

Classic Iterative Methods. *

① Description of the methods.

The classic iterative methods to solve the matrix equation $Ax = b$, use a splitting of the matrix A , namely,

$$A = N - P$$

and use it to generate a sequence $\{x^n\}_{n \geq 0}$ that should converge to the solution x . The sequence is generated as follows:

$$(I) \quad \left\{ \begin{array}{l} \textcircled{1} \text{ Pick } x^0 = x_0. \\ \textcircled{2} \text{ If } r^k = Ax^k - b = 0 \text{ stop, otherwise compute } x^{k+1} \text{ by solving} \\ N x^{k+1} = P x^k + b. \end{array} \right.$$

Note that if $e^k = x - x^k$, we get that the error equation is

$$N e^{k+1} = P e^k \quad k \geq 0.$$

The matrix $M = N^{-1}P$ is called the iteration matrix. The convergence of the method (I) depends only on properties of M .

* From: Introd. à l'analyse numérique et à l'optimisation, T.G. Cioclet

the classic iterative methods are described in the table below; we use the fact that D is the diagonal part of A , U is the upper part and L the lower part

Name	iteration matrix
Jacobi	$-D^{-1}(A - D)$
Gauss-Seidel	$-(D + L)^{-1}U$
SOR	$(D + \omega L)^{-1}(D - \omega(A - L)) = M_{SOR}(L)$
SSOR	$M_{SOR}(U) M_{SOR}(L)$

Note that:

$$\textcircled{1} \quad M_{SOR}(L) \Big|_{\omega=1} = M_{GS}.$$

\textcircled{2} the SOR can be rewritten as:

$$\begin{cases} D \tilde{x}^{k+1} = -L x^{k+1} - U x + b \\ x^{k+1} = (1-\omega) x^k + \omega \tilde{x}^{k+1} \end{cases}$$

\textcircled{3} M_{SSOR} is symmetric if A is symmetric.
 M_{SOR} is not symmetric even if A is.

② Strictly diagonal dominant matrices.

A matrix A is said to be strictly diagonal dominant if

$$\rho = \max_{1 \leq i \leq N} \sum_{\substack{j=1 \\ j \neq i}}^N \frac{|a_{ij}|}{|a_{ii}|} < 1.$$

For these matrices, the error of both the Jacobi and the Gauss-Seidel methods satisfy

$$\|e^n\|_{\ell^\infty} \leq \rho^n \|e^0\|_{\ell^\infty}.$$

③ Symmetric and positive definite matrices.

For these matrices, we have that

(i) the Jacobi method converges if $a_{ij} \leq 0$ for $i \neq j$.

(ii) the SOR method converges if $\omega \in (0, 2)$.
It diverges if $\omega \in (-\infty, 0) \cup (2, \infty)$.

To show this result, we prove the following theorem.

Theorem. Let A be symmetric and positive definite. Suppose that

$$Q = N + N^T - A$$

is positive definite. Then $\rho(M) < 1$.

Proof. Let λ be an eigenvalue of M and let u be its eigenvector. Then

$$\begin{aligned} Mu &= \lambda u \\ \Rightarrow (I - N^T A)u &= \lambda u \\ \Rightarrow Au &= (1-\lambda)Nu \end{aligned}$$

Since A is positive definite $Au \neq 0$ and hence $\lambda \neq 1$.

Next, let us exploit the fact that Q is positive definite:

$$\begin{aligned} 0 &< \overline{u^T Q u} \\ &= \overline{u^T (N + N^T - A)u} \\ &= \overline{u^T Nu} + \overline{u^T N^T u} - \overline{u^T Au} \\ &= \frac{\overline{u^T Au}}{1-\lambda} + \frac{\overline{u^T Au}}{1-\bar{\lambda}} - \overline{u^T Au} \\ &= \frac{1-\bar{\lambda}}{(1-\lambda)(1-\bar{\lambda})} \quad \Rightarrow |\lambda| < 1. \\ &\Rightarrow \rho(M) < 1. \end{aligned}$$

Let us prove the claims about the convergence of the Jacobi and SOR methods. We assume that A is real.

Let us begin with the Jacobi method. In this case $N = D$ and

$$Q = N + N^T - A = 2D - A = |A|.$$

Since A is positive definite, so is $|A|$. Hence Q is positive definite.

Let us now consider the SOR method. In this case, we have

$$N = \frac{1}{\omega} (D + \omega L).$$

Hence

$$\begin{aligned} Q &= \frac{1}{\omega} (D + \omega L + D + \omega U) - A \\ &= \left(\frac{2}{\omega} - 1\right) D, \end{aligned}$$

which is positive definite if $0 < \omega < 2$, as claimed. Now, note that

$$\begin{aligned} \det M_{SOR} &= \det ((D + \omega L)^{-1} (D - \omega(D + U))) \\ &= \det ((I + \omega D^{-1}L)^{-1} (Id - \omega(Id + D^{-1}U))) \\ &= (1 - \omega)^N \end{aligned}$$

and so, $P(M_{SOR}) \geq |1-\omega|$. This implies that the SOR method diverges if $\omega \notin [0, 2]$.

