

MA 7007: Numerical Solution of Differential Equations I

Elliptic Partial Differential Equations



Suh-Yuh Yang (楊肅煜)

Department of Mathematics, National Central University
Jhongli District, Taoyuan City 32001, Taiwan

E-mail: syyang@math.ncu.edu.tw

Website: <http://www.math.ncu.edu.tw/~syyang/>

Introduction

In two space dimensions a constant-coefficient elliptic equation has the form

$$a_1 u_{xx}(x, y) + a_2 u_{xy}(x, y) + a_3 u_{yy}(x, y) + a_4 u_x(x, y) + a_5 u_y(x, y) + a_6 u(x, y) = f(x, y),$$

for all $(x, y) \in \Omega$, where $\Omega \subseteq \mathbb{R}^2$ is typically an open bounded domain and the coefficients a_1, a_2, a_3 satisfy

$$a_2^2 - 4a_1a_3 < 0.$$

This equation must be complemented with some boundary condition on the boundary $\partial\Omega$ such as the Dirichlet boundary condition

$$u(x, y) = g(x, y) \quad \text{for all } (x, y) \in \partial\Omega.$$

Steady-state heat conduction

Heat conduction problem in two space dimensions:

$$\begin{cases} u_t = (\kappa u_x)_x + (\kappa u_y)_y + \psi, & t \in (0, T), (x, y) \in \Omega, \\ \text{"Initial and boundary conditions."} \end{cases}$$

where $\kappa(x, y) > 0$ is the diffusivity and $\psi(t, x, y)$ is a source function. If the boundary conditions and the source term are independent of time t , then we expect a steady state to exist,

$$(\kappa u_x)_x + (\kappa u_y)_y = -\psi := f \quad \text{in } \Omega + \text{"boundary conditions."}$$

Let $\kappa(x, y) \equiv 1$ for all $(x, y) \in \Omega$.

- 1 Poisson equation: $u_{xx} + u_{yy} = f$.
- 2 Laplace equation: $u_{xx} + u_{yy} = 0$.

Solutions to the Laplace equation are called harmonic functions.

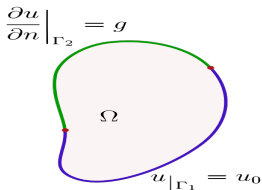
Notation and boundary conditions

Notation: $\nabla := [\partial_x, \partial_y]^\top$.

- 1 gradient operator: $\nabla u = [u_x, u_y]^\top$.
- 2 divergence operator: $\nabla \cdot [u, v]^\top = u_x + v_y$.
- 3 Laplacian operator: $\nabla^2 u := \nabla \cdot \nabla u = u_{xx} + u_{yy} := \Delta u$.

Boundary conditions:

- 1 Dirichlet BC: $u(x, y) = g(x, y), \quad \forall (x, y) \in \partial\Omega$
- 2 Neumann BC: $\frac{\partial u}{\partial \mathbf{n}}(x, y) \left(:= \nabla u(x, y) \cdot \mathbf{n}(x, y) \right) = g(x, y), \quad \forall (x, y) \in \partial\Omega$
- 3 Robin BC: $au(x, y) + b\frac{\partial u}{\partial \mathbf{n}}(x, y) = g(x, y), \quad \forall (x, y) \in \partial\Omega$
- 4 Mixed BC:



Centered difference scheme

For example, we consider the Poisson equation with the Dirichlet BC:

$$\begin{aligned}\nabla^2 u &= f \quad \text{in } \Omega := (0, 1) \times (0, 1), \\ u &= g \quad \text{on } \partial\Omega.\end{aligned}$$

We will use the uniform Cartesian grid: (x_i, y_j) , where $x_i = i\Delta x$ and $y_j = j\Delta y$, Δx and Δy are the grid sizes in x - and y - directions.

Let u_{ij} represent an approximation to $u(x_i, y_j)$ and $f_{ij} := f(x_i, y_j)$.

