

## Sec. 2.5 Sensitivity analysis of linear system

Goal: 1. Condition number

2. Sensitivity analysis

1. Condition number.

eg. Consider  $\underbrace{\begin{pmatrix} 0.780 & 0.563 \\ 0.913 & 0.659 \end{pmatrix}}_A \underbrace{\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}}_{\vec{x}} = \underbrace{\begin{pmatrix} 0.217 \\ 0.254 \end{pmatrix}}_{\vec{b}}.$

approximate solutions:  $\vec{x}^{(1)} = \begin{pmatrix} 0.341 \\ -0.087 \end{pmatrix}, \quad \vec{x}^{(2)} = \begin{pmatrix} 0.999 \\ 1.00 \end{pmatrix}.$

Which one is preferred?

$$A\vec{x}^{(1)} - \vec{b} = \begin{pmatrix} 0.000001 \\ 0 \end{pmatrix}, \quad A\vec{x}^{(2)} - \vec{b} = \begin{pmatrix} 0.000780 \\ 0.000913 \end{pmatrix}.$$

If we want small residual, then  $\vec{x}^{(1)}$  is preferred.

However, exact solution  $\vec{x} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ .  $\|\vec{x}^{(1)} - \vec{x}\| \ll \|\vec{x}^{(2)} - \vec{x}\|$ ,  $\vec{x}^{(2)}$  is more accurate.

Some understanding behind the problem:

$$\tilde{A} = \begin{pmatrix} 0.780 & 0.5630001095 \dots \\ 0.913 & 0.659 \end{pmatrix} \text{ is singular}$$

So  $A$  is nearly singular.

A  $0(10^{-6})$  perturbation of the data will render the problem  $A\vec{x} = \vec{b}$  insolvable.

Def: Condition number of  $A$  is  $K_p(A) = \|A\|_p \cdot \|A^{-1}\|_p.$

Recall: For  $\vec{x} = (x_1, \dots, x_n)^T$ ,  $\|\vec{x}\|_p = \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$ ,  $1 \leq p < \infty$ ,  $\|\vec{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$

$$\text{For } A \in \mathbb{C}^{n \times n}, \quad \|A\|_p = \sup_{\|\vec{x}\|=1} \|A\vec{x}\|_p, \quad 1 \leq p < \infty$$

$$\|A\|_1 = \max_{1 \leq j \leq n} \|A(:, j)\|_1, \quad \|A\|_\infty = \max_{1 \leq i \leq n} \|A(i, :)\|_1, \quad \|A\|_2 \leq \|A\|_F := \left(\sum_{i,j=1}^n |a_{ij}|^2\right)^{1/2}$$

eg.  $A = \begin{pmatrix} 0.780 & 0.563 \\ 0.913 & 0.659 \end{pmatrix} \Rightarrow A^{-1} = 10^6 \begin{pmatrix} 0.659 & -0.563 \\ -0.913 & 0.780 \end{pmatrix}$

$$\Rightarrow \kappa_1(A) = \|A\|_1 \cdot \|A^{-1}\|_1 \approx 10^6$$

Theorem 1: 1.  $\kappa(A) \geq 1$ . (If  $A$  is unitary,  $\kappa(A) = 1$ )

2.  $\kappa(A)$  is large if  $A$  is close to singular

3.  $\kappa(\alpha A) = \kappa(A)$ , (scaling invariant).

4.  $\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$ , where  $\sigma_{\max}(A)$  and  $\sigma_{\min}(A)$  are the maximal and minimal singular values of  $A$ , i.e.  $\sigma(A) = \sqrt{\lambda(A^*A)}$

Remark: 1. Small eigenvalues or determinants does NOT imply "near singular".

eg.  $\alpha I_n$ , all eigenvalues =  $\alpha$ ,  $\det(\alpha I_n) = \alpha^n$ .

If  $\alpha = 10^{-10}$ , small, but  $\kappa(\alpha I_n) = \kappa(I_n) = 1$ .

