

## 4.4 (continued)

### Condition number

Question: when solving  $A\bar{x} = \bar{b}$ , if  $\bar{b}$  is perturbed, how does it affect the sol'n?

Suppose  $\bar{x}$  is the sol'n to  $A\bar{x} = \bar{b}$ , and  $\tilde{x}$  is the sol'n to

$$A\tilde{x} = \tilde{b}, \text{ then}$$

$$A(\bar{x} - \tilde{x}) = \bar{b} - \tilde{b}$$

$$\bar{x} - \tilde{x} = A^{-1}(\bar{b} - \tilde{b})$$

$$\|\bar{x} - \tilde{x}\| \leq \|A^{-1}\| \cdot \|\bar{b} - \tilde{b}\| \quad \text{for some vector/matrix norm.}$$

To measure in relative error,

$$\|\bar{b}\| \leq \|A\| \cdot \|\bar{x}\|$$

$$\frac{1}{\|\bar{x}\|} \leq \|A\| \cdot \frac{1}{\|\bar{b}\|}$$

$$\Rightarrow \frac{\|\bar{x} - \tilde{x}\|}{\|\bar{x}\|} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{:= k(A)} \cdot \frac{\|\bar{b} - \tilde{b}\|}{\|\bar{b}\|}$$

$:= k(A)$  "condition number of  $A$ "

$$b_0 = 1.23456 \times 10^0$$

$$x_0 = 2.31478 \times 10^0 \quad \text{if } k(A) \approx 10^3$$

- A matrix  $A$  is ill-conditioned if  $k(A)$  is "large". Otherwise it's well-conditioned

Ex  $A = \begin{pmatrix} 1 & 1+\varepsilon \\ 1-\varepsilon & 1 \end{pmatrix}$   $A^{-1} = \varepsilon^{-2} \begin{pmatrix} 1 & -1-\varepsilon \\ -1+\varepsilon & 1 \end{pmatrix}$

In  $\ell^1$ -norm,  $\|A\| = 2+\varepsilon$   $\|A^{-1}\| = \varepsilon^{-2}(2+\varepsilon)$

$\kappa(A) = \varepsilon^{-2}(2+\varepsilon)^2$

Thm Denote  $\vec{e} = \vec{x} - \vec{\tilde{x}}$ ,  $\vec{r} = \vec{b} - \vec{\tilde{b}}$

$A\vec{x} = \vec{b}$

$\vec{\tilde{x}} = A^{-1}\vec{\tilde{b}}$

$\frac{1}{\kappa(A)} \frac{\|\vec{r}\|}{\|\vec{\tilde{b}}\|} \leq \frac{\|\vec{e}\|}{\|\vec{\tilde{x}}\|} \leq \kappa(A) \frac{\|\vec{r}\|}{\|\vec{\tilde{b}}\|}$

4.5 Neumann series, iterative refinement

Let  $\|\cdot\|$  be fixed vector/matrix norm.

Thm If  $A$  is  $n \times n$  w/  $\|A\| < 1$ , then  $I-A$  is invertible and

$(I-A)^{-1} = \sum_{k=0}^{\infty} A^k$

$\frac{1}{1-x} = 1+x+x^2+\dots$   
 $|x| < 1$

$\hookrightarrow$  means  $\lim_{m \rightarrow \infty} \sum_{k=0}^m A^k$  in the sense of  $\|\cdot\|$

i.e.,  $\lim_{m \rightarrow \infty} \left\| (I-A)^{-1} - \sum_{k=0}^m A^k \right\| = 0$ .

Proof ① Prove that  $\lim_{m \rightarrow \infty} \sum_{k=0}^m A^k$  exists.

$(a+b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}$

$(A+B)^2 = A^2 + \underbrace{AB}_{\sim} + \underbrace{BA}_{\sim} + B^2$

For  $m_1 > m_2 > M$ ,

$\left\| \sum_{k=0}^{m_1} A^k - \sum_{k=0}^{m_2} A^k \right\| = \left\| \sum_{k=m_2+1}^{m_1} A^k \right\| \leq \sum_{k=m_2+1}^{m_1} \|A^k\|$

$\leq \sum_{k=m_2+1}^{m_1} \|A\|^k \leq \sum_{k=M}^{\infty} \|A\|^k = \|A\|^M \cdot \frac{1}{1-\|A\|}$

is arbitrarily small for large  $M$

$\Rightarrow$  partial sums are a Cauchy sequence.

