

Lecture 3 Iterative methods for solving linear system

Weinan E^{1,2} and Tiejun Li²

¹Department of Mathematics,
Princeton University,
weinan@princeton.edu

²School of Mathematical Sciences,
Peking University,
tieli@pku.edu.cn
No.1 Science Building, 1575

Outline

Iterations

- ▶ Iterative methods

Object: construct sequence $\{\mathbf{x}^k\}_{k=1}^{\infty}$, such that \mathbf{x}^k converge to a fixed vector \mathbf{x}^* , and \mathbf{x}^* is the solution of the linear system.

- ▶ General iteration idea:

If we want to solve equations

$$\mathbf{g}(\mathbf{x}) = \mathbf{0},$$

and the equation $\mathbf{x} = \mathbf{f}(\mathbf{x})$ has the same solution as it, then construct

$$\mathbf{x}^{k+1} = \mathbf{f}(\mathbf{x}^k).$$

If $\mathbf{x}^k \rightarrow \mathbf{x}^*$, then $\mathbf{x}^* = \mathbf{f}(\mathbf{x}^*)$, thus the root of $\mathbf{g}(\mathbf{x})$ is obtained.

Matrix form for Jacobi iterations

- ▶ Decompose $A = D - L - U$, where

$$L = \begin{pmatrix} 0 & & & \\ -a_{21} & 0 & & \\ \vdots & \vdots & \ddots & \\ -a_{n1} & -a_{n2} & \cdots & 0 \end{pmatrix}, \quad U = \begin{pmatrix} 0 & -a_{12} & \cdots & -a_{1n} \\ & 0 & \cdots & -a_{2n} \\ & & \ddots & \vdots \\ & & & 0 \end{pmatrix}$$

and $D = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$.

Transformation: $Dx = b + (L + U)x$, thus

$$x = D^{-1}(b + (L + U)x) := Bx + g.$$

Here $B = D^{-1}(L + U)$ and $g = D^{-1}b$.

- ▶ Define the iterations

$$x^{k+1} = Bx^k + g$$

This is Jacobi iteration.

Component form of Jacobi iterations

- ▶ Jacobi iterations

$$x_1^{k+1} = \left(b_1 - \sum_{j \neq 1} a_{1j} x_j^k \right) / a_{11}$$

$$x_2^{k+1} = \left(b_2 - \sum_{j \neq 2} a_{2j} x_j^k \right) / a_{22}$$

... ..

$$x_n^{k+1} = \left(b_n - \sum_{j \neq n} a_{nj} x_j^k \right) / a_{nn}$$

- ▶ Update each component of x^k according to each equation simultaneously and independently by freezing the other variables at former step.
- ▶ Parallel in nature!

Jacobi iterations

- ▶ Example 1: Solving the linear system $Ax = b$ with

$$A = \begin{pmatrix} 1 & 2 & -2 \\ 1 & 1 & 1 \\ 2 & 2 & 1 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}.$$

- ▶ Example 2: Solving the linear system $Ax = b$ with

$$A = \begin{pmatrix} 2 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & -2 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}.$$

Jacobi iterations

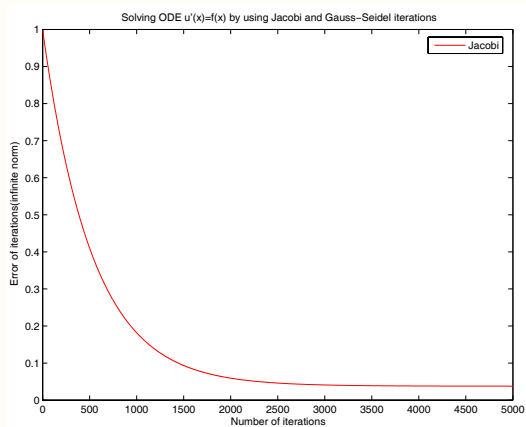
- ▶ Example 3: Solving the system $\mathbf{Ax} = \mathbf{b}$ ($n = 30$)

$$\mathbf{A} = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 \end{pmatrix}$$

and \mathbf{b} is discretized from $-\pi^2 \sin \pi x$, $x \in [0, 1]$.