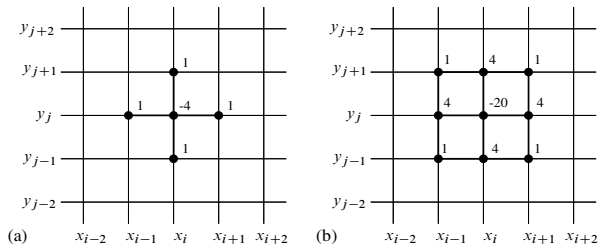
$$\frac{1}{(\Delta x)^2} (u_{i-1,j} - 2u_{ij} + u_{i+1,j}) + \frac{1}{(\Delta y)^2} (u_{i,j-1} - 2u_{ij} + u_{i,j+1}) = f_{ij}.$$

For simplicity, we set $\Delta x = \Delta y = h$. Then we have

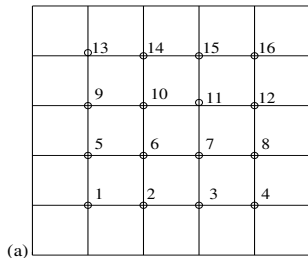
$$\nabla_5^2 u_{ij} := \frac{1}{h^2} (u_{i-1,j} + u_{i+1,j} + u_{i,j-1} + u_{i,j+1} - 4u_{ij}) = f_{ij}.$$

Let $0 = x_0 < x_1 < \cdots < x_m < x_{m+1} = 1$ and $0 = y_0 < y_1 < \cdots < y_m < y_{m+1} = 1$ be the partitions. Then $h = 1/(m+1)$. From the above equations, we have an $m^2 \times m^2$ linear system $\mathbf{A}\mathbf{u} = \mathbf{F}$ of m^2 unknowns u_{ij} for $1 \leq i \leq m, 1 \leq j \leq m$, where \mathbf{A} is sparse (Roughly speaking, at least $\frac{2}{3} \uparrow$ zeros).

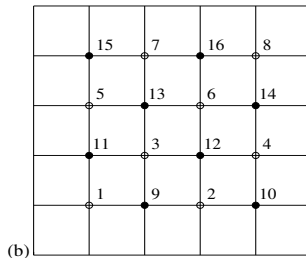
Computational grid: 5-point stencil and 9-point stencil



Ordering the unknowns and equations



(a) The rowwise ordering.



(b) The red-black ordering.

The rowwise ordering

Let

$$\mathbf{u} = \begin{bmatrix} u^{[1]} \\ u^{[2]} \\ \vdots \\ u^{[m]} \end{bmatrix}, \quad u^{[j]} = \begin{bmatrix} u_{1j} \\ u_{2j} \\ \vdots \\ u_{mj} \end{bmatrix}, \quad \mathbf{F} = \begin{bmatrix} f^{[1]} \\ f^{[2]} \\ \vdots \\ f^{[m]} \end{bmatrix} + BV, \quad f^{[j]} = \begin{bmatrix} f_{1j} \\ f_{2j} \\ \vdots \\ f_{mj} \end{bmatrix}.$$

Then

$$A = \frac{1}{h^2} \begin{bmatrix} T & I & & & & \\ I & T & I & & & \\ & & \ddots & \ddots & \ddots & \\ & & & I & T & I \\ & & & & I & T \end{bmatrix}, \quad T = \begin{bmatrix} -4 & 1 & & & & \\ 1 & -4 & 1 & & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & -4 & 1 \\ & & & & 1 & -4 \end{bmatrix}.$$

Accuracy and stability

The local truncation error τ_{ij} at the grid point (i, j) is defined by

$$\tau_{ij} := \frac{1}{h^2} \left(u(x_{i-1}, y_j) + u(x_{i+1}, y_j) + u(x_i, y_{j-1}) + u(x_i, y_{j+1}) - 4u(x_i, y_j) \right) - f(x_i, y_j).$$

By the Taylor expansion, we have

$$\tau_{ij} = \frac{1}{12} h^2 \left(u_{xxxx}(x_i, y_j) + u_{yyyy}(x_i, y_j) \right) + O(h^4)$$

and

$$A\mathbf{u}^{exact} = \mathbf{F} + \boldsymbol{\tau},$$

where A is the discretization matrix corresponding to the rowwise ordering. Letting the global error $E_{ij} := u_{ij} - u(x_i, y_j)$ and noting that $A\mathbf{u} = \mathbf{F}$, we obtain

$$A\mathbf{E} = -\boldsymbol{\tau} \implies \mathbf{E} = A^{-1}(-\boldsymbol{\tau}).$$

The method will be globally second order accurate **in some grid function norm** provided that $\|A^{-1}\|$ is uniformly bounded as $h \rightarrow 0^+$.

