

## Lecture 18

Last time we introduced norms on  $\mathbb{R}^n$ . We also studied matrix norms

Recall a norm on  $\mathbb{R}^n$  is a function  $\|\cdot\|$  which associates to each vector  $\bar{a}$  in  $\mathbb{R}^n$ , a real number  $\|\bar{a}\| \geq 0$

$$1) \|\bar{a}\| \geq 0 \quad \forall \bar{a} \in \mathbb{R}^n$$

$$\|\bar{a}\| = 0 \quad \text{iff} \quad \bar{a} = \mathbf{0}$$

$$2) \text{ for } \alpha \in \mathbb{R},$$

$$\|\alpha \bar{a}\| = |\alpha| \|\bar{a}\|$$

$$3) \text{ (triangle inequality)}$$

$$\|\bar{a} + \bar{b}\| \leq \|\bar{a}\| + \|\bar{b}\|$$

---

Examples of norms

$$\bar{a} = (a_1, \dots, a_n)$$

$$\|\bar{a}\|_2 = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2}$$

$$\|a\|_1 = |a_1| + |a_2| + \dots + |a_n|$$

$$\|a\|_\infty = \max_{1 \leq i \leq n} |a_i|$$

## Matrix Norms

a matrix norm is a real valued function defined over the set of  $n \times n$  matrices such that for any  $n \times n$  matrices  $A, B$  and real number  $\alpha$  we have

$$(1) \quad \|A\| \geq 0 \quad \forall A$$

and  $\|A\| = 0 \iff A = 0$

$$(2) \quad \|\alpha A\| = |\alpha| \|A\|$$

$$(3) \quad \|A+B\| \leq \|A\| + \|B\|$$

$$(4) \quad \|AB\| \leq \|A\| \|B\|$$

Theorem Let  $\|\cdot\|$  be a norm in  $\mathbb{R}^n$

Then for a  $n \times n$  matrix  $A$

$$\|A\| = \max_{\|x\|=1} \|Ax\|$$

defines a <sup>matrix</sup> norm on the set of  $n \times n$  matrices

Remarks

$$\begin{aligned} \|A\| &= \max_{\|x\|=1} \|Ax\| \\ &= \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} \end{aligned}$$

Thus we get  $\|Ax\| \leq \|A\| \|x\|$   
for all  $x \in \mathbb{R}^n$ .

---

Examples

$$1) \|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$$

2-norm on matrices

"difficult to compute".

$$(2) \|A\|_{\infty} = \max_{\|x\|_{\infty}=1} \|Ax\|_{\infty}$$

Thus

$$\|A\|_{\infty} = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$$

$$(3) \|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1$$

is

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|$$

### Remark

If the matrix  $A$  is invertible then

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

$$\text{So } \|\bar{x}\| = \|A^{-1}(Ax)\| \leq \|A^{-1}\| \|Ax\|$$

Thus

$$\frac{\|\bar{x}\|}{\|A^{-1}\|} \leq \|Ax\| \leq \|A\| \|x\|$$

## Reason for studying norm

We want to study errors in computing solutions to linear equation

Solving  $Ax = b$

we get an approximate sol<sup>n</sup>  $\hat{x}$

error  $\bar{e} = \bar{x} - \hat{x}$

This error is unknown to us

"residual error"

$$\bar{r} = A\bar{e} = Ax - A\hat{x} = b - A\hat{x}$$

$$\bar{r} = b - A\hat{x}$$

This we can compute

---

last time I gave example which showed that the size of the residual  $\bar{r}$  is not always a reliable indicator of the size of error  $\bar{e}$ .

It depends on the size of  $A$  and  $A^{-1}$

error  $\bar{e} = x - \hat{x}$

residual  $\bar{r} = A\bar{e} = b - A\hat{x}$

$$\therefore \bar{e} = A^{-1}\bar{r}$$

Recall we have proved that for an invertible matrix  $B$

$$\frac{\|B^{-1}u\|}{\|B^{-1}\|} \leq \|Bu\| \leq \|B\| \|u\|$$

apply this to  $B = A^{-1}$  and  $u = \bar{r}$

$$(*) \quad \frac{\|\bar{r}\|}{\|A\|} \leq \|e\| = \|A^{-1}\bar{r}\| \leq \|A^{-1}\| \|\bar{r}\|$$

