

# Méthodes numériques matricielles avancées: analyse et expérimentation

**Marc Bonnet, Luiz Faria**

Propagation des Ondes: Etudes Mathématiques et Simulation (POEMS)

UMR 7231 CNRS-INRIA-ENSTA

Unité de Mathématiques Appliquées

ENSTA Paris

*mbonnet@ensta.fr*

*<https://perso.ensta-paris.fr/~mbonnet/enseignement.html>*

ANN 203, ENSTA PARIS, 2023-24

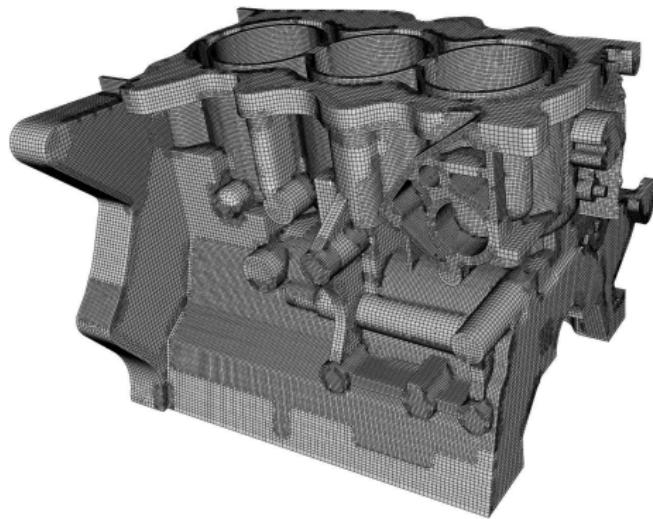
**Partie 1:** Généralités

**Partie 2:** Méthodes directes

**Partie 3:** Méthodes itératives

**Partie 4:** Problèmes aux valeurs et vecteurs propres

**Partie 5:** Systèmes linéaires mal conditionnés



- Linear statics:  $KU = F$
- Free vibrations:  $(K - \omega^2 M)U = 0$
- Forced vibrations:  $(K - \omega^2 M)U = F$

## Solving non-linear equations

E.g. mechanical structure involving nonlinear material properties (or large strains, or contact, or...)

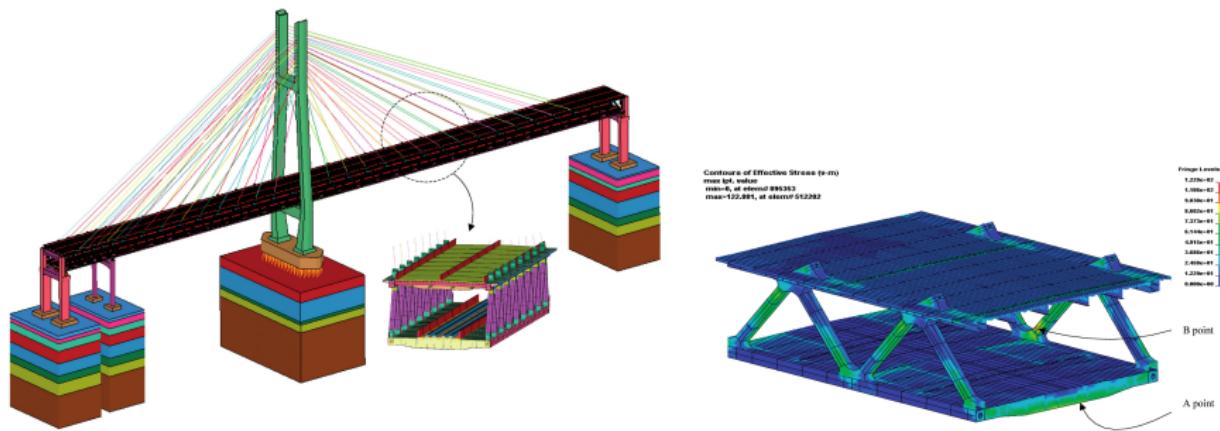
$$\mathcal{F}(U) = 0$$

Newton-Raphson:

$$\mathcal{F}(U + \Delta U) = 0 \implies \mathcal{F}(U) + \mathcal{F}'(U)\Delta U = 0$$