Accuracy and stability (continued)

We consider the 2-norm for the discretization matrix A . By further computations, one can show that for $p, q = 1, 2, \dots, m$, the eigenvector $\mathbf{u}^{p,q}$ has the m^2 elements,

$$u_{ij}^{p,q} = \sin(p\pi ih) \sin(q\pi ih)$$

and the corresponding eigenvalue is

$$\lambda_{p,q} = \frac{2}{h^2} \left((\cos(p\pi h) - 1) + (\cos(q\pi h) - 1) \right) < 0. \quad \left(\text{Note that } h = \frac{1}{m+1} \right)$$

Thus, the one closest to origin is $\lambda_{1,1} = -2\pi^2 + O(h^2)$. (Hint: By Taylor expansion: $\cos(x) = 1 - x^2/2! + x^4/4! - \dots$) The spectral radius of A^{-1} is

$$\rho(A^{-1}) = \frac{1}{|\lambda_{1,1}|} \approx \frac{1}{2\pi^2} \quad \text{as } h \rightarrow 0^+,$$

and then as $h \rightarrow 0^+$,

$$\|A^{-1}\|_2 = \sqrt{\rho(A^{-\top}A^{-1})} = \sqrt{\rho((A^{-1})^2)} = \sqrt{(\rho(A^{-1}))^2} = \rho(A^{-1}) \approx \frac{1}{2\pi^2},$$

which is uniformly bounded.

Accuracy and stability (continued)

From the centered difference scheme with uniform mesh size, $\nabla_5^2 u_{ij} = f_{ij}$, we obtain an $m^2 \times m^2$ linear system of m^2 unknowns u_{ij} for $1 \leq i \leq m, 1 \leq j \leq m$,

$$Au = F,$$

(or more precisely, $A^h u^h = F^h$). Now suppose source term F is perturbed by a small vector p (say, $\|p\|_2 < \delta$ for a small $\delta > 0$) and the corresponding solution is denoted by \tilde{u} . Then we have

$$A\tilde{u} = F + p,$$

and

$$A(\tilde{u} - u) = p,$$

which implies

$$\tilde{u} - u = A^{-1}p \implies \|\tilde{u} - u\|_2 \leq \|A^{-1}\|_2 \|p\|_2 < \|A^{-1}\|_2 \delta,$$

where $\|\tilde{u} - u\|_2$ and $\|p\|_2$ are **grid function norms**. Since $\|A^{-1}\|_2$ is uniformly bounded, we have $\|\tilde{u} - u\|_2 \leq C\delta$. Hence, the centered difference scheme for the Poisson problem is *stable*.

Condition number

The 2-norm condition number of the discretization matrix A is defined by

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

Notice that as $h \rightarrow 0^+$,

$$\|A\|_2 = \rho(A) = \max_{1 \leq p, q \leq m} |\lambda_{p,q}| = |\lambda_{m,m}| = \frac{4}{h^2} \left| \cos\left(\frac{m}{m+1}\pi\right) - 1 \right| \approx \frac{4}{h^2} | -2 | = \frac{8}{h^2}.$$

Therefore

$$\kappa(A) \approx \frac{4}{\pi^2 h^2} = O(h^{-2}) \quad \text{as } h \rightarrow 0^+.$$

The discretization matrix A is very ill-conditioned as we refine the grid.