(2) Prove that  $(I-A) \sum_{k=0}^{\infty} A^k = I$ .

$$\begin{aligned} (I-A) \sum_{k=0}^{\infty} A^k &= (I-A) \lim_{m \rightarrow \infty} \sum_{k=0}^m A^k = \lim_{m \rightarrow \infty} (I-A) \sum_{k=0}^m A^k \\ &= \lim_{m \rightarrow \infty} (I - A^{m+1}) = I. \end{aligned} \quad \|A^{m+1}\| \leq \|A\|^{m+1} \rightarrow 0$$

• It follows that  $\|(I-A)^{-1}\| \leq \frac{1}{1-\|A\|}$

if you solve  $(I-A)\vec{x} = \vec{b}$  w/  $\|A\|$  small, then it's well-conditioned.

### Iterative refinement

Suppose we solve  $A\vec{x} = \vec{b}$  by a numerical method and get an approximate sol'n  $\vec{x}^{(0)}$

Then, denote  $\vec{e}^{(0)} = \vec{x} - \vec{x}^{(0)}$

$$\Rightarrow A\vec{e}^{(0)} = A\vec{x} - A\vec{x}^{(0)} = \vec{b} - A\vec{x}^{(0)} = \vec{r}^{(0)}$$

Then we solve for  $\vec{e}^{(0)}$  from  $A\vec{e}^{(0)} = \vec{r}^{(0)}$  numerically, get  $\tilde{\vec{e}}^{(0)}$

$$\text{define } \vec{x}^{(1)} = \vec{x}^{(0)} + \tilde{\vec{e}}^{(0)}$$

This can be applied iteratively.

• To make it work, one needs to use higher precision for calculating  $\vec{r}^{(0)} = \vec{b} - A\vec{x}^{(0)}$  because it's a difference of two close numbers.

Iterative refinement for Gaussian elimination is cheap ( $O(n^2)$  each step if LU-decomp. for  $A$  is saved).

• Error analysis: Suppose the numerical solver for  $A\vec{x} = \vec{b}$  produces  $\vec{x} = B\vec{b}$  where  $B$  is an approximation of  $A^{-1}$ .

Then  $\vec{x}^{(0)} = B \vec{b}$

$$\begin{aligned}\vec{x}^{(1)} &= \vec{x}^{(0)} + B(\vec{b} - A\vec{x}^{(0)}) = B\vec{b} + B(\vec{b} - AB\vec{b}) \\ &= B(\vec{b} + \vec{b} - AB\vec{b}) = B(\mathbf{I} + (\mathbf{I} - AB))\vec{b}\end{aligned}$$

$$\begin{aligned}\vec{x}^{(2)} &= \vec{x}^{(1)} + B(\vec{b} - A\vec{x}^{(1)}) \\ &= B(\mathbf{I} + (\mathbf{I} - AB))\vec{b} + B(\vec{b} - AB(\mathbf{I} + (\mathbf{I} - AB))\vec{b}) \\ &= B\left((\mathbf{I} + (\mathbf{I} - AB))\vec{b} + \vec{b} - AB(\mathbf{I} + (\mathbf{I} - AB))\vec{b}\right) \\ &= B\left(\mathbf{I} + \underbrace{(\mathbf{I} - AB)} + \underbrace{\mathbf{I} - AB} - \underbrace{AB(\mathbf{I} - AB)}\right)\vec{b} \\ &= B\left(\mathbf{I} + (\mathbf{I} - AB) + (\mathbf{I} - AB)^2\right)\vec{b}\end{aligned}$$

...

$$\vec{x}^{(m)} = B \sum_{k=0}^m (\mathbf{I} - AB)^k \vec{b} \quad (\text{can be proved by induction})$$

If  $\|\mathbf{I} - AB\| < 1$  then  $\sum_{k=0}^m (\mathbf{I} - AB)^k \rightarrow (\mathbf{I} - (\mathbf{I} - AB))^{-1} = (AB)^{-1} = B^{-1}A^{-1}$

$$\Rightarrow \vec{x}^{(m)} \rightarrow A^{-1}\vec{b}$$