Iterations:

$$K(U_k)\Delta U_k + \mathcal{F}(U_k) = 0, \quad U_{k+1} := U_k + \Delta U_k \quad k = 1, 2, 3, \dots$$



## PDE discretization, boundary elements

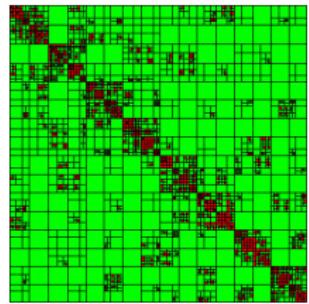
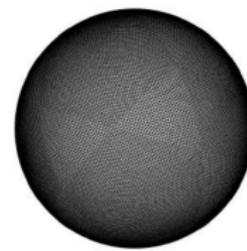
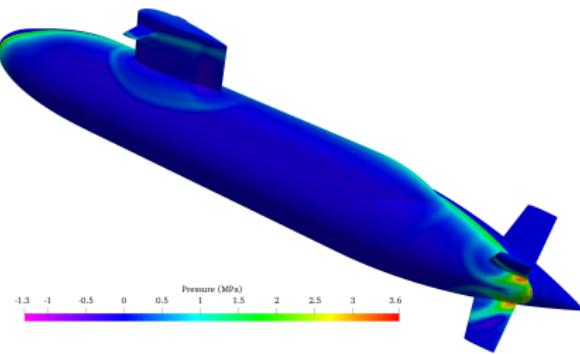
Boundary integral equation

$$\int_S G(x, y) u(y) dS(y) = f(x) \quad \text{for all } x \in S \quad S: \text{surface}$$

Often used for wave propagation in large/unbounded media.

After boundary element discretization:

$$GU = F$$



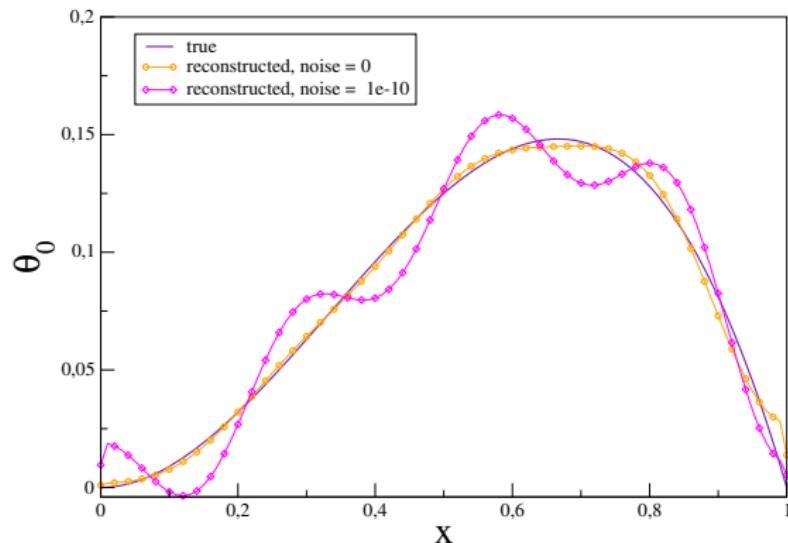
## Inverse problems

Example: backward heat conduction problem (BHCP): quantify initial temperature field  $\Theta(x, 0)$  using later measurement

$$\underbrace{\Theta(\cdot, T)}_{\text{measurement at time } T} = \underbrace{\mathcal{A}([0, T])}_{\text{heat diffusion eq.}} \underbrace{\Theta(\cdot, 0)}_{\text{unknown}}$$



numerical solution of 1D BHCP



## Image restoration

$$\hat{f}(x) = \int_Y k(x-y) f(y) dy, \quad x \in Y$$

blurred      blurring image

e.g.  $k(z) = Ce^{-|z|^2/2\sigma^2}$  (atmospheric blur)

