

Numerical Methods I: Iterative solvers for $Ax = b$

Georg Stadler
Courant Institute, NYU
stadler@cims.nyu.edu

December 7, 2017

Iterative solution of (symmetric) linear systems

Target problems: very **large** ($n = 10^5, 10^6, \dots$), A is usually **sparse** and has specific **properties**.

To solve

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

we construct a sequence

$$\mathbf{x}_1, \mathbf{x}_2, \dots$$

of iterates that converges as fast as possible to the solution \mathbf{x} , where \mathbf{x}_{k+1} can be computed from $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ with as little cost as possible (e.g., one matrix-vector multiplication).

Iterative solution of (symmetric) linear systems

Let Q be invertible, then $Q \in \mathbb{R}^{n \times n}$, invertible

$$Ax = b \Leftrightarrow Q^{-1}(b - Ax) = 0$$

$$\Leftrightarrow (I - Q^{-1}A)x + Q^{-1}b = x$$

$$\Leftrightarrow Gx + c = x$$

Fixed point method:

$$x_{k+1} = Gx_k + c,$$

$$k = 0, 1, 2, \dots$$

$x_0 \in \mathbb{R}^n$ initialization

Iterative solution of (symmetric) linear systems

Theorem: The fixed point method $\mathbf{x}_{k+1} = G\mathbf{x}_k + \mathbf{c}$ with an invertible G converges for each starting point \mathbf{x}_0 if and only if

$$\rho(G) < 1,$$

where $\rho(G)$ is the largest eigenvalue of G (i.e., the spectral radius).

Proof: G spd. $\exists Q$ orthogonal s.t. $\boxed{\mathbf{x}^* = \bar{A}^{-1}\mathbf{b}}$

$$Q G Q^T = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

Since $|\lambda_i| \leq \rho(G) < 1 \implies \Lambda^k \rightarrow 0$ as $k \rightarrow \infty$

$$\implies (Q \Lambda Q^T)^k = G^k = Q \Lambda^k Q^T \rightarrow 0 \text{ as } k \rightarrow \infty$$

Since $\rho(G) \leq \|G\|$ for any matrix norm induced by vector norm.

$$\mathbf{x}_k - \mathbf{x}^* = G\mathbf{x}_{k-1} + \mathbf{c} - (G\mathbf{x}^* + \mathbf{c}) = G(\mathbf{x}_{k-1} - \mathbf{x}^*) = \dots = G^k(\mathbf{x}_0 - \mathbf{x}^*)$$

$$\implies \|\mathbf{x}_k - \mathbf{x}^*\| \leq \|G\|^k \|\mathbf{x}_0 - \mathbf{x}^*\| \rightarrow 0 \text{ if } \|G\| < 1.$$

Iterative solution of (symmetric) linear systems

Choices for Q :

- ▶ Choose $Q = I \dots$ **Richardson method**

$$G = I - Q^T A = I - A, \quad x_{k+1} = x_k - A x_k + b$$

$$A \text{ spd}, \quad \rho(G) = \rho(I - A) = \max \{ |1 - \lambda_{\min}(A)|, |1 - \lambda_{\max}(A)| \}$$

$$\implies \rho(G) < 1 \text{ if } \lambda_{\max}(A) < 2$$

For more choices, consider $A = L + D + U$, where D is diagonal, L and U are lower and upper triangular with zero diagonal.

The diagram illustrates the decomposition of a matrix A into three components: L , D , and U . On the left, a square box contains the letter A . To its right is an equals sign, followed by three matrices separated by plus signs. The first matrix, labeled L below it, is a lower triangular matrix with a diagonal line from top-left to bottom-right, an asterisk in the bottom-left cell, and zeros in the top-right cell and along the diagonal. The second matrix, labeled D below it, is a diagonal matrix with an asterisk in the top-left cell, an asterisk in the bottom-right cell, and zeros in the top-right cell and bottom-left cell. The third matrix, labeled U below it, is an upper triangular matrix with an asterisk in the top-right cell, an asterisk in the bottom-right cell, and zeros in the top-left cell and along the diagonal.

