

Numerical analysis

System of linear equations

Kazuaki Tanaka

$$I = \|I\| = \|A \cdot A^{-1}\|$$

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

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

1 Condition number

$n \times n$ matrix

Definition 1.1. The condition number $\text{cond}(A)$ of nonsingular matrix A is defined by

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

Note: Condition numbers depend on the norm imposed on spaces (e.g., l^1 , l^2 or l^∞).

1.1 Impact on the stability of linear equations

$$Ax = b$$

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

small error (pointing to Δb)
small? (pointing to Δx)

$$\|b\| = \|Ax\| \leq \|A\| \|x\|$$

$$\Leftrightarrow \|x\| \geq \|A\|^{-1} \|b\| \dots \textcircled{1}$$

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

$$\Rightarrow \|\Delta x\| \leq \|A^{-1}\| \|\Delta b\| \dots \textcircled{2}$$

By $\frac{\textcircled{2}}{\textcircled{1}}$

$$\frac{\|\Delta x\|}{\|x\|} \leq \underbrace{\|A\| \|A^{-1}\|}_{\substack{\text{cond}(A) \\ \gg 1}} \underbrace{\frac{\|\Delta b\|}{\|b\|}}_{\text{small}}$$

Can be large

2 Stationary Iterative methods

$x^0 \xrightarrow{\text{give}} x^1 \xrightarrow{\text{relation}} x^2 \xrightarrow{\text{stop}} \dots \xrightarrow{\text{stop}} x^n \sim x^* (\text{exact})$

$Ax = b$
 $A = D + U + L$
 $D = \begin{pmatrix} a_{11} & 0 & \dots & 0 \\ 0 & a_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_{nn} \end{pmatrix}$
 $U = \begin{pmatrix} 0 & a_{12} & \dots & a_{1n} \\ 0 & 0 & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}$
 $L = \begin{pmatrix} 0 & 0 & \dots & 0 \\ a_{21} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix}$

$$(D + U + L)x = b$$

$$\Leftrightarrow Dx = -(U + L)x + b$$

$$\Leftrightarrow x^{k+1} = -D^{-1}(U + L)x^k + D^{-1}b$$

Jacobi

$$(D + L)x = -Ux + b$$

$$\Leftrightarrow x^{k+1} = -(D + L)^{-1}Ux^k + (D + L)^{-1}b$$

Gauss-Seidel

Jacobi method

$$x^{k+1} = -D^{-1}(L + U)x^k + D^{-1}b, \quad k = 0, 1, 2, \dots \quad (2)$$

This can be written in the component form

$x^k = \begin{pmatrix} x_1^k \\ \vdots \\ x_n^k \end{pmatrix}$

$$x_i^{k+1} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1, j \neq i}^n a_{ij}x_j^k \right), \quad i = 1, 2, \dots, n, \quad k = 0, 1, 2, \dots \quad (3)$$

Gauss-Seidel method

$$x^{k+1} = -(L + D)^{-1}Ux^k + (L + D)^{-1}b, \quad k = 0, 1, 2, \dots \quad (4)$$

This can be written in the component form

$$x_i^{k+1} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{k+1} - \sum_{j=i+1}^n a_{ij}x_j^k \right), \quad i = 1, 2, \dots, n, \quad k = 0, 1, 2, \dots \quad (5)$$

2.1 Stopping criterion

For small $\varepsilon > 0$,

- $\frac{\|x^{k+1} - x^k\|}{\|x^k\|} < \varepsilon$ then stop

Very easy, but it sometimes does not converge enough.

- $\frac{\|Ax^k - b\|}{\|b\|} < \varepsilon$ then stop

Stable, but computational cost is higher than the above.

2.2 Convergence analysis

The spectral radius of matrices is important for analyzing convergence property of iterative methods.

Definition 2.1. The spectral radius $\rho(A)$ of an $n \times n$ matrix A is defined by

$$\rho(A) := \max_{1 \leq i \leq n} |\lambda_i|, \quad (6)$$

where λ_i is the i -th eigenvalue of A .

Lemma 2.2. For any matrix norm $\|A\|$ induced from by the vector norm, we have

$$\rho(A) \leq \|A\|. \quad (7)$$

Proof. $Ax = \lambda x, x \neq 0$

Therefore

$$|\lambda| \|x\| = \|\lambda x\| = \|Ax\| \leq \|A\| \|x\|$$

So, we have

$$|\lambda| \leq \|A\| \text{ for all eigenvalues } \lambda.$$

□

Both Jacobi and Gauss-Seidel method can be written in the general form

$$x^{k+1} = Mx^k + c, \quad k = 0, 1, 2, \dots \quad (8)$$

for finding a solution of

$$x = Mx + c, \quad (9)$$

where

$$\text{[Jacobi method]} \quad M = -D^{-1}(L + U), \quad c = D^{-1}b,$$

$$\text{[Gauss-Seidel method]} \quad M = -(L + D)^{-1}U, \quad c = (L + D)^{-1}b.$$

Theorem 2.3. Equation (9) has a unique solution and the sequence $\{x^k\}$ defined by (8) converges to the unique solution from any starting point x^0 if and only if $\rho(M) < 1$.

Proof. See, for example, Section 8.2 in [J. Stoer, R. Bulirsch, Introduction to Numerical Analysis, Springer, 2002]. □

Corollary 2.4. If $\|M\| < 1$, then the sequence $\{x^k\}$ defined by (8) converges to the unique solution of (9) from any starting point x^0 .

