

# Numerical Calculus

Pătrulescu Flavius

Technical University of Cluj-Napoca

2024

# References

1. R.L. Burden, D.J. Faires, A.M. Burden, *Numerical Analysis*, 10th Ed., Cengage Learning, Boston, 2022.
2. A. Quarteroni, R. Sacco, F. Saleri, *Numerical Mathematics*, Springer-Verlag, New-York, 2000.

# Iterative techniques for linear systems (Burden *et al.*, 2022)

Iterative techniques are seldom used for solving linear systems of small dimension since the time required for sufficient accuracy exceeds that required for direct techniques, such as Gaussian elimination. For large systems with a high percentage of 0 entries, these techniques are efficient in terms of both computer storage and computation. Systems of this type arise frequently in circuit analysis and in the numerical solution of boundary-value problems.

# Iterative techniques for linear systems (Burden *et al.*, 2022; Quarteroni *et al.* 2000)

In general, *iterative techniques* to solve linear systems involve a process that converts the system

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

into an equivalent system of the form

$$\mathbf{x} = R\mathbf{x} + \mathbf{c}$$

for some fixed matrix (*iteration matrix*)  $R$  and vector  $\mathbf{c}$ . After the initial vector  $\mathbf{x}^{(0)}$  is selected, the sequence of approximate solution vectors is generated by computing

$$\mathbf{x}^{(k+1)} = R\mathbf{x}^{(k)} + \mathbf{c}$$

for each  $k = 1, 2, \dots$

# Iterative techniques for linear systems (Burden *et al.*, 2022; Quarteroni *et al.* 2000)

We consider a decomposition  $A = M - K$  with  $\det M \neq 0$ .  
We obtain

$$A\mathbf{x} = \mathbf{b} \Leftrightarrow (M - K)\mathbf{x} = \mathbf{b} \Leftrightarrow M\mathbf{x} = K\mathbf{x} + \mathbf{b}$$

$$\mathbf{x} = \underbrace{M^{-1}K}_{R:=} \mathbf{x} + \underbrace{M^{-1}\mathbf{b}}_{\mathbf{c}:=}$$

Find a splitting  $A = M - K$  s.t.

$R$  and  $\mathbf{c}$  are easy to evaluate

$\rho(R)$  is *small*.

# Iterative techniques for linear systems (Burden *et al.*, 2022)

$$A = D - L - U$$

$D$  is the diagonal matrix of the diagonal entries of  $A$

$-L$  is the strictly lower-triangular part of  $A$

$$L = (l_{ij}) = \begin{cases} -a_{ij}, & i > j \\ 0, & i \leq j \end{cases}$$

$-U$  is the strictly upper-triangular part of  $A$

$$U = (u_{ij}) = \begin{cases} -a_{ij}, & i < j \\ 0, & i \geq j \end{cases}$$

## Jacobi method (Burden *et al.*, 2022)

$$A\mathbf{x} = \mathbf{b} \Leftrightarrow (D - L - U)\mathbf{x} = \mathbf{b} \Leftrightarrow D\mathbf{x} = (L + U)\mathbf{x} + \mathbf{b}$$

$$\mathbf{x} = \underbrace{D^{-1}(L + U)}_{R_J :=} \mathbf{x} + \underbrace{D^{-1}\mathbf{b}}_{\mathbf{c}_J :=} \Leftrightarrow \mathbf{x} = R_J\mathbf{x} + \mathbf{c}_J$$

$$\mathbf{x}^{(k+1)} = R_J\mathbf{x}^{(k)} + \mathbf{c}_J, k \geq 0$$

Jacobi method is obtained by solving the  $i$ th equation in  $A\mathbf{x} = \mathbf{b}$  for  $x_i$  (provided that  $a_{ii} \neq 0$ ). For each  $k \geq 0$ , Jacobi method generates the components  $x_i^{(k+1)}$  of  $\mathbf{x}^{(k+1)}$  from the components of  $\mathbf{x}^{(k)}$  by

