

CS 450: Numerical Analysis  
Lecture 28  
Chapter 11 Partial Differential Equations  
Solving Sparse Linear Systems

Edgar Solomonik

Department of Computer Science  
University of Illinois at Urbana-Champaign

April 29, 2018

## Sparse Linear Systems

- ▶ Finite-difference and finite-element methods for time-independent PDEs give rise to sparse linear systems:

- ▶ *typified by the 2D Laplace equation, where for both finite differences and FEM,*

$$\underbrace{(\mathbf{I} \otimes \mathbf{T} + \mathbf{T} \otimes \mathbf{I})}_{\mathbf{A}} \mathbf{x} = \mathbf{b} \quad \text{where} \quad \mathbf{T} = \frac{1}{h^2} \begin{bmatrix} -2 & 1 & & \\ 1 & \ddots & \ddots & \\ & & \ddots & \\ & & & \ddots \end{bmatrix}$$

- ▶ *often have  $O(1)$  nonzeros per row/column of the matrix,*
    - ▶ *for simple/regular problems matrices are near-**Toeplitz** (same entry along each (sub/super)diagonal), permitting fast solvers via e.g. FFT.*
  - ▶ **Direct methods** apply LU or other factorization to  $\mathbf{A}$ , while **iterative methods** refine  $\mathbf{x}$  by minimizing  $\mathbf{r} = \mathbf{A}\mathbf{x} - \mathbf{b}$ , e.g. via Krylov subspace methods.
    - ▶ *Direct methods provide a high-accuracy solution, but may not be effective at leveraging sparsity to reduce cost.*
    - ▶ *Iterative methods effectively leverage sparsity by computing matrix-vector products with  $\mathbf{A}$ , but may require many iterations to achieve high-accuracy.*

## Direct Methods for Sparse Linear Systems

- ▶ It helps to think of  $A$  as the adjacency matrix of graph  $G = (V, E)$  where  $V = \{1, \dots, n\}$  and  $a_{ij} \neq 0$  if and only if  $(i, j) \in E$ :

*The graph is invariant under permutation of vertices, i.e. for any permutation matrix  $P$ , reordering  $V$  accordingly transforms the adjacency matrix of the graph into  $P^T A P$ . Note that such reorderings of variables essentially do not change the linear system,*

$$P^T A P \underbrace{P^T x}_{\hat{x}} = b$$

- ▶ Factorizing the  $l$ th row/column in Gaussian elimination corresponds to removing node  $i$ , with nonzeros (new edges) introduced for each  $k, l$  such that  $(i, k)$  and  $(i, l)$  are in the graph.
  - ▶ *creates clique (fully connected subgraph) among neighbors of vertex  $i$ ,*
  - ▶ *different orderings of vertices can result in radically different amounts of fill,*
  - ▶ *finding optimal ordering to reduce fill is NP complete.*

## Vertex Orderings for Sparse Direct Methods

- ▶ Select the node of minimum degree at each step of factorization:  
*Each step minimizes work, but not necessarily amount of fill at that step (depends on whether neighbors are already connected).*
- ▶ Graph partitioning also serves to bound fill, remove vertex separator  $S \subset V$  so that  $V \setminus S = V_1 \cup \dots \cup V_k$  become disconnected, then order  $V_1, \dots, V_k, S$ :  
*Matrix takes on the form*

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & & & \mathbf{A}_{1S} \\ & \ddots & & \vdots \\ & & \mathbf{A}_{kk} & \mathbf{A}_{kS} \\ \mathbf{A}_{S1} & \cdots & \mathbf{A}_{Sk} & \mathbf{A}_{SS} \end{bmatrix}$$

where each  $\mathbf{A}_{ii}$  for  $i \in \{1, \dots, k\}$  can be factored independently.

- ▶ *Nested dissection* ordering partitions graph into halves recursively, ordering each separator last.

## Sparse Iterative Methods

- ▶ Direct sparse factorization is ineffective in memory usage and/or cost for many typical sparsity matrices, motivating iterative methods:

*Many such basic iterative methods corresponds to linear fixed point iterations defined by linear system:*

$$\mathbf{M}\mathbf{x}_{k+1} = \mathbf{N}\mathbf{x}_k + \mathbf{b}$$

*which can be written in the form  $\mathbf{g}(\mathbf{x}) = \mathbf{M}^{-1}\mathbf{N}\mathbf{x} + \mathbf{M}^{-1}\mathbf{b}$ . We desire to have a fixed point whenever  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , which implies*

$$\mathbf{M}\mathbf{x} = \mathbf{N}\mathbf{x} + \mathbf{A}\mathbf{x},$$

*so generally  $\mathbf{M}$  and  $\mathbf{N}$  are chosen so that  $\mathbf{M} - \mathbf{N} = \mathbf{A}$ . Moreover, to achieve convergence we need  $\rho(\mathbf{g}) = \rho(\mathbf{M}^{-1}\mathbf{N}) < 1$ .*