Iterative solution of (symmetric) linear systems

$$A = L + D + U$$

- Choose $Q = D$... **Jacobi method**

$$G = I - Q^{-1}A = I - D^{-1}A = I - D^{-1}(L+D+U) = -D^{-1}(L+U)$$
$$\boxed{x_{k+1} = D^{-1}(L+U)x_k + D^{-1}b}$$

Theorem: The Jacobi method converges for any starting point x_0 to the solution of $Ax = b$ if A is strictly diagonal dominant, i.e.,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad \text{for } i = 1, \dots, n.$$

Proof:

$$\rho(G) = \rho(-D^{-1}(L+U)) \leq \|D^{-1}(L+U)\|_{\infty} = \max_i \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|} < 1 \quad \text{since } A \text{ is strictly diagonal dom.}$$

Iterative solution of (symmetric) linear systems

$$A = L + D + U$$

- ▶ Choose $Q = D + L$... **Gauss-Seidel method**

$$G = I - Q^{-1}A = I - (D+L)^{-1}A = I - (D+L)^{-1}(L+D+U)$$
$$= - (D+L)^{-1}U$$
$$\boxed{x_{k+1} = (D+L)^{-1}U x_k + (D+L)^{-1}b}$$

Theorem: The Gauss-Seidel method converges for any starting point x_0 if A is spd.

Iterative solution of (symmetric) linear systems

Relaxation methods: Use linear combination between new and previous iterate:

$$G_\omega = \omega G + (1-\omega)I$$

$$\mathbf{x}_{k+1} = \omega(G\mathbf{x}_k + \mathbf{c}) + (1-\omega)\mathbf{x}_k = G_\omega\mathbf{x}_k + \omega\mathbf{c},$$

where $\omega \in [0, 1]$ is a **damping/relaxation parameter** (sometimes, $\omega > 1$ is used, leading to overrelaxation). Target is to choose ω such that $\rho(G_\omega)$ is as small as possible.

Iterative solution of (symmetric) linear systems

Def: A fixed point method $\mathbf{x}_{k+1} = G\mathbf{x}_k + \mathbf{c}$ with $G = G(A)$ is called *symmetrizable* if for any spd matrix A , $I - G$ is similar to an spd matrix. That is, $\exists W \in \mathbb{R}^{n \times n}$ invertible such that $W(I - G)W^{-1}$ is spd.

Examples: Richardson $G = I - A$, so $I - G = A$ is spd
[$W = I$]

Jacobi: $G = I - D^{-1}A$, $W = D^{-\frac{1}{2}}$
 $D^{-\frac{1}{2}}(I - G)D^{\frac{1}{2}} = D^{-\frac{1}{2}}(I - I + D^{-1}A)D^{\frac{1}{2}} = D^{\frac{1}{2}}AD^{\frac{1}{2}}$ spd if A spd.

Iterative solution of (symmetric) linear systems

Let the fixed point method be symmetrizable, and A an spd matrix. Then all eigenvalues of G are real and less than 1.

Proof: method is symmetrizable $\Rightarrow I-G$ is similar to an spd matrix \Rightarrow eigenvalues of $I-G$ are real and positive \Rightarrow eigenvalues of G are real and < 1 .



Iterative solution of (symmetric) linear systems

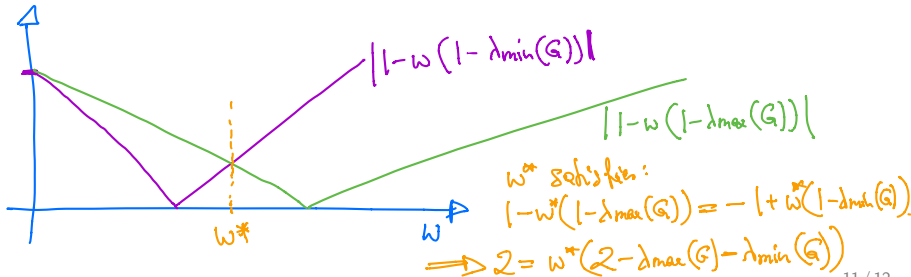
Finding the optimal damping parameter:

Symmetrizable method, $\lambda_{\min} \leq \lambda_{\max} < 1$ extreme eigenvalues of G .

Eigenvalues of $G_w = wG + (1-w)I$:

$$\lambda_i(G_w) = w\lambda_i(G) + (1-w) = 1 - w(1 - \lambda_i(G)) < 1$$

$$\rho(G_w) = \max \left\{ |1 - w(1 - \lambda_{\min}(G))|, |1 - w(1 - \lambda_{\max}(G))| \right\}$$



Iterative solution of (symmetric) linear systems

Krylov methods:

Idea: Build a basis for the Krylov subspace $\{\mathbf{r}_0, A\mathbf{r}_0, A^2\mathbf{r}_0 \dots\}$ and reduce residual optimally in that space.

- ▶ spd matrices: **Conjugate gradient (CG)** method
- ▶ symmetric matrices: **Minimal residual method (MINRES)**
- ▶ general matrices: **Generalized residual method (GMRES), BiCG, BiCGSTAB**

Iterative solution of (symmetric) linear systems

Krylov methods:

Idea: Build a basis for the Krylov subspace $\{\mathbf{r}_0, A\mathbf{r}_0, A^2\mathbf{r}_0 \dots\}$ and reduce residual optimally in that space.