Jacobi iterations

- ▶ Convergence rate for example 3: (error vs. iteration times)



Gauss-Seidel iterations

- **Idea:** Update each component of $x^{(k)}$ sequentially with the most updated information and freeze the left components.

$$\begin{aligned}x_1^{k+1} &= \left(b_1 - \sum_{j>1} a_{1j}x_j^k\right)/a_{11} \\x_2^{k+1} &= \left(b_2 - \sum_{j<2} a_{2j}x_j^{k+1} - \sum_{j>2} a_{2j}x_j^k\right)/a_{22} \\&\dots \quad \dots \quad \dots \\x_n^{k+1} &= \left(b_n - \sum_{j<n} a_{nj}x_j^{k+1}\right)/a_{nn}\end{aligned}$$

- NOT parallel directly.

Matrix form for Gauss-Seidel iteration

- ▶ Decompose $A = D - L - U$, and perform transformation $(D - L)x = b + Ux$, thus

$$x = (D - L)^{-1}(b + Ux) := Bx + g$$

Here $B = (D - L)^{-1}U$ and $g = (D - L)^{-1}b$.

- ▶ Define the iterations

$$x^{k+1} = D^{-1}Lx^{k+1} + D^{-1}Ux^k + D^{-1}b$$

This is Gauss-Seidel iteration.

Gauss-Seidel iterations

- ▶ Example 1: Solving the linear system $Ax = b$ with

$$A = \begin{pmatrix} 1 & 2 & -2 \\ 1 & 1 & 1 \\ 2 & 2 & 1 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}.$$

- ▶ Example 2: Solving the linear system $Ax = b$ with

$$A = \begin{pmatrix} 2 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & -2 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}.$$

Gauss-Seidel iterations

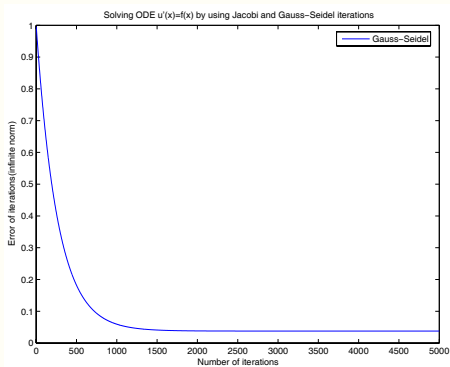
- ▶ Example 3: Solving the system $\mathbf{Ax} = \mathbf{b}$ ($n = 30$)

$$\mathbf{A} = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix}$$

and \mathbf{b} is discretized from $-\pi^2 \sin \pi x$, $x \in [0, 1]$.

Gauss-Seidel iterations

- ▶ Convergence rate for example 3: (error vs. iteration times)



Successive relaxation and SOR

- ▶ Gauss-Seidel

$$\mathbf{x}^{k+1} = D^{-1}L\mathbf{x}^{k+1} + D^{-1}U\mathbf{x}^k + D^{-1}\mathbf{b}.$$

Define

$$\Delta\mathbf{x} = \mathbf{x}^{k+1} - \mathbf{x}^k$$

then $\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta\mathbf{x}$.

- ▶ Successive relaxation is to add a weight ω

$$\begin{aligned}\mathbf{x}^{k+1} &= \mathbf{x}^k + \omega\Delta\mathbf{x} \\ &= (1 - \omega)\mathbf{x}^k + \omega(D^{-1}L\mathbf{x}^{k+1} + D^{-1}U\mathbf{x}^k + D^{-1}\mathbf{b}).\end{aligned}$$

If $0 < \omega < 1$, it is called under-relaxation; if $1 < \omega < 2$, it is called over-relaxation;

Convergence theorem

Theorem (Necessary condition)

The necessary condition for the convergence of SR method is

$$0 < \omega < 2$$

Theorem (SPD matrix)

If A is symmetric positive definite, then

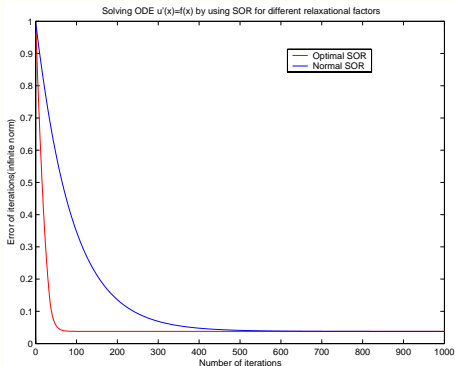
$$0 < \omega < 2$$

is the sufficient condition for convergence.

Examples

SOR method for the ODE problem ($n=30$).

- ▶ $\omega = \frac{2}{1 + \sqrt{1 - (\rho(J))^2}}$, $\mathbf{x}^0 = \mathbf{0}$
- ▶ $\omega = \text{another one}$, $\mathbf{x}^0 = \mathbf{0}$



Outline

Puzzle?