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

## Gauss-Seidel method(Burden *et al.*, 2022)

$$A\mathbf{x} = \mathbf{b} \Leftrightarrow (D - L - U)\mathbf{x} = \mathbf{b} \Leftrightarrow (D - L)\mathbf{x} = U\mathbf{x} + \mathbf{b}$$

$$\mathbf{x} = \underbrace{(D - L)^{-1}U\mathbf{x}}_{R_{GS}} + \underbrace{(D - L)^{-1}\mathbf{b}}_{\mathbf{c}_{GS}} \Leftrightarrow \mathbf{x} = R_{GS}\mathbf{x} + \mathbf{c}_{GS}$$

$$\mathbf{x}^{(k+1)} = R_{GS}\mathbf{x}^{(k)} + \mathbf{c}_{GS}, k \geq 0$$

For each  $k \geq 0$ , Gauss-Seidel method generates the components  $x_i^{(k+1)}$  of  $\mathbf{x}^{(k+1)}$  from the components of  $\mathbf{x}^{(k)}$  by

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



# Successive Over-Relaxation (S.O.R.) method (Burden et al., 2022)

Let  $\omega > 0$ .

$$\omega A\mathbf{x} = \omega \mathbf{b} \Leftrightarrow \omega(D - L - U)\mathbf{x} = \omega \mathbf{b} \Leftrightarrow \omega(D - L)\mathbf{x} = \omega U\mathbf{x} + \omega \mathbf{b}$$

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

$$\Leftrightarrow \mathbf{x} = \underbrace{(D - \omega L)^{-1}((1 - \omega)D + \omega U)}_{R(\omega) :=} \mathbf{x} + \underbrace{\omega(D - \omega L)^{-1}\mathbf{b}}_{\mathbf{c}(\omega) :=} \Leftrightarrow$$

$$\Leftrightarrow \mathbf{x} = R(\omega)\mathbf{x} + \mathbf{c}(\omega)$$

## S.O.R. (Burden *et al.*, 2022)

$$\mathbf{x}^{(k+1)} = R(\omega)\mathbf{x}^{(k)} + \mathbf{c}(\omega), \quad k \geq 0$$

For each  $k \geq 0$ , S.O.R. method generates the components  $x_i^{(k+1)}$  of  $\mathbf{x}^{(k+1)}$  from the components of  $\mathbf{x}^{(k)}$  by

$$x_i^{(k+1)} = (1 - \omega)x_i^{(k)} + \omega \underbrace{\frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} \right)}_{\text{the component generated by Gauss-Seidel method}}, \quad i = \overline{1, n}.$$

For  $\omega = 1$  S.O.R. coincides with Gauss-Seidel method. In particular, for  $\omega \in (0, 1)$ , the method is called *under-relaxation*.

## Backward Gauss-Seidel (B.G.S.) method (Quarteroni et al., 2000)

An analogue of the Gauss-Seidel method can be constructed by simply exchanging  $L$  with  $U$

$$A\mathbf{x} = \mathbf{b} \Leftrightarrow (D - L - U)\mathbf{x} = \mathbf{b} \Leftrightarrow (D - U)\mathbf{x} = L\mathbf{x} + \mathbf{b}$$

$$\mathbf{x} = \underbrace{(D - U)^{-1}L}_{R_{BGS}}\mathbf{x} + \underbrace{(D - U)^{-1}\mathbf{b}}_{\mathbf{c}_{BGS}} \Leftrightarrow \mathbf{x} = R_{BGS}\mathbf{x} + \mathbf{c}_{BGS}$$

$$\mathbf{x}^{(k+1)} = R_{BGS}\mathbf{x}^{(k)} + \mathbf{c}_{BGS}, \quad k \geq 0.$$

For each  $k \geq 0$ , backward Gauss-Seidel method generates the components  $x_i^{(k+1)}$  of  $\mathbf{x}^{(k+1)}$  from the components of  $\mathbf{x}^{(k)}$  by

