

## [6] Indirect solvers : iterative methods for sparse linear systems

→ We are going to look at a couple of simple & basic methods first, including:

- (i) Jacobi,
- (ii) Gauss-Seidel,
- (iii) Successive overrelaxation (SOR),

and later we also study more effective ones:

- (iv) Conjugate gradient (CG),
- (v) descent method,
- (vi) preconditioner
- (v) multigrid methods, etc.

## §6.1. Stationary iterative methods

→ Before we begin studying iterative methods, let's first take a look at "stationary iterative methods" by splitting the matrix  $\underline{A}$ .

→ Consider we are interested in solving a linear system

$$\boxed{\underline{A}U = f} \quad \dots \textcircled{1}$$

as a result of discretizing the elliptic eqn:

$$U_{xx} + U_{yy} = f.$$

→ The simplest type of iterative method for solving  $\textcircled{1}$  has the form

$$\boxed{U^{(k+1)} = \underline{G}U^{(k)} + c} \quad \dots \textcircled{2}$$

where the matrix  $\underline{G}$  and the vector  $c$  are chosen so that a fixed pt. of the equation

$$\boxed{U = \underline{G}U + c} \quad \dots \textcircled{3}$$

is a soln to  $\textcircled{1}$ ,

→  $\textcircled{3}$  is called stationary if  $\underline{G}$  and  $c$  are constant over all iterations.

→ One way to obtain  $\underline{G}$  is by splitting  $\underline{A}$  into

$$\boxed{\underline{A} = \underline{M} - \underline{N}}, \quad \underline{M}: \text{nonsingular.} \quad \dots \textcircled{4}$$

$$\rightarrow \underline{A}U = f \Leftrightarrow \underline{M}U - \underline{N}U = f$$

$$\Leftrightarrow U - \underline{M}^{-1}\underline{N}U = \underline{M}^{-1}f$$

$\Leftrightarrow$  Letting  $\underline{G} = \underline{M}^{-1}\underline{N}$ , &  $c = \underline{M}^{-1}f$ , we get the iteration scheme as

$$U^{(k+1)} = \underline{G}U^{(k)} + c; \quad \dots \textcircled{2}$$

which can be implemented as

$$\boxed{\underline{M}U^{(k+1)} = \underline{N}U^{(k)} + f} \quad \dots \textcircled{3}$$

$\rightarrow$  In  $\textcircled{3}$  we solve a linear system with  $\underline{M}$  at each  $k$ , instead of solving the original system with  $\underline{A}$ .

$\rightarrow$  From  $\textcircled{2}$ , we see that the iteration scheme converges if

$$\boxed{\rho(\underline{G}) = \rho(\underline{M}^{-1}\underline{N}) < 1}, \text{ where}$$

$$\rho(\underline{G}) = \max \{ |\lambda| \mid \lambda : \text{eigenvalues of } \underline{G} \}$$

is the spectral radius.

$\rightarrow$  The smaller the  $\rho(\underline{G})$ , the faster the convergence. For rapid convergence, we should choose  $\underline{M}$  &  $\underline{N}$  s.t.  $\rho(\underline{M}^{-1}\underline{N})$  is as small as possible.

$\rightarrow$  Note that the cost of solving  $\textcircled{3}$  is determined by the linear system with  $\underline{M}$ .

→ As an extreme case, if  $\underline{M} = \underline{A}$ , then the scheme converges in a single iteration,

→ In practice,  $\underline{M}$  is chosen to approximate  $\underline{A}$  in some sense, but is usually constrained to have a simple form, such as diagonal or triangular, so that the linear system (5) is easy to solve.

## Ex 2.1. Jacobi Method.

→ The simplest choice for  $\underline{M} = \text{diag}(A)$ ,  $\underline{A} = \underline{M} - \underline{N}$ .

Let's have the followings:

(i)  $\underline{D} \equiv \text{diag}(A)$ ,

(ii)  $\underline{L} \equiv$  strict lower matrix of  $A$ ,

(iii)  $\underline{U} \equiv$  strict upper matrix of  $A$ ,

then

(iv)  $\underline{M} = \underline{D}$  --- ①

(v)  $A = \underline{M} - \underline{N} = \underline{D} - \underline{N}$

$\searrow$   
 $\underline{L} + \underline{D} + \underline{U}$

$\Rightarrow \underline{N} = -(\underline{L} + \underline{U})$  --- ②

→ ① & ② give the splitting of  $A$ .

