



Math 673

Multigrid Methods: A Mostly Matrix-Based Approach

Chapter 01: Classical Iterative Methods

Abner J. Salgado and Steven M. Wise

asalgad1@utk.edu swise1@utk.edu
University of Tennessee

Fall 2024



Chapter 01, Part 2 of 2

Classical Iterative Methods





Classical Splitting Methods

In this section we define a couple of the classical iterative methods that one meets in an elementary course on numerical analysis, namely, the Jacobi, Gauss-Seidel, and Richardson's methods. We analyze them using the spectral theory introduced in the last section.

The Jacobi and Gauss-Seidel methods are based on the canonical (diagonal) splitting of A .

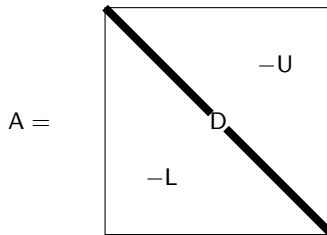
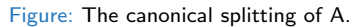
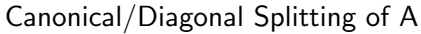


Figure: The canonical splitting of A.



Let $A \in \mathbb{R}^{n \times n}$ be given. Consider

where $D \in \mathbb{R}^{n \times n}$ is the diagonal of A , $-L \in \mathbb{R}^{n \times n}$ is the strictly lower triangular part of A , and $-U \in \mathbb{R}^{n \times n}$ is the strictly upper triangular part of A , respectively, as illustrated in the figure on the previous slide. The decomposition (1) is called the **canonical splitting of A** or **diagonal splitting of A** .





Suppose that

$$A = \begin{bmatrix} 10 & -2 & -10 & 5 & -10 \\ 10 & 9 & 7 & -2 & -5 \\ 0 & 6 & 9 & 3 & -10 \\ 6 & 10 & 4 & -7 & -8 \\ -8 & 3 & 5 & 4 & 7 \end{bmatrix}.$$

Then

$$D = \begin{bmatrix} 10 & 0 & 0 & 0 & 0 \\ 0 & 9 & 0 & 0 & 0 \\ 0 & 0 & 9 & 0 & 0 \\ 0 & 0 & 0 & -7 & 0 \\ 0 & 0 & 0 & 0 & 7 \end{bmatrix},$$

$$U = \begin{bmatrix} 0 & 2 & 10 & -5 & 10 \\ 0 & 0 & -7 & 2 & 5 \\ 0 & 0 & 0 & -3 & 10 \\ 0 & 0 & 0 & 0 & 8 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \quad \text{and} \quad L = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ -10 & 0 & 0 & 0 & 0 \\ 0 & -6 & 0 & 0 & 0 \\ -6 & -10 & -4 & 0 & 0 \\ 8 & -3 & -5 & -4 & 0 \end{bmatrix}.$$



Proposition

If A is SPD, then D has only positive diagonal elements, and $U = L^T$.

Proof.

Exercise.







Definition (The Jacobi Method)

$$B = B_I := D^{-1}$$
$$T = T_I := I - D^{-1}A.$$
$$\mathbf{u}^{k+1} = \mathbf{u}^k + \mathbf{D}^{-1} \left(\mathbf{f} - \mathbf{A} \mathbf{u}^k \right). \quad (2)$$



The error transfer matrix for the Jacobi method can be expressed as

$$\begin{aligned} \mathbf{T}_J &= \mathbf{D}^{-1}\mathbf{D} - \mathbf{D}^{-1}\mathbf{A} \\ &= \mathbf{D}^{-1}(\mathbf{D} - \mathbf{D} + \mathbf{U} + \mathbf{L}) \\ &= \mathbf{D}^{-1}(\mathbf{U} + \mathbf{L}). \end{aligned} \quad (3)$$

In component form, the Jacobi method can be written as

$$u_i^{k+1} = \frac{f_i - \sum_{j=1}^{i-1} a_{i,j} u_j^k - \sum_{j=i+1}^n a_{i,j} u_j^k}{a_{i,i}}. \quad (4)$$

The order in which we obtain the updated components is completely immaterial.



We say that a matrix $A \in \mathbb{R}^{n \times n}$ is **diagonally dominant** or **row-wise diagonally dominant** iff, for each $i \in \{1, \dots, n\}$,

$$\sum_{\substack{j=1 \\ j \neq i}}^n |a_{i,j}| < |a_{i,i}|.$$

If $A \in \mathbb{R}^{n \times n}$ is diagonally dominant, then it is invertible and none of its diagonal elements is zero.

Exercise.





If $A \in \mathbb{R}^{n \times n}$ is diagonally dominant, then the Jacobi method is well defined and unconditionally convergent.

We will show that $\|T_J\|_\infty < 1$. Since A is diagonally dominant, its diagonal elements are all non-zero. Thus the method is well defined because D is invertible. Now, observe that

$$\mathbf{T}_J = \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix} - \begin{bmatrix} \frac{1}{a_{1,1}} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{a_{n,n}} \end{bmatrix} \begin{bmatrix} a_{1,1} & \cdots & a_{1,n} \\ \vdots & \ddots & \vdots \\ a_{n,1} & \cdots & a_{n,n} \end{bmatrix}.$$

Thus

$$t_{i,j} = \begin{cases} 0, & \text{if } i = j, \\ -\frac{a_{i,j}}{a_{i,i}}, & \text{if } i \neq j. \end{cases}$$



Proof (Cont.)

So,

$$\|T_J\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |t_{i,j}| = \max_{1 \leq i \leq n} \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{i,j}}{a_{i,i}} \right|.$$

Since A is diagonally dominant, there is some $\delta > 0$, such that

$$|a_{i,i}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{i,j}| + \delta,$$

for each $i = 1, \dots, n$.



Proof (Cont.)

