

Chapter 7

Finite Difference (FD) Methods for Boundary Value Problems (BVPs)

Definition 7.1. A *partial differential equation* (PDE) is an equation involving an unknown function of two or more variables and some of its partial derivatives.

Definition 7.2. *Laplace equation* is a second-order PDE of the form

$$\Delta u(\mathbf{x}) = 0, \quad (7.1)$$

where the unknown is a function $u : \bar{\Omega} \rightarrow \mathbb{R}$, $\bar{\Omega}$ is the closure of an open set $\Omega \subset \mathbb{R}^n$, and the *Laplacian operator* $\Delta : \mathcal{C}^2(\Omega) \rightarrow \mathcal{C}(\Omega)$ is

$$\Delta := \sum_{i=1}^n \frac{\partial^2}{\partial x_i^2}. \quad (7.2)$$

Example 7.3. *Potential flow* is a special type of flow where the velocity field \mathbf{u} can be expressed as the gradient of a scalar function:

$$\mathbf{u} = \nabla \varphi,$$

where φ is called the *velocity potential*. For incompressible fluids with $\nabla \cdot \mathbf{u} = 0$, the velocity potential satisfies a Laplace equation $\Delta \varphi = 0$.

Definition 7.4. *Poisson's equation* is a second-order PDE of the form

$$-\Delta u(\mathbf{x}) = f(\mathbf{x}), \quad (7.3)$$

where the unknown is a function $u : \bar{\Omega} \rightarrow \mathbb{R}$, $\bar{\Omega}$ is the closure of an open set $\Omega \subset \mathbb{R}^n$, and the RHS function $f : \Omega \rightarrow \mathbb{R}$ is given a priori.

Definition 7.5. A *boundary value problem* (BVP) is a differential equation together with a set of additional constraints, called the *boundary conditions*, that hold only on the domain boundary.

Definition 7.6. Common types of boundary conditions for a one-dimensional interval $\Omega = (a, b)$ are

- *Dirichlet conditions:* $u(a) = \alpha$ and $u(b) = \beta$;
- *Mixed conditions:* $u(a) = \alpha$ and $\left. \frac{\partial u}{\partial x} \right|_b = \beta$;
- *Neumann conditions:* $\left. \frac{\partial u}{\partial x} \right|_a = \alpha$ and $\left. \frac{\partial u}{\partial x} \right|_b = \beta$.

Theorem 7.7. Suppose f and g are two sufficiently smooth functions. Then there exists a unique solution (up to an additive constant) for the Neumann BVP

$$\Delta \phi = f \quad \text{in } \Omega; \quad (7.4a)$$

$$\mathbf{n} \cdot \nabla \phi = g \quad \text{on } \partial\Omega \quad (7.4b)$$

if and only if

$$\int_{\Omega} f \, dV = \int_{\partial\Omega} g \, dA. \quad (7.5)$$

Proof. See [Taylor, 2011, page 409]. \square

Example 7.8. The fundamental theorem of vector calculus (Theorem 13.31) states that a continuously differentiable vector field \mathbf{v}^* can be uniquely decomposed into a divergence-free part and a curl-free part:

$$\begin{cases} \mathbf{v}^* = \mathbf{v} + \nabla \phi, \\ \nabla \cdot \mathbf{v} = 0, \quad \nabla \times \nabla \phi = \mathbf{0}. \end{cases} \quad (7.6)$$

When a given boundary condition of \mathbf{v} satisfies $\oint_{\partial\Omega} \mathbf{v} \cdot \mathbf{n} = 0$, the decomposition is realized by solving the Neumann BVP

$$\Delta \phi = \nabla \cdot \mathbf{v}^* \quad \text{in } \Omega, \quad (7.7a)$$

$$\mathbf{n} \cdot \nabla \phi = \mathbf{n} \cdot (\mathbf{v}^* - \mathbf{v}) \quad \text{on } \partial\Omega, \quad (7.7b)$$

for which the existence of the unique solution is guaranteed by Theorem 7.7 and

$$\int_{\Omega} \nabla \cdot \mathbf{v}^* \, dV = \int_{\Omega} \nabla \cdot (\mathbf{v}^* - \mathbf{v}) \, dV = \int_{\partial\Omega} \mathbf{n} \cdot (\mathbf{v}^* - \mathbf{v}) \, dA.$$

7.1 The FD discretization

Formula 7.9. In solving a linear BVP, the general procedures of an FD method are as follows.

(FD-1) Discretize the problem domain by a grid.

(FD-2) Approximate each spatial derivative in the PDE with some finite difference formula at every grid point to get a system of linear equations $\mathbf{A}\mathbf{U} = \mathbf{F}$

where the vector \mathbf{U} approximates the unknown variable on the grid while the vector \mathbf{F} contains given conditions of the BVP such as boundary conditions and derivatives of the unknown function.

(FD-3) Solve the system of algebraic equations.