Note that

$$P(M_{SSOR}) \leq P^2(M_{SOR}).$$

and so, the SSOR method converges for $\omega \in (0, 2)$. An argument similar to that used for the SOR method shows that for $\omega \notin [0, 2]$, the SSOR method diverges.

- ④ Symmetric, positive definite matrices that are block-tridiagonal.

For these matrices, we have the following result:

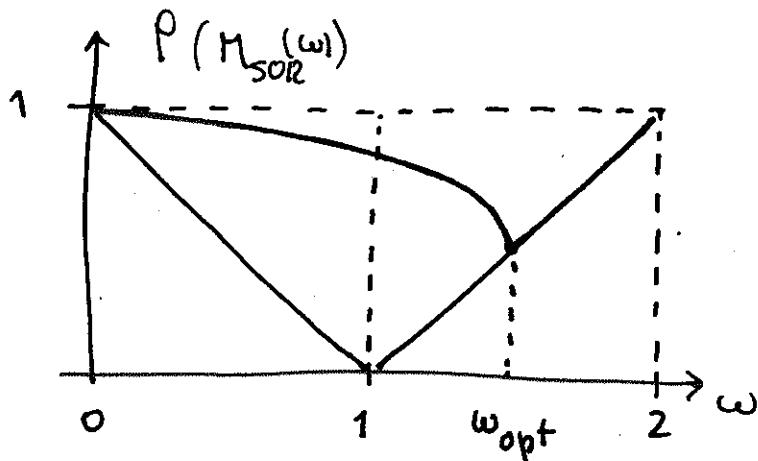
the Jacobi method, the Gauss-Seidel method and the SOR with $\omega \in (0, 2)$ converge. Moreover, we have:

$$\begin{aligned} P(M_{SOR}(\omega_{opt})) &= \inf_{\omega \in \mathbb{R}} P(M_{SOR}(\omega)) \\ &= \omega_{opt} - 1 \\ &< P(M_{GS}) \\ &= P^2(M_J) < P(M_J), \end{aligned}$$

where

$$\omega_{\text{opt}} = \frac{2}{1 + \sqrt{1 - P^*(M_j)}}.$$

Note that



To prove this result, we use a fundamental fact for block-tridiagonal matrices. Consider the matrix

$$A(\mu) = \begin{bmatrix} D_1 & \bar{\mu} U_1 & 0 \\ \mu L_1 & D_2 & \bar{\mu} U_2 \\ 0 & \mu L_2 & D_3 & \ddots \\ & & \ddots & \ddots & \bar{\mu} U_{N-1} \\ & & & & \mu U_N & D_N \end{bmatrix}$$

$$= D + \mu L + \bar{\mu} U.$$

then

$$\det A(\mu) = \det A(s)$$

or, equivalently,

$$(*) \quad \det(D + \mu L + \bar{\mu} U) = \det(D + L + U).$$

To prove this result, it is enough to realize that

$$A(\mu) = \Phi(\mu) A(s) \Phi(\mu)^{-1}$$

where

$$\Phi(\mu) = \begin{bmatrix} \mu I_1 & & & 0 \\ & \bar{\mu}^2 I_2 & & \\ & & \ddots & \\ 0 & & & \bar{\mu}^n I_n \end{bmatrix}.$$

Now, to prove the convergence result, we proceed as follows. First, let us consider the following polynomial

$$p_j(\lambda) = \det(-\bar{J}^{-1}(U+L) - \lambda Id).$$

By definition, the zeroes of p_j are the eigenvalues of the iteration matrix for the Jacobi method. Then

$$p_j(\lambda) = \det(-\bar{J}) \cdot \det(U+L+\lambda D).$$

Now, let us consider the SOR method:

$$\begin{aligned}
 p_{\text{SOR}}(\lambda) &= \det \left((\bar{D} + \omega \bar{L})^{-1} ((1-\omega) D - \omega U) - \lambda \text{Id} \right) \\
 &= \det \left(\left(\frac{\bar{D}}{\omega} + \bar{L} \right)^{-1} \left(\frac{1-\omega}{\omega} D - U \right) - \lambda \text{Id} \right) \\
 &= \det \left(\frac{\bar{D}}{\omega} + \bar{L} \right)^{-1} \cdot \det \left(\frac{1-\omega}{\omega} D - U - \frac{\lambda}{\omega} D - \lambda L \right) \\
 &= \det \left(\frac{\bar{D}}{\omega} + \bar{L} \right)^{-1} \cdot \det \left(\frac{1-\lambda\omega}{\omega} D - U - \lambda L \right) \\
 &= \det \left(\frac{\bar{D}}{\omega} + \bar{L} \right)^{-1} \cdot (-\sqrt{\lambda})^n \cdot \det \left(\frac{\lambda+\omega-1}{\omega\sqrt{\lambda}} D + \frac{1}{\sqrt{\lambda}} U + \sqrt{\lambda} L \right)
 \end{aligned}$$

by the property (*),

$$p_{\text{SOR}}(\lambda) = \det \left(\frac{\bar{D}}{\omega} + \bar{L} \right)^{-1} (-\sqrt{\lambda})^n \cdot \det \left(\frac{\lambda+\omega-1}{\omega\sqrt{\lambda}} D + U + L \right),$$