## Sparse Iterative Methods

- ▶ The *Jacobi method* is the simplest iterative solver:

*We split up  $A = D + L + U$  where  $D$  is diagonal while  $L$  and  $U^T$  are strictly lower triangular. Jacobi iteration uses a fixed point scheme with  $M = D$  and  $N = -(L + U)$ , yielding a diagonal system of equations,*

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

*The cost of each iteration of Jacobi is proportional to SpMV with  $A$ .*

- ▶ The Jacobi method converges if  $A$  is strictly row-diagonally-dominant:

*A strictly row-diagonally-dominant  $A$  satisfies  $|a_{ii}| < \sum_{j \neq i} |a_{ij}|$ , which implies that for  $B = D^{-1}(L + U)$ ,*

$$\left| \sum_j b_{ij} \right| = \left| \sum_{j \neq i} a_{ij}/a_{ii} \right| < 1,$$

*so that  $\rho(M^{-1}N) = \rho(B) < 1$ . For the Laplace problem, this condition holds, and the method corresponds to averaging neighbors. However, the coefficient in linear convergence is  $\cos(\pi h)$  (goes to 1 with decreasing  $h$ ).*

## Gauss-Seidel Method

- ▶ The Jacobi method takes weighted sums of  $\mathbf{x}^{(k)}$  to produce each entry of  $\mathbf{x}^{(k+1)}$ , while Gauss-Seidel uses the latest available values, i.e. to compute  $x_i^{(k+1)}$  it uses a weighted sum of

$$x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_i^{(k)}, \dots, x_n^{(k)}.$$

*We can define the method by the splitting  $\mathbf{M} = \mathbf{D} + \mathbf{L}$  so that we have*

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

*The Gauss-Seidel method performs an in-order traversal of the directed acyclic adjacency graph induced by the vertex ordering, and updates each vertex by taking newly computed values from incoming edges and values from the previous iteration from outgoing edges.*

- ▶ Gauss-Seidel provides somewhat better convergence than Jacobi:  
*Convergence and efficiency depend on vertex ordering and connectivity:*
  - ▶ *for 2-D Poisson, spectral radius is  $\cos^2(\pi h)$ ,*
  - ▶ *computational cost is same as Jacobi, but less parallelism available.*

## Successive Over-Relaxation

- ▶ The *successive over-relaxation* (SOR) method seeks to improve the spectral radius achieved by Gauss-Seidel, by choosing

$$M = \frac{1}{\omega}D + L, \quad N = \left(\frac{1}{\omega} - 1\right)D - U$$

*In the resulting iterative scheme, we have*

$$\left(\frac{1}{\omega}D + L\right)\mathbf{x}^{(k+1)} = ((1 - \omega)D + U)\mathbf{x}^{(k)}$$

*so that if  $\mathbf{x}_{GS}^{(k+1)}$  is the iterate produced by Gauss-Seidel, SOR instead produces*

$$\mathbf{x}^{(k+1)} = (1 - \omega)\mathbf{x}^{(k)} + \omega\mathbf{x}_{GS}^{(k+1)}.$$

- ▶ The parameter  $\omega$  in SOR controls the ‘step-size’ of the iterative method:
  - ▶ *over-relaxation corresponds to  $\omega > 1$ ,*
  - ▶ *under-relaxation corresponds to  $\omega < 1$ ,*
  - ▶ *generally best choice of  $\omega \in (0, 2)$  is hard to determine.*



## Conjugate Gradient

- ▶ The solution to  $\mathbf{Ax} = \mathbf{b}$  is a minima of the quadratic optimization problem,

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

*We can leverage this and employ optimization methods such as conjugate gradient (CG) in the case when  $\mathbf{A}$  is SPD.*

- ▶ Conjugate gradient works by picking  $\mathbf{A}$ -orthogonal descent directions  
*Ensures search directions make progress and converge in at most  $n$  iterations.*

- ▶ The convergence rate of CG is linear with coefficient  $\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$  where  $\kappa = \text{cond}(\mathbf{A})$ :

*This convergence rate motivates techniques to improve the conditioning of  $\mathbf{A}$  to accelerate convergence.*

## Preconditioning

- ▶ Preconditioning techniques choose matrix  $M \approx A$  and solve the linear system

$$M^{-1}Ax = M^{-1}b$$

*For example at each step of CG, we replace multiplication with  $A$  by multiplication by  $M^{-1}A$ :*

- ▶ *cost of iteration depends on difficulty of applying  $M^{-1}$ ,*
  - ▶ *convergence rate depends on how close  $M$  is to  $A$ .*
- ▶  $M$  is usually chosen to be an effective approximation to  $A$  with a simple structure:
  - ▶ *Jacobi preconditioning takes  $M = D$  where  $D$  is the diagonal of  $A$ ,*
  - ▶ *incomplete factorization (ILU) uses  $A \approx LU$  where the sparsity pattern of  $L$  and  $U$  is restricted to that of  $A$  (factorization is then generally inexact), and employs  $M = LU$ .*