Example 7.10 (An FD method for Poisson's equation in a unit interval). Consider the one-dimensional BVP

$$-u''(x) = f(x) \text{ in } \Omega := (0, 1) \quad (7.8)$$

with Dirichlet boundary conditions

$$u(0) = \alpha, \quad u(1) = \beta. \quad (7.9)$$

The general procedures of an FD method based on the central difference are as follows.

(a) Discretize Ω by a Cartesian grid with uniform spacing,

$$x_j = jh, \quad h = \frac{1}{m+1}, \quad j = 0, 1, \dots, m+1.$$

Set $U_0 = \alpha$, $U_{m+1} = \beta$ and we will compute m values U_1, \dots, U_m where each U_j approximates $u(x_j)$.

(b) Approximate the second derivative u'' with a centered difference

$$u''(x_j) = \frac{u(x_{j-1}) - 2u(x_j) + u(x_{j+1}))}{h^2} + O(h^2) \quad (7.10)$$

and we get the following system of linear equations:

$$\begin{aligned} -\frac{\alpha - 2U_1 + U_2}{h^2} &= f(x_1), \\ -\frac{U_{j-1} - 2U_j + U_{j+1}}{h^2} &= f(x_j), \quad j = 2, \dots, m-1, \\ -\frac{U_{m-1} - 2U_m + \beta}{h^2} &= f(x_m). \end{aligned}$$

These equations are written in the form

$$A\mathbf{U} = \mathbf{F}, \quad (7.11)$$

where

$$\mathbf{U} = \begin{bmatrix} U_1 \\ U_2 \\ U_3 \\ \vdots \\ U_{m-1} \\ U_m \end{bmatrix}, \quad \mathbf{F} = \begin{bmatrix} f(x_1) + \frac{\alpha}{h^2} \\ f(x_2) \\ f(x_3) \\ \vdots \\ f(x_{m-1}) \\ f(x_m) + \frac{\beta}{h^2} \end{bmatrix}, \quad (7.12)$$

$$A = \frac{1}{h^2} \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & -1 & 2 & -1 & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 \end{bmatrix}. \quad (7.13)$$

(c) Solve the linear system (7.11).

7.2 Errors and consistency

Definition 7.11. The *global error* or *solution error* of an FD method in Formula 7.9 is

$$\mathbf{E} = \mathbf{U} - \hat{\mathbf{U}}, \quad (7.14)$$

where $\hat{\mathbf{U}} = [u(x_1), u(x_2), \dots, u(x_m)]^T$ is the vector of true values and \mathbf{U} the computed solution.

Definition 7.12. A *grid function* is a function $\mathbf{g} : \mathbf{X} \rightarrow \mathbb{R}$ on a discrete grid \mathbf{X} that contains a finite number of points.

Definition 7.13. The *q-norm* of a grid function \mathbf{g} on a one-dimensional grid $\mathbf{X} := \{x_1, x_2, \dots, x_N\}$ is

$$\|\mathbf{g}\|_q = \left(h \sum_{i=1}^N |g_i|^q \right)^{\frac{1}{q}}, \quad (7.15)$$

where $\mathbf{g} = (g_1, g_2, \dots, g_N)$. In particular, the *1-norm* is

$$\|\mathbf{g}\|_1 = h \sum_{i=1}^N |g_i| \quad (7.16)$$

and the *max-norm* is

$$\|\mathbf{g}\|_\infty = \max_{1 \leq i \leq N} |g_i|. \quad (7.17)$$

Exercise 7.14. Suppose a grid function $\mathbf{g} : \mathbf{X} \rightarrow \mathbb{R}$ has $\mathbf{X} := \{x_1, x_2, \dots, x_N\}$, $g_1 = O(h)$, $g_N = O(h)$, and $g_j = O(h^2)$ for all $j = 2, \dots, N-1$. Show that

$$\|\mathbf{g}\|_\infty = O(h), \quad \|\mathbf{g}\|_1 = O(h^2), \quad \|\mathbf{g}\|_2 = O(h^{\frac{3}{2}}). \quad (7.18)$$

As the main point of this exercise, the differences in the max-norm, 1-norm, and 2-norm of a grid function often reveal the percentage of components with large magnitude.

Definition 7.15. The *local truncation error* (LTE) of an FD method in Formula 7.9 is the error caused by replacing a continuous derivative with an FD formula.

Example 7.16. When we approximate Δu with

$$D^2 u(x_j) := \frac{u(x_{j-1}) - 2u(x_j) + u(x_{j+1}))}{h^2}, \quad (7.19)$$

the LTE of the FD method in Example 7.10 is

$$\tau_j = -D^2 u(x_j) - (-u''(x_j)) = -\frac{h^2}{12} u''''(x_j) + O(h^4).$$

Lemma 7.17. Let $A\mathbf{U} = \mathbf{F}$ be the linear system obtained by applying Formula 7.9 to a linear BVP $\mathcal{L}u = f(x)$ in $(0, 1)$ with Dirichlet conditions. Then the LTE of this FD method is the error of calculating the RHS function \mathbf{F} by replacing U_j with the exact solution $u(x_j)$, i.e.,

$$\boldsymbol{\tau} = A\hat{\mathbf{U}} - \mathbf{F}, \quad (7.20)$$

where $\hat{\mathbf{U}}$ is the vector of true solution values. In particular, this holds for the FD method in Example 7.10.