Therefore, for each $i = 1, \dots, n$,

$$\sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{i,j}}{a_{i,i}} \right| \leq 1 - \frac{\delta}{|a_{i,i}|} \leq 1 - \frac{\delta}{\alpha},$$

where

$$\alpha = \max_{1 \leq i \leq n} |a_{i,i}|.$$

Hence,

$$\|\mathbf{T}_J\|_{\infty} \leq 1 - \frac{\delta}{\alpha} < 1.$$





Definition

Suppose that $A = D - L - U$ is the canonical splitting of $A \in \mathbb{R}^{n \times n}$ and $\omega \in (0, 1]$. Assume D is invertible. The **weighted Jacobi method** is defined by the iteration

$$\begin{aligned} \mathbf{z} &= \mathbf{u}^k + D^{-1}(\mathbf{f} - A\mathbf{u}^k) \\ \mathbf{u}^{k+1} &= \omega \mathbf{z} + (1 - \omega)\mathbf{u}^k. \end{aligned}$$

where the starting value $\mathbf{u}^0 \in \mathbb{R}^n$ is given.

Thus the weighted Jacobi method can be viewed as doing one step of the Jacobi method, followed by a weighting step. We can eliminate \mathbf{z} in the second equation to obtain the following.



Proposition

With the same assumptions as in the last definition, the weighted Jacobi method is a GLIS with iterator

$$\mathbf{B} = \mathbf{B}_{J,\omega} := \omega \mathbf{D}^{-1}$$

and error transfer matrix

$$\mathbf{T} = \mathbf{T}_{J,\omega} := \mathbf{I} - \omega \mathbf{D}^{-1} \mathbf{A}.$$

The iteration sequence can be expressed as

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \omega \mathbf{D}^{-1} (\mathbf{f} - \mathbf{A} \mathbf{u}^k). \quad (5)$$

Proof.

Exercise. □

We will come back to the weighted Jacobi method in a later chapter.

The Gauss-Seidel Method



Suppose that $A = D - L - U$ is the canonical splitting of $A \in \mathbb{R}^{n \times n}$. Then

$$(D - L)\mathbf{u} = U\mathbf{u} + \mathbf{f}.$$

The Gauss-Seidel method is simply stated as

$$(D - L)\mathbf{u}^{k+1} = U\mathbf{u}^k + \mathbf{f}.$$

Assuming $D - L$ is invertible, some manipulations give

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



Definition (Forward Gauss-Seidel Method)

Suppose that $A = D - L - U$ is the canonical splitting of $A \in \mathbb{R}^{n \times n}$ and $D - L$ is invertible. The **Gauss-Seidel method** – also called the **forward Gauss-Seidel method** – is a GLIS with iterator matrix

$$B = B_{GS} := (D - L)^{-1}$$

and error transfer matrix

$$T = T_{GS} := I - (D - L)^{-1}A.$$

The iteration sequence can be expressed as

$$\mathbf{u}^{k+1} = \mathbf{u}^k + (D - L)^{-1}(\mathbf{f} - A\mathbf{u}^k). \quad (6)$$



Remark

The error transfer matrix for the Gauss-Seidel method can be expressed equivalently as

$$\begin{aligned}
 T_{GS} &= I - (D - L)^{-1}A \\
 &= (D - L)^{-1}(D - L) - (D - L)^{-1}(D - L - U) \\
 &= (D - L)^{-1}U.
 \end{aligned} \tag{7}$$

The Component Form of Gauss-Seidel



The Gauss-Seidel method is simply stated as

$$(D - L)\mathbf{u}^{k+1} - U\mathbf{u}^k = \mathbf{f}.$$

In component form, we can express the forward Gauss-Seidel method as

$$\sum_{j=1}^i a_{i,j} u_j^{k+1} + \sum_{j=i+1}^n a_{i,j} u_j^k = f_i, \quad i = 1, \dots, n.$$

The Component Form of Gauss-Seidel



To compute the update \mathbf{u}^{k+1} , we proceed from $i = 1$, in order, to $i = n$. So, the first equation we solve is

$$u_1^{k+1} = \frac{f_1 - \sum_{j=2}^n a_{1,j} u_j^k}{a_{1,1}}.$$

Next, we have

$$u_2^{k+1} = \frac{f_2 - a_{2,1} u_1^{k+1} - \sum_{j=3}^n a_{2,j} u_j^k}{a_{2,2}},$$

and

$$u_3^{k+1} = \frac{f_3 - \sum_{j=1}^2 a_{3,j} u_j^{k+1} - \sum_{j=4}^n a_{3,j} u_j^k}{a_{3,3}}.$$

We continue, in the same way, with $i = 4$, in order, to $i = n - 1$, and we end with

$$u_n^{k+1} = \frac{f_n - \sum_{j=1}^{n-1} a_{n,j} u_j^{k+1}}{a_{n,n}}.$$

The Component Form of Gauss-Seidel



The idea with the method is simple. Once we obtain an updated component u_i^{k+1} , we immediately use it to obtain u_{i+1}^{k+1} . One full pass, from $i = 1$ to $i = n$ is called a **forward sweep**. The generic i^{th} step is

$$u_i^{k+1} = \frac{f_i - \sum_{j=1}^{i-1} a_{i,j} u_j^{k+1} - \sum_{j=i+1}^n a_{i,j} u_j^k}{a_{i,i}}. \quad (8)$$



Theorem (Convergence of the forward Gauss-Seidel method)

Suppose that $A = D - L - U$ is the canonical splitting of $A \in \mathbb{R}^{n \times n}$. If A is diagonally dominant, then the forward Gauss-Seidel method is well defined, that is, $D - L$ is invertible, and the method is unconditionally convergent.

Proof.

If A is diagonally dominant, then D has only nonzero diagonal entries. Hence, $D - L$ is invertible.

For convergence, we again want to show that, if A is diagonally dominant, then $\|T_{GS}\|_{\infty} < 1$. This will establish the unconditional convergence of the method. Since A is diagonally dominant, it follows that, for some $\delta > 0$,

$$|a_{i,i}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{i,j}| + \delta = \sum_{j>i} |a_{i,j}| + \sum_{j<i} |a_{i,j}| + \delta.$$



Proof (Cont.)