Pixel discretization:

$$KF = \hat{F}$$

( $K$  dense, ill-conditioned, numerically rank-deficient)



reference



blurred



restored

### Basic problem:

- Population ( $j = 1, \dots, m$  individuals)  
 $x_{ij}$  (explanatory variables),  $y_j$
- Seek best affine model  $y = a^T x + b$  for an individual ( $a \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ ):

$$\sum_{j=1}^m |y_j - a^T x_j - b|^2 \rightarrow \min$$

- Regression, least squares...

**And much more:** optimization, data analysis, machine learning...

Floating-point representation of numbers:

$$x = s m b^e$$

s: sign (1 bit), x: mantissa (52 bits), b: basis (normally 2), e: exponent (11 bits)

Observations:

- Numbers subject to **roundoff error** (relative to orders of magnitude)
- Numbers not spaced evenly
- Estimating **relative errors** usually makes better sense.

IEEE 754 norm: floating point number representation ensuring

- |  |   |
|--|---|
| (a) for all $x \in \mathbb{R}$ , exists $\varepsilon$ , $ \varepsilon  < \varepsilon_{\text{mach}}$    | $\text{fl}(x) = x(1 + \varepsilon)$   |
| (b) for all $x, y \in \mathbb{F}$ ,  | $x \circledast y = \text{fl}(x \star y)$ ( $\star$ one of $+, -, \times, /, \sqrt{\phantom{x}}$ ) |
| (c) for all $x, y \in \mathbb{F}$ , exists $\varepsilon$ , $ \varepsilon  < \varepsilon_{\text{mach}}$ | $\text{fl}(x \star y) = (x \star y)(1 + \varepsilon)$   |

## Vector spaces, matrices: a few reminders

**Vector spaces** A set  $E$  is a vector space over  $\mathbb{K} = \mathbb{R}$  or  $\mathbb{C}$  provided:

1. Addition of vectors is commutative:

$$xy = y + x$$

2. Addition of vectors is associative:

$$x + (y + z) = (x + y) + z$$

3. There exists a zero vector:

$$x + 0 = x$$

4. Each vector has an opposite vector:

$$x + (-x) = 0$$

5. Multiplication of scalars and with a vector are compatible:

$$(\alpha\beta)x = \alpha(\beta x)$$

6. Scalar multiplication has a unit element:

$$1x = x$$

7. Scalar multiplication is distributive w.r.t. scalar addition:

$$(\alpha + \beta)x = \alpha x + \beta x$$

8. Scalar multiplication is distributive w.r.t. vector addition:

$$\alpha(x + y) = \alpha x + \alpha y$$

**Matrices:** represent action of linear mappings  $A : \mathbb{K}^n \rightarrow \mathbb{K}^m$  (relative to bases chosen *a priori*):

$$y_i = \sum_{i=1}^n a_{ij} x_j, \quad \text{or} \quad \begin{Bmatrix} y_1 \\ \vdots \\ y_m \end{Bmatrix} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \begin{Bmatrix} x_1 \\ \vdots \\ x_n \end{Bmatrix} \quad \text{i.e. } \boxed{y = Ax}$$

### Terminology:

- $A \in \mathbb{K}^{n \times n}$  is called *square* (otherwise: *rectangular*)
- $A \in \mathbb{K}^{m \times n}$  with zeros in most entries is called *sparse* (otherwise: *dense* or *full*).
- *Transpose*  $A^T \in \mathbb{K}^{n \times m}$  of  $A \in \mathbb{K}^{m \times n}$ :  $(A^T)_{ij} = A_{ji}$ .
- *Conjugate transpose*  $A^H \in \mathbb{C}^{n \times m}$  of  $A \in \mathbb{C}^{m \times n}$ :  $(A^H)_{ij} = \overline{A_{ji}}$ , that is,  $A^H = \overline{A^T}$ .
- $A \in \mathbb{R}^{n \times n}$  verifying  $A^T = A$ ,  $a_{ji} = a_{jj}$  is called *symmetric*.
- $A \in \mathbb{C}^{n \times n}$  verifying  $A^H = A$ ,  $a_{ji} = \overline{a_{jj}}$  is called *Hermitian*.
- $A \in \mathbb{K}^{n \times n}$  Hermitian with  $x^H A x > 0$  for all  $x \neq 0$  is called *symmetric positive definite* (SPD).