Proof. We only prove the case in Example 7.10 since other linear BVPs can be proven following similar arguments.

By (7.8) and (7.11), we have, $\forall j = 2, \dots, m-1$,

$$(A\hat{U} - \mathbf{F})_j = -\frac{u(x_{j-1}) - 2u(x_j) + u(x_{j+1}))}{h^2} + \Delta u(x_j),$$

which also holds for $j = 1$ and $j = m$ since the boundary conditions yield $u(x_0) = u(0) = \alpha$ and $u(x_{m+1}) = u(1) = \beta$. Then the proof is completed by Definition 7.15. \square

Lemma 7.18. The LTE and the global error are related as

$$AE = -\tau. \quad (7.21)$$

Proof. $AE = A(\mathbf{U} - \hat{\mathbf{U}}) = \mathbf{F} - (\mathbf{F} + \tau) = -\tau$. \square

Definition 7.19. An FD method in Formula 7.9 is said to be *consistent* with the BVP if

$$\lim_{h \rightarrow 0} \|\tau^h\| = 0, \quad (7.22)$$

where τ^h is the LTE.

7.3 Stability and convergence

Definition 7.20. An FD method is *convergent* if

$$\lim_{h \rightarrow 0} \|\mathbf{E}^h\| = 0, \quad (7.23)$$

where \mathbf{E}^h is the solution error in Definition 7.11 and $\|\cdot\|$ is a q -norm in Definition 7.13.

Definition 7.21. An FD method in Formula 7.9 is *stable* if

- (a) $\exists h_0 \in \mathbb{R}^+$ s.t. $\forall h \in (0, h_0)$, $\det(A) \neq 0$, where A is the matrix of the linear system for the grid size h ;
- (b) $\lim_{h \rightarrow 0} \|A^{-1}\| = O(1)$.

Theorem 7.22. A consistent and stable FD method is convergent.

Proof. Lemma 7.18 and Definitions 7.19 and 7.21 yield

$$\lim \|\mathbf{E}^h\| \leq \lim \|(A^h)^{-1}\| \lim \|\tau^h\| \leq C \lim \|\tau^h\| = 0,$$

where a norm is either a vector norm or an induced matrix norm. Then (7.23) follows from the observation that, by Definition 7.13, the q -norm $\|\cdot\|_q$ and the corresponding vector norm $\|\cdot\|$ are related by $\|\mathbf{E}\|_q = h^{\frac{1}{q}} \|\mathbf{E}\|$. \square

7.3.1 Convergence in the 2-norm

Definition 7.23 (Matrix norms induced by vector norms). The *norm* of a matrix $A \in \mathbb{R}^{n \times n}$ is defined by

$$\begin{aligned} \|A\| &= \sup \left\{ \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} : \mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq \mathbf{0} \right\} \\ &= \sup \left\{ \|A\mathbf{x}\| : \|\mathbf{x}\| = 1 \right\}. \end{aligned}$$

Example 7.24. Commonly used matrix norms include

$$\begin{cases} \|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|, \\ \|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|, \\ \|A\|_2 = \sqrt{\rho(A^T A)}, \end{cases} \quad (7.24)$$

where $\rho(B) := \max_i |\lambda_i(B)|$ is the spectral radius of the matrix B , i.e. the maximum modulus of eigenvalues of B .

Lemma 7.25. The eigenvalues λ_k and eigenvectors \mathbf{w}_k of the matrix A in (7.13) are

$$\lambda_k(A) = \frac{4}{h^2} \sin^2 \frac{k\pi}{2(m+1)}, \quad (7.25)$$

$$w_{k,j} = \sin \frac{jk\pi}{m+1}, \quad (7.26)$$

where $j, k = 1, 2, \dots, m$.

Proof. It is straightforward to verify the conclusions using the trigonometric identity

$$\sin \alpha + \sin \beta = 2 \sin \frac{\alpha + \beta}{2} \cos \frac{\alpha - \beta}{2}. \quad \square$$

Theorem 7.26. The FD method in Example 7.10 is second-order convergent in the 2-norm.

Proof. The symmetry of A gives $\|A\|_2 = \rho(A)$. Then Lemma 7.25 implies the stability, i.e., condition (b) in Definition 7.21:

$$\lim_{h \rightarrow 0} \|A^{-1}\|_2 = \lim_{h \rightarrow 0} \frac{1}{\min |\lambda_k(A)|} = \lim_{h \rightarrow 0} \frac{h^2}{4 \sin^2 \frac{\pi h}{2}} = \frac{1}{\pi^2}.$$

The rest of the proof follows from Example 7.16, Definition 7.19, and Theorem 7.22. \square

7.3.2 Green's function

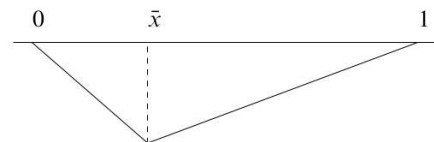
Definition 7.27. For any fixed $\bar{x} \in [0, 1]$, the *Green's function* $G(x; \bar{x})$ is the function of x that solves the BVP

$$\begin{cases} u''(x) = \delta(x - \bar{x}); \\ u(0) = u(1) = 0, \end{cases} \quad (7.27)$$

where $\delta(x - \bar{x})$ is the Dirac delta function in Definition 5.41.