and hence

$$p_{\text{SOR}}(\lambda) = C(\lambda) \cdot p_J \left(\frac{\lambda+\omega-1}{\omega\sqrt{\lambda}} \right),$$

$$C(\lambda) = \det \left(\frac{\bar{D}}{\omega} + \bar{L} \right)^{-1} \cdot \det \bar{D} \cdot \lambda^{n_2}.$$

this means that

$\alpha \in \text{spectrum of } M_J$

$$\Leftrightarrow \lambda: \frac{\lambda+\omega-1}{\omega\sqrt{\lambda}} = \alpha \in \text{spectrum of } H_{\text{SOR}}^{(\omega)}!$$

If $\omega=1$, $\sqrt{\lambda}=\alpha$ and so

$$\rho(M_{SQR}(\omega=1)) = \rho(M_{GS}) = \rho^2(M_J).$$

To prove the remainder of the result, let us show that $\alpha \in \text{spectrum of } M_J$ is real. Indeed, if we set that

$$-\bar{J}(U+L)V = \alpha V,$$

then

$$\begin{aligned} O &= (U+L + D\alpha) V \\ &= (A + (\alpha-1)D) V \end{aligned}$$

and so

$$\bar{V}^T A V^T = (1-\alpha) \bar{V}^T D V,$$

which proves the claim since A is symmetric and positive definite. Moreover, $\alpha \in [0, 1]$.

As a consequence, we must study the number

$$\rho(M_{SQR}(\omega)) = \max_{\lambda \in \text{spectrum } M_J} \left\{ |\lambda| : \frac{\lambda + \omega - 1}{\sqrt{\omega}} = \alpha \right\}.$$

assuming that $\text{spectrum } M_J \subset (0, 1)$.

the theorem follows from such study.

(5) the model problem.

Let Ω be the unit square and let U_{xj} denote the classic finite difference approximations to the exact solutions of

$$(P) \quad -\Delta u = f \text{ in } \Omega; \quad u=0 \text{ on } \partial\Omega,$$

given by

$$\left\{ \begin{array}{l} -\frac{1}{h^2} (U_{i+1,j} + U_{i-1,j} + U_{i,j+1} + U_{i,j-1} - 4U_{ij}) = f_{ij} \\ i, j = 1, \dots, N-1 \\ U_{ij} = 0 \text{ if } i=0, j=0, i=N \text{ or } j=N. \end{array} \right.$$

It is not difficult to show that if A denotes the matrix associated to (P_h) , it is positive definite and symmetric. Moreover,

$$AU = \lambda U$$

$$\Leftrightarrow U = U^{nm}, \quad \lambda = \lambda^{nm}$$

$$U_{ij}^{nm} = (\sin n\pi x_i) (\sin m\pi y_j)$$

$$\lambda^{nm} = \frac{4}{h^2} \left(\sin^2 \frac{n\pi h}{2} + \sin^2 \frac{m\pi h}{2} \right)$$

while $x_i = ih$, $y_j = jh$.

In this case, we can write

$$e^k = \sum_{m,n=1}^{N-1} c_{mn}^k v^{mn},$$

where $e^k = x - x^k$ and x^k is the k -th iterate of the Jacobi method. Then, since

$$\begin{aligned} e^{k+1} &= -\bar{D}(U+L)e^k \\ &= \bar{D}^{-1}(D-A)e^k \\ &= (Id - \bar{D}^{-1}A)e^k \end{aligned}$$

we get that

$$c_{mn}^{k+1} = \left(1 - \frac{h^2}{4} \lambda^{mn}\right) c_{mn}^k.$$

This implies that

$$c_{mn}^{k+1} = \left(1 - \sin^2 \frac{m\pi h}{2} - \sin^2 \frac{n\pi h}{2}\right) c_{mn}^k.$$

If $m=n=N/2$,

$$\begin{aligned} c_{mn}^{k+1} &= \left(1 - 2 \sin^2 \frac{N/2 \pi}{4}\right) c_{mn}^k \\ &= 0, \end{aligned}$$

and this implies that the components of the error associated with v^{mn} , $m, n \approx N/2$ are damped extremely fast.

If $m=n=1$,

$$\begin{aligned} C_{mn}^{k+1} &= \left(1 - 2 \sin^2 \frac{\pi h}{2}\right) C_{mn}^k \\ &= \left(1 - \frac{\pi^2 h^2}{2} + O(h^4)\right) C_{mn}^k, \end{aligned}$$

and if $m=n=N-1$,

$$\begin{aligned} C_{mn}^{k+1} &= \left(1 - 2 \sin^2 \frac{(N-1)\pi h}{2}\right) C_{mn}^k \\ &= -\left(1 - 2 \sin^2 \frac{\pi h}{2}\right) C_{mn}^k \\ &= \left(-1 + \frac{\pi^2 h^2}{2} + O(h^4)\right) C_{mn}^k, \end{aligned}$$

which means that for $m \approx n = 1$ and $m \approx n \approx N-1$, the corresponding frequencies are barely damped.