- ▶ Arbitrary decomposition (F is invertible)

$$A = F - G$$

then the linear system is transformed into

$$x = F^{-1}Gx + F^{-1}b.$$

and construct the iteration

$$x^{k+1} = F^{-1}Gx^k + F^{-1}b.$$

How about the result?

- ▶ A simpler case (1D case) $ax = b$, $a = f - g$ then construct

$$x^{k+1} = f^{-1}gx^k + f^{-1}b.$$

How about the convergence?

Analysis for simple case

- ▶ Linear system $ax = b$

Exact solution $x^* = f^{-1}gx^* + f^{-1}b$

Iteration $x^{k+1} = f^{-1}gx^k + f^{-1}b.$

- ▶ Subtract two expressions and define error $e^k = x^k - x^*$ thus

$$e^{k+1} = f^{-1}ge^k$$

Then $|e^k| = |f^{-1}g|^k |e^0|$. In order $|e^k| \rightarrow 0$, we must have

$$|f^{-1}g| < 1$$

That's the convergence condition for 1D case.

A fundamental theorem

- ▶ Spectral radius

The spectral radius of a matrix \mathbf{A} is defined as

$$\rho(\mathbf{A}) = \max_{\lambda} |\lambda(\mathbf{A})|$$

Theorem (Spectral radius and iterative convergence)

If $\mathbf{A} \in \mathbb{R}^n$, then

$$\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{0} \iff \rho(\mathbf{A}) < 1$$

Analysis for general case

- ▶ Linear system $Ax = b$

Exact solution $x^* = Bx^* + g$

Iteration $x^{k+1} = Bx^k + g$

- ▶ Subtract two expressions and define error $e^k = x^k - x^*$ thus

$$e^{k+1} = Be^k$$

By induction we have

$$e^k = B^k e^0.$$

Thus we have

$$e^k \rightarrow 0 \iff B^k \rightarrow 0 \iff \rho(B) < 1$$

That's the convergence criterion for general iterations.

Examples

- ▶ Analysis for example 1.
Spectral radius for Jacobi:
Spectral radius for Gauss-Seidel:
- ▶ Analysis for example 2.
Spectral radius for Jacobi:
Spectral radius for Gauss-Seidel:

Rate of convergence

Theorem (Rate of convergence)

If $\|B\| = q < 1$, then the iterated solution has the following convergence rate:

$$\|e^k\| \leq \frac{q^k}{1-q} \|\mathbf{x}^1 - \mathbf{x}^0\|$$

Loosely speaking, if we ask a tolerance ε_{tol} and $\|e^k\| \leq \varepsilon_{tol}$, we need the number of iterations

$$k \geq \frac{\ln(\varepsilon_{tol}/C)}{\ln q}.$$

where $C = \frac{\|\mathbf{x}^1 - \mathbf{x}^0\|}{1-q}$.

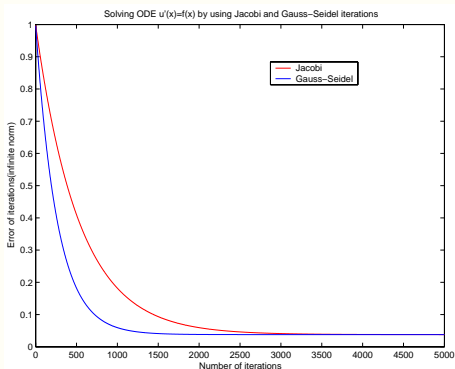
Remark (Linear convergence)

This type of convergence is called linear convergence, which means

$$\lim_{k \rightarrow \infty} \frac{\|e^{k+1}\|}{\|e^k\|} = C \neq 0.$$

Examples

- ▶ Convergence rate for example 3
Spectral radius for Jacobi:
Spectral radius for Gauss-Seidel:
- ▶ Comparison of numerical convergence rate:



Diagonally dominant matrix and its properties

Definition (Diagonally dominant matrix (DDM))

Suppose $A \in \mathbb{R}^{n \times n}$, if

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

A is called a DDM.

- ▶ DDM Example

$$A = \begin{pmatrix} 3 & -1 & 1 \\ 1 & 3 & 1 \\ 0 & 1 & -4 \end{pmatrix}$$

- ▶ Example 3 does NOT belong to this class.