### Notation conventions (used throughout):

- Column vectors (e.g.  $x \in \mathbb{K}^{n,1}$ ), (conjugate) transpose are row vectors. Consistent with  $y = Ax$  (matrix-vector product),  $(x, y) = x^H y$  (scalar product).
- Vectors (matrices): lowercase (uppercase) letters, e.g.  $x$  (generic entry  $x_i$ ),  $A$  (generic entry  $a_{ij}$ ).
- MATLAB-like colon ":" to define submatrices by index ranges, e.g.  
 $A_{k:\ell, p:q} := [a_{ij}]_{k \leq i \leq \ell, p \leq j \leq q}$  (rectangular submatrix of  $A$ ),  
 $A_{k:\ell, p} := [a_{ip}]_{k \leq i \leq \ell}$  (part of  $p$ -th column of  $A$ )

## Vector norms

- Measuring “smallness/largeness” of vectors/matrices is essential  
(e.g. convergence of an algorithm: solution errors becoming “increasingly small”).
- Magnitudes measured using (vector, matrix) norms. Defining requirements:

zero norm:

$$\|x\| = 0 \text{ if and only if } x = 0,$$

positive homogeneity:

$$\|\lambda x\| = |\lambda| \|x\| \quad \text{for any } \lambda \in \mathbb{K}$$

for all  $x \in \mathbb{K}^n$

triangle inequality:

$$\|x+y\| \leq \|x\| + \|y\|$$

- Common vector norms ( $p=2$  is Euclidean 2-norm):

$$\|x\|_1 := \sum_{i=1}^n |x_i|, \quad \|x\|_2 := \left( \sum_{i=1}^n |x_i|^2 \right)^{1/2}, \quad \|x\|_p := \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}, \quad \|x\|_\infty := \max_{1 \leq i \leq n} |x_i|$$

- All vector norms for finite-dimensional spaces are equivalent (fails as  $n \rightarrow \infty$ ):

$$C_1 \|x\|_\alpha \leq \|x\|_\beta \leq C_2 \|x\|_\alpha, \quad \text{e.g.} \quad \begin{cases} \|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2 \\ \|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty \\ \|x\|_\infty \leq \|x\|_1 \leq n \|x\|_\infty \end{cases} \quad \text{for all } x \in \mathbb{K}^n.$$

- Classical inequalities:

$$|x^H y| \leq \|x\|_2 \|y\|_2 \quad (\text{Cauchy-Schwarz}), \quad |x^H y| \leq \|x\|_p \|y\|_q \quad (\text{Hölder}, \frac{1}{p} + \frac{1}{q} = 1)$$

## Matrix norms

- Matrix norms induced by vector norms:

$$\|A\|_p := \max_{x \in \mathbb{K}^n, x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} = \max_{\|x\|_p=1} \|Ax\|_p,$$

Provides best upper bound on matrix-vector products: for any  $A \in \mathbb{K}^{m \times n}$  and  $x \in \mathbb{K}^m$ ,

$$\|Ax\|_p \leq \|A\|_p \|x\|_p \quad (\text{with equality for at least one } x),$$

(infinite-dimensional extension: **operator norm**, see e.g. MA102)

- Another norm: Frobenius (**not** an induced norm)

$$\|A\|_F := \left( \sum_{i,j} |a_{ij}|^2 \right)^{1/2} = \sqrt{\text{Tr}(AA^H)}$$

**Sub-multiplicativity:** All induced matrix norms (and Frobenius norm) verify:  $\|AB\| \leq \|A\| \|B\|$

Very important property for deriving (e.g. error) estimates.