2. If  $A$  is normal, i.e.  $A^*A = A \cdot A^*$ , then  $\kappa(A) = \frac{|\lambda_{\max}(A)|}{|\lambda_{\min}(A)|}$

3. All eigenvalues are 1 does NOT imply small condition number.

eg,  $A = \begin{pmatrix} 1 & & \alpha \\ & \ddots & \\ & & 1 \end{pmatrix}$ ,  $A^{-1} = \begin{pmatrix} 1 & & -\alpha \\ & \ddots & \\ & & 1 \end{pmatrix}$ ,  $\kappa_{\infty}(A) = \|A\|_{\infty} \|A^{-1}\|_{\infty} = (1 + |\alpha|)^2$

If  $\alpha = 10^8$ ,  $\kappa_{\infty}(A) \approx 10^{16}$

## 2. Sensitivity analysis

linear system:  $A\vec{x} = \vec{b}$

perturbed linear system:  $\hat{A}\hat{x} = \hat{b}$

$\|\hat{x} - \vec{x}\|$ : absolute error,  $\frac{\|\hat{x} - \vec{x}\|}{\|\vec{x}\|}$ : relative error ✓

(The sensitivity of the linear system has nothing to do with the numerical algorithm used.)

Theorem 2: Suppose  $A \in \mathbb{R}^{n \times n}$  is nonsingular and  $A\vec{x} = \vec{b}$ . If eps.  $\kappa(A) < 1$ , then the stored linear system  $\hat{A}\hat{x} = \hat{b}$  is nonsingular, and

$$\frac{\|\hat{x} - \vec{x}\|}{\|\vec{x}\|} \leq \frac{2 \cdot \text{eps} \cdot \kappa(A)}{1 - \text{eps} \cdot \kappa(A)}$$

To prove Theorem 2, we need the following Lemmas.

Lemma 1: If  $\hat{A}$  is the stored version of any  $A \in \mathbb{R}^{m \times n}$ , then  $\hat{A} = A + E$ , where  $E \in \mathbb{R}^{m \times n}$  and  $\|E\| \leq \text{eps} \cdot \|A\|$ . eps: machine precision

proof: Let  $\hat{A} = (\hat{a}_{ij})$ . Then  $\hat{a}_{ij} = a_{ij} (1 + \varepsilon_{ij})$ , where  $|\varepsilon_{ij}| \leq \text{eps}$

$$\|E\|_1 = \|\hat{A} - A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |\hat{a}_{ij} - a_{ij}| = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij} \cdot \varepsilon_{ij}| \leq \text{eps} \cdot \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| = \text{eps} \cdot \|A\|_1$$

The proof for  $\|\cdot\|_\infty$  is similar.

Lemma 2: Suppose  $A\vec{x} = \vec{b}$  is perturbed to  $(A+E)\hat{x} = (\vec{b} + \vec{e})$ .

$$\text{Then } \frac{\|\hat{x} - \vec{x}\|}{\|\vec{x}\|} \leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|E\|} \left( \frac{\|\vec{e}\|}{\|\vec{b}\|} + \frac{\|E\|}{\|A\|} \right)$$

$$\begin{aligned} \text{proof: } A(\hat{x} - \vec{x}) &= A\hat{x} - A\vec{x} = (\vec{b} + \vec{e} - E\hat{x}) - \vec{b} = \vec{e} - E\hat{x} = (\vec{e} + E\vec{x}) - E(\hat{x} - \vec{x}) \\ \hat{x} - \vec{x} &= A^{-1} [(\vec{e} + E\vec{x}) - E(\hat{x} - \vec{x})] \end{aligned}$$