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

# Symmetric Gauss-Seidel method (Quarteroni *et al.*, 2000)

The symmetric Gauss-Seidel method is obtained by combining an iteration of the Gauss-Seidel method with an iteration of the backward Gauss-Seidel method.

$$(D - L)\mathbf{x}^{(k+\frac{1}{2})} = U\mathbf{x}^{(k)} + \mathbf{b},$$

$$(D - U)\mathbf{x}^{(k+1)} = L\mathbf{x}^{(k+\frac{1}{2})} + \mathbf{b}.$$

Eliminating  $\mathbf{x}^{(k+\frac{1}{2})}$ , the following scheme is obtained

$$\mathbf{x}^{(k+1)} = R_{SGS}\mathbf{x}^{(k)} + \mathbf{c}_{SGS}, \quad k \geq 0$$

$$\begin{cases} R_{SGS} = (D - U)^{-1}L(D - L)^{-1}U \\ \mathbf{c}_{SGS} = (D - U)^{-1}[L(D - L)^{-1} + I_n]\mathbf{b}. \end{cases}$$

# General Iterations Methods (Burden *et al.*, 2022)

For any  $\mathbf{x}^{(0)} \in \mathbb{R}^n$  the sequence  $(\mathbf{x}^{(k)})_k$  defined by

$$\mathbf{x}^{(k+1)} = R\mathbf{x}^{(k)} + \mathbf{c}, \quad k \geq 0$$

converges to the unique solution of  $\mathbf{x} = R\mathbf{x} + \mathbf{c}$  if and only if

$$\rho(R) < 1 \text{ (spectral radius).}$$

Moreover, if  $\|R\| < 1$  for any induced norm then the sequence  $(\mathbf{x}^{(k)})_k$  converges to  $\mathbf{x}$  and the following error bounds hold

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|R\|^k \|\mathbf{x} - \mathbf{x}^{(0)}\|$$

and

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \frac{\|R\|^k}{1 - \|R\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|.$$

# Convergence Results (Burden *et al.*, 2022)

If  $A$  is strictly diagonally dominant, then for any choice of  $\mathbf{x}^{(0)}$  both the Jacobi and Gauss-Seidel methods give sequences  $(\mathbf{x}^{(k)})_k$  that converge to the unique solution of  $A\mathbf{x} = \mathbf{b}$ .

*Stein-Rosenberg*: If  $a_{ij} < 0$  for each  $i \neq j$  and  $a_{ii} > 0$  for each  $i = \overline{1, n}$  then one and only one of the following statements holds:

1.  $0 \leq \rho(R_{GS}) < \rho(R_J) < 1$
2.  $1 \leq \rho(R_J) < \rho(R_{GS})$
3.  $\rho(R_J) = \rho(R_{GS}) = 0$
4.  $\rho(R_J) = \rho(R_{GS}) = 1$

## S.O.R. (Burden *et al.*, 2022; Quarteroni *et al.*, 2000)

*Kahan*: If  $a_{ii} \neq 0$ ,  $i = \overline{1, n}$  then  $\rho(R(\omega)) > |\omega - 1|$ . This implies that the S.O.R. method can converge only if  $0 < \omega < 2$  (necessary condition).

*Ostrowski-Reich*: If  $A$  is symmetric and positive definite, then the S.O.R. method is convergent if and only if  $0 < \omega < 2$  (condition becomes also sufficient for convergence).

If  $A$  is strictly diagonally dominant by rows, the S.O.R. method converges if  $0 < \omega \leq 1$ .

If  $A$  is symmetric, positive definite and tridiagonal then  $\rho(R_{GS}) = (\rho(R_J))^2 < 1$  and the optimal choice of  $\omega$  for the S.O.R. method is

$$\omega = \frac{2}{1 + \sqrt{1 - (\rho(R_J))^2}}$$

With this choice of  $\omega$ , we have  $\rho(R(\omega)) = \omega - 1$ .