The 9-point Laplacian

By the Taylor expansion, we have

$$\begin{aligned}\nabla_5^2 u(x_i, y_j) &= \nabla^2 u(x_i, y_j) + \frac{1}{12} h^2 \frac{\partial^4 u}{\partial x^4}(x_i, y_j) + \frac{1}{12} h^2 \frac{\partial^4 u}{\partial y^4}(x_i, y_j) + O(h^4) \\ \implies \nabla_5^2 u(x_i, y_j) &+ \frac{2}{12} h^2 \frac{\partial^4 u}{\partial y^2 \partial x^2}(x_i, y_j) \\ &= \nabla^2 u(x_i, y_j) + \frac{1}{12} h^2 \left\{ \frac{\partial^4 u}{\partial x^4} + 2 \frac{\partial^4 u}{\partial x^2 \partial y^2} + \frac{\partial^4 u}{\partial y^4} \right\} (x_i, y_j) + O(h^4) \\ &= \nabla^2 u(x_i, y_j) + \frac{1}{12} h^2 \nabla^2 f(x_i, y_j) + O(h^4) \\ \implies \nabla_5^2 u(x_i, y_j) &+ \frac{2}{12} h^2 \frac{\partial^4 u}{\partial y^2 \partial x^2}(x_i, y_j) - \frac{1}{12} h^2 \nabla^2 f(x_i, y_j) = \nabla^2 u(x_i, y_j) + O(h^4).\end{aligned}$$

The 9-point Laplacian (continued)

$$\nabla_5^2 u(x_i, y_j) + \frac{2}{12} h^2 \frac{\partial^4 u}{\partial y^2 \partial x^2}(x_i, y_j) - \frac{1}{12} h^2 \nabla^2 f(x_i, y_j) = \nabla^2 u(x_i, y_j) + O(h^4)$$

$$\begin{aligned} \implies \nabla_5^2 u(x_i, y_j) + \frac{h^2}{6h^4} \{ & u(x_{i-1}, y_{j-1}) - 2u(x_{i-1}, y_j) + u(x_{i-1}, y_{j+1}) \\ & - 2u(x_i, y_{j-1}) + 4u(x_i, y_j) - 2u(x_i, y_{j+1}) \\ & + u(x_{i+1}, y_{j-1}) - 2u(x_{i+1}, y_j) + u(x_{i+1}, y_{j+1}) \} + O(h^4) \\ & - \frac{1}{12} h^2 \nabla^2 f(x_i, y_j) = \nabla^2 u(x_i, y_j) + O(h^4). \end{aligned}$$

$$\begin{aligned} \therefore \nabla_9^2 u_{ij} := \frac{1}{6h^2} \{ & 4u_{i-1,j} + 4u_{i+1,j} + 4u_{i,j-1} + 4u_{i,j+1} + u_{i-1,j-1} + u_{i-1,j+1} \\ & + u_{i+1,j-1} + u_{i+1,j+1} - 20u_{ij} \} = f_{ij} + \frac{1}{12} h^2 \nabla^2 f(x_i, y_j) \end{aligned}$$

is a finite difference scheme for the Poisson problem with local truncation error $O(h^4)$. The term $\frac{1}{12} h^2 \nabla^2 f(x_i, y_j)$ can be exactly computed or approximated by $\frac{1}{12} h^2 \nabla_5^2 f(x_i, y_j)$.

Estimates from the true solution

Suppose we know the true solution. Let $E(h)$ denote the error function of grid size h , i.e., $E(h) = \|U(h) - \hat{U}(h)\|$, where $U(h)$ is the numerical solution vector and $\hat{U}(h)$ is the true solution evaluated on the same grid.

If the method is p -th order accurate, i.e., $E(h) = Ch^p + O(h^{p+1})$ as $h \rightarrow 0$, then for $0 < h_2 < h_1$ sufficiently small, we expect $E(h_1) \approx Ch_1^p$ and $E(h_2) \approx Ch_2^p$. The order of convergence can be estimated using

$$p \approx \frac{\log(E(h_1)/E(h_2))}{\log(h_1/h_2)},$$

this is because

$$\log \frac{E(h_1)}{E(h_2)} \approx \log \frac{Ch_1^p}{Ch_2^p} = \log \left(\frac{h_1}{h_2} \right)^p = p \log \frac{h_1}{h_2}.$$

Estimates from a fine-grid solution

Now suppose we don't know the exact solution but that we can afford to run the problem on a very fine grid, say \bar{h} , and use the numerical solution $U(\bar{h})$ as a reference solution.