Lemma 7.28. The Green's function $G(x; \bar{x})$ that solves (7.27) is

$$G(x; \bar{x}) = \begin{cases} (\bar{x} - 1)x, & x \in [0, \bar{x}], \\ \bar{x}(x - 1), & x \in [\bar{x}, 1]. \end{cases} \quad (7.28)$$



Proof. For any fixed ϵ , we have from (7.27) and (5.41b) that

$$\begin{aligned} \int_{x_0-\epsilon}^{x_0+\epsilon} G'''(x)dx &= \int_{x_0-\epsilon}^{x_0+\epsilon} \delta(x-\bar{x})dx \\ &= \begin{cases} 0, & \bar{x} \notin (x_0-\epsilon, x_0+\epsilon), \\ 1, & \bar{x} \in (x_0-\epsilon, x_0+\epsilon). \end{cases} \end{aligned}$$

Take limit $\epsilon \rightarrow 0$ of the above equation and we deduce from the second fundamental theorem of calculus (Theorem C.74) that

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} G'(x_0+\epsilon) - \lim_{\epsilon \rightarrow 0} G'(x_0-\epsilon) \\ = \begin{cases} 0 & \text{if } x_0 \in (0, \bar{x}) \cup (\bar{x}, 1), \\ 1 & \text{if } x_0 = \bar{x}. \end{cases} \end{aligned} \quad (7.29)$$

Substitute

$$G(x; \bar{x}) = \begin{cases} ax + b, & x \in [0, \bar{x}], \\ cx + d, & x \in [\bar{x}, 1] \end{cases}$$

into (7.29) and (7.27) and the continuity of $G(x; \bar{x})$ yields

$$\begin{cases} c = a + 1 \\ b = 0 \\ c + d = 0 \\ a\bar{x} + b = c\bar{x} + d \end{cases} \Rightarrow \begin{cases} a = \bar{x} - 1 \\ b = 0 \\ c = \bar{x} \\ d = -\bar{x} \end{cases},$$

which completes the proof. \square

Corollary 7.29. The solution to the linear BVP

$$\begin{cases} u''(x) = c\delta(x-\bar{x}), \\ u(0) = u(1) = 0 \end{cases}$$

is

$$u(x) = cG(x; \bar{x}).$$

Proof. This follows directly from Lemma 7.28. \square

7.3.3 Convergence in the max-norm

Lemma 7.30. For the matrix A in (7.13), any element of its inverse $B = A^{-1}$ is

$$b_{ij} = -hG(x_i; x_j) = \begin{cases} -h(x_j - 1)x_i, & i \leq j, \\ -hx_j(x_i - 1), & i \geq j. \end{cases} \quad (7.30)$$

More explicitly, the matrix B is

$$B = -h \begin{bmatrix} x_1(x_1 - 1) & x_1(x_2 - 1) & \cdots & x_1(x_m - 1) \\ x_1(x_2 - 1) & x_2(x_2 - 1) & \cdots & x_2(x_m - 1) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(x_m - 1) & x_2(x_m - 1) & \cdots & x_m(x_m - 1) \end{bmatrix}.$$

Proof. To verify that B is indeed the inverse of A , it suffices to multiply the i th row of h^2A and the j th column of $-\frac{1}{h}B$,

$$[0, \dots, 0, -1, 2, -1, 0, \dots, 0] \begin{bmatrix} x_1(x_j - 1) \\ \vdots \\ x_{j-1}(x_j - 1) \\ x_j(x_j - 1) \\ x_j(x_{j+1} - 1) \\ \vdots \\ x_j(x_m - 1) \end{bmatrix},$$

the only nonzero case is when $i = j$:

$$2x_j(x_j - 1) - x_{j-1}(x_j - 1) - x_j(x_{j+1} - 1) = -h. \quad \square$$

Theorem 7.31. The max-norm of $B = A^{-1}$ satisfies

$$\|B\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^m |b_{ij}| \leq 1. \quad (7.31)$$

Proof. Lemma 7.30 yields

$$\begin{aligned} \sum_{j=1}^m |b_{ij}| &= \sum_{j=1}^i hx_j|x_i - 1| + \sum_{j=i+1}^m hx_j|x_j - 1| \\ &\leq \sum_{j=1}^i h \left(\frac{m}{m+1} \right)^2 + \sum_{j=i+1}^m h \left(\frac{m}{m+1} \right)^2 \\ &= mh \left(\frac{m}{m+1} \right)^2 = \left(\frac{m}{m+1} \right)^3 \leq 1. \quad \square \end{aligned}$$