Thus

$$|a_{i,i}| - \sum_{j<i} |a_{i,j}| \geq \sum_{j>i} |a_{i,j}| + \delta > \sum_{j>i} |a_{i,j}|,$$

which implies

$$\gamma := \max_{1 \leq i \leq n} \left\{ \frac{\sum_{j>i} |a_{i,j}|}{|a_{i,i}| - \sum_{j<i} |a_{i,j}|} \right\} < 1.$$

Now, we will show $\|T_{GS}\|_{\infty} \leq \gamma$.

Let $\mathbf{x} \in \mathbb{R}^{n \times n}$ and $\mathbf{y} = T_{GS}\mathbf{x}$, that is,

$$\mathbf{y} = T_{GS}\mathbf{x} = (D - L)^{-1}U\mathbf{x}.$$

Let k be the index such that $\|\mathbf{y}\|_{\infty} = |y_k|$. Then we have

$$|[(D - L)\mathbf{y}]_k| = |[U\mathbf{x}]_k| = \left| \sum_{j>k} a_{k,j}x_j \right| \leq \sum_{j>k} |a_{k,j}| |x_j| \leq \sum_{j>k} |a_{k,j}| \|\mathbf{x}\|_{\infty}.$$



Proof (Cont.)

Now notice that

$$\begin{aligned}
 |[(D - L)y]_k| &= \left| \sum_{j < k} a_{k,j} y_j + a_{k,k} y_k \right| \\
 &\geq |a_{k,k} y_k| - \left| \sum_{j < k} a_{k,j} y_j \right| \\
 &= |a_{k,k}| \|y\|_\infty - \left| \sum_{j < k} a_{k,j} y_j \right| \\
 &\geq |a_{k,k}| \|y\|_\infty - \sum_{j < k} |a_{k,j}| \|y\|_\infty.
 \end{aligned}$$



Proof (Cont.)

Therefore, we have

$$|a_{k,k}| \|y\|_{\infty} - \sum_{j < k} |a_{k,j}| \|y\|_{\infty} \leq \sum_{j > k} |a_{k,j}| \|x\|_{\infty},$$

which implies

$$\|y\|_{\infty} \leq \frac{\sum_{j > k} |a_{k,j}|}{|a_{k,k}| - \sum_{j < k} |a_{k,j}|} \|x\|_{\infty}.$$

So,

$$\|T_{GS}x\|_{\infty} \leq \gamma \|x\|_{\infty},$$

which implies

$$\|T_{GS}\|_{\infty} \leq \gamma < 1.$$





Definition (Backward Gauss-Seidel method)

Suppose that $A \in \mathbb{R}^{n \times n}$ is invertible and $A = D - L - U$ is the canonical splitting of A . Assume $D - U$ is invertible. The **backward Gauss-Seidel method** is a GLIS with iterator matrix

$$B = B_{BGS} := (D - U)^{-1}$$

and error transfer matrix

$$T = T_{BGS} := I - (D - U)^{-1}A.$$

The iteration scheme is precisely

$$\mathbf{u}^{k+1} = \mathbf{u}^k + (D - U)^{-1}(\mathbf{f} - A\mathbf{u}^k). \quad (9)$$



Remark

The error transfer matrix for the backward Gauss-Seidel method can be re-expressed as

$$\begin{aligned}
 T_{BGS} &= I - (D - U)^{-1}A \\
 &= I - (D - U)^{-1}(D - L - U) \\
 &= (D - U)^{-1}L.
 \end{aligned} \tag{10}$$



It should not be hard to modify the proof of the last theorem to obtain the following:

Theorem (Convergence of the backward Gauss-Seidel method)

Suppose that $A = D - L - U$ is the canonical splitting of $A \in \mathbb{R}^{n \times n}$. If A is diagonally dominant, then the backward Gauss-Seidel method is well defined, that is, $D - U$ is invertible, and it is unconditionally convergent.

Component Form of Backward Gauss-Seidel



The Gauss-Seidel method is simply stated as

$$-\mathbf{L}\mathbf{u}^k + (\mathbf{D} - \mathbf{U})\mathbf{u}^{k+1} = \mathbf{f}.$$

In component form, we can express the backward Gauss-Seidel method as

$$\sum_{j=1}^{i-1} a_{i,j} u_j^k + \sum_{j=i}^n a_{i,j} u_j^{k+1} = f_i, \quad i = 1, \dots, n.$$

Component Form of Backward Gauss-Seidel



To compute the update \mathbf{u}^{k+1} , we proceed from $i = n$, in reverse order, to $i = 1$. We find the n^{th} component first:

$$u_n^{k+1} = \frac{f_n - \sum_{j=1}^{n-1} a_{n,j} u_j^k}{a_{n,n}}.$$

Next, we have

$$u_{n-1}^{k+1} = \frac{f_{n-1} - \sum_{j=1}^{n-2} a_{n-1,j} u_j^k - a_{n-1,n} u_n^{k+1}}{a_{n-1,n-1}},$$

and then

$$u_{n-2}^{k+1} = \frac{f_{n-2} - \sum_{j=1}^{n-3} a_{n-2,j} u_j^k - \sum_{j=n-1}^n a_{n-2,j} u_j^{k+1}}{a_{n-2,n-2}}.$$

We continue, in the same way, with $i = n - 3$, in reverse order, to $i = 2$, and we end with

$$u_1^{k+1} = \frac{f_1 - \sum_{j=1}^{n-1} a_{1,j} u_j^{k+1}}{a_{1,1}}.$$

Component Form of Backward Gauss-Seidel



One full pass, from $i = n$, proceeding in reverse order, to $i = 1$ is called a **backward sweep**. The generic i^{th} step is

$$u_i^{k+1} = \frac{f_i - \sum_{j=1}^{i-1} a_{i,j} u_j^k - \sum_{j=i+1}^n a_{i,j} u_j^{k+1}}{a_{i,i}}. \quad (11)$$

Richardson's Method



Suppose $\omega \in \mathbb{R}_* := \mathbb{R} \setminus \{0\}$, and consider the splitting

$$A = \omega I + A - \omega I.$$

Suppose

$$A\mathbf{u} = \mathbf{f}.$$

Then,

$$\omega \mathbf{u} = (\omega I - A)\mathbf{u} + \mathbf{f}.$$

Richardson's method is, essentially,

$$\omega \mathbf{u}^{k+1} = (\omega I - A)\mathbf{u}^k + \mathbf{f}.$$

Notice that this method is not based on the canonical diagonal splitting, but it is still quite useful and interesting.



