

## 15. Problem condition

- condition of a mathematical problem
- matrix norm
- condition number

# Sources of error in numerical computation

**Example:** evaluate a function  $f : \mathbf{R} \rightarrow \mathbf{R}$  at a given  $x$

sources of error in the result:

- $x$  is not exactly known
  - measurement errors
  - errors in previous computations
  - how sensitive is  $f(x)$  to errors in  $x$ ?
- the algorithm for computing  $f(x)$  is not exact
  - discretization (e.g., algorithm uses a table to look up function values)
  - truncation (e.g., function is evaluated by truncating a Taylor series)
  - rounding error during the computation
  - how large is the error introduced by the algorithm?

# Condition (conditioning) of a problem

describes sensitivity of the solution to changes in the problem data

- **well-conditioned problem:**

small changes in the data produce small changes in the solution

- **ill-conditioned (badly conditioned) problem:**

small changes in the data can produce large changes in the solution

a rigorous definition depends on what 'large error' means

- absolute or relative error, which norm is used, ...
- the informal definition is sufficient for our purposes

## Example: function evaluation

here the problem is: given  $x$ , evaluate  $y = f(x)$

- if  $x$  is changed to  $x + \Delta x$ , solution changes to

$$y + \Delta y = f(x + \Delta x)$$

- condition with respect to absolute error in  $x$  and  $y$

$$|\Delta y| \approx |f'(x)| |\Delta x|$$

problem is ill-conditioned with respect to absolute error if  $|f'(x)|$  is very large

- condition with respect to relative errors in  $x$  and  $y$

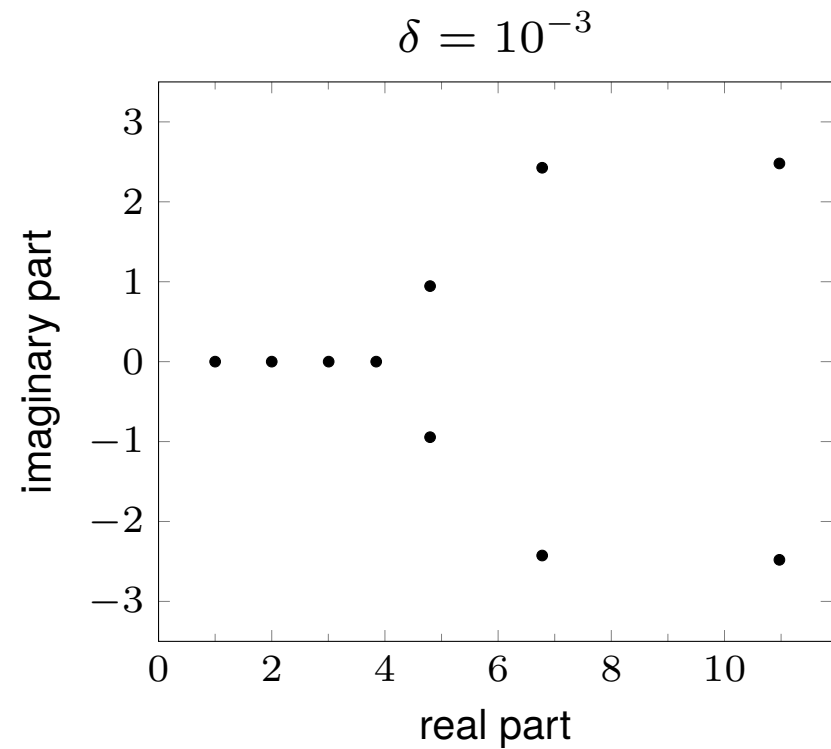
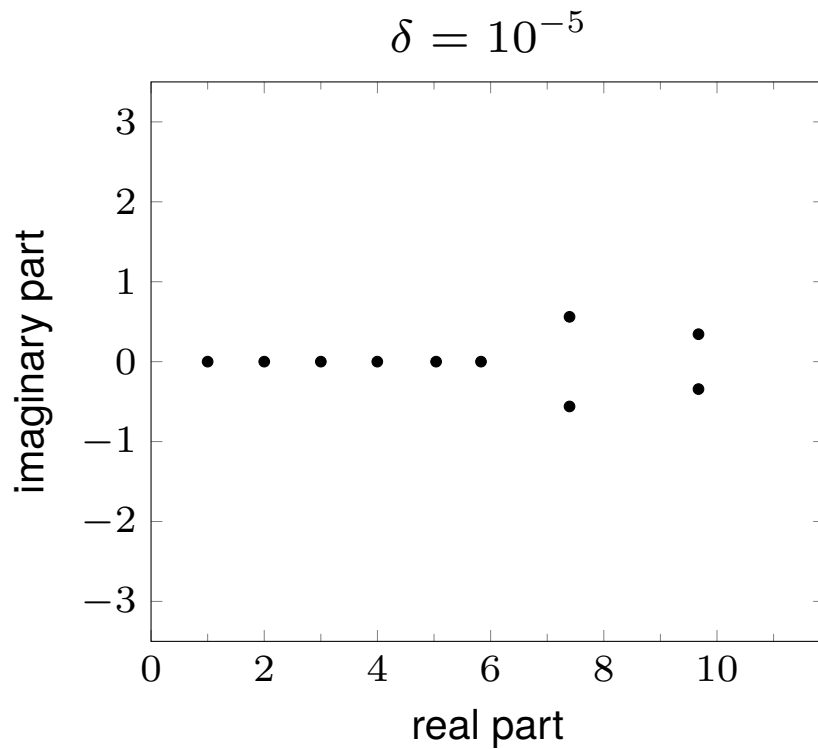
$$\frac{|\Delta y|}{|y|} \approx \frac{|f'(x)| |x|}{|f(x)|} \frac{|\Delta x|}{|x|}$$

