

# Iterative methods for solving systems of linear equations

Note Title

Let  $A$  be a  $n \times n$  matrix. Let  $Ax = b$  be a linear system having a unique solution.

Defn: A splitting of  $A$  is a decomposition  $A = S - T$  where  $S$  is non-singular.

A splitting yields an iterative method as follows-

$$Ax = (S - T)x = b$$

$$Sx = Tx + b$$

$$x = \underbrace{S^{-1}Tx}_{R} + S^{-1}b.$$

$$x = Rx + c, \quad R = S^{-1}T, \quad c = S^{-1}b.$$

This suggests the foll. iterative method :  $x_{k+1} = Rx_k + c$ .

$R$  is called the matrix of the iterative method.

Lemma 1: Let  $\|\cdot\|$  denote any induced matrix norm.

If  $\|R\| < 1$ , then  $x_{k+1} = Rx_k + c$  converges for any  $x_0$ .

Proof: Consider

$$\begin{aligned} x_{k+1} &= Rx_k + c \\ - \quad x &= Rx + c \\ x_{k+1} - x &= R(x_k - x). \end{aligned}$$

$$\text{Now } \|x_{k+1} - x\| = \|R(x_k - x)\| \leq \|R\| \cdot \underbrace{\|x_k - x\|}_{\|x_{k+1} - x\|}$$

$$\leq \|R\| \cdot \|R\| \cdot \|x_{k-1} - x\|$$

$$\text{i.e. } \|x_{k+1} - x\| \leq \|R\|^{k+1} \|x_0 - x\|. \quad = \|R\|^{k+1} \|x_{k-1} - x\| \leq \dots \leq \|R\|^{k+1} \|x_0 - x\|$$

If  $\|R\| < 1$ , then  $\|R\|^{k+1} \rightarrow 0$  as  $k \rightarrow \infty$ ,  $\therefore \|x_{k+1} - x\| \rightarrow 0$  as  $k \rightarrow \infty$  i.e.  $\{x_k\} \xrightarrow[k \rightarrow \infty]{} x$ .

Lemma 2: Let  $\rho(R)$  denote the spectral radius of  $n \times n$  matrix  $R$ .

(i.e.  $\rho(R) = \text{largest eigenvalue of } R \text{ in abs. value.}$ )

①  $\rho(R) \leq \|R\|$ , for every induced matrix norm.

② For every  $\varepsilon > 0$ ,  $\exists \|\cdot\|_*$  such that  $\|R\|_* \leq \rho(R) + \varepsilon$ .  
 ↓    ↑  
 (depending on  $R$  &  $\varepsilon$ )    an induced norm.

Proof: ① Let  $\lambda$  be s.t.  $\rho(R) = |\lambda|$  & let  $x$  be the corr. eigenvector so that  
 $Rx = \lambda x$ .

$$\|R\| = \max_{y \neq 0} \frac{\|Ry\|}{\|y\|} \geq \frac{\|Rx\|}{\|x\|} = \frac{\|\lambda x\|}{\|x\|} = |\lambda| = \delta(R).$$

② Refer to Lemma 6.5. (Demmel).

Theorem: The iterative method  $x_{k+1} = R x_k + c$  converges to the solution of  $Ax = b$ ,  $\forall x_0 \neq b$

$$\Leftrightarrow \delta(R) < 1.$$

Proof: If  $\delta(R) \geq 1$ , let  $|\lambda| = \delta(R)$  & let  $Ry = \lambda y$ .

Suppose  $y + x = x_0$  i.e.  $y = x_0 - x$ .

then  $Ry = R(x_0 - x) = \lambda(x_0 - x)$ . ✓

$$x_{k+1} - x = R(x_k - x) = R^2(x_{k-1} - x) = \dots = R^{k+1}(x_0 - x) = \lambda^{k+1}(x_0 - x)$$

a contradiction to the assumption.  $\not\rightarrow 0$  as  $k \rightarrow \infty$

Conversely, if  $\delta(R) < 1$ , choose the norm  $\|\cdot\|_\infty$  such that

$\|R\|_\infty < 1$ , then by lemma 1, the iterative method converges.

