

# CS 450: Numerical Analysis<sup>1</sup>

## Partial Differential Equations

University of Illinois at Urbana-Champaign

---

<sup>1</sup> *These slides have been drafted by Edgar Solomonik as lecture templates and supplementary material for the book “Scientific Computing: An Introductory Survey” by Michael T. Heath ([slides](#)).*

# Partial Differential Equations

- ▶ *Partial differential equations (PDEs)* describe physical laws and other continuous phenomena:
  - ▶ *They contain partial derivatives in multiple variables.*
  - ▶ *Examples include: electromagnetism, fluid flow, quantum mechanics, and general relativity.*
- ▶ The *advection PDE* describes basic phenomena in fluid flow,

$$u_t = -a(t, x)u_x$$

where  $u_t = \partial u / \partial t$  and  $u_x = \partial u / \partial x$ .

- ▶ *Generally, we impose an initial condition with respect to  $t$ , i.e.,  $u(0, x) = u_0(x)$ .*
- ▶ *Boundary conditions* are also often exposed on boundary of domain.
- ▶ When  $a(t, x) = c$  this is the *Cauchy problem* with solution,

$$u(t, x) = u_0(x - ct).$$

# Types of PDEs

- ▶ Some of the most important PDEs are *second order*:

- ▶ Heat equation (diffusion),  $u_t = u_{xx}$
- ▶ Wave equation (oscillation),  $u_{tt} = u_{xx}$
- ▶ Laplace equation (steady-state),  $u_{xx} + u_{yy} = 0$

Any PDE of the form

$$au_{xx} + bu_{xy} + cu_{yy} + du_x + eu_y + fu + g = 0$$

behaves (has the character) of one of the above equations.

- ▶ The *discriminant* determines the canonical form of second-order PDEs:

The discriminant is  $r = b^2 - 4ac$  and linear PDEs are classified as

$$\begin{cases} r > 0 & : \text{hyperbolic, wave-equation-like} \\ r = 0 & : \text{parabolic, heat-equation-like} \\ r < 0 & : \text{elliptic, Laplace-equation-like} \end{cases}$$

When coefficients are varying, a PDE may exhibit different behavior in different parts of the domain.

# Characteristic Curves

- ▶ A *characteristic* of a PDE is a level curve in the solution:
  - ▶ For the Cauchy form of the advection equation,  $u_t = -cu_x$ , a characteristic  $\hat{x}(t)$  satisfies  $u(t, \hat{x}(t)) = \text{const.}$
  - ▶ One of their uses is identifying where boundary conditions must be imposed, e.g. for the Cauchy equation, it tells us whether boundary conditions are needed on the left or the right (in terms of  $x$ ) of the domain.
- ▶ More generally, characteristic curves describe curves in the solution field  $u(t, x)$  that correspond to solutions of ODEs, e.g. for  $u_t = -a(t, x)u_x$  with  $u(0, x) = u_0(x)$ ,

$$\hat{x}'(t) = -a(t, \hat{x}(t)) \text{ with initial condition } \hat{x}(0) = x_0.$$

*These characteristic curves give solutions  $\hat{u}(t, x) = u_0(\hat{x}(t))$ , which satisfy the advection PDE  $u_t = -a(t, x)u_x$ , since*

$$\hat{u}_t = \hat{x}'(t)u'_0(\hat{x}(t)) = -a(t, \hat{x}(t))u'_0(\hat{x}(t)) = -a(t, x)\hat{u}_x.$$

## Method of Lines

- ▶ *Semidiscrete methods* obtain an approximation to the PDE by solving a system of ODEs. Consider the heat equation,

$$u_t = cu_{xx} \text{ on } 0 \leq x \leq 1, \quad u(0, x) = f(x), u(t, 0) = u(t, 1) = 0.$$

- ▶ *We discretize over  $x$  and use finite differences to approximate*

$$u_{xx} \approx \frac{u(t, x_{i+1}) - 2u(t, x_i) + u(t, x_{i-1}))}{(\Delta x)^2}.$$

- ▶ *This approximation yields a system of ODE IVPs, where  $y_i(t) \approx u(t, x_i)$ ,*

$$y'_i(t) = \frac{c}{(\Delta x)^2} (y_{i+1}(t) - 2y_i(t) + y_{i-1}(t)), \quad y_i(0) = f(x_i).$$

- ▶ *In vector form, we obtain a linear constant-coefficient ODE,  $\mathbf{y}'(t) = \mathbf{A}\mathbf{y}(t)$ , where  $\mathbf{A}$  is tridiagonal.*

- ▶ This *method of lines* often yields a stiff ODE:

- ▶ *The eigenvalues of  $\mathbf{A}$  are in  $[-4c/(\Delta x)^2, 0]$ , which gives stiff ODE for  $\Delta x \ll 1$ .*
- ▶ *Heat dissipates very rapidly when concentrated, but slowly overall.*

## Semidiscrete Collocation

- Instead of finite-differences, we can express  $u(t, x)$  in a spatial basis  $\phi_1(x), \dots, \phi_n(x)$  with time-dependent coefficients  $\alpha_1(t), \dots, \alpha_n(t)$ :

$$u(t, x) \approx v(t, x, \boldsymbol{\alpha}(t)) = \sum_{j=1}^n \alpha_j(t) \phi_j(x)$$

*Semidiscrete collocation* methods then ensure the approximation is exact on  $x_1, \dots, x_n$ , yielding system of ODEs.