Definition (Richardson's Method)

Suppose that $A \in \mathbb{R}^{n \times n}$ is invertible and $\omega \in \mathbb{R}_*$. **Richardson's method** is a GLIS with iterator matrix

$$B = B_R := \omega^{-1}I$$

and error transfer matrix

$$T = T_R := I - \omega^{-1}A.$$

The iteration scheme is

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \omega^{-1}(\mathbf{f} - A\mathbf{u}^k). \quad (12)$$

We make a slight departure from the now familiar approach and analyze this method using the 2-norm, $\|\cdot\|_2$, under the assumption that A is SPD.



Theorem (Convergence of Richardson's Method)

Let $A \in \mathbb{R}^{n \times n}$ be SPD and $\sigma(A) = \{\lambda_1, \dots, \lambda_n\}$, with $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. Richardson's method converges unconditionally iff $\omega \in (0, 2/\lambda_n)$. In this case, we have the estimate

$$\|\mathbf{e}^k\|_2 \leq \rho^k \|\mathbf{e}^0\|_2, \quad \rho = \rho(\omega) = \max\{|1 - \omega\lambda_n|, |1 - \omega\lambda_1|\}.$$

From this, it follows that setting

$$\omega = \omega_{\text{opt}} := \frac{2}{\lambda_1 + \lambda_n},$$

one obtains the smallest possible value of ρ , ρ_{opt} , and

$$\rho_{\text{opt}} := \rho(\omega_{\text{opt}}) = \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} = \frac{\kappa_2(A) - 1}{\kappa_2(A) + 1},$$

where

$$\kappa_2(A) = \frac{\lambda_n}{\lambda_1}.$$



Proof of Convergence of Richardson's Method.

(\Leftarrow): Since A is SPD, we know that the eigenvalues of A are positive real numbers and

$$\lambda_n = \|A\|_2.$$

Notice also that $T_R = I_n - \omega A = T_R^T$, which implies that the eigenvalues of T_R are real. Observe that $(\lambda_i, \mathbf{w}_i)$ is an eigenpair of A iff $(\nu_i = 1 - \omega\lambda_i, \mathbf{w}_i)$ is an eigenpair of T_R . Assume that $0 < \omega < 2/\lambda_n$. Then

$$0 < \lambda_i \omega < 2 \frac{\lambda_i}{\lambda_n}, \quad i = 1, \dots, n,$$

which implies

$$1 > 1 - \lambda_i \omega > 1 - 2 \frac{\lambda_i}{\lambda_n} \geq -1, \quad i = 1, \dots, n.$$

It follows that

$$1 > \nu_1 \geq \dots \geq \nu_n > -1, \quad \nu_i = 1 - \omega\lambda_i.$$

This guarantees that $\|T_R\|_2 = \rho(T_R) < 1$, which implies convergence.



Proof (Cont.)

(\Rightarrow): Conversely, if $\omega \notin (0, 2/\lambda_n)$, then $\rho(\mathbf{T}_R) \geq 1$, and the method cannot converge unconditionally.

(Error estimate): By consistency,

$$\|\mathbf{e}^k\|_2 = \|\mathbf{T}_R^k \mathbf{e}^0\|_2 \leq \rho^k \|\mathbf{e}^0\|_2.$$

Of course, it is easy to see that

$$\rho = \rho(\mathbf{T}_R) = \max\{|\nu_1|, |\nu_n|\} = \max\{|1 - \omega\lambda_n|, |1 - \omega\lambda_1|\}.$$



Proof (Cont.)

(Optimality): Finally, showing optimality amounts to minimizing ρ . See the figure on the next slide. From this we see that the minimum of ρ is attained when

$$|1 - \omega\lambda_1| = |1 - \omega\lambda_n|$$

or

$$1 - \omega\lambda_n = \omega\lambda_1 - 1,$$

which implies

$$\omega_{\text{opt}} = \frac{2}{\lambda_1 + \lambda_n}.$$

Therefore

$$\rho_{\text{opt}} = 1 - \omega_{\text{opt}}\lambda_1 = \frac{\lambda_1 + \lambda_n - 2\lambda_1}{\lambda_1 + \lambda_n} = \frac{\lambda_n - \lambda_1}{\lambda_1 + \lambda_n} = \frac{\kappa_2(A) - 1}{\kappa_2(A) + 1}.$$



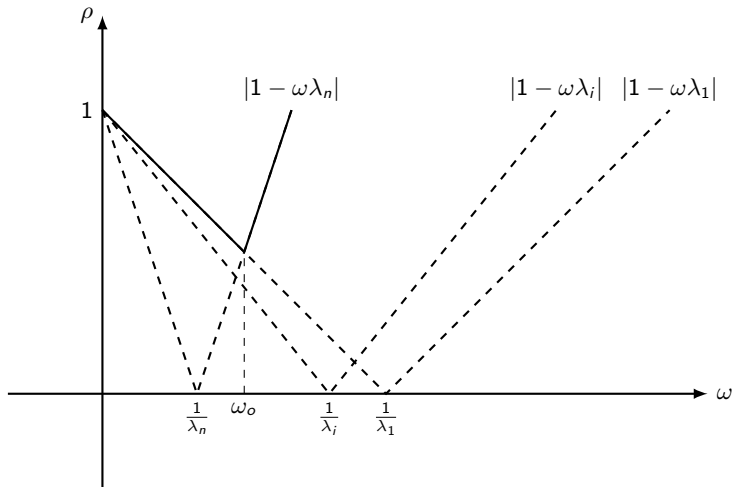


Figure: The curve $\rho(T_R)$ (in solid black) as a function of ω .