Questions to be considered when working with iterative methods -

1) Is the given method convergent? I.e. does there exist a matrix norm such that  $\|R\| < 1$ ? Or is  $\delta(R) < 1$ ?

2) Given 2 iterative methods, which converges faster?

(i.e. which matrix has a smaller spectral radius?)

(smaller spectral radius  $\Rightarrow$  faster convergence)

We will study 3 iterative methods for  $Ax = b$ :

- Jacobi
- Gauss-Seidel
- SOR

(I) Jacobi's method: Assume  $a_{ii} \neq 0 \ \forall i$ .

Split A as  $A = \underbrace{D}_{S} - \underbrace{E}_{-F} - F$ , where  $D = \text{diag}(A)$

$$\begin{bmatrix} D \\ -E \end{bmatrix}$$

$-E = \text{lower } \Delta \text{ of } A$

$-F = \text{upper } \Delta \text{ of } A$

The iterative method is  $x_{k+1} = D^{-1}(E+F)x_k + D^{-1}b$ .

i.e.  $x_{k+1} = Jx_k + c$ , where  $c = D^{-1}b$   
 $J$  is called the Jacobi matrix of  $A$ .

$$\begin{aligned} J &= D^{-1}(E+F) = D^{-1}E + D^{-1}F = I - I + D^{-1}E + D^{-1}F \\ &= I - D^{-1}(D - E - F) \\ &= I - D^{-1}A \end{aligned}$$

### Implementing Jacobi method

$$Ax = b : \quad a_{11}x_1 + \dots + a_{1n}x_n = b_1 \\ \vdots + a_{21}x_2 + \dots = b_2$$

$$a_{n1}x_1 + \dots + a_{nn}x_n = b_n$$

$$x_0 = \begin{pmatrix} x_1^0 \\ x_2^0 \\ \vdots \\ x_n^0 \end{pmatrix}$$

$$x_1 = \frac{1}{a_{11}}(b_1 - a_{12}x_2 - \dots - a_{1n}x_n)$$

$$x_2 = \frac{1}{a_{22}}(b_2 - a_{21}x_1 - \dots - a_{2n}x_n)$$

$$\vdots \quad x_n = \frac{1}{a_{nn}}(b_n - a_{n1}x_1 - \dots - a_{n,n-1}x_{n-1})$$

} \*

- ① Choose an appropriate initial approximation  $x^0 = (x_1^0, \dots, x_n^0)$   
& substitute in the RHS of \*, to get the first approximation  $x^1 = (x_1^1, \dots, x_n^1)$

- ② Substitute  $x^1$  in the RHS of \* to get the second approximation  
 $x^2 = (x_1^2, \dots, x_n^2)$  & so on.

To see what happens at the  $(k+1)^{\text{th}}$  step,

$$\text{consider } J = I - D^{-1}A$$

$$\text{so } x^{k+1} = Jx^k + c$$

$$= (I - D^{-1}A)x^k + c$$

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

$$= x^k - D^{-1}(Ax^k - b)$$

$$\boxed{\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}}$$

$$\left\{ \begin{array}{l} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 = b_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 = b_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 = b_3. \end{array} \right.$$

$$\rightarrow x_1^1 = (b_1 - a_{12}x_2^0 - a_{13}x_3^0)/a_{11}$$

$$x_2^1 = (b_2 - a_{21}x_1^0 - a_{23}x_3^0)/a_{22}$$

$$x_3^1 = (b_3 - a_{31}x_1^0 - a_{32}x_2^0)/a_{33}$$

$$\boxed{x^0 = (x_1^0, x_2^0, x_3^0)}$$

$$\begin{pmatrix} x_1^{k+1} \\ x_2^{k+1} \\ x_3^{k+1} \end{pmatrix}$$

$$\rightarrow x_i^{k+1} = x_i^k - \frac{1}{a_{ii}}(a_{i1}x_1^k - \dots - a_{in}x_n^k - b_i)$$

$$\begin{aligned}
 &= (a_{ii}x_i^k - a_{i1}x_1^k - \dots - a_{in}x_n^k - \dots - a_{in}x_n^k + b_i) / a_{ii} \\
 &= \frac{1}{a_{ii}} (b_i - a_{i1}x_1^k - \dots - \widehat{a_{ii}x_i^k} - \dots - a_{in}x_n^k)
 \end{aligned}$$