- For the heat equation  $u_t = cu_{xx}$ , we obtain a linear constant-coefficient vector ODE:

$$\sum_{j=1}^n \frac{\partial \alpha_j}{\partial t}(t) \underbrace{\phi_j(x_i)}_{m_{ij}} = c \sum_{j=1}^n \alpha_j(t) \underbrace{\frac{\partial^2 \phi_j}{\partial x^2}(x_i)}_{n_{ij}}$$

written in matrix form,

$$\boldsymbol{\alpha}'(t) = c\mathbf{M}^{-1}\mathbf{N}\boldsymbol{\alpha}(t)$$

# Fully Discrete Methods

- ▶ Generally, both time and space dimensions are discretized, either by applying an ODE solver to a semidiscrete method or using finite differences.
  - ▶ *Again consider the heat equation  $u_t = cu_{xx}$  and discretize so  $u_i^{(k)} \approx u(t_k, x_i)$ ,*

$$\frac{u_i^{(k+1)} - u_i^{(k)}}{\Delta t} = c \frac{u_{i+1}^{(k)} - 2u_i^{(k)} + u_{i-1}^{(k)}}{(\Delta x)^2}.$$

- ▶ *This iterative scheme corresponds to a 3-point **stencil**,*

$$u_i^{(k+1)} = u_i^{(k)} + c\Delta t \frac{u_{i+1}^{(k)} - 2u_i^{(k)} + u_{i-1}^{(k)}}{(\Delta x)^2}.$$

- ▶ *The same scheme can be derived by applying Euler's method to the ODE given by the method of lines.*

## Implicit Fully Discrete Methods

- ▶ Using Euler's method for the heat equation, stability requirement is

$$\Delta t = O((\Delta x)^2)$$

- ▶ This step-size restriction on stability can be circumvented by use of implicit time-stepper, such as backward Euler,

$$u_i^{(k+1)} = u_i^{(k)} + c\Delta t \frac{u_{i+1}^{(k+1)} - 2u_i^{(k+1)} + u_{i-1}^{(k+1)}}{(\Delta x)^2}.$$

*This scheme requires for a tridiagonal matrix system to be solved at each time-step, but obtains unconditional stability, albeit only first-order accuracy.*

- ▶ Using the trapezoid method to solve the ODE we obtain the second-order *Crank-Nicolson method*,

$$u_i^{(k+1)} = u_i^{(k)} + c\Delta t \frac{u_{i+1}^{(k+1)} - 2u_i^{(k+1)} + u_{i-1}^{(k+1)} + u_{i+1}^{(k)} - 2u_i^{(k)} + u_{i-1}^{(k)}}{2(\Delta x)^2}.$$



# Convergence and Stability

- ▶ *Lax Equivalence Theorem*: consistency + stability = convergence
  - ▶ Consistency means that the local truncation error goes to zero, and is easy to verify by Taylor expansions.
  - ▶ Stability implies that the approximate solution at any time  $t$  must remain bounded.
  - ▶ Together these conditions are necessary and sufficient for convergence.
- ▶ Stability can be ascertained by spectral or Fourier analysis:
  - ▶ In the method of lines, we saw that the eigenvalues of the resulting ODE define the stability region.
  - ▶ Fourier analysis decomposes the solution into a sum of harmonic functions and bounds their amplitudes.

## CFL Condition

- ▶ The domain of dependence of a PDE for a given point  $(t, x)$  is the portion of the problem domain influencing this point through the PDE:
  - ▶ *Generally determined by characteristics of PDE.*
  - ▶ *For a stencil method, the numerical solution depends on the set of mesh-points influencing the mesh point at  $(t, x)$ .*
- ▶ *The Courant, Friedrichs, and Lewy (CFL) condition* states that for an explicit finite-differencing scheme to be stable for a hyperbolic PDE, it is necessary that the domain of the dependence of the PDE must be contained in the domain of dependence of the scheme:

*Intuitively, we can then achieve stability for fixed  $\Delta t$  in two ways,*

  - ▶ *by choosing a sufficiently large grid spacing  $\Delta x$ , or*
  - ▶ *including more mesh points in our stencil.*

## Time-Independent PDEs

- ▶ We now turn our focus to time-independent PDEs as exemplified by the *Helmholtz equation*:

$$u_{xx} + u_{yy} + \lambda u = f(x, y)$$

- ▶  $\lambda = 0$  yields the Poisson equation.
  - ▶  $\lambda = 0$  and  $f = 0$  yields the Laplace equation.
  - ▶ Boundary conditions (e.g. Dirichlet or Neumann or mixed) are imposed on domain surface.
- ▶ We discretize as before, but no longer perform time stepping:  
For example given a domain  $[0, 1]^2$ , we can
  - ▶ tile it using  $n \times n$  mesh points,
  - ▶ setup finite-difference equations on interior,
  - ▶ setup boundary-condition equations on perimeter.