Matrix norms are all equivalent. In particular:

$$\frac{1}{\sqrt{m}} \|A\|_1 \leq \|A\|_2 \leq \sqrt{n} \|A\|_1$$

$$\frac{1}{\sqrt{n}} \|A\|_\infty \leq \|A\|_2 \leq \sqrt{m} \|A\|_\infty$$

$$\frac{1}{\sqrt{\min(m,n)}} \|A\|_F \leq \|A\|_2 \leq \|A\|_F$$

for all  $A \in \mathbb{K}^{m \times n}$

**Convention:** Generic symbol  $\|A\|$  always denotes an **induced** norm ( $\|A\|_F$  for Frobenius).

- Many scientific computing tasks boil down to:

apply “function”  $\mathcal{F}$  to data  $x \in \mathcal{X}$ , obtain  $y = \mathcal{F}(x) \in \mathcal{Y}$

Example: solve  $y - f(y) = 0$  by **fixed-point iterations** from initial guess  $x$ :

$$\mathcal{F}(x) := \lim_{n \rightarrow \infty} f^n(x) \quad (\text{assuming } f \text{ to be contracting!})$$

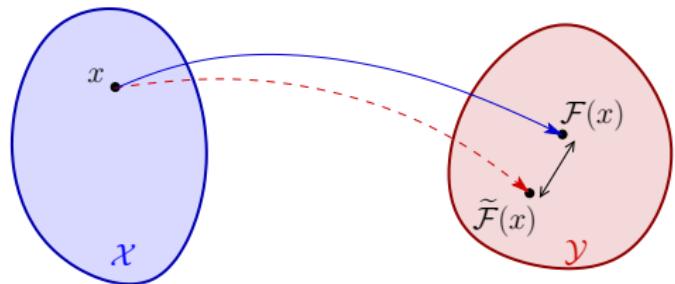
- In practice, data round-off, imperfect implementation of  $\mathcal{F}$ :  $\tilde{y} := \tilde{\mathcal{F}}(x)$

Example (no data round-off):  $\tilde{\mathcal{F}}(x) := \tilde{f}^{N+1}(x)$  with  $N$  such that  $|\tilde{f}^{N+1}(x) - \tilde{f}^N(x)| < \varepsilon$

How close to  $y = \mathcal{F}(x)$  is the approximation  $\tilde{y} := \tilde{\mathcal{F}}(x)$ ?

- Relative solution accuracy:**

$$e_{\text{rel}} := \frac{\|\tilde{\mathcal{F}}(x) - \mathcal{F}(x)\|}{\|\mathcal{F}(x)\|}$$



Best conceivable accuracy:  $e_{\text{rel}} = O(\varepsilon_{\text{mach}})$  (achieved by individual floating-point operations). Requirement  $e_{\text{rel}} \approx \varepsilon_{\text{mach}}$  overly demanding (large-scale and/or ill-conditioned problems).

- **Stability:** a more-appropriate aim:

for each  $x \in \mathcal{X}$ :

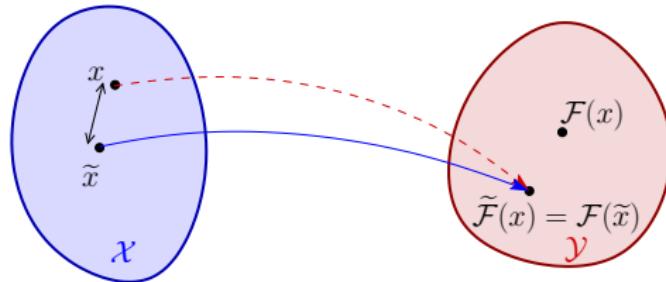
$$\frac{\|\tilde{\mathcal{F}}(x) - \mathcal{F}(\tilde{x})\|}{\|\mathcal{F}(\tilde{x})\|} = O(\varepsilon_{\text{mach}}) \quad \text{for some } \tilde{x} \text{ with } \frac{\|x - \tilde{x}\|}{\|x\|} = O(\varepsilon_{\text{mach}})$$

"A stable algorithm yields nearly the right answer if given a nearly correct data."