Suppose $A \in \mathbb{R}^{n \times n}$ is invertible and $f \in \mathbb{R}^n$ is given. Consider a family of iterator matrices, $B_i \in \mathbb{R}^{n \times n}$, for $i = 1, \dots, r$. An iterative scheme is called **additive with respect to** $\{B_i\}_{i=1}^r$ iff

$$\begin{aligned} \mathbf{z}_1 &= \mathbf{u}^k + \mathbf{B}_1 \left(\mathbf{f} - \mathbf{A} \mathbf{u}^k \right), \\ \mathbf{z}_2 &= \mathbf{z}_1 + \mathbf{B}_2 \left(\mathbf{f} - \mathbf{A} \mathbf{u}^k \right), \\ \mathbf{z}_3 &= \mathbf{z}_2 + \mathbf{B}_3 \left(\mathbf{f} - \mathbf{A} \mathbf{u}^k \right), \\ &\vdots \\ \mathbf{z}_{r-1} &= \mathbf{z}_{r-2} + \mathbf{B}_{r-1} \left(\mathbf{f} - \mathbf{A} \mathbf{u}^k \right), \\ \mathbf{u}^{k+1} &= \mathbf{z}_{r-1} + \mathbf{B}_r \left(\mathbf{f} - \mathbf{A} \mathbf{u}^k \right). \end{aligned} \tag{13}$$



An iterative scheme is called **multiplicative with respect to** $\{B_i\}_{i=1}^r$ iff

$$\begin{aligned} \mathbf{z}_1 &= \mathbf{u}^k + \mathbf{B}_1 (\mathbf{f} - \mathbf{A}\mathbf{u}^k), \\ \mathbf{z}_2 &= \mathbf{z}_1 + \mathbf{B}_2 (\mathbf{f} - \mathbf{A}\mathbf{z}_1), \\ \mathbf{z}_3 &= \mathbf{z}_2 + \mathbf{B}_3 (\mathbf{f} - \mathbf{A}\mathbf{z}_2), \\ &\vdots \\ \mathbf{z}_{r-1} &= \mathbf{z}_{r-2} + \mathbf{B}_{r-1} (\mathbf{f} - \mathbf{A}\mathbf{z}_{r-2}), \\ \mathbf{u}^{k+1} &= \mathbf{z}_{r-1} + \mathbf{B}_r (\mathbf{f} - \mathbf{A}\mathbf{z}_{r-1}). \end{aligned} \tag{14}$$



Suppose $A \in \mathbb{R}^{n \times n}$ is invertible and $f \in \mathbb{R}^n$ is given. Consider a family of iterator matrices, $B_i \in \mathbb{R}^{n \times n}$, for $i = 1, \dots, r$. An additive iterative scheme with respect to $\{B_i\}_{i=1}^r$ is a GLIS with iterator

$$B = \sum_{i=1}^r B_i.$$

A multiplicative iterative scheme with respect to $\{B_i\}_{i=1}^r$ is a GLIS with the following recursively-defined iterator:

$$\mathbb{B} = \tilde{\mathbb{B}}_r,$$

where $\tilde{B}_1 = B_1$, and, for $2 \leq i \leq r$,

$$\tilde{B}_i = \tilde{B}_{i-1} + B_i - B_i A \tilde{B}_{i-1}.$$

Proof.

Exercise.





Proposition

Suppose $A \in \mathbb{R}^{n \times n}$ is invertible, $f \in \mathbb{R}^n$ is given, and $\{B_i\}_{i=1}^r \subset \mathbb{R}^{n \times n}$ is a family of iterator matrices. The error transfer matrix for the additive GLIS with respect to $\{B_i\}_{i=1}^r$ can be expressed as

$$T = I - \left(\sum_{i=1}^r B_i \right) A.$$

The error transfer matrix for the multiplicative GLIS with respect to $\{B_i\}_{i=1}^r$ can be expressed as

$$T = \prod_{i=1}^r (I - B_i A).$$

Proof.

Exercise. □



Definition (Symmetrized Multiplicative GLIS)

Suppose $A \in \mathbb{R}^{n \times n}$ is invertible and $\mathbf{f} \in \mathbb{R}^n$ is given. Consider a generic GLIS with iterator $B \in \mathbb{R}^{n \times n}$. The **symmetrized multiplicative GLIS (SMGLIS) with respect to B** is defined as follows: given \mathbf{u}^0 , find $\mathbf{u}^1, \mathbf{u}^2, \dots$ via

$$\mathbf{z} = \mathbf{u}^k + B(\mathbf{f} - A\mathbf{u}^k), \quad (15)$$

$$\mathbf{u}^{k+1} = \mathbf{z} + B^T(\mathbf{f} - A\mathbf{z}). \quad (16)$$



Suppose $A \in \mathbb{R}^{n \times n}$ is invertible and $f \in \mathbb{R}^n$ is given. Assume that B is invertible. The SMGLIS with respect to B can be written as

where

In other words, the SMGLIS with respect to B is a GLIS with iterator B_{SM} . If B is invertible, then the iterator may be expressed as

If A is symmetric, then B_{SM} is as well.