Let $U(h)$ be the numerical solution on a coarser grid h , and $\bar{U}(h)$ be the restriction of $U(\bar{h})$ to the h -grid. Define the approximate error and the true error as

$$\bar{E}(h) = \|U(h) - \bar{U}(h)\| \quad \text{and} \quad E(h) = \|U(h) - \hat{U}(h)\|,$$

respectively. Then consider

$$U(h) - \bar{U}(h) = (U(h) - \hat{U}(h)) + (\hat{U}(h) - \bar{U}(h)).$$

If the method is supposed to be p -th order accurate and $\bar{h}^p \ll h^p$, then we will have $U(h) - \bar{U}(h) \approx U(h) - \hat{U}(h)$ since the second term $\hat{U}(h) - \bar{U}(h)$ should be negligible compared to the first term $U(h) - \hat{U}(h)$. In this case, the approximate error $\bar{E}(h)$ can be used as a good estimate of the true error $E(h)$.

L^p -norm and discrete L^p -norm for grid functions, $1 \leq p \leq \infty$

- 1 **L^p -norm:** Let $U(x)$ be an approximate solution of $u(x)$ on $\bar{\Omega} = [a, b]$ and let $e(x) := U(x) - u(x)$, where $U(x)$ and $u(x)$ are smooth enough. Then

$$\|e\|_{L^\infty(\Omega)} := \max_{a \leq x \leq b} |e(x)| \quad \text{and} \quad \|e\|_{L^p(\Omega)} := \left(\int_a^b |e(x)|^p dx \right)^{1/p}, \quad p \geq 1.$$

- 2 **Discrete L^p -norm of grid function e :** Let $U_i \approx u(x_i)$, $1 \leq i \leq N$. Let $e_i = U_i - u(x_i)$ and $e = (e_1, \dots, e_N)^\top$. Then

$$\|e\|_\infty := \max_{1 \leq i \leq N} |e_i| \quad \text{and} \quad \|e\|_p := \left(h \sum_{i=1}^N |e_i|^p \right)^{1/p}, \quad p \geq 1.$$

- 3 **2-D discrete L^p -norm of grid function e :**

$$\|e\|_\infty := \max_{1 \leq i, j \leq N} |e_{ij}| \quad \text{and} \quad \|e\|_p := \left(h^2 \sum_i \sum_j |e_{ij}|^p \right)^{1/p}, \quad p \geq 1.$$

Review: Vector norm

Let V be a vector space over \mathbb{R} , e.g., $V = \mathbb{R}^n$. A norm is a real-valued function $\|\cdot\| : V \rightarrow \mathbb{R}$ that satisfies

- 1 $\|x\| \geq 0, \forall x \in V$, and $\|x\| = 0$ if and only if $x = \mathbf{0}$;
- 2 $\|\lambda x\| = |\lambda| \|x\|, \forall x \in V$ and $\lambda \in \mathbb{R}$;
- 3 $\|x + y\| \leq \|x\| + \|y\|, \forall x, y \in V$ (triangle inequality).

Note: $\|x\|$ is called the norm of x , the length or magnitude of x .

Some vector norms on \mathbb{R}^n

Let $\mathbf{x} = (x_1, x_2, \dots, x_n)^\top \in \mathbb{R}^n$:

- ① The 2-norm (Euclidean norm, or ℓ^2 norm):

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}.$$

- ② The infinity norm (ℓ^∞ -norm):

$$\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|.$$

- ③ The 1-norm (ℓ^1 -norm):

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|.$$

The difference between the above norms

- ① Take three vectors $x = (4, 4, -4, 4)^\top$, $v = (0, 5, 5, 5)^\top$, $w = (6, 0, 0, 0)^\top$:

	$\ \cdot\ _1$	$\ \cdot\ _2$	$\ \cdot\ _\infty$
x	16	8	4
v	15	8.66	5
w	6	6	6

- ② What is the unit ball $\{x \in \mathbb{R}^2 : \|x\| \leq 1\}$ for the three norms above?
- 2-norm: a circle;
 - ∞ -norm: a square;
 - 1-norm: a diamond.