By Cauchy inequality and triangle inequality,

$$\begin{aligned} \|\hat{x} - \vec{x}\| &\leq \|A^{-1}\| \cdot \|\vec{e} + E\vec{x} - E(\hat{x} - \vec{x})\| \leq \|A^{-1}\| \cdot (\|\vec{e}\| + \|E\vec{x}\| + \|E(\hat{x} - \vec{x})\|) \\ &\leq \|A^{-1}\| (\|\vec{e}\| + \|E\| \cdot \|\vec{x}\| + \|E\| \cdot \|\hat{x} - \vec{x}\|) \end{aligned}$$

$$\begin{aligned} \Rightarrow (1 - \|A^{-1}\| \cdot \|E\|) (\|\hat{x} - \vec{x}\|) &\leq \|A^{-1}\| \cdot (\|\vec{e}\| + \|E\| \cdot \|\vec{x}\|) \\ \|\hat{x} - \vec{x}\| &\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|E\|} (\|\vec{e}\| + \|E\| \cdot \|\vec{x}\|) \end{aligned}$$

$$\frac{\|\hat{x} - \vec{x}\|}{\|\vec{x}\|} \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|E\|} \left( \frac{\|\vec{e}\|}{\|\vec{x}\|} + \|E\| \right) = \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|E\|} \cdot \left( \frac{\|\vec{e}\|}{\|A\| \cdot \|\vec{x}\|} + \frac{\|E\|}{\|A\|} \right)$$

$$\text{Note that } \|\vec{b}\| = \|A\vec{x}\| \leq \|A\| \cdot \|\vec{x}\| \Rightarrow \frac{1}{\|A\| \cdot \|\vec{x}\|} \leq \frac{1}{\|\vec{b}\|}.$$

$$\frac{\|\hat{x} - \vec{x}\|}{\|\vec{x}\|} \leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|E\|} \left( \frac{\|\vec{e}\|}{\|\vec{b}\|} + \frac{\|E\|}{\|A\|} \right)$$

proof of Theorem 2:

① If  $\hat{A}$  is singular, then there is a nonzero vector  $\vec{z}$  such that  $\hat{A}\vec{z} = \vec{0}$ .

$$\hat{A} = A + E \Rightarrow (A + E)\vec{z} = \vec{0} \Rightarrow A\vec{z} = -E\vec{z} \Rightarrow \vec{z} = -A^{-1}E\vec{z}$$

$$\|\vec{z}\| \leq \|A^{-1}\| \cdot \|E\| \cdot \|\vec{z}\| \leq \|A^{-1}\| \cdot \text{eps} \cdot \|A\| \cdot \|\vec{z}\| = \underbrace{\text{eps} \cdot \kappa(A)}_{< 1} \cdot \|\vec{z}\| < \|\vec{z}\|, \text{ contradiction!}$$

So  $\hat{A}$  is nonsingular.

② By Lemma 1,  $\|E\| \leq \text{eps} \cdot \|A\|$ ,  $\|\vec{e}\| \leq \text{eps} \cdot \|\vec{b}\|$ ,

Using Lemma 2, we get

$$\frac{\|\hat{x} - \vec{x}\|}{\|\vec{x}\|} \leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|E\|} \left( \frac{\|\vec{e}\|}{\|\vec{b}\|} + \frac{\|E\|}{\|A\|} \right) \leq \frac{\kappa(A)}{1 - \|A^{-1}\| \cdot \|E\|} (\text{eps} + \text{eps}) \leq \frac{2 \cdot \text{eps} \cdot \kappa(A)}{1 - \|A^{-1}\| \cdot \|E\|}$$

$$\|E\| \leq \text{eps} \cdot \|A\| \Rightarrow -\|A^{-1}\| \cdot \|E\| \geq -\|A^{-1}\| \cdot \text{eps} \cdot \|A\| = -\text{eps} \cdot \kappa(A)$$

$$\Rightarrow 1 - \|A^{-1}\| \cdot \|E\| \geq 1 - \text{eps} \cdot \kappa(A) \Rightarrow \frac{1}{1 - \|A^{-1}\| \cdot \|E\|} \leq \frac{1}{1 - \text{eps} \cdot \kappa(A)}$$

$$\text{So } \frac{\|\hat{x} - \vec{x}\|}{\|\vec{x}\|} \leq \frac{2 \cdot \text{eps} \cdot \kappa(A)}{1 - \text{eps} \cdot \kappa(A)}.$$