7.4 A solution via Green's function

Lemma 7.32. Suppose \mathcal{L} is an invertible linear differential operator that satisfies

$$\mathcal{L}u(x) = f(x). \quad (7.32)$$

Then we have

$$u(x) = \int G(x; \bar{x})f(\bar{x})d\bar{x}, \quad (7.33)$$

where G is the Green's function satisfying

$$\mathcal{L}G(x; \bar{x}) = \delta(x - \bar{x}). \quad (7.34)$$

Proof. Multiply (7.34) by $f(\bar{x})$, integrate w.r.t. \bar{x} , and we have

$$\int \mathcal{L}G(x; \bar{x})f(\bar{x})d\bar{x} = \int \delta(x - \bar{x})f(\bar{x})d\bar{x} = f(x),$$

where the second equality follows from the sifting property of the Dirac delta function (Lemma 5.43). Therefore

$$\mathcal{L} \int G(x; \bar{x})f(\bar{x})d\bar{x} = f(x),$$

which further implies (7.33) since \mathcal{L} is invertible. \square

Theorem 7.33. The Dirichlet BVP

$$\begin{cases} u''(x) = f(x), \\ u(0) = \alpha, u(1) = \beta \end{cases} \quad (7.35)$$

is solved by

$$u(x) = \alpha G_0(x) + \beta G_1(x) + \hat{U}(x), \quad (7.36)$$

where $G_0(x)$, $G_1(x)$ and $\hat{U}(x)$ are defined by BVPs as follows

$$\begin{cases} G_0''(x) = 0, \\ G_0(0) = 1, G_0(1) = 0 \end{cases} \Rightarrow G_0(x) = 1 - x, \quad (7.37a)$$

$$\begin{cases} G_1''(x) = 0, \\ G_1(0) = 0, G_1(1) = 1 \end{cases} \Rightarrow G_1(x) = x, \quad (7.37b)$$

$$\begin{cases} \hat{U}''(x) = f(x), \\ \hat{U}(0) = 0, \hat{U}(1) = 0 \end{cases} \Rightarrow \hat{U}(x) = \int_0^1 f(\bar{x})G(x; \bar{x})d\bar{x}. \quad (7.37c)$$

Proof. This follows from the linearity of the BVP (7.35). \square

7.5 Other boundary conditions

Example 7.34. Consider the second-order BVP

$$u''(x) = f(x) \quad \text{in } (0, 1) \quad (7.38)$$

with mixed boundary conditions

$$u'(0) = \sigma, \quad u(1) = \beta. \quad (7.39)$$

As the crucial difference between this BVP and the BVP with pure Dirichlet conditions in Example 7.10, the value of $u(x)$ at $x = 0$ becomes an unknown to be solved for.

The first approach is to use a one-sided expression

$$\frac{U_1 - U_0}{h} = \sigma \quad (7.40)$$

to arrive at

$$A_E \mathbf{U}_E = \mathbf{F}_E \quad (7.41)$$

where

$$A_E = \frac{1}{h^2} \begin{bmatrix} -h & h & & & & & \\ 1 & -2 & 1 & & & & \\ & 1 & -2 & 1 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & 1 & -2 & 1 & \\ & & & & 1 & -2 & 1 \\ & & & & & 0 & h^2 \end{bmatrix},$$

$$\mathbf{U}_E = \begin{bmatrix} U_0 \\ U_1 \\ U_2 \\ \vdots \\ U_{m-1} \\ U_m \\ U_{m+1} \end{bmatrix}, \quad \mathbf{F}_E = \begin{bmatrix} \sigma \\ f(x_1) \\ f(x_2) \\ \vdots \\ f(x_{m-1}) \\ f(x_m) \\ \beta \end{bmatrix}.$$

The LTE at $x_0 = 0$ is

$$\begin{aligned} \tau_0 &= \frac{1}{h^2}(hu(x_1) - hu(x_0)) - \sigma \\ &= u'(x_0) + \frac{1}{2}hu''(x_0) + O(h^2) - \sigma \\ &= \frac{1}{2}hu''(x_0) + O(h^2), \end{aligned}$$

which is only first order accurate.

The second approach is to extend the domain with a *ghost cell* $x_{-1} = -h$ and use a central difference to obtain

$$\frac{U_1 - U_{-1}}{2h} = \sigma \quad (7.42)$$

that is second-order accurate for the LTE. We do not have any information for U_{-1} , so we want to eliminate it by

$$\frac{1}{h^2}(U_{-1} - 2U_0 + U_1) = f(x_0). \quad (7.43)$$

(7.42) and (7.43) yield

$$\frac{1}{h}(-U_0 + U_1) = \sigma + \frac{h}{2}f(x_0). \quad (7.44)$$

The resulting matrix is the same as that of the first approach in (7.41) except that the first component of \mathbf{F}_E has an additional term $\frac{h}{2}f(x_0)$.