Matrix norm

Let A be an $n \times n$ real matrix. If $\|\cdot\|$ is any norm on \mathbb{R}^n , then

$$\|A\| := \sup\{\|Ax\| : x \in \mathbb{R}^n, \|x\| = 1\} \left(\iff \|A\| := \sup\left\{ \frac{\|Ax\|}{\|x\|} : x \in \mathbb{R}^n, x \neq \mathbf{0} \right\} \right)$$

defines a norm on the vector space of all $n \times n$ real matrices.

(This is called the matrix norm associated with the given vector norm)

Proof:

- $\because \|Ax\| \geq 0 \forall x \in \mathbb{R}^n, \|x\| = 1. \therefore \|A\| \geq 0.$
Moreover, one can check that $\|A\| = 0$ if and only if $A = \mathbf{0}$.
- $\|\lambda A\| = \sup\{\|\lambda Ax\| : \|x\| = 1\} = \sup\{|\lambda| \|Ax\| : \|x\| = 1\}$
 $= |\lambda| \sup\{\|Ax\| : \|x\| = 1\} = |\lambda| \|A\|.$
- $\|A + B\| = \sup\{\|(A + B)x\| : \|x\| = 1\} \leq \sup\{\|Ax\| + \|Bx\| : \|x\| = 1\}$
 $\leq \sup\{\|Ax\| : \|x\| = 1\} + \sup\{\|Bx\| : \|x\| = 1\} = \|A\| + \|B\|.$

Some additional properties

① $\|Ax\| \leq \|A\|\|x\|, \forall x \in \mathbb{R}^n.$

Proof:

Let $x \neq \mathbf{0}$. Then $v = \frac{x}{\|x\|}$ is of norm 1. $\therefore \|A\| \geq \|Av\| = \frac{\|Ax\|}{\|x\|}.$

② $\|I\| = 1.$

③ $\|AB\| \leq \|A\|\|B\|.$

Proof:

$$\begin{aligned} \|AB\| &:= \sup\{\|(AB)x\| : x \in \mathbb{R}^n, \|x\| = 1\} \\ &\leq \sup\{\|A\|\|Bx\| : x \in \mathbb{R}^n, \|x\| = 1\} \\ &\leq \sup\{\|A\|\|B\|\|x\| : x \in \mathbb{R}^n, \|x\| = 1\} = \|A\|\|B\|. \end{aligned}$$

Some matrix norms

Let $A_{n \times n} = (a_{ij})$ be an $n \times n$ real matrix. Then

① The ∞ -matrix norm:

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

② The 1-matrix norm:

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

③ The 2-matrix norm:

$$\|A\|_2 = \sup_{\|x\|_2=1} \|Ax\|_2.$$

The 2-matrix norm

- 1 $\|A\|_2$ is not easy to compute.
- 2 Since $A^\top A$ is symmetric, $A^\top A$ has n real eigenvalues, $\lambda_1, \lambda_2, \dots, \lambda_n \in \mathbb{R}$. Moreover, one can prove that they are all nonnegative. Then

$$\rho(A^\top A) := \max_{1 \leq i \leq n} \{\lambda_i\} \geq 0.$$

is called the spectral radius of $A^\top A$.

- 3 Then the 2-matrix norm of A is given by

$$\|A\|_2 = \sqrt{\rho(A^\top A)}.$$

- 4 The 2-matrix norm is also called the *spectral norm*.

Some error analysis

- 1 Suppose that we want to solve the linear system $Ax = b$, but b is somehow perturbed to \tilde{b} (this may happen when we convert a real b to a floating-point b).
- 2 Then actual solution would satisfy a slightly different linear system

$$A\tilde{x} = \tilde{b}.$$

- 3 *Question:* Is \tilde{x} very different from the desired solution x of the original system?
- 4 Of course, the answer should depend on *how good the matrix A is*.
- 5 Let $\|\cdot\|$ be a vector norm, we consider two types of errors:
 - absolute error: $\|x - \tilde{x}\|$?
 - relative error: $\|x - \tilde{x}\|/\|x\|$?