Proof. This follows from Lemma 2.2 and Theorem 2.3. □

$$\rho(M) \leq \|M\| < 1$$

Definition 2.5. An $n \times n$ matrix A is called **strictly diagonally dominant** if

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}| \quad (10)$$

for all $i = 1, 2, \dots, n$.

$$\begin{pmatrix} 10 & -1 & 2 \\ -2 & 10 & 3 \\ 1 & -5 & 10 \end{pmatrix} \rightarrow \begin{matrix} 10 > 3 & 0 < \\ 10 > 5 & 0 < \\ 10 > 6 & 0 < \end{matrix} \quad \left| \begin{pmatrix} 10 & 1 & 0 & 0 & 0 \\ 0 & 10 & 1 & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 10 \end{pmatrix} \right.$$

Theorem 2.6. If A is strictly diagonally dominant, both Jacobi and Gauss-Seidel methods generate sequences $\{x^k\}$ that converge to unique solution of $Ax = b$ for any starting point x^0 .

Proof. [Jacobi] We want to prove $\|M\|_{\infty} = \|D^{-1}(L+U)\|_{\infty} < 1$.

$$M = \begin{pmatrix} 0 & \frac{a_{12}}{a_{11}} & \dots & \frac{a_{1n}}{a_{11}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{a_{n1}}{a_{nn}} & \dots & \dots & 0 \end{pmatrix}, \quad \|M\|_{\infty} = \max_{1 \leq i \leq n} \sum_{j \neq i}^n \left| \frac{a_{ij}}{a_{ii}} \right| = \max_{1 \leq i \leq n} \left(\frac{1}{|a_{ii}|} \sum_{j \neq i}^n |a_{ij}| \right) < 1$$

[Gauss-Seidel] Ex 6.1

Hind: proving by induction

$$\|M\|_{\infty} = \|(L+D)^{-1}U\|_{\infty} = \max_{\|z\|_{\infty}=1} \|\underbrace{(L+D)^{-1}Uz}_{Mz} \|_{\infty} < 1 \quad \text{we want}$$

Prove, by induction, that

(i) $|z_1| = ? < 1$

(ii) Prove $|z_{i+1}| < 1$ when $|z_i| < 1$

$$\begin{aligned} (x, z) &= x \cdot z \\ &:= x^T A z = (x, z)_A \\ \text{when } A = I \\ (x, z) &= x^T z \end{aligned}$$

Remark 2.7. In fact, if A is a symmetric positive-definite matrix ($A = A^T$ and $x^T A x > 0$ for all $x \neq 0$), Gauss-Seidel method generates sequences $\{x^k\}$ that converge to unique solution of $Ax = b$ for any starting point x^0 .

$$\begin{aligned} Ax &= b \\ \underbrace{A^T A}_{A'} x &= \underbrace{A^T b}_{b'}, \quad \underbrace{A'}_{\text{symmetric positive-definite}} x = b' \end{aligned} \quad \left| \quad \text{cond}(A') = \text{cond}(A)^2 \right.$$

3 Nonstationary iterative methods — Conjugate gradient method (CG method)

One of the top 10 algorithms in the 20th century for solving linear equation $Ax = b$, where $b \in \mathbb{R}^n$ and $A \in \mathbb{R}^n \times \mathbb{R}^n$ is a symmetric positive-definite matrix ($A = A^T$ and $x^T Ax > 0$ for all $x \neq 0$).

Algorithm

Compute $p_0 = r_0 = b - Ax_0$.

For $k = 0, 1, 2, \dots$

If r_k is small enough, break.

$$\alpha_k = \frac{(r_k, p_k)}{(Ap_k, p_k)} \left(= \frac{\|r_k\|_2^2}{(Ap_k, p_k)} \right)$$

$$x_{k+1} = x_k + \alpha_k p_k$$

$$r_{k+1} = r_k - \alpha_k Ap_k$$

$$\beta_k = -\frac{(Ap_k, r_{k+1})}{(Ap_k, p_k)} \left(= \frac{\|r_{k+1}\|_2^2}{\|r_k\|_2^2} \right)$$

$$p_{k+1} = r_{k+1} + \beta_k p_k$$

End

Theorem 3.1. *The sequence generated by the conjugate gradient method converges to the solution of $Ax = b$ at most n steps.*

Theorem 3.2. *Let*

$$\kappa := \frac{\sqrt{\rho} - 1}{\sqrt{\rho} + 1}, \quad \rho = \text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2$$

Then, we have

$$\|x^* - x_k\|_2 \leq \frac{2\rho\kappa^k}{1 + \kappa^{2k}} \|x^* - x_0\|_2, \quad k = 0, 1, 2, \dots, n.$$

3.1 Preconditioned conjugate gradient method (PCG method)

Theorem 3.3. *If A is a symmetric positive-definite matrix, then there exists a unique real lower triangular matrix L with positive diagonals such that*

$$A = LL^T = \begin{pmatrix} \times & & \\ & \times & \\ & & \times \end{pmatrix} \begin{pmatrix} \times & & \\ & \times & \\ & & \times \end{pmatrix} \quad (11)$$

This factorization is called Cholesky's factorization or Cholesky's decomposition.

$$Ax = b, \quad \text{cond}(A) \gg 1 \text{ very large}$$

$$\downarrow \underbrace{(L^{-1} A L^{-T})}_{A'} \underbrace{(L^T x)}_{x'} = \underbrace{L^{-1} b}_{b'}$$

$$A' x' = b' \quad (\Leftrightarrow Ax = b)$$

$$A' = L^{-1} A L^{-T} = L^{-1} L L^T L^{-T} = I$$

$$A = L L^T + \underset{\text{error}}{E} \Rightarrow A' \sim I$$

Incomplete cholesky's factorization
 $A \sim L L^T$ to avoid fill-in