Algorithm: (one step of Jacobi's method)

for  $j = 1$  to  $n$

$$x_j^{k+1} = \frac{1}{a_{jj}} \left( b_j - \sum_{\substack{i=1 \\ i \neq j}}^n a_{ji} x_i^k \right)$$

end.

## (II) Gauss-Seidel method

The number of allocations required in Jacobi's method can be reduced using the following trick: use  $x_1^{k+1}$  to calculate  $x_2^{k+1}$  and so on.

$$\text{Then } x_1^{k+1} = \frac{1}{a_{11}} (b_1 - a_{12}x_2^k - \dots - a_{1n}x_n^k)$$

$$x_2^{k+1} = \frac{1}{a_{22}} (b_2 - a_{21}x_1^{k+1} - a_{23}x_3^k - \dots - a_{2n}x_n^k)$$

$$x_3^{k+1} = \frac{1}{a_{33}} (b_3 - a_{31}x_1^{k+1} - a_{32}x_2^{k+1} - a_{34}x_4^k - \dots - a_{3n}x_n^k)$$

& so on.

$$E = \begin{pmatrix} 0 & 0 & \dots \\ a_{21} & 0 & \dots \\ a_{31} & a_{32} & 0 & \dots \end{pmatrix}$$

This corresponds to the iterative method -

$$\rightarrow X^{k+1} = D^{-1}(E X^{k+1} + F X^k + b)$$

$$\text{i.e. } D X^{k+1} = E X^{k+1} + F X^k + b$$

$$\text{i.e. } (D - E) X^{k+1} = F X^k + b$$

$$\text{i.e. } X^{k+1} = \underbrace{(D - E)^{-1} F X^k}_{\text{denote by } L_1} + \underbrace{(D - E)^{-1} b}_{c}.$$

$$X^{k+1} = L_1 X^k + c.$$

$L_1$  is the matrix of this method & is called the Gauss-Seidel matrix of  $A$ .

Splitting is  
 $A = D - E - F$ .  
Let  $S = D - E$ ,  $T = -F$ .  
 $S$  is invertible since  $D$  is.

Algorithm (one step of G-S)

for  $j = 1$  to  $n$

$$x_j^{k+1} = \frac{1}{a_{jj}} \left( b_j - \sum_{i=1}^{j-1} a_{ji} x_i^{k+1} - \sum_{i=j+1}^n a_{ji} x_i^k \right)$$

↓  
updated  
 $x_i$ 's

↑  
old  $x_i$ 's

end.

(III) Successive over-relaxation:

$\omega$  : relaxation parameter ( $\omega \neq 0$ ).

Idea: to improve on G-S by taking appropriate weighted average of  $x_j^{k+1}$  &  $x_i^k$ 's:

$$\text{in SOR } x_j^{k+1} = (1-\omega)x_j^k + \omega(x_i^{k+1}'s \text{ & } x_i^k's)_{i \neq j}$$

For this, the splitting of  $A$  is defined as-

$$A = \left( \frac{D}{\omega} - E \right) - \left( \frac{1-\omega}{\omega} D + F \right) \quad \begin{matrix} \text{(this simplifies} \\ \text{to } D - E - F \end{matrix}$$

The associated iterative method is -

$$x_{k+1} = \left( \frac{D}{\omega} - E \right)^{-1} \left( \frac{1-\omega}{\omega} D + F \right) x_k + c, \text{ where}$$

$$(D - \omega E) x^{k+1} = ((1-\omega)D + \omega F) x^k + b = (D - \omega D + \omega F) x^k + b \quad c = \left( \frac{D}{\omega} - E \right)^{-1} b.$$

The matrix of the method is -

$$L_\omega = \left( \frac{D}{\omega} - E \right)^{-1} \left( \frac{1-\omega}{\omega} D + F \right)$$

$$\text{OR } (D - \omega E)^{-1} ((1-\omega)D + \omega F)$$

If  $\omega > 1$ , the method is called "over-relaxation"

$\omega < 1$ ,

"under-relaxation".