→  $\underline{D}$  is invertible if  $A$  has no zero entries in diagonal.

→ In this case, we obtain the iterative scheme, known as the Jacobi method:

$$\underline{U}^{(k+1)} = -\underline{D}^{-1}(\underline{L} + \underline{U}) \underline{U}^{(k)} + \underline{D}^{-1} \underline{f} \quad \text{--- ③}$$

→ To write ③ in the componentwise form, we revisit Eqn ⑤ (5a ~ 5e) in ③ ordering the unknowns & eqns.

→ We reorganize (5) into :

$$(5b) (5c) = (5e) - ((5a) + (5c)), \text{ or}$$

$$\underbrace{\begin{bmatrix} -4 & 1 & & \\ 1 & -4 & & \\ & & \ddots & \\ & & & -4 \end{bmatrix}}_{(= \tilde{A})} \underbrace{\begin{bmatrix} U_{1,j} \\ U_{2,j} \\ \vdots \\ U_{i,j} \\ \vdots \\ U_{m,j} \end{bmatrix}}_{(= U_j)} = \frac{1}{h^2} \underbrace{\begin{bmatrix} f_{1,j} \\ f_{2,j} \\ \vdots \\ f_{i,j} \\ \vdots \\ f_{m,j} \end{bmatrix}}_{(= f_j)} - \underbrace{\begin{bmatrix} U_{1,j-1} + U_{1,j+1} \\ U_{2,j-1} + U_{2,j+1} \\ \vdots \\ U_{i,j-1} + U_{i,j+1} \\ \vdots \\ U_{m,j-1} + U_{m,j+1} \end{bmatrix}}_{(= f_j)} \leftarrow (i\text{th})$$

$$\rightarrow U_{ij}^{(k+1)} = \frac{1}{4} \left( U_{i-1,j}^{(k)} + U_{i+1,j}^{(k)} + U_{i,j-1}^{(k)} + U_{i,j+1}^{(k)} \right) - \frac{h^2}{4} f_{ij} \quad \text{--- (4)}$$

→ (4) can be thought that the new approximate  $U_{ij}^{(k+1)}$  is simply the average of the previous soln components at the four neighboring grid pts.

Rank

① The Jacobi method does NOT always converge, but it is guaranteed to converge under conditions that are often satisfied in practice (e.g., if the matrix is diagonally dominant by rows,

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|, \quad \forall i$$

② Unfortunately, the convergence rate of the Jacobi method is in general unacceptably slow.

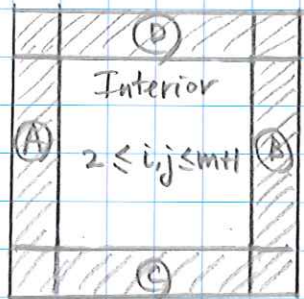
Rule. A pseudo-code for Jacobi looks like:

```

for iter = 0 : maxIter
  for j = 2 : m+1
    for i = 2 : m+1
      
$$u_{\text{new}}(i,j) = 0.25 * (u(i-1,j) + u(i+1,j) + u(i,j-1) + u(i,j+1)) - \frac{h^2}{4} * f(i,j)$$

    end for
  end for
  u = unew
end for
  
```

→ BCs are stored in those four GCs band regions



GCs

- (A) →  $u(1, :)$
- (B) →  $u(m+2, :)$
- (C) →  $u(:, 1)$
- (D) →  $u(:, m+2)$

→ Two arrays  $u$  &  $u_{\text{new}}$  need to be used.

### §6.3. Gauss-Seidel Method.

→ There may be some who wish to save the u<sup>new</sup> array space from the Jacobi method, and try the following:

```
for iter = 0 : maxIter
  for j = 2 : m+1
    for i = 2 : m+1
      u(i,j) = 0.25 * (u(i-1,j) + u(i+1,j) + u(i,j-1) + u(i,j+1))
              - h^2 * f(i,j)
    end for
  end for
end for
```

→ This will give a different result though, which is known as the Gauss-Seidel method, given as

$$u_{ij}^{(k+1)} = \frac{1}{4} \left( u_{i-1,j}^{(k+1)} + u_{i+1,j}^{(k)} + u_{i,j-1}^{(k+1)} + u_{i,j+1}^{(k)} \right) - \frac{h^2}{4} f_{ij} \quad \textcircled{1}$$

→ This can be considered as a lucky coding error!  
However, the method does work!

