue of multiplicity 7 and the spectral projector norm is large, so the eigenvalue is ill-conditioned in this sense. On the other hand, the same eigenvalue is well-conditioned in multiplicity support sensitivity. Interestingly, this is not a contradiction at all. The conflicting sensitivity measures imply that the cluster mean is not accurate for approximating the eigenvalue but the pseudo-eigenvalue is, and Algorithm EIGENITERATION converges to the defective eigenvalue with all the digits correct.

**Example 4** Let

$$A = \begin{bmatrix} 3.006 & 2 & 1.005 & -1.001 & -0.002 & -0.001 & -0.001 & -1 \\ 5 & 2 & 5 & -1 & -2 & -1 & -1 & 0 \\ -5.006 & -3 & -3.005 & 2.001 & 3.002 & 2.001 & 0.001 & 2 \\ -6 & -1 & -6 & 3 & 5 & 3 & 0 & 1 \\ -5 & -1 & -5 & 1 & 6 & 3 & 0 & 1 \\ 1 & 0 & 1 & 0 & -1 & 1 & 0 & 0 \\ -4 & -2 & -4 & 1 & 3 & 2 & 2 & 2 \\ 5 & 0 & 5 & -1 & -2 & -1 & -1 & 2 \end{bmatrix}$$

with two distinct eigenvalues in exact sense: A simple eigenvalue  $\lambda_1 = 2.001$ , and a defective eigenvalue  $\lambda_2 = 2$  with a nonzero Segre characteristic  $\{5, 2\}$  (i.e. multiplicity support is  $2 \times 2$ ). Let  $P_2$  be the spectral projector associated with  $\lambda_2 = 2$ . The defective

eigenvalue  $\lambda_2$  is *both* highly ill-conditioned in spectral projector norm and almost perfectly conditioned measured by its  $2 \times 2$  condition number with a sharp contrast:

$$\frac{1}{m} \|P_2\|_2 \approx 4.05 \times 10^{14} \quad \text{while} \quad \tau_{A,2 \times 2}(\lambda_2) \leq 19.95.$$

This may seem to be a contradiction except it is not. Both conditions accurately measure the sensitivities of same end (finding the defective eigenvalue) through different means (cluster mean versus pseudo-eigenvalue). The Francis QR algorithm implemented in Matlab produces computed eigenvalues

```

530          2.003667055821394,          2.001912473859015 + 0.002992156370408i,
531          1.996674198110247,          2.001912473859015 - 0.002992156370408i,
532          2.000000046670435,          1.998416899175164 + 0.002994143122392i,
533          1.999999953329568,          1.998416899175164 - 0.002994143122392i.
```

There is no apparent way to group 7 computed eigenvalues to use the cluster mean for the defective eigenvalue even if we know the multiplicity is 7. Out of all 8 possible groups of 7 eigenvalues, the best approximation to  $\lambda_2 = 2.0$  by the average is 2.000142850475652 with a substantial error  $1.4 \times 10^{-4}$  predicted by the spectral projector norm. In contrast, Algorithm EIGENITERATION accurately converges to  $\lambda_2 = 2.0$  with an error below the unit round off  $2.2 \times 10^{-16}$  using the correct multiplicity support  $2 \times 2$  that can easily be identified using the method described in §9, as accurately predicted by the multiplicity support condition number.

This seemingly contradicting sensitivities can be explained by the fact that there are infinitely many matrices nearby possessing a single eigenvalue of nonzero Segre characteristic  $\{6, 2\}$  within 2-norm distances of  $5.2 \times 10^{-5}$ . Namely, such a small perturbation increases the multiplicity from 7 to 8 but can not increase the multiplicity support  $2 \times 2$ . Using the publicly available Matlab functionality NUMERICALJORDANFORM on the matrix  $A$  with error tolerance  $10^{-5}$  in the software package NACLAB<sup>1</sup> for numerical algebraic computation, we obtain approximately nearest matrix  $B$  with a single eigenvalue associated with Jordan blocks sizes 6 and 2 with first 14 digits of its entries given as

```

550  3.0059955942886   1.9999978851470   1.0049959180573  -1.0010020728471  -0.0020046893569  -0.0010002300301  -0.0010132897111  -0.9999977586058
551  4.9999998736434   1.9999937661529   5.0000001193777  -1.0000000065301  -2.0000000129428  -0.9999999934070  -0.9999999926252  -0.0000169379637
552  -5.0060008381014  -3.0000021146845  -3.0050076499797  2.0009979267360  3.0019953102172  2.0009997699688  0.0009867094117  2.0000022421372
553  -5.9999927405774  -1.0000021677309  -6.0000074892946  2.9999962015789  4.9999997701478  2.999999999775  -0.0000002324249  0.9999978331627
554  -4.999995930006  -0.9999961877596  -5.0000095377366  0.9999880178349  5.9999870716625  3.0000002295144  -0.0000150335883  0.999994536709
555  0.9999971940837  -0.0000010574036  1.0000006545259  -0.0000047736356  -1.0000023807994  0.9999966612522  -0.0000036987224  -0.0000054161918
556  -4.0000092166543  -1.999997569827  -3.9999765043908  1.0000142841290  3.0000142782479  2.0000000005386  2.0000249853630  2.0000002431865
557  4.9999998983338  0.0000026655939  5.0000001062663  -0.9999999958672  -1.999999894366  -1.0000000065916  -1.0000000120790  2.0000133696815
```

The spectrum of  $B$  consists of a single eigenvalue  $\lambda = 2.00125$ . This lurking nearby matrix indicates that the multiplicity 7 of  $\lambda_2 = 2.0 \in \text{eig}(A)$  can be increased to 8 with a small perturbation  $\|A - B\|$ , which is exactly the kind of cases where spectral projectors have large norms as elaborated by Kahan [9] and grouping method fails. However, those nearby defective matrices have the same multiplicity support  $2 \times 2$ , implying a small perturbation does not increase either the geometric multiplicity or the Segre characteristic anchor. As a result, the multiplicity support condition number is benign, and computing the defective eigenvalue via pseudo-eigenvalue is stable.