Implementing SOR :

$$x^{k+1} = [D - \omega E]^{-1} \left[ D x^k - \omega D x^k + \omega F x^k + b \right]$$

$$x^{k+1} = L_\omega x^k + c$$

This corresponds to equations -

$$x_1^{k+1} = \frac{1}{a_{11}} \left( a_{11} x_1^k - \omega (a_{11} x_1^k + \dots + a_{1n} x_n^k - b_1) \right)$$

$$(1-\omega) x_1^k + \frac{\omega}{a_{11}} (b_1 - a_{12} x_2^k - \dots - a_{1n} x_n^k)$$

$$x_2^{k+1} = \frac{1}{a_{22}} \left( a_{22} x_2^k - \omega (a_{21} x_1^{k+1} + \dots + a_{2n} x_n^k - b_2) \right)$$

$$(1-\omega) x_2^k + \frac{\omega}{a_{22}} (b_2 - a_{21} x_1^{k+1} - \dots - a_{2n} x_n^k)$$

$$x_i^{k+1} = \frac{1}{a_{ii}} \left( a_{ii} x_i^k - \omega (a_{in} x_1^{k+1} + \dots + a_{i,i-1} x_{i-1}^{k+1} + \dots + a_{in} x_n^k - b_i) \right)$$

$$(1-\omega) x_i^k + \frac{\omega}{a_{ii}} (b_i - a_{i1} x_1^{k+1} - \dots - a_{i,i-1} x_{i-1}^{k+1} - \dots - a_{in} x_n^k)$$

By choosing appropriate  $\omega$ , it is possible to attain a higher rate of convergence than C-S.

Algorithm (one step of SOR)

for  $j=1$  to  $n$

$$x_j^{k+1} = (1-\omega) x_j^k + \frac{\omega}{a_{jj}} \left[ b_j - \sum_{i=1}^{j-1} a_{ji} x_i^{k+1} - \sum_{i=j+1}^n a_{ji} x_i^k \right]$$

end.

Defn: The rate of convergence of  $x_{k+1} = R x_k + c$  is

See sections

6.5.3 & 6.5.4  
of Demmel

$$\gamma(R) = -\log_{10} \beta(R).$$

Compare  $\beta(I)$ ,  $\beta(L_1)$ ,  $\beta(L_\omega)$  for some values  $\omega$

so smaller the  $\beta(R)$ , higher the rate of convergence i.e. more number of correct decimal places computed per iteration.

Key facts about Jacobi's method -

- ① requires non-zero diagonal entries (can be accomplished by permuting rows/columns if not already true)
- ② requires  $2n$  memory allocations ( $\dim A = n \times n$ ).
- ③ components do not depend on one another, so they can be computed simultaneously.
- ④ does not always converge; converges for sure when  $A$  is SDD. i.e.  $|a_{ii}| > \sum_{\substack{j \neq i \\ j=1}}^n |a_{ij}|$

Key facts about G-S method.

- ① requires non-zero diagonal entries.
- ② requires  $n$  memory allocations at each step.
- ③ each component depends on previous ones, so must be computed successively.
- ④ converges if  $A$  is HPD (weaker than SDD)
- ⑤ when they converge together, G-S convergence twice as fast as Jacobi.

Key facts/questions about SOR:

- ① requires non-zero diagonal entries.
- ② 2 questions : (i) for what  $\omega$  is  $S(L\omega) < 1$  ?

Is there an interval  $I \subset \mathbb{R}$  such that

$$S(L\omega) < 1 \quad \forall \omega \in I.$$

(ii) Is there an optimal  $\omega_0$  such that

$$S(L\omega_0) = \inf_{\omega \in I} S(L\omega).$$