Diagonally dominant matrix and its properties

Theorem (Convergence for DDM)

If A is a DDM, A is nonsingular and the Jacobi and Gauss-Seidel method for $Ax = b$ is convergent.

Outline

Symmetric positive definite (SPD) matrix and Quadratic function

- ▶ Linear system

$$\mathbf{Ax} = \mathbf{b}$$

where \mathbf{A} is a SPD matrix.

- ▶ Quadratic function

$$\varphi(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Ax} - \mathbf{b}^T \mathbf{x}$$

- ▶ Connection

$$\text{Solving } \mathbf{Ax} = \mathbf{b} \iff \min \varphi(\mathbf{x})$$

Symmetric positive definite (SPD) matrix and Quadratic function

- ▶ Basic observation

$$\frac{\partial \varphi}{\partial x_i} = \sum_{j=1}^n a_{ij} x_j - x_i b_i$$

Define gradient vector

$$\text{grad} \varphi = \left(\frac{\partial \varphi}{\partial x_i} \right)_{i=1}^n = \mathbf{Ax} - \mathbf{b}$$

The minimizer \mathbf{x}^* of φ satisfies

$$\text{grad} \varphi = 0, \text{ i.e. } \mathbf{Ax}^* = \mathbf{b}$$

- ▶ The necessity is skipped.

Symmetric positive definite (SPD) matrix and Quadratic function

The connection between linear system and quadratic function minimization tells us **if we have an algorithm to deal with quadratic function minimization we have an algorithm for solving the linear system.**

Steepest decent for quadratic minimization

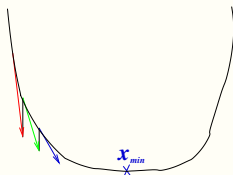
- ▶ Basic idea: Find a series of decent directions \mathbf{p}^k and corresponding stepsize α_k such that the iterations

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{p}^k$$

and

$$\varphi(\mathbf{x}^{k+1}) \leq \varphi(\mathbf{x}^k).$$

- ▶ Schematics for Steepest decent



- ▶ How to find \mathbf{p}^k and α_k ?

Steepest decent method (SDM) for quadratic minimization

- ▶ The negative gradient direction $-\text{grad}\varphi$ is the “steepest” decent direction, so choose

$$\mathbf{p}^k = \mathbf{r}^k := -\text{grad}\varphi(\mathbf{x}^k)$$

and choose α_k such that

$$\min_{\alpha} \varphi(\mathbf{x}^k + \alpha\mathbf{p}^k)$$

- ▶ Define $f(\alpha) = \varphi(\mathbf{x}^k + \alpha\mathbf{p}^k)$ then

$$\frac{\partial f(\alpha)}{\partial \alpha} = 0.$$

We obtain

$$\alpha_k = \frac{(\mathbf{r}^k)^T \mathbf{p}^k}{(\mathbf{p}^k)^T \mathbf{A} \mathbf{p}^k}$$

where $\mathbf{r}^k = -\text{grad}\varphi(\mathbf{x}^k) = \mathbf{b} - \mathbf{A}\mathbf{x}^k$.

- ▶ This is the so called steepest decent method.

Example

- ▶ Example 1:

$$\mathbf{A} = \begin{pmatrix} 4 & 1 & 1 & 0 \\ 1 & 4 & 1 & 1 \\ 1 & 1 & 4 & 1 \\ 0 & 1 & 1 & 4 \end{pmatrix} \quad \mathbf{b} = \begin{pmatrix} 6 \\ 7 \\ 7 \\ 6 \end{pmatrix}$$

- ▶ The ODE example ($n=30$):
Convergence rate.

Convergence of steepest decent method

Theorem (Convergence of steepest decent method)

Suppose the eigenvalues of \mathbf{A} are $0 < \lambda_1 \leq \dots \leq \lambda_n$, then the iterating sequence $\{\mathbf{x}^k\}_{k=1}^{\infty}$ by steepest decent method has the following convergence property

$$\|\mathbf{x}^k - \mathbf{x}^*\|_{\mathbf{A}} \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^k \|\mathbf{x}^0 - \mathbf{x}^*\|_{\mathbf{A}}$$

where $\mathbf{x}^* = \mathbf{A}^{-1}\mathbf{b}$, $\|\mathbf{x}\|_{\mathbf{A}} := \sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}}$.