Plugging (15) into (16) yields

$$\begin{aligned} \mathbf{u}^{k+1} &= \mathbf{u}^k + \mathbf{B}(\mathbf{f} - \mathbf{A}\mathbf{u}^k) + \mathbf{B}^T \left(\mathbf{f} - \mathbf{A} \left[\mathbf{u}^k + \mathbf{B}(\mathbf{f} - \mathbf{A}\mathbf{u}^k) \right] \right) \\ &= \mathbf{u}^k + \mathbf{B}(\mathbf{f} - \mathbf{A}\mathbf{u}^k) + \mathbf{B}^T \left(\mathbf{f} - \mathbf{A}\mathbf{u}^k - \mathbf{A}\mathbf{B}(\mathbf{f} - \mathbf{A}\mathbf{u}^k) \right) \\ &= \mathbf{u}^k + \mathbf{B}(\mathbf{f} - \mathbf{A}\mathbf{u}^k) + \mathbf{B}^T \left(\mathbf{f} - \mathbf{A}\mathbf{u}^k \right) - \mathbf{B}^T \mathbf{A} \mathbf{B} (\mathbf{f} - \mathbf{A}\mathbf{u}^k) \\ &= \mathbf{u}^k + \left(\mathbf{B} + \mathbf{B}^T - \mathbf{B}^T \mathbf{A} \mathbf{B} \right) (\mathbf{f} - \mathbf{A}\mathbf{u}^k), \end{aligned}$$

which shows the SMGLIS with respect to B is a GLIS with iterator

$$\mathbf{B} = \mathbf{B}_{SM} = \mathbf{B} + \mathbf{B}^T - \mathbf{B}^T \mathbf{A} \mathbf{B}.$$



Proof (Cont.)

If B is invertible, then

$$\begin{aligned}
 B_{SM} &= B + B^T - B^T A B \\
 &= (I + B^T B^{-1} - B^T A) B \\
 &= (B^T B^{-T} + B^T B^{-1} - B^T A) B \\
 &= B^T (B^{-T} + B^{-1} - A) B.
 \end{aligned}$$

If A is symmetric, then B_{SM} is symmetric, since

$$\begin{aligned}
 B_{SM}^T &= B^T + B - B^T A^T B \\
 &= B^T + B - B^T A B \\
 &= B_{SM}.
 \end{aligned}$$





Remark

We could also consider a symmetrized additive method:

$$\mathbf{z} = \mathbf{u}^k + \mathbf{B}(\mathbf{f} - \mathbf{A}\mathbf{u}^k), \quad (20)$$

$$\mathbf{u}^{k+1} = \mathbf{z} + \mathbf{B}^T(\mathbf{f} - \mathbf{A}\mathbf{u}^k). \quad (21)$$

This is a GLIS with iterator matrix

$$B = B_{SA} := B + B^T.$$

But this, it turns out, is not as useful to us in the multigrid setting.



The forward Gauss-Seidel method is the best example of a GLIS where the iterator is non-symmetric, even when A is symmetric. However, we can use the methodology presented in the last section to symmetrize it.

Definition (The Symmetric Gauss-Seidel Method)

Suppose that $A \in \mathbb{R}^{n \times n}$ is invertible and $A = D - L - U$ is the canonical splitting of A . Assume $D - L$ is invertible, and consider the Gauss-Seidel method, which is characterized by the iterator

$$B_{GS} = (D - L)^{-1}.$$

The **symmetric Gauss Seidel method** is the SMGLIS with respect to B_{GS} and is, therefore, a GLIS with iterator matrix

$$\mathbf{B}_{SGS} = \mathbf{B}_{GS} + \mathbf{B}_{GS}^T - \mathbf{B}_{GS}^T \mathbf{A} \mathbf{B}_{GS}. \quad (22)$$



Suppose that $A \in \mathbb{R}^{n \times n}$ is symmetric and invertible and $A = D - L - U$ is the canonical splitting of A . Assume that $D - L$ is invertible. The iterator for the symmetric Gauss Seidel method can be written as

$$\mathbf{B}_{SGS} = (\mathbf{D} - \mathbf{U})^{-1} \mathbf{D} (\mathbf{D} - \mathbf{L})^{-1}, \quad (23)$$

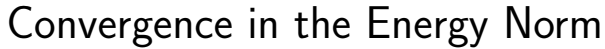
and it is symmetric.

Proof.

Using Equation (18), we find

$$\begin{aligned} \mathbf{B}_{SGS} &= \mathbf{B}_{GS}^T (\mathbf{B}_{GS}^{-T} + \mathbf{B}_{GS}^{-1} - \mathbf{A}) \mathbf{B}_{GS} \\ &= (\mathbf{D}^T - \mathbf{L}^T)^{-1} (\mathbf{D}^T - \mathbf{L}^T + \mathbf{D} - \mathbf{L} - \mathbf{A}) (\mathbf{D} - \mathbf{L})^{-1} \\ &= (\mathbf{D} - \mathbf{U})^{-1} (\mathbf{D} - \mathbf{U} + \mathbf{D} - \mathbf{L} - \mathbf{A}) (\mathbf{D} - \mathbf{L})^{-1} \\ &= (\mathbf{D} - \mathbf{U})^{-1} \mathbf{D} (\mathbf{D} - \mathbf{L})^{-1}. \end{aligned}$$







If the coefficient matrix A is SPD, we can use it to construct a useful norm, one which will be the basis of measuring convergence in a new way.

Definition (Energy Norm)

Suppose that A is SPD. The **energy norm** associated with A is

$$\|\mathbf{u}\|_A = \sqrt{(\mathbf{u}, \mathbf{u})_A}.$$

The energy norm is a natural metric for the convergence of the GLIS, when A is SPD, as the next few results show.



Suppose $A \in \mathbb{R}^{n \times n}$ is SPD and $B \in \mathbb{R}^{n \times n}$ is symmetric and invertible. Define

$$Q := B^{-1} - \frac{1}{2}A.$$

If Q is SPD, then the GLIS

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \mathbf{B}(\mathbf{f} - \mathbf{A}\mathbf{u}^k) \quad (24)$$

is unconditionally convergent with respect to the A-norm, that is,

$$\left\| \mathbf{e}^k \right\|_A \rightarrow 0, \quad \text{as } k \rightarrow \infty.$$

Moreover, the convergence is always monotonic in the sense that

$$\left\| \mathbf{e}^{k+1} \right\|_{\mathbf{A}}^2 \leq \left\| \mathbf{e}^k \right\|_{\mathbf{A}}^2, \quad \forall k \in \mathbb{N}_0 := \mathbb{N} \cup \{0\}.$$



Proof.