- ▶ spd matrices: **Conjugate gradient (CG)** method
- ▶ symmetric matrices: **Minimal residual method (MINRES)**
- ▶ general matrices: **Generalized residual method (GMRES), BiCG, BiCGSTAB**

Properties:

Do not require eigenvalue estimates; require usually one matrix-vector multiplication per iteration; convergence depends on eigenvalue structure of matrix (clustering of eigenvalues aids convergence). Availability of a good preconditioner is often important. Some methods require storage of iteration vectors.

The conjugate gradient method (CG, pcg in Matlab)

Solve $Ax=b$, A spd, $x, b \in \mathbb{R}^n$,
 $A \in \mathbb{R}^{n \times n}$ unique solution

Solving $Ax=b \iff \min_{x \in \mathbb{R}^n} \frac{1}{2} x^T A x - b^T x$

A -weighted norm

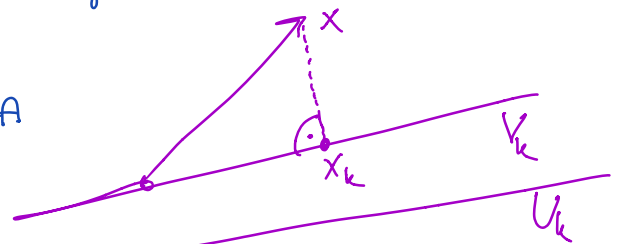
$\|y\|_A = \sqrt{(y, Ay)}$ is a norm

$x = A^{-1}b$, Search approximations $x_k \in \mathbb{R}^n$ of x
 in an affine subspace $V_k = x_0 + U_k$,
 $U_k \subset \mathbb{R}^n$ subspace

$$x_k = \arg \min_{y \in V_k} \|y - x\|_A$$



$$(x - x_k, u)_A := (A(x - x_k), u) = 0 \quad \text{for all } u \in U_k$$



Why choose A -weighted norm?

$r_k = b - Ax_k$ residual

$$(x - x_k, u)_A = (A(x - x_k), u) = (r_k, u)$$

does not include x ,
 which is
 not known.

i.e.: residual $\perp u$ in Euclidean inner product

\Leftrightarrow error $x-x_k \perp$ to u in A -weighted inner product.

Let p_1, \dots, p_k A -orthogonal basis in U_k , i.e.

orthog. proj on U_k

$$(p_i, p_j)_A = \delta_{ij}$$

A -orthogonal
also called A -conjugate

$$x_k = P_k x = x_0 + \sum_{j=1}^k \frac{(p_j, x-x_0)_A}{(p_j, p_j)_A} p_j =$$

$$= x_0 + \sum_{j=1}^k \frac{(p_j, \overbrace{x-x_0}^{r_0})_A}{(p_j, p_j)_A} p_j$$

$$= x_0 + \sum_{j=1}^k \underbrace{\frac{(p_j, r_0)_A}{(p_j, p_j)_A}}_{\alpha_k} p_j \leftarrow \text{indep. of } x$$

$$x_k = x_{k-1} + \alpha_k p_k$$

$$r_k = A(x-x_k) = A(x-x_{k-1} - \alpha_k p_k) = r_{k-1} - A\alpha_k p_k$$

Specific choices: $V_k = x_0 + U_k$

$$U_k = \text{span} \{ r_0, A r_0, A^2 r_0, \dots, A^{k-1} r_0 \}$$

Krylov spaces

$$p_{k+1} = r_k - \frac{(r_k, p_k)_A}{(p_k, p_k)_A} p_k$$

β_{k+1}

- We have:
- Method to choose V_k , and compute A-orthog. basis
 - method to compute x_k that does not need x
 - recurrence for r_k

Algorithm: A spd, x_0 starting value

$$p_1 = r_0 = b - Ax_0$$

for $k=1, 2, \dots$

$$\alpha_k = \frac{(r_{k-1}, r_{k-1})}{(p_k, A p_k)}$$

$$x_k := x_{k-1} + \alpha_k p_k$$

if converged stop

$$r_k := r_{k-1} - \alpha_k A p_k$$

$$\beta_{k+1} := \frac{(r_k, r_k)}{(r_{k-1}, r_{k-1})}$$

$$p_{k+1} := r_k + \beta_k p_k$$

end

- Converges after at most n iterations because

$$x_k = \arg \min_{y \in V_k} \|x - y\|_A, \quad V_k = \mathbb{R}^k \text{ after } k \text{ iterations}$$

Cost:

1 application of A

2 inner products

- Converges much faster depending on the eigenvalues of A , in particular if eigenvalue of A are clustered, i.e.