→ The major difference from the Jacobi method is to use the already updated data at  $(i-1,j)$  &  $(i,j-1)$  in order to update the data at  $(i,j)$ .



→ Mathematically, this idea of using the latest updated information to update the new grid data makes the Gauss-Seidel scheme faster in convergence than the Jacobi method.

→ In general, one can write the GS method in matrix terms using the splitting of the form:

$$\begin{cases} \underline{M} = \underline{D} + \underline{L} \\ \underline{N} = -\underline{U} \end{cases}, \quad \begin{array}{l} \leftarrow \text{updated information} \\ \end{array}$$

which corresponds to the matrix form:

$$\begin{aligned} \underline{U}^{(k+1)} &= (\underline{D} + \underline{L})^{-1} (-\underline{U} \underline{U}^{(k)}) + (\underline{D} + \underline{L})^{-1} \underline{f} \\ &= (\underline{D} + \underline{L})^{-1} (\underline{f} - \underline{U} \underline{U}^{(k)}). \end{aligned}$$

Rank ① GS is faster than Jacobi in convergence.

② GS only requires one vector array  $\underline{u}^{(k)}$ .

③ GS needs the updates successively, whereas Jacobi can update its soln in any order.

Hence, Jacobi is easier for parallelization.

④ GS does NOT always converge, but can converge under a weaker condition than Jacobi (e.g., if the matrix is SPD), but still too slow.

SPD: symmetric positive definite

## §6.4. Successive Overrelaxation.

→ The convergence rate of GS can be accelerated by a technique called "successive overrelaxation (SOR)".

→ The idea is based on the fact that

- (i) GS update moves  $U_{ij}^{(k)}$  in the right direction, but
- (ii) the direction is too conservative

→ So let's be bit more aggressive by pushing the GS update bit further (i.e., overrelaxing, or extrapolation based on the GS update).

→ Step 1:  $U_{ij}^{(k+1), GS} = \frac{1}{4}(U_{ij}^{(k+1)} + U_{i+1,j}^{(k)} + U_{i,j-1}^{(k+1)} + U_{i,j+1}^{(k)}) - \frac{h^2}{4}f_{ij}$   
(GS update)

Step 2:  $U_{ij}^{(k+1)} = U_{ij}^{(k)} + \omega(U_{ij}^{(k+1), GS} - U_{ij}^{(k)})$   
(relaxation)

→ Note that Step 2 is equivalent to

$$U_{ij}^{(k+1)} = (1-\omega)U_{ij}^{(k)} + \omega U_{ij}^{(k+1), GS} \quad \dots \textcircled{1}$$

→  $\textcircled{1}$  is then a weighted average of

- (i) the current  $U_{ij}^{(k)}$ , &
- (ii) the updated GS soln  $U_{ij}^{(k+1), GS}$ .

→ We note that from  $\Phi$ ,

(i)  $\omega < 1 \Rightarrow$  interpolation, or under-relaxation,

(ii)  $\omega > 1 \Rightarrow$  extrapolation, or over-relaxation.

(iii)  $\omega = 1 \Rightarrow$  GS method.

(iv)  $\omega = 2 \Rightarrow$  Not converge!

→ We always have  $0 < \omega < 2$  for convergence.

→ Choosing a specific value of  $\omega$  to attain the best possible convergence rate is a difficult problem in general.

→ One also can write SOR in the matrix form using

$$(i) \underline{M} = \frac{1}{\omega} \underline{D} + \underline{L},$$

$$(ii) \underline{N} = \left(\frac{1}{\omega} - 1\right) \underline{D} - \underline{U},$$

resulting in the corresponding form of:

$$\underline{U}^{(k+1)} = (\underline{D} + \omega \underline{L})^{-1} \left[ (1-\omega) \underline{D} - \omega \underline{U} \right] \underline{U}^{(k)} + \omega (\underline{D} + \omega \underline{L})^{-1} \underline{f},$$