Let \mathbf{u} be the exact solution and \mathbf{u}^k be the GLIS approximation at the k^{th} step. Then

$$\mathbf{e}^{k+1} = \mathbf{T}\mathbf{e}^k = \mathbf{e}^k - \mathbf{B}\mathbf{A}\mathbf{e}^k.$$

Set $\mathbf{v}^{k+1} := \mathbf{e}^{k+1} - \mathbf{e}^k$. Then

$$\mathbf{B}^{-1}\mathbf{v}^{k+1} + \mathbf{A}\mathbf{e}^k = \mathbf{0}. \quad (25)$$

Taking the Euclidean inner product of the last equation with \mathbf{v}^{k+1} gives

$$\left(\mathbf{B}^{-1}\mathbf{v}^{k+1}, \mathbf{v}^{k+1}\right) + \left(\mathbf{A}\mathbf{e}^k, \mathbf{v}^{k+1}\right) = 0.$$

Now, observe that

$$\mathbf{e}^k = \frac{1}{2}(\mathbf{e}^{k+1} + \mathbf{e}^k) - \frac{1}{2}(\mathbf{e}^{k+1} - \mathbf{e}^k) = \frac{1}{2}(\mathbf{e}^{k+1} + \mathbf{e}^k) - \frac{1}{2}\mathbf{v}^{k+1}.$$



Proof (Cont.)

Then,

$$\begin{aligned}
 0 &= (\mathbf{B}^{-1} \mathbf{v}^{k+1}, \mathbf{v}^{k+1}) + (\mathbf{A} \mathbf{e}^k, \mathbf{v}^{k+1}) \\
 &= (\mathbf{B}^{-1} \mathbf{v}^{k+1}, \mathbf{v}^{k+1}) + \frac{1}{2} (\mathbf{A} (\mathbf{e}^{k+1} + \mathbf{e}^k), \mathbf{v}^{k+1}) - \frac{1}{2} (\mathbf{A} \mathbf{v}^{k+1}, \mathbf{v}^{k+1}) \\
 &= \left(\left(\mathbf{B}^{-1} - \frac{1}{2} \mathbf{A} \right) \mathbf{v}^{k+1}, \mathbf{v}^{k+1} \right) + \frac{1}{2} (\mathbf{A} (\mathbf{e}^{k+1} + \mathbf{e}^k), \mathbf{v}^{k+1}) \\
 &= \left(\left(\mathbf{B}^{-1} - \frac{1}{2} \mathbf{A} \right) \mathbf{v}^{k+1}, \mathbf{v}^{k+1} \right) + \frac{1}{2} (\mathbf{A} (\mathbf{e}^{k+1} + \mathbf{e}^k), \mathbf{e}^{k+1} - \mathbf{e}^k) \\
 &= \left(\left(\mathbf{B}^{-1} - \frac{1}{2} \mathbf{A} \right) \mathbf{v}^{k+1}, \mathbf{v}^{k+1} \right) + \frac{1}{2} \left(\|\mathbf{e}^{k+1}\|_{\mathbf{A}}^2 - \|\mathbf{e}^k\|_{\mathbf{A}}^2 \right) \\
 &= \|\mathbf{v}^{k+1}\|_{\mathbf{Q}}^2 + \frac{1}{2} \left(\|\mathbf{e}^{k+1}\|_{\mathbf{A}}^2 - \|\mathbf{e}^k\|_{\mathbf{A}}^2 \right). \tag{26}
 \end{aligned}$$

Proof (Cont.)

By assumption, $Q := B^{-1} - \frac{1}{2}A$ is SPD. Hence

$$\|e^{k+1}\|_A^2 \leq \|e^k\|_A^2.$$

and

$$\|e^{k+1}\|_A \leq \|e^k\|_A.$$

Thus $\{\|e^k\|_A\}$ is a decreasing sequence of nonnegative numbers. By the Monotone Convergence Theorem, there is some $\gamma \geq 0$, such that

$$\|e^k\|_A \rightarrow \gamma \geq 0, \quad \text{as } k \rightarrow \infty.$$

Passing to the limit $k \rightarrow \infty$ in (26), we get

$$\lim_{k \rightarrow \infty} \|v^{k+1}\|_Q^2 = 0.$$

Using (25),

$$e^k = -A^{-1}B^{-1}v^{k+1},$$

we must have $\|e^k\|_A \rightarrow 0$.





We can generalize the last result. But first, we need the following:

Lemma

Suppose that $Q \in \mathbb{R}^{n \times n}$ is positive definite in the sense that

$$(Q\mathbf{y}, \mathbf{y}) > 0, \quad \forall \mathbf{y} \in \mathbb{R}_*^n,$$

but Q is not necessarily symmetric. Then

$$\|\mathbf{w}\|_Q = \sqrt{(Q\mathbf{w}, \mathbf{w})}, \quad \forall \mathbf{w} \in \mathbb{R}^n,$$

defines a norm.



Proof.

Suppose that Q is not symmetric, to avoid the simple case. Then

$$Q = Q_S + Q_A,$$

where

$$Q_S := \frac{1}{2} (Q + Q^T), \quad Q_A := \frac{1}{2} (Q - Q^T),$$

are the symmetric and anti-symmetric parts, respectively. Observe that $Q_S^T = Q_S$ and $Q_A^T = -Q_A$. It follows that $(Q_A y, y)$ is zero, for any $y \in \mathbb{R}^n$, because

$$(Q_A y, y) = y^T Q_A y = y^T Q_A^T y = -y^T Q_A y = -(Q_A y, y).$$

Therefore, for all $y \in \mathbb{R}_*^n$,

$$0 < (Q y, y) = (Q_S y, y) + (Q_A y, y) = (Q_S y, y).$$

Thus, Q_S is SPD. Since

$$\|w\|_{Q_S} = \sqrt{(Q_S w, w)}, \quad \forall w \in \mathbb{R}^n,$$

defines a norm, the result follows. □



Corollary

Suppose A is SPD and B is invertible. Define

$$Q := B^{-1} - \frac{1}{2}A,$$

and assume that Q is positive definite in the sense that

$$(Qy, y) > 0, \quad \forall y \in \mathbb{R}_*^n,$$

but is not necessarily symmetric. Then the GLIS is convergent with respect to the A -norm, that is,

$$\|e^k\|_A \rightarrow 0, \quad \text{as } k \rightarrow \infty.$$

Moreover, the convergence is always monotonic:

$$\|e^{k+1}\|_A^2 \leq \|e^k\|_A^2, \quad \forall k \in \mathbb{N}_0 := \mathbb{N} \cup \{0\}.$$

Proof.