ill-conditioned with respect to relative error if  $|f'(x)| |x| / |f(x)|$  is very large

# Roots of a polynomial

$$p(x) = (x - 1)(x - 2) \cdots (x - 10) + \delta \cdot x^{10}$$

roots of  $p$  computed by MATLAB for two values of  $\delta$



roots are very sensitive to errors in the coefficients

# Condition of a set of linear equations

- assume  $A$  is nonsingular and  $Ax = b$
- if we change  $b$  to  $b + \Delta b$ , the new solution is  $x + \Delta x$  with

$$A(x + \Delta x) = b + \Delta b$$

- the change in  $x$  is

$$\Delta x = A^{-1} \Delta b$$

## Condition

- the equations are *well-conditioned* if small  $\Delta b$  results in small  $\Delta x$
- the equations are *ill-conditioned* if small  $\Delta b$  can result in large  $\Delta x$

## Example of ill-conditioned equations

$$A = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 + 10^{-10} & 1 - 10^{-10} \end{bmatrix}, \quad A^{-1} = \begin{bmatrix} 1 - 10^{10} & 10^{10} \\ 1 + 10^{10} & -10^{10} \end{bmatrix}$$

- solution for  $b = (1, 1)$  is  $x = (1, 1)$
- change in  $x$  if we change  $b$  to  $b + \Delta b$ :

$$\Delta x = A^{-1} \Delta b = \begin{bmatrix} \Delta b_1 - 10^{10}(\Delta b_1 - \Delta b_2) \\ \Delta b_1 + 10^{10}(\Delta b_1 - \Delta b_2) \end{bmatrix}$$

small  $\Delta b$  can lead to extremely large  $\Delta x$

# Outline

- condition of a mathematical problem
- **matrix norm**
- condition number



# Matrix norms

the **Frobenius norm** of an  $m \times n$  matrix  $A$  is defined as

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2}$$

- denoted  $\|A\|$  in the textbook
- in MATLAB: `norm(A, 'fro')`

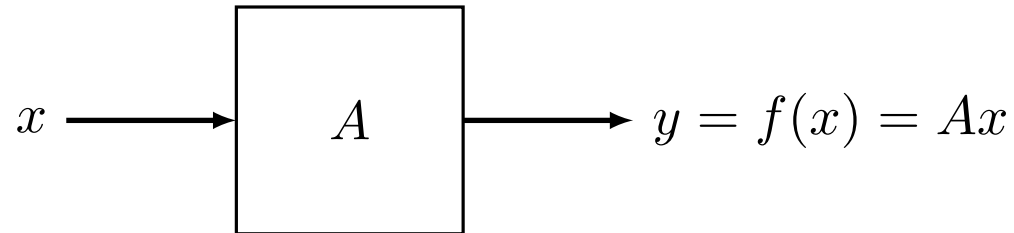
the **2-norm** or **spectral norm** is defined as