The absolute error

For the absolute error, we have

$$\|x - \tilde{x}\| = \|A^{-1}b - A^{-1}\tilde{b}\| = \|A^{-1}(b - \tilde{b})\| \leq \|A^{-1}\| \|b - \tilde{b}\|.$$

Therefore, the absolute error of x depends on two factors: the absolute error of b and the matrix norm of A^{-1} .

The relative error

For the relative error, we have

$$\begin{aligned}\|x - \tilde{x}\| &= \|A^{-1}b - A^{-1}\tilde{b}\| = \|A^{-1}(b - \tilde{b})\| \\ &\leq \|A^{-1}\| \|b - \tilde{b}\| = \|A^{-1}\| \|Ax\| \frac{\|b - \tilde{b}\|}{\|b\|} \\ &\leq \|A^{-1}\| \|A\| \|x\| \frac{\|b - \tilde{b}\|}{\|b\|}.\end{aligned}$$

That is

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

Therefore, the relative error of x depends on two factors: the relative error of b and $\|A\| \|A^{-1}\|$.

Condition number

- 1 Therefore, we define a condition number of the matrix A as

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

$\kappa(A)$ measures how good the matrix A is.

- 2 Example: Let $\varepsilon > 0$ and

$$A = \begin{bmatrix} 1 & 1 + \varepsilon \\ 1 - \varepsilon & 1 \end{bmatrix} \implies A^{-1} = \varepsilon^{-2} \begin{bmatrix} 1 & -1 - \varepsilon \\ -1 + \varepsilon & 1 \end{bmatrix}.$$

Then $\|A\|_{\infty} = 2 + \varepsilon$, $\|A^{-1}\|_{\infty} = \varepsilon^{-2}(2 + \varepsilon)$, and $\kappa(A) = \left(\frac{2 + \varepsilon}{\varepsilon}\right)^2 \geq \frac{4}{\varepsilon^2}$.

Condition number (continued)

- 1 For example, if $\varepsilon = 0.01$, then $\kappa(A) \geq 40000$.
- 2 What does this mean?
It means that the relative error in x can be 40000 times greater than the relative error in b .
- 3 If $\kappa(A)$ is large, we say that A is *ill-conditioned*, otherwise A is *well-conditioned*.
- 4 In the ill-conditioned case, the solution is very sensitive to the small changes in the right-hand vector b (higher precision in b may be needed).

Another way to measure the error

Consider the linear system $Ax = b (\neq 0)$. Let \tilde{x} be a computed solution (an approximation to x).

- ① Residual vector:

$$r = b - A\tilde{x}.$$

- ② Error vector:

$$e = x - \tilde{x}.$$

- ③ They satisfy

$$Ae = Ax - A\tilde{x} = b - A\tilde{x} = r.$$

- ④ Moreover, we have

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

(Theorem on bounds involving condition number)

Proof of the Theorem

$$\because \mathbf{Ae} = \mathbf{r}.$$

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

$$\therefore \|\mathbf{e}\|\|\mathbf{b}\| = \|\mathbf{A}^{-1}\mathbf{r}\|\|\mathbf{Ax}\| \leq \|\mathbf{A}^{-1}\|\|\mathbf{r}\|\|\mathbf{A}\|\|\mathbf{x}\|.$$

$$\therefore \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \kappa(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.$$

On the other hand, we have $\|\mathbf{r}\|\|\mathbf{x}\| = \|\mathbf{Ae}\|\|\mathbf{A}^{-1}\mathbf{b}\| \leq \|\mathbf{A}\|\|\mathbf{e}\|\|\mathbf{A}^{-1}\|\|\mathbf{b}\|.$

$$\therefore \frac{1}{\kappa(\mathbf{A})} \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|}.$$