Exercise.





Theorem (Householder–John)

Suppose that $A \in \mathbb{R}^{n \times n}$ is nonsingular and symmetric and $B \in \mathbb{R}^{n \times n}$ is invertible. Assume that

$$Q = B^{-1} + B^{-T} - A$$

is SPD. Then the GLIS, with error transfer matrix $T = I_n - BA$, converges unconditionally iff A is SPD. Moreover, if A is SPD, the convergence is always monotonic in the energy norm:

$$\|e^{k+1}\|_A^2 \leq \|e^k\|_A^2, \quad \forall k \in \mathbb{N}_0 := \mathbb{N} \cup \{0\}.$$

Proof.

(\Leftarrow) : Suppose that A is SPD. Recall that $\|\cdot\|_A$ defines a norm on \mathbb{R}^n . The error equation is precisely

$$e^{k+1} = (I - BA)e^k.$$


$$\begin{aligned}
 \left\| \mathbf{e}^{k+1} \right\|_A^2 &= \left((\mathbf{I} - \mathbf{B}\mathbf{A}) \mathbf{e}^k \right)^T \mathbf{A} \left((\mathbf{I} - \mathbf{B}\mathbf{A}) \mathbf{e}^k \right) \\
 &= \left(\mathbf{e}^{kT} (\mathbf{I} - \mathbf{A}\mathbf{B}^T) \right) \mathbf{A} \left((\mathbf{I} - \mathbf{B}\mathbf{A}) \mathbf{e}^k \right) \\
 &= \left(\mathbf{e}^{kT} - \mathbf{e}^{kT} \mathbf{A}\mathbf{B}^T \right) \left(\mathbf{A}\mathbf{e}^k - \mathbf{A}\mathbf{B}\mathbf{A}\mathbf{e}^k \right) \\
 &= \mathbf{e}^{kT} \mathbf{A}\mathbf{e}^k - \mathbf{e}^{kT} \mathbf{A}\mathbf{B}^T \mathbf{A}\mathbf{e}^k - \mathbf{e}^{kT} \mathbf{A}\mathbf{B}\mathbf{A}\mathbf{e}^k + \mathbf{e}^{kT} \mathbf{A}\mathbf{B}^T \mathbf{A}\mathbf{B}\mathbf{A}\mathbf{e}^k \\
 &= \left\| \mathbf{e}^k \right\|_A^2 - \mathbf{e}^{kT} \mathbf{A} \left(\mathbf{B}^T + \mathbf{B} - \mathbf{B}^T \mathbf{A}\mathbf{B} \right) \mathbf{A}\mathbf{e}^k \\
 &= \left\| \mathbf{e}^k \right\|_A^2 - \mathbf{e}^{kT} \mathbf{A}\mathbf{B}^T \left(\mathbf{B}^{-1} + \mathbf{B}^{-T} - \mathbf{A} \right) \mathbf{B}\mathbf{A}\mathbf{e}^k \\
 &= \left\| \mathbf{e}^k \right\|_A^2 - \left(\mathbf{B}\mathbf{A}\mathbf{e}^k \right)^T \left(\mathbf{B}^{-1} + \mathbf{B}^{-T} - \mathbf{A} \right) \mathbf{B}\mathbf{A}\mathbf{e}^k.
 \end{aligned}$$



11. *Journal of the American Medical Association*, 2000; 283: 2686-2692.



The energy methods described earlier prove to be a more useful tool for establishing the convergence of the Gauss-Seidel Methods than the spectral approach when the coefficient matrix is invertible.

Theorem

Suppose that $A \in \mathbb{R}^{n \times n}$ is SPD. Then the forward, backward and symmetric Gauss-Seidel methods are well defined, and the methods converge unconditionally.

Proof.

(Forward): For the forward Gauss-Seidel method, consider

$$\begin{aligned} Q_1 &= B_{GS}^{-1} + B_{GS}^{-T} - A \\ &= D - L + D - U - D + L + U \\ &= D. \end{aligned}$$

Since A is SPD, its diagonal elements are positive. Thus, D is SPD. Now we can apply the Householder-John Theorem to see that the Gauss-Seidel method converges unconditionally.



We can also use our previous corollary. For that, set

$$\begin{aligned} Q_2 &= B_{GS}^{-1} - \frac{1}{2}A. \\ &= D - L - \frac{1}{2}D + \frac{1}{2}L + \frac{1}{2}U \\ &= \frac{1}{2}D - \frac{1}{2}L + \frac{1}{2}U. \end{aligned}$$

Now, let us consider the symmetric part of Q_2 :

$$\begin{aligned} Q_{2,S} &= \frac{1}{2} \left(Q_2 + Q_2^T \right) \\ &= \frac{1}{2} \left(\frac{1}{2}D - \frac{1}{2}L + \frac{1}{2}U + \frac{1}{2}D - \frac{1}{2}U + \frac{1}{2}L \right) \\ &= \frac{1}{2}D. \end{aligned}$$

Since $Q_{2,5}$ is SPD, the corollary implies that the Gauss-Seidel method converges.



(Backward): The same techniques show that the backward Gauss-Seidel method converges.

$$B_{SGS} = (D - U)^{-1} D (D - L)^{-1}.$$

Observe that that the matrix

$$Q = B_{SGS}^{-1} + B_{SGS}^{-T} - A = A + 2U^T D^{-1} U$$

is SPD, since A is SPD. Now one can apply the Householder-John Theorem to see that the symmetric Gauss-Seidel method converges. \square