## Finite-Differencing for Poisson

- ▶ Consider the Poisson equation with equispaced mesh-points on  $[0, 1]$ :
  - ▶ If  $\mathbf{u}$  is a vector containing the mesh-points, we have that

$$\mathcal{D}_x \mathbf{u} + \mathcal{D}_y \mathbf{u} = \mathbf{b} \text{ where } b_i = f(u_i)$$

where  $\mathcal{D}_x$  and  $\mathcal{D}_y$  are finite-difference operators along  $x$  and  $y$  dimensions, respectively.

- ▶ Given a differencing matrix  $\mathbf{D}$  (e.g. tridiagonal with  $1, -2, 1$ ), we obtain the matrix equation,

$$(\mathbf{I} \otimes \mathbf{D} + \mathbf{D} \otimes \mathbf{I})\mathbf{u} = \mathbf{b}$$

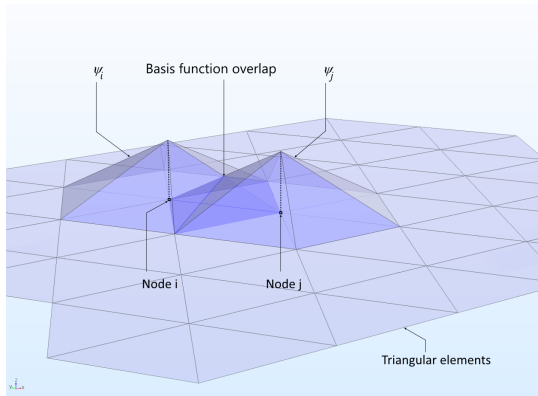
where the *Kronecker product* is defined as

$$\mathbf{C} = \mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots \\ a_{21}\mathbf{B} & \ddots & \\ \vdots & & \end{bmatrix},$$

and the elements of  $\mathbf{b}$  contain the mesh elements in column-major (or row-major in this example) order.

# Multidimensional Finite Elements

- There are many ways to define localized basis functions, for example in the 2D FEM method<sup>2</sup>:



*We partition the domain into triangles (elements) and define linear basis functions that are 1 at the intersection of three or more elements (nodes).*

<sup>2</sup>Source: Comsol Multiphysics Cyclopedia <https://www.comsol.com/multiphysics/finite-element-method>

## 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,*

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

- ▶ *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  $\mathbf{r} = \mathbf{Ax} - \mathbf{b}$ , 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 to permutation of vertices by permutation matrix  $P$ .*
  - ▶ *Reordering  $V$  accordingly transforms the adjacency matrix into  $P^T A P$ .*
  - ▶ *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) introduces 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 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

- ▶ Sparse iterative methods avoid overhead of fill in sparse direct factorization.  
*Matrix splitting* methods provide the most basic iterative methods:

- ▶ *These are linear fixed point iterations that solve the linear system:*

$$Mx_{k+1} = Nx_k + b$$

- ▶ *The fixed point function is*

$$g(x) = M^{-1}Nx + M^{-1}b.$$

- ▶ *We desire to have a fixed point whenever  $Ax = b$ , which implies*

$$Mx = Nx + Ax.$$

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

## 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,

$$Dx^{(k+1)} = -(L + U)x^{(k)} + b.$$

- ▶ 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  $B = D^{-1}(L + U)$ ,

$$\left| \sum_j b_{ij} \right| = \left| \sum_{j \neq i} a_{ij}/a_{ii} \right| < 1 \Rightarrow \rho(M^{-1}N) = \rho(B) < 1.$$

- ▶ For the 2D Laplace problem, this condition holds.
  - ▶ However, the coefficient in linear convergence  $\cos(\pi h) \rightarrow 1$  as  $h \rightarrow 0$ .

## 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*

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

- ▶ *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*

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

- ▶ *If  $x_{GS}^{(k+1)}$  is the iterate produced by Gauss-Seidel, SOR instead produces*

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

- ▶ 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  $\mathbf{Ax} = \mathbf{b}$  when  $\mathbf{A}$  is symmetric positive definite is the minima of the quadratic optimization problem,

$$\min_{\mathbf{x}} \mathbf{x}^T \mathbf{Ax} - \mathbf{x}^T \mathbf{b}$$

*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  $\frac{\sqrt{\kappa(\mathbf{A})}-1}{\sqrt{\kappa(\mathbf{A})}+1}$ :

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

## Preconditioning

- ▶ Preconditioning techniques choose matrix  $M \approx A$  that is easy to invert and solve a modified linear system with an equivalent solution to  $Ax = b$ ,

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

*We can then use Krylov subspace methods that build a Krylov subspace of  $M^{-1}A$  rather than  $A$  with actually forming the matrix  $M^{-1}A$ :*

- ▶ *cost of iteration depends on difficulty of applying  $M^{-1}$ ,*
  - ▶ *convergence rate depends on how close  $M$  is to  $A$ .*
- ▶  $M$  is 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  $M = LU$ .*