(\*\*) gives an upper and a lower bound of the relative error  $\frac{\|e\|}{\|x\|}$  in terms of

relative residual  $\frac{\|\bar{r}\|}{\|b\|}$

$$\frac{\|b\|}{\|A\| \|x\|} \frac{\|\tilde{b}\|}{\|b\|} \leq \frac{\|e\|}{\|x\|} \leq \frac{\|A^{-1}\| \|b\|}{\|x\|} \cdot \frac{\|\lambda\|}{\|b\|}$$

(\*) for  $\hat{x} = 0$  gives

$$\frac{\|b\|}{\|A\|} \leq \|x\| \leq \|A^{-1}\| \|b\|$$

$$\frac{1}{\|A^{-1}\| \|A\|} \frac{\|\lambda\|}{\|b\|} \leq \frac{\|e\|}{\|x\|} \leq \|A^{-1}\| \|A\| \frac{\|\lambda\|}{\|b\|}$$

$\text{Cond}(A) = \ A\  \ A^{-1}\ $	Condition number of A
-------------------------------------	-----------------------------

$$\frac{1}{\text{cond}(A)} \frac{\|\lambda\|}{\|b\|} \leq \frac{\|e\|}{\|x\|} \leq \text{cond}(A) \frac{\|\lambda\|}{\|b\|}$$

example

$$A = \begin{bmatrix} 1.01 & 0.99 \\ 0.99 & 1.01 \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} \frac{101}{4} & -\frac{99}{4} \\ -\frac{99}{4} & \frac{101}{4} \end{bmatrix}$$

$$\|A\|_{\infty} = 2$$

$$\|A^{-1}\|_{\infty} = 50$$

$$\text{So } \text{cond}(A) = 100$$

-----x-----x-----x-----

Today we do Backward error  
Analysis and Iterative improvement.

Theorem Let  $A$  be a  $n \times n$  invertible matrix.

$$\frac{1}{\text{cond}(A)} = \min \left\{ \frac{\|A-B\|}{\|A\|} \mid B \text{ is not invertible} \right\}$$

Pf We only prove  $\leq$

So we prove



$$\frac{1}{\text{cond}(A)} \leq \frac{\|A-B\|}{\|A\|} \quad \text{for any non-invertible matrix } B$$

$$\text{cond}(A) = \|A\| \|A^{-1}\|$$

So we have to prove

$$\frac{1}{\|A^{-1}\|} \leq \|A-B\| \quad B \text{ non-invertible}$$

$B$  non-invertible. So  $\exists x \neq 0$  s.t.  
 $Bx = 0$

$$1) \| (A-B)x \| \leq \|A-B\| \|x\|$$

$$\|(A-B)x\| = \|Ax - Bx\| = \|Ax\| \geq \frac{\|x\|}{\|A^{-1}\|}$$

$$\hookrightarrow \|A-B\| \|x\| \leq \frac{\|x\|}{\|A^{-1}\|}$$

$$\text{as } \|x\| \neq 0 \text{ we get } \frac{1}{\|A^{-1}\|} \leq \|A-B\|$$

Corollary :- If  $A$  is invertible and  $B$  is a matrix such that

$$\|A - B\| < \frac{1}{\|A^{-1}\|}$$

then  $B$  is invertible

---

We use the inequality

$$\frac{1}{\|A^{-1}\|} \leq \|A - B\|$$

where  $B$  is singular

to estimate  $\|A^{-1}\|$  without directly computing  $A^{-1}$

Example ①  $A = \begin{pmatrix} 1.01 & 0.99 \\ 0.99 & 1.01 \end{pmatrix}$

$$B = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

$B$  is singular

$$A - B = \begin{pmatrix} 0.01 & -0.01 \\ -0.01 & 0.01 \end{pmatrix}$$

$$\|A - B\|_{\infty} = 0.02$$

$$\text{so } \|A^{-1}\|_{\infty} \geq \frac{1}{0.02} = 50$$

$$\text{So } \kappa(A) = \|A\| \|A^{-1}\| \geq 100$$

In this case it is exact.

Example 2

$A$  is an invertible  
upper triangular matrix

$$\text{Then } \text{cond}(A) \geq \frac{\|A\|_{\infty}}{\min_i |a_{ii}|}$$

pf  $A$  invertible & upper triangular.