Remarks on steepest decent method

Remark

1. Consider $\mathbf{p}^k = -\text{grad}\varphi(\mathbf{x}^k)$, $\mathbf{p}^{k+1} = -\text{grad}\varphi(\mathbf{x}^{k+1})$, $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{p}^k$, we have

$$\left. \frac{d}{d\alpha} \varphi(\mathbf{x}^k + \alpha \mathbf{p}^k) \right|_{\alpha=\alpha_k} = 0$$

i.e.

$$(\mathbf{p}^{k+1})^T \mathbf{p}^k = 0$$

the decent directions are orthogonal for two neighboring steps. This means steepest decent is locally steepest decent, not globally steepest decent!

This effect becomes quite severe when \mathbf{x}^k approaches \mathbf{x}^ .*

2. If $\lambda_n \gg \lambda_1$, then $\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \approx 1$, the convergence rate will be very slow!
3. Though the steepest decent method is easily implemented, it is very rarely used for the drawback stated above.
4. It can be slightly modified to obtain the famous **conjugate gradient method**.

Conjugate gradient method (CGM)

- ▶ Only modify SDM a little bit to obtain CGM.
- ▶ Suppose we have $\mathbf{x}^k, \mathbf{p}^{k-1}, \mathbf{r}^k$. The strategy of CGM is not to choose $\mathbf{p}^k = \mathbf{r}^k$, but choose \mathbf{p}^k in 2D plane

$$\pi_2 := \{\xi \mathbf{r}^k + \eta \mathbf{p}^{k-1}, \xi, \eta \in \mathbb{R}\}$$

Define $\psi(\xi, \eta) = \varphi(\mathbf{x}_k + \xi \mathbf{r}^k + \eta \mathbf{p}^{k-1})$, then take

$$\frac{\partial \psi}{\partial \xi} = 0, \quad \frac{\partial \psi}{\partial \eta} = 0$$

then

$$\begin{cases} \xi_0 (\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k + \eta_0 (\mathbf{r}^k)^T \mathbf{A} \mathbf{p}^{k-1} & = (\mathbf{r}^k)^T \mathbf{r}^k \\ \xi_0 (\mathbf{r}^k)^T \mathbf{A} \mathbf{p}^{k-1} + \eta_0 (\mathbf{p}^{k-1})^T \mathbf{A} \mathbf{p}^{k-1} & = \mathbf{0} \end{cases}$$

i.e. we can take

$$\mathbf{p}^k = \mathbf{r}_k + \frac{\eta_0}{\xi_0} \mathbf{p}^{k-1}$$

- ▶ Same as SDM to obtain α_k and \mathbf{x}^{k+1} .

Conjugate gradient method (CGM)

Final formulation of CGM:

1. Initial step: $\mathbf{x}^0, \mathbf{p}^0 = \mathbf{r}^0 = \mathbf{b} - \mathbf{A}\mathbf{x}^0$
2. Suppose we have $\mathbf{x}^k, \mathbf{r}^k, \mathbf{p}^k$, the CGM step

2.1 Search the optimal α_k along \mathbf{p}^k ;

$$\alpha_k = \frac{(\mathbf{r}^k)^T \mathbf{p}^k}{(\mathbf{p}^k)^T \mathbf{A} \mathbf{p}^k}$$

2.2 Update \mathbf{x}^k and gradient direction \mathbf{r}^k ;

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{p}^k, \quad \mathbf{r}^{k+1} = \mathbf{b} - \mathbf{A} \mathbf{x}^{k+1}$$

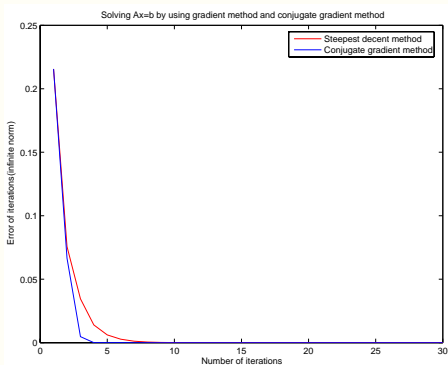
2.3 According to the calculation before to form new search direction \mathbf{p}^{k+1}

$$\beta_k = -\frac{(\mathbf{r}^{k+1})^T \mathbf{A} \mathbf{p}^k}{(\mathbf{p}^k)^T \mathbf{A} \mathbf{p}^k}, \quad \mathbf{p}^{k+1} = \mathbf{r}^{k+1} + \beta_k \mathbf{p}^k$$