$$\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

- the norms  $\|Ax\|$  and  $\|x\|$  are Euclidean norms of vectors
- no simple explicit expression, except for special  $A$
- readily computed numerically (in MATLAB: `norm(A)`)

# Interpretation of 2-norm

the  $m \times n$  matrix  $A$  defines a linear function  $f(x) = Ax$



- $\|Ax\|/\|x\|$  gives the *amplification factor* or *gain* for input  $x$
- the gain only depends on the direction of  $x$
- the 2-norm of  $A$  is the maximum gain over all directions:

$$\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \|Ax\|$$

# Computing the 2-norm of a matrix

**Simple matrices:** sometimes it is easy to maximize  $\|Ax\|/\|x\|$

- zero matrix:  $\|0\|_2 = 0$
- identity matrix:  $\|I\|_2 = 1$
- diagonal matrix:

$$A = \begin{bmatrix} A_{11} & 0 & \cdots & 0 \\ 0 & A_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & A_{nn} \end{bmatrix}, \quad \|A\|_2 = \max_{i=1,\dots,n} |A_{ii}|$$

- matrix with orthonormal columns:  $\|A\|_2 = 1$

**General matrices:**  $\|A\|_2$  must be computed by numerical algorithms

# Properties of the matrix norm

## Properties satisfied by all matrix norms

- *nonnegative*:  $\|A\|_2 \geq 0$  for all  $A$
- *positive definiteness*:  $\|A\|_2 = 0$  only if  $A = 0$
- *homogeneity*:  $\|\beta A\|_2 = |\beta| \|A\|_2$
- *triangle inequality*:  $\|A + B\|_2 \leq \|A\|_2 + \|B\|_2$

## Additional properties satisfied by the 2-norm

- $\|Ax\| \leq \|A\|_2 \|x\|$  if the product  $Ax$  exists
- $\|AB\|_2 \leq \|A\|_2 \|B\|_2$  if the product  $AB$  exists
- if  $A$  is nonsingular:  $\|A\|_2 \|A^{-1}\|_2 \geq 1$
- if  $A$  is nonsingular:  $1/\|A^{-1}\|_2 = \min_{x \neq 0} (\|Ax\|_2 / \|x\|)$
- $\|A^T\|_2 = \|A\|_2$

# Outline

- condition of a mathematical problem
- matrix norm
- **condition number**

## Bound on absolute error

suppose  $A$  is nonsingular and define

$$x = A^{-1}b, \quad \Delta x = A^{-1}\Delta b$$

**Upper bound** on  $\|\Delta x\|$ :

$$\|\Delta x\| \leq \|A^{-1}\|_2 \|\Delta b\|$$

- follows from property 4 on page 15-11
- small  $\|A^{-1}\|_2$  means that  $\|\Delta x\|$  is small when  $\|\Delta b\|$  is small
- large  $\|A^{-1}\|_2$  means that  $\|\Delta x\|$  can be large, even when  $\|\Delta b\|$  is small
- for every  $A$ , there exists  $\Delta b$  such that  $\|\Delta x\| = \|A^{-1}\|_2 \|\Delta b\|$  (no proof)

## Bound on relative error

suppose in addition that  $b \neq 0$ ; hence  $x \neq 0$

**Upper bound** on  $\|\Delta x\|/\|x\|$ :

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A\|_2 \|A^{-1}\|_2 \frac{\|\Delta b\|}{\|b\|} \quad (1)$$