- Stronger requirement (replacing  $O(\varepsilon_{\text{mach}})$  with zero): **backward stability:**

for each  $x \in \mathcal{X}$ :

$$\tilde{\mathcal{F}}(x) = \mathcal{F}(\tilde{x}) \quad \text{for some } \tilde{x} \text{ with } \frac{\|x - \tilde{x}\|}{\|x\|} = O(\varepsilon_{\text{mach}})$$

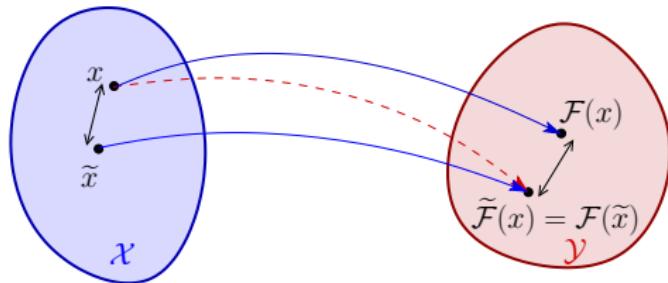


"A backward-stable algorithm yields the exact answer for some nearly correct data."

## Conditioning, condition number

- Connect forward and backward errors by relative sensitivity:

$$\varrho = \varrho(\mathcal{F}; x, \tilde{x}) := \frac{\|\tilde{\mathcal{F}}(x) - \mathcal{F}(x)\|}{\|\mathcal{F}(x)\|} \left( \frac{\|\tilde{x} - x\|}{\|x\|} \right)^{-1} = \frac{\|\mathcal{F}(\tilde{x}) - \mathcal{F}(x)\|}{\|\mathcal{F}(x)\|} \frac{\|x\|}{\|\tilde{x} - x\|}$$



- Condition number (of  $\mathcal{F}$  at  $x$ ): limiting value of  $\varrho$  for  $\|\tilde{x} - x\|$  small:

$$\kappa(\mathcal{F}; x) := \lim_{\delta \rightarrow 0} \sup_{\|\tilde{x} - x\| \leq \delta} \varrho(\mathcal{F}; x, \tilde{x})$$

Explicit formula if  $\mathcal{F}$  regular enough:

$$\kappa(\mathcal{F}, x) = \frac{\|\mathcal{F}'(x)\| \|x\|}{\|\mathcal{F}(x)\|},$$

- $\kappa(\mathcal{F}; x)$ : dimensionless number;
- A solution process  $\mathcal{F}$  is *well-conditioned* (*ill-conditioned*) if  $\kappa = O(1)$  ( $\kappa \gg 1$ )

## Condition number of linear systems

Solution of linear system  $Ay = b$ : sensitivity to data  $A, b$  ( $A \in \mathbb{K}^{n \times n}$  invertible)

- Perturbation  $z$  of solution  $y = A^{-1}b$  satisfies  $(A+E)(y+z) = b+f$ , i.e.

$$(A+E)z = f - Ey.$$

- If  $\|A^{-1}\| \|E\| < 1$  (perturbation of  $A$  small enough),  $(A+E)^{-1} = A^{-1}(I + EA^{-1})^{-1}$  exists.

$$\|(A+E)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|EA^{-1}\|} \leq \frac{\|A^{-1}\|}{1 - \|E\| \|A^{-1}\|}, \quad (\text{using submultiplicativity})$$

- Solution error estimate:

$$\|z\| = \|(A+E)^{-1}(f - Ey)\| \implies \|z\| \leq \frac{\|A^{-1}\|}{1 - \|E\| \|A^{-1}\|} (\|f\| + \|E\| \|y\|).$$