The third approach is to use U_0 , U_1 , and U_2 to approximate $u'(0)$ and we can get a second-order FD formula, c.f. Example 6.38,

$$-\frac{1}{h} \left(\frac{3}{2}U_0 - 2U_1 + \frac{1}{2}U_2 \right) = \sigma + O(h^2). \quad (7.45)$$

This results in the linear system

$$A_F \mathbf{U}_E = \mathbf{F}_E \quad (7.46)$$

where

$$A_F = \frac{1}{h^2} \begin{bmatrix} -\frac{3}{2}h & 2h & -\frac{1}{2}h & & & & \\ 1 & -2 & 1 & & & & \\ & 1 & -2 & 1 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & 1 & -2 & 1 & \\ & & & & 1 & -2 & 1 \\ & & & & & 0 & h^2 \end{bmatrix}.$$

Exercise 7.35. Show that all elements of the first column of $B_E = A_E^{-1}$ are $O(1)$.

Example 7.36. Consider the second-order BVP

$$u''(x) = f(x) \quad \text{in } (0, 1), \quad (7.47)$$

with pure Neumann conditions

$$u'(0) = \sigma_0, \quad u'(1) = \sigma_1. \quad (7.48)$$

To ensure the existence of a solution, the following compatibility condition on $f(x)$, σ_0 , and σ_1 must be satisfied:

$$\int_0^1 f(x)dx = \int_0^1 u''(x)dx = u'(1) - u'(0) = \sigma_1 - \sigma_0. \quad (7.49)$$

In fact, if (7.49) holds, there are an infinite number of solutions: if v is a solution of (7.47), $v + \mathbb{R}$ are also solutions.

Using procedures similar to those in Example 7.34, we can discretize (7.47) and (7.48) as

$$A_F \mathbf{U}_E = \mathbf{F}_F, \quad (7.50)$$

where

$$A_F = \frac{1}{h^2} \begin{bmatrix} -h & h & & & & & \\ 1 & -2 & 1 & & & & \\ & 1 & -2 & 1 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & 1 & -2 & 1 & \\ & & & & 1 & -2 & 1 \\ & & & & & h & -h \end{bmatrix}, \quad (7.51)$$

$$\mathbf{F}_F = \begin{bmatrix} \sigma_0 + \frac{h}{2}f(x_0) \\ f(x_1) \\ f(x_2) \\ \vdots \\ f(x_{m-1}) \\ f(x_m) \\ -\sigma_1 + \frac{h}{2}f(x_{m+1}) \end{bmatrix}.$$

Lemma 7.37. The matrix A_F in (7.51) satisfies

$$\dim \mathcal{N}(A_F) = 1. \quad (7.52)$$

Proof. Clearly, $\mathbf{e} = (1, 1, \dots, 1)^T$ is in the null space of A_F . The rest follows from the well-posedness of the BVP with mixed conditions. \square

Theorem 7.38 (Solvability condition). The linear system (7.50) has a solution if and only if

$$\frac{h}{2}f(x_0) + h \sum_{i=1}^m f(x_i) + \frac{h}{2}f(x_{m+1}) = \sigma_1 - \sigma_0. \quad (7.53)$$

Proof. The fundamental theorem of linear algebra (Theorem B.89) implies

$$\mathbb{R}^{m+2} = \mathcal{R}(A_F) \oplus \mathcal{N}(A_F^T) \quad (7.54)$$

and $\dim \mathcal{N}(A_F^T) = \dim \mathcal{N}(A_F)$. Lemma 7.37 further yields $\dim \mathcal{N}(A_F^T) = 1$. Then it is readily verified that

$$\mathcal{N}(A_F^T) = \text{span} \{ (1, h, h, \dots, h, 1)^T \}.$$

For sufficiency, the above equation and (7.53) imply that \mathbf{F}_F is orthogonal to $\mathcal{N}(A_F^T)$ and thus (7.54) yields $\mathbf{F}_F \in \mathcal{R}(A_F)$. Hence (7.50) must have a solution. As for necessity, the existence of a solution of (7.50) implies $\mathbf{F}_F \in \mathcal{R}(A_F)$, which, together with the above two equations, implies (7.53). \square

7.6 BVPs in two dimensions

Example 7.39 (An FD method for Poisson's equation in a unit square). Consider the two-dimensional BVP

$$-\frac{\partial^2}{\partial x^2}u(x, y) - \frac{\partial^2}{\partial y^2}u(x, y) = f(x, y) \quad (7.55)$$

in $\Omega := (0, 1)^2$ with homogeneous Dirichlet conditions

$$u(x, y)|_{\partial\Omega} = 0. \quad (7.56)$$

A uniform Cartesian grid can be generated with

$$x_i = ih, \quad y_j = jh, \quad i, j = 0, 1, \dots, m, m+1, \quad (7.57)$$

where $h = \Delta x = \Delta y = \frac{1}{m+1}$ is the uniform grid size.