- follows from  $\|\Delta x\| \leq \|A^{-1}\|_2 \|\Delta b\|$  and  $\|b\| \leq \|A\|_2 \|x\|$
- $\|A\|_2 \|A^{-1}\|_2$  small means  $\|\Delta x\|/\|x\|$  is small when  $\|\Delta b\|/\|b\|$  is small
- $\|A\|_2 \|A^{-1}\|_2$  large means  $\|\Delta x\|/\|x\|$  can be much larger than  $\|\Delta b\|/\|b\|$
- for every  $A$ , there exist  $b, \Delta b$  such that equality holds in (1) (no proof)

# Condition number

**Definition:** the condition number of a nonsingular matrix  $A$  is

$$\kappa(A) = \|A\|_2 \|A^{-1}\|_2$$

## Properties

- $\kappa(A) \geq 1$  for all  $A$  (last property on page page 15-11)
- $A$  is a *well-conditioned* matrix if  $\kappa(A)$  is small (close to 1):  
the relative error in  $x$  is not much larger than the relative error in  $b$
- $A$  is *badly conditioned* or *ill-conditioned* if  $\kappa(A)$  is large:  
the relative error in  $x$  can be much larger than the relative error in  $b$



# Example

- $A$  is blurring matrix, nonsingular with condition number  $\approx 10^9$
- we apply  $A$  to image  $x$



blurred image

$$y_1 = Ax$$



blurred and noisy image

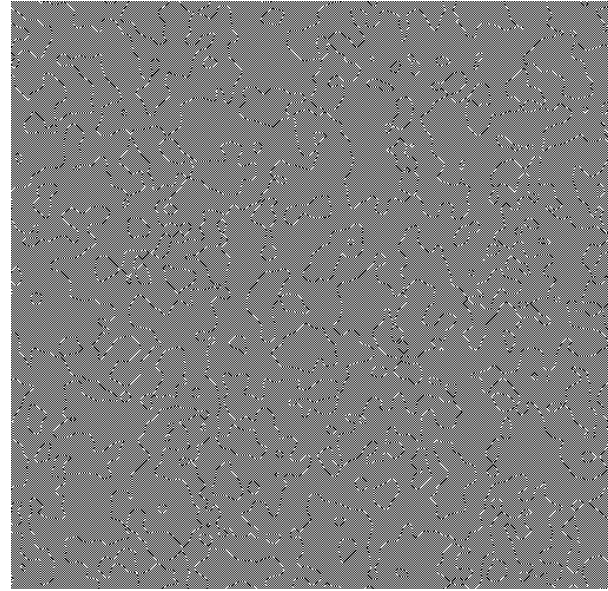
$$y_2 = Ax + \text{small noise}$$

# Example

we solve  $Ax = y$  for the two blurred images



$$A^{-1}y_1$$



$$A^{-1}y_2$$

- illustrates ill conditioning of  $A$
- explains need for regularization in deblurring algorithms

# Exercises

## Exercise 1

$$A = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1+a & 1-a \end{bmatrix}, \quad A^{-1} = \frac{1}{a} \begin{bmatrix} a-1 & 1 \\ a+1 & -1 \end{bmatrix}$$

$a$  is small and nonzero ( $a = 10^{-10}$  on page 15-7); show that  $\kappa(A) \geq 1/|a|$

## Exercise 2

suppose  $A = UBV$  with  $U, V$  orthogonal, and  $B$  nonsingular; show that

$$\kappa(A) = \kappa(B)$$

## Exercise 3

suppose  $A = uv^T$  where  $u$  and  $v$  are vectors; show that  $\|A\|_2 = \|u\| \|v\|$

# Exercises

## Exercise 4 (ex. 15.3)

- let  $u$  be a vector; show that

$$\|u\| = \max_{v \neq 0} \frac{v^T u}{\|v\|}$$

- let  $A$  be a matrix; show that

$$\|A\|_2 = \max_{y \neq 0, x \neq 0} \frac{y^T A x}{\|x\| \|y\|}$$

therefore  $\|A\|_2 = \|A^T\|_2$