Formulate using relative errors:

$$\frac{\|z\|}{\|y\|} \leq \frac{\|A^{-1}\| \|A\|}{1 - \|E\| \|A^{-1}\|} \left( \frac{\|f\|}{\|b\|} \frac{\|b\|}{\|A\| \|y\|} + \frac{\|E\|}{\|A\|} \right) \leq \frac{\|A^{-1}\| \|A\|}{1 - \|E\| \|A^{-1}\|} \left( \frac{\|f\|}{\|b\|} + \frac{\|E\|}{\|A\|} \right),$$

**Relative sensitivity of solution w.r.t. data:**

$$\frac{\|z\|/\|y\|}{\|f\|/\|b\| + \|E\|/\|A\|} \leq \frac{\kappa(A)}{1 - \kappa(A) \|E\|/\|A\|} = \kappa(A) + O(\|E\|/\|A\|) \quad \text{with } \kappa(A) := \|A^{-1}\| \|A\|.$$

Relative sensitivity of solution w.r.t. data:

$$\frac{\|z\|/\|y\|}{\|f\|/\|b\| + \|E\|/\|A\|} \leq \kappa(A) + O(\|E\|/\|A\|) \quad \text{with } \kappa(A) := \|A^{-1}\| \|A\|.$$

Condition number  $\kappa(A) = \|A^{-1}\| \|A\|$  of  $A$ : upper bound of condition number for solving  $Ay = b$ .

Properties of  $\kappa(A)$ :

- Always  $\kappa(A) \geq 1$  ( $\|A^{-1}\| \|A\| \geq \|A^{-1}A\| = \|I\| = 1$  for any induced norm).
- $\kappa(A)$  depends on choice of (matrix) norm.
- If  $A$  normal ( $AA^H = A^H A$ ), we have  $A = Q\Lambda Q^H$  for some  $Q$  unitary. Then:  
$$\|A\|_2 = |\lambda_{\max}|, \quad \|A^{-1}\|_2 = 1/|\lambda_{\min}|, \quad \text{and hence} \quad \boxed{\kappa_2(A) = |\lambda_{\max}|/|\lambda_{\min}|}.$$
- $\|Q\|_2 = 1$  and  $\|Q^{-1}\|_2 = 1$  if  $Q$  orthogonal or unitary. Consequently,  $\boxed{\kappa_2(Q) = 1}$ .
- For arbitrary  $A \in \mathbb{K}^{m \times n}$ ,  $\kappa_2(A)$  given in terms of either singular values or pseudo-inverse of  $A$  (see Part 3).

## A simple numerical example

- Example (exact matrix inverse):

$$A = \begin{bmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{bmatrix} \implies A^{-1} = \begin{bmatrix} 25 & 41 & 10 & -6 \\ -41 & 68 & -17 & 10 \\ 10 & -17 & 5 & -3 \\ -6 & 10 & -3 & 2 \end{bmatrix}$$

Note:  $AA^T = A^T A$  (i.e.  $A$  is normal)

- Effect of perturbations of  $A$  or  $b$  on solution of  $x$  of  $Ax = b$ :

$$b = [32 \ 23 \ 33 \ 31]^T \implies x = [1 \ 1 \ 1 \ 1]^T$$

$$\delta b = [0.1 \ -0.1 \ 0.1 \ -0.1]^T \implies x = [9.2 \ -12.6 \ 4.5 \ -1.1]^T$$

$$\delta A_{23} = 0.1 \implies x \approx [-4.86 \ -10.7 \ -1.43 \ -2.43]^T$$

- Eigenvalues of  $A$ :

$$\Lambda \approx \text{Diag}[ 30.29 \ 3.858 \ 0.8431 \ 0.01015 ], \quad \kappa_2(A) \approx 3 \cdot 10^3$$

$A$  is a rather ill-conditioned  $4 \times 4$  matrix.