Approximate $\frac{\partial^2 u}{\partial x^2}$ and $\frac{\partial^2 u}{\partial y^2}$ separately and we have, $\forall i, j = 1, 2, \dots, m$,

$$-\frac{U_{i-1,j} - 2U_{ij} + U_{i+1,j}}{h^2} - \frac{U_{i,j-1} - 2U_{ij} + U_{i,j+1}}{h^2} = f_{ij}. \quad (7.58)$$

These $m \times m$ equations organize into a single system

$$A_{2D} \mathbf{U} = \mathbf{F}. \quad (7.59)$$

Exercise 7.40. Show that the LTE τ of the FD method in Example 7.39 is

$$\tau_{i,j} = -\frac{1}{12}h^2 \left(\frac{\partial^4 u}{\partial x^4} + \frac{\partial^4 u}{\partial y^4} \right) \Big|_{(x_i, y_j)} + O(h^4). \quad (7.60)$$

		7	8	9
		4	5	6
		1	2	3

Example 7.41. For $m = 3$ with ordering as shown above, we have

$$A_{2D} = \frac{1}{h^2} \begin{bmatrix} +4 & -1 & & & & & & \\ -1 & +4 & -1 & & & & & \\ & -1 & +4 & & & & & \\ -1 & & & +4 & -1 & & & \\ & -1 & & -1 & +4 & -1 & & \\ & & -1 & & -1 & +4 & & \\ & & & -1 & & -1 & +4 & \\ & & & & & -1 & & -1 & +4 \end{bmatrix}$$

$$= \frac{1}{h^2} \begin{bmatrix} T & -I & \\ -I & T & -I \\ & -I & T \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} U_{11} \\ U_{21} \\ U_{31} \\ U_{12} \\ U_{22} \\ U_{32} \\ U_{13} \\ U_{23} \\ U_{33} \end{bmatrix},$$

where

$$T = \begin{bmatrix} +4 & -1 & 0 \\ -1 & +4 & -1 \\ 0 & -1 & +4 \end{bmatrix}. \quad (7.61)$$

where \mathbf{U} is obtained by stacking the columns on top of each other.

Lemma 7.42. Let $\mathbf{E} = \mathbf{U} - \hat{\mathbf{U}}$ denote the global error of the linear system (7.59). Then the LTE (7.60) satisfies

$$A_{2D} \mathbf{E} = -\boldsymbol{\tau}. \quad (7.62)$$

Proof. The proof is the same as that of Lemma 7.18. \square

7.6.1 Kronecker product

Definition 7.43. The *Kronecker product* of two matrices $A \in \mathbb{C}^{m \times n}$ and $B \in \mathbb{C}^{p \times q}$ is another matrix $A \otimes B \in \mathbb{C}^{mp \times nq}$ given by

$$A \otimes B := \begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{bmatrix}, \quad (7.63)$$

where a_{ij} is the (i, j) th element of A .

Example 7.44.

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \otimes \begin{bmatrix} 0 & 5 \\ 6 & 7 \end{bmatrix} = \begin{bmatrix} 0 & 5 & 0 & 10 \\ 6 & 7 & 12 & 14 \\ 0 & 15 & 0 & 20 \\ 18 & 21 & 24 & 28 \end{bmatrix}.$$

Definition 7.45. For $X \in \mathbb{C}^{m \times n}$, $\text{vec}(X)$ is defined to be a column vector of size mn made of the columns of X stacked on top of one another from left to right.

Lemma 7.46. Any $A \in \mathbb{C}^{m \times m}$, $B \in \mathbb{C}^{n \times n}$, and $X \in \mathbb{C}^{m \times n}$ satisfy

$$\text{vec}(AX) = (I_n \otimes A)\text{vec}(X), \quad (7.64)$$

$$\text{vec}(XB) = (B^T \otimes I_m)\text{vec}(X). \quad (7.65)$$

Proof. We have

$$\begin{aligned} \text{vec}(AX) &= \text{vec}([AX_1, AX_2, \dots, AX_n]) \\ &= \begin{bmatrix} AX_1 \\ AX_2 \\ \vdots \\ AX_n \end{bmatrix} = \begin{bmatrix} A & & & \\ & A & & \\ & & \ddots & \\ & & & A \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \\ &= (I_n \otimes A) \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}. \end{aligned}$$

Let $Y = XB$, then

$$Y_j = Xb_j \Rightarrow y_{kj} = \sum_{i=1}^n x_{ki}b_{ij}. \quad (7.66)$$

Let $C = B^T \otimes I_m$, then the (i, j) -th sub-block of C is

$$C_{ij} = b_{ji}I_m. \quad (7.67)$$

Let $D = C\text{vec}(X)$, then the j -th block of D is

$$D_j = \sum_{i=1}^n C_{ji}X_i = \sum_{i=1}^n b_{ij}I_mX_i = \sum_{i=1}^n b_{ij}X_i, \quad (7.68)$$

and the (k, j) -th entry of D is (Here we also use (k, j) to denote the scalar index corresponding to the multi-index (k, j) .)

$$d_{(k,j)} = \sum_{i=1}^n b_{ij}x_{ki} = \sum_{i=1}^n x_{ki}b_{ij}. \quad (7.69)$$

Combining (7.66) and (7.69) yields (7.65). \square

7.6.2 Convergence in the 2-norm

Lemma 7.47. The linear system (7.59) is equivalent to

$$AU_{m \times m} + U_{m \times m}A = F_{m \times m}, \quad (7.70)$$

where the (i, j) th element of $U_{m \times m}$ is the computed solution at the (i, j) th grid point, the (i, j) th element of $F_{m \times m}$ is

$$(F_{m \times m})_{ij} = f(ih, jh),$$

and A is the 1D discrete Laplacian in (7.13).