Interestingly, even though the matrix  $B$  is only known via the above empirical data, the spectral projector associated with its eigenvalue 2.00125 is known to be identity since there

<sup>1</sup><http://homepages.neiu.edu/~naclab>

568 is only one distinct eigenvalue. Consequently, the mean of all approximate eigenvalues  
569 computed by Francis QR algorithm is 2.000124999999987 with 14 digits accuracy, same as  
570 the empirical data. Algorithm EIGENITERATION produces the  $2 \times 2$  pseudo-eigenvalue  
571 2.000125000000078 with the same number of correct digits due to a small  $2 \times 2$  condition  
572 number about 14.47.

573 The software NUMERICALJORDANFORM accurately produces the Jordan Canonical Forms  
574 of both matrices  $A$  and  $B$ . □

## 575 References

- 576 [1] J. V. Burke and M. L. Overton. Stable perturbations of nonsymmetric matrices. *Linear*  
577 *Algebra and Its Applications*, 171:249–273, 1992.
- 578 [2] F. Chaitin-Chatelin and V. Frayssé. *Lectures on Finite Precision Computations*. SIAM,  
579 Philadelphia, 1996.
- 580 [3] F. Chatelin. Ill conditioned eigenproblems. In J. Cullum and R. A. Willoughby, editors,  
581 *Large Scale Eigenvalue Problems*, North-Holland, Amsterdam, 1986. Elsevier Science  
582 Publishers B. V.
- 583 [4] J. W. Demmel. A numerical analyst’s jordan canonical form. Ph.D. Dissertation,  
584 Computer Science Department, University of California, 1983.
- 585 [5] J. W. Demmel. Computing stable eigendecompositions of matrices. *Lin. Alg. and*  
586 *Appl.*, 79:163–193, 1986.
- 587 [6] G. H. Golub and J. H. Wilkinson. Ill-conditioned eigensystems and the computation  
588 of the Jordan canonical form. *SIAM Review*, 18:578–619, 1976.
- 589 [7] B. Kågström and A. Ruhe. Algorithm 560: JNF, an algorithm for numerical compu-  
590 tation of the Jordan Normal Form of a complex matrix. *ACM Trans. Math. Software*,  
591 6:437–443, 1980.
- 592 [8] B. Kågström and A. Ruhe. An algorithm for numerical computation of the Jordan  
593 normal form of a complex matrix. *ACM Trans. Math. Software*, 6:398–419, 1980.
- 594 [9] W. Kahan. Conserving confluence curbs ill-condition. Technical Report 6, Computer  
595 Science, University of California, Berkeley, 1972.
- 596 [10] T. Kato. *Perturbation Theory for Linear Operators*. Springer, Berlin, Heidelberg, New  
597 York, 1966 & 1980.
- 598 [11] V. N. Kublanovskaya. On a method of solving the complete eigenvalue problem for a  
599 degenerate matrix. *USSR Computational Math. and Math. Phys.*, 6:1–14, 1968.
- 600 [12] T.-Y. Li and Z. Zeng. A rank-revealing method with updating, downdating and appli-  
601 cations. *SIAM J. Matrix Anal. Appl.*, 26:918–946, 2005.



- 602 [13] V. Lidskii. Perturbation theory of non-conjugate operators. *U.S.S.R. Comput. Math.*  
603 *and Math. Phys.*, 6:73–85, 1966.
- 604 [14] R. A. Lippert and A. Edelman. The computation and sensitivity of double eigenvalues.  
605 In *Advances in computational mathematics, Lecture Notes in Pure and Appl. Math.*  
606 *202*, pages 353–393, New York, 1999. Dekker.
- 607 [15] J. Moro, J. V. Burke, and M. L. Overton. On the Lidskii-Vishik-Lyusternik perturba-  
608 tion theory for eigenvalues of matrices with arbitrary Jordan structure. *SIAM J. Matrix*  
609 *Anal. Appl.*, 18:793–817, 1997.
- 610 [16] A. Ruhe. An algorithm for numerical determination of the structure of a general matrix.  
611 *BIT*, 10:196–216, 1970.
- 612 [17] A. Ruhe. Perturbation bounds for means of eigenvalues and invariant subspaces. *BIT*,  
613 10:343–354, 1970.
- 614 [18] S. Rump. Computational error bounds for multiple or nearly multiple eigenvalues.  
615 *Linear Algebra and its Applications*, 324:209–226, 2001.
- 616 [19] S. Rump. Eigenvalues pseudospectrum and structured perturbations. *Linear Algebra*  
617 *and its Applications*, 413:567–593, 2006.
- 618 [20] B. Sridhar and D. Jordan. An algorithm for calculation of the Jordan Canonical Form  
619 of a matrix. *Comput. & Elect. Engng.*, 1:239–254, 1973.
- 620 [21] L. N. Trefethen and M. Ebbree. *Spectra and Pseudospectra*. Princeton University Press,  
621 Princeton and Oxford, 2005.
- 622 [22] J. H. Wilkinson. Sensitivity of eigenvalues. *Utilitas Mathematica*, 25:5–76, 1984.
- 623 [23] J. H. Wilkinson. Sensitivity of eigenvalues, II. *Utilitas Mathematica*, 30:243–286, 1986.