Example

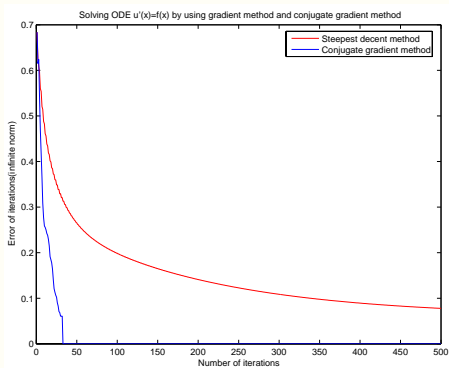
► Example 1:

$$\mathbf{A} = \begin{pmatrix} 4 & 1 & 1 & 0 \\ 1 & 4 & 1 & 1 \\ 1 & 1 & 4 & 1 \\ 0 & 1 & 1 & 4 \end{pmatrix} \quad \mathbf{b} = \begin{pmatrix} 6 \\ 7 \\ 7 \\ 6 \end{pmatrix}$$



Example

- ▶ The ODE example ($n=30$):



Properties of CGM

Theorem

The vectors generated by CGM have the properties

1. $(\mathbf{r}^i)^T \mathbf{r}^j = 0, \quad i \neq j, \quad 0 \leq i < j \leq k;$

2. $(\mathbf{p}^i)^T \mathbf{A} \mathbf{p}^j = 0, \quad i \neq j, \quad 0 \leq i < j \leq k;$

- ▶ The property $(\mathbf{p}^i)^T \mathbf{A} \mathbf{p}^j = 0$ makes $\{\mathbf{p}^j\}_{j=1,2,\dots}$ are called conjugate directions of \mathbf{A} . From property 1, we have

Theoretically, CGM will find the minimum (or solution) within n steps.

This is not the case in real computations. CGM is applied as an iteration.

- ▶ **Virtue of CGM:** Suitable for sparse matrix \mathbf{A} , without adjustable parameters as SOR.

Convergence of CGM

Theorem

The sequence $\{\mathbf{x}^k\}$ has the error estimate

$$\|\mathbf{x}^k - \mathbf{x}^*\|_{\mathbf{A}} \leq 2 \left(\frac{\sqrt{\kappa_2} - 1}{\sqrt{\kappa_2} + 1} \right) \|\mathbf{x}^0 - \mathbf{x}^*\|_{\mathbf{A}}$$

where $\kappa_2 := \text{Cond}_2(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2$.

This theorem tells us if the matrix is in good condition ($\text{Cond}_2(\mathbf{A}) \approx 1$), the convergence is very fast. The real computations are even faster than the estimate.

Preconditioning conjugate gradient method (PCG)

- ▶ If the conditioning of \mathbf{A} is bad, we need the so-called preconditioning technique.
- ▶ Formally transform $\mathbf{A}\mathbf{x} = \mathbf{b}$ into

$$\hat{\mathbf{A}}\hat{\mathbf{x}} = \hat{\mathbf{b}}$$

where $\hat{\mathbf{A}} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}^{-1}$, $\hat{\mathbf{x}} = \mathbf{P}\mathbf{x}$, $\hat{\mathbf{b}} = \mathbf{P}^{-1}\mathbf{b}$, \mathbf{P} is a SPD matrix.

- ▶ Apply CGM to the modified system to obtain $\hat{\mathbf{x}}^k, \hat{\mathbf{r}}^k, \hat{\mathbf{p}}^k$.
- ▶ The matrix $\mathbf{M} := \mathbf{P}^2$ is called preconditioner. Usually \mathbf{M} is chosen a sparse SPD matrix, and the eigenvalues of $\mathbf{M}^{-1}\mathbf{A}$ concentrate on some value, and the equation $\mathbf{M}\mathbf{y} = \mathbf{r}$ is easy to be solved.

Extension to general linear system

- ▶ Technique 1: Regularization. Transform the linear system

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A} \mathbf{b}$$

If the conditioning of \mathbf{A} is bad, the convergence is very slow.

- ▶ Generalized Minimum RESidual method (GMRES):
Compute the minimal residual problem

$$\min\{\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2, \mathbf{x} \in \mathbf{x}_0 + \text{some subspace}\}$$

This formulation is well-posed for non-symmetric matrix.

- ▶ References:

Y. Saad, Iterative methods for sparse linear systems, PWS Publishing Company, 1996.

Homework assignment 3

1. Using Gauss-Seidel and conjugate gradient method to solve the second order ODEs($n=30, 50, 100$). Plot the figure for the convergence rate.