So all diagonal entries of  $A$  are non-zero.

without loss of generality say

$$|a_{11}| = \min_i |a_{ii}|$$

$$B = A - \begin{pmatrix} a_{11} & 0 & \dots & 0 \\ 0 & 0 & & 0 \\ 0 & 0 & \dots & 0 \end{pmatrix}$$

note  $B$  is singular ( $\because B$  is upper triangular &  $b_{11} = 0$ )

$$\|A - B\|_{\infty} = |a_{11}|$$

$$\text{so } \|A^{-1}\| \geq \frac{1}{|a_{11}|}$$

## Perturbations of linear systems of equations

If the linear system  $Ax = b$  derives from a practical problem, we must expect the coefficients of this system to be subject to error either because they result from other calculations, or from physical measurement.

Hence assuming for example the RHS is accurate we are in effect solving

the system  $\hat{A}\hat{x} = b$

instead of  $Ax = b$

where  $A = \hat{A} + E$ , the matrix

$E$  contains the errors in the coefficients,

Even if all calculations are carried out exactly we only have a solution of  $\hat{A}\hat{x} = b$  rather than  $Ax = b$

$$\begin{aligned}\text{Now } x &= A^{-1}b \\ &= A^{-1}\hat{A}\hat{x} \\ &= A^{-1}(A + \hat{A} - A)\hat{x} \\ &= \hat{x} + A^{-1}(\hat{A} - A)\hat{x}\end{aligned}$$

$$\text{Now } \hat{A} - A = -E$$

$$\text{So } x = \hat{x} - A^{-1}E\hat{x}$$

$$\begin{aligned}\text{So } \|x - \hat{x}\| &\leq \|A^{-1}\| \|E\| \|\hat{x}\| \\ &= \|A^{-1}\| \|A\| \frac{\|E\|}{\|A\|} \|\hat{x}\|\end{aligned}$$

$$\text{Thus } \boxed{\frac{\|x - \hat{x}\|}{\|\hat{x}\|} \leq \text{cond}(A) \frac{\|E\|}{\|A\|}}$$

Thus if the coefficients of the linear system  $Ax=b$  are known to be accurate only to about  $10^{-5}$  (relative to size of  $A$ ) and  $\text{cond}(A) \approx 10^t$

Then there is no point in calculating the solution to a relative accuracy of  $10^{t-5}$ .

---

Quite loosely we say the linear system  $Ax=b$  is "ill-conditioned" if  $\text{cond}(A)$  is "large".

## Iterative improvement of solution

Let  $e = x - \hat{x}^{(1)}$  be the (unknown) error in the approximate solution  $\hat{x}^{(1)}$  for  $Ax = b$

$$Ae = r = b - A\hat{x}^{(1)} \quad (*)$$

Let  $\hat{e}^{(1)}$  be approximate solution of  $(*)$ . Then  $\hat{e}^{(1)}$  need not equal  $e$  but at the very least  $\hat{e}^{(1)}$  gives an indication of the size of  $e$

If  $\frac{\|\hat{e}^{(1)}\|}{\|\hat{x}^{(1)}\|} \approx 10^{-3}$  then we

conclude that the first 3 decimals



of  $\hat{x}^{(1)}$  agree with that of the exact answer  $x$ .

We would then also expect

$\hat{e}^{(1)}$  to be an accurate approximation to  $e$ .

So  $\hat{x}^{(2)} = \hat{x}^{(1)} + \hat{e}^{(1)}$  to be a better approximation to  $x$  than  $\hat{x}^{(1)}$

We can now, if necessary compute the new residual

$$r = b - A \hat{x}^{(2)}$$

and solve  $Ae = r$  to obtain a

new correction  $\hat{e}^{(2)}$  and a new approximation  $\hat{x}^{(3)} = \hat{x}^{(2)} + \hat{e}^{(2)}$

The number of places in agreement in the successive approximations  $\hat{x}^{(1)}, \hat{x}^{(2)}, \dots$  as well as an examination of the successive residuals should give an indication of the accuracy of these approximate solutions.

One normally carries out this iteration until  $\frac{\|\hat{e}^{(k)}\|}{\|\hat{x}^{(k)}\|} \approx 10^{-t}$  if  $t$  decimal places are carried during the calculation.

# Iteration steps increase with  $\text{cond}(A)$

If  $\text{cond}(A)$  is "very large" the correction  $\hat{e}^{(1)}, \hat{e}^{(2)}, \dots$ , may never decrease in size

## Remarks

For the success of iterative improvement it is absolutely mandatory that residuals be computed as accurately as possible.

If, as is usual, floating pt arithmetic is used, the residual should always be calculated in double-precision arithmetic.