Proof. A direct computation gives

$$\begin{cases} (AU_{m \times m})_{ij} = \frac{1}{h^2}(-U_{i-1,j} + 2U_{ij} - U_{i+1,j}), \\ (U_{m \times m}A)_{ij} = \frac{1}{h^2}(-U_{i,j-1} + 2U_{ij} - U_{i,j+1}), \end{cases} \quad (7.71)$$

and the homogeneous Dirichlet condition yields

$$AU_{m \times m} + U_{m \times m}A = F_{m \times m}. \quad \square$$

Lemma 7.48. The 1D discrete Laplacian A in (7.13) satisfies

$$\text{vec}(AU_{m \times m} + U_{m \times m}A) = (I_m \otimes A + A \otimes I_m)\text{vec}(U_{m \times m}).$$

Proof. By Lemma 7.46, we have

$$\text{vec}(AU_{m \times m}) = (I_m \otimes A)\text{vec}(U_{m \times m}),$$

and

$$\text{vec}(U_{m \times m}A) = (A^T \otimes I_m)\text{vec}(U_{m \times m}) = (A \otimes I_m)\text{vec}(U_{m \times m}),$$

where the second equality follows from the symmetry of A . Adding these two equations gives the desired result. \square

Theorem 7.49. With matrix ordering, the linear system (7.59) can be written as

$$A_{2D} = I_m \otimes A + A \otimes I_m, \quad \mathbf{U} = \text{vec}(U_{m \times m}), \quad \mathbf{F} = \text{vec}(F_{m \times m}).$$

Proof. This follows from Lemma 7.47 and Lemma 7.48. \square

Definition 7.50. The *discrete Laplacian* in n -dimensional space analogous to the 1D discrete Laplacian (7.13) is

$$A_{nD} = \sum_{j=0}^{n-1} \underbrace{I_m \otimes \dots \otimes I_m}_{\#I_m=j} \otimes A \otimes \underbrace{I_m \otimes \dots \otimes I_m}_{\#I_m=n-j-1}. \quad (7.72)$$

Example 7.51. For $n = 3$, we have

$$A_{3D} = A \otimes I_m \otimes I_m + I_m \otimes A \otimes I_m + I_m \otimes I_m \otimes A.$$

Theorem 7.52. The eigenpairs of A_{2D} are

$$\lambda_{ij} = \lambda_i + \lambda_j, \quad \mathbf{W}_{ij} = \text{vec}(\mathbf{w}_i \mathbf{w}_j^T), \quad (7.73)$$

where $i, j = 1, 2, \dots, m$ and $(\lambda_i, \mathbf{w}_i)$ is an eigenpair of A in Lemma 7.25.

Proof. By Lemma 7.25, we have

$$\begin{aligned} A\mathbf{w}_i \mathbf{w}_j^T + \mathbf{w}_i \mathbf{w}_j^T A &= \lambda_i \mathbf{w}_i \mathbf{w}_j^T + \mathbf{w}_i \mathbf{w}_j^T A^T \\ &= \lambda_i \mathbf{w}_i \mathbf{w}_j^T + \mathbf{w}_i (A\mathbf{w}_j)^T \\ &= \lambda_i \mathbf{w}_i \mathbf{w}_j^T + \lambda_j \mathbf{w}_i \mathbf{w}_j^T \\ &= (\lambda_i + \lambda_j) \mathbf{w}_i \mathbf{w}_j^T. \end{aligned}$$

Then Theorem 7.49 and Lemma 7.48 yield

$$\begin{aligned} A_{2D} \text{vec}(\mathbf{w}_i \mathbf{w}_j^T) &= (I_m \otimes A + A \otimes I_m) \text{vec}(\mathbf{w}_i \mathbf{w}_j^T) \\ &= \text{vec}(A(\mathbf{w}_i \mathbf{w}_j^T) + (\mathbf{w}_i \mathbf{w}_j^T)A) \\ &= (\lambda_i + \lambda_j) \text{vec}(\mathbf{w}_i \mathbf{w}_j^T), \end{aligned}$$

and hence $\lambda_i + \lambda_j$ is an eigenvalue of A_{2D} with corresponding eigenvector \mathbf{W}_{ij} . \square

Theorem 7.53. The FD method in Example 7.39 is second-order convergent in the 2-norm.

Proof. We have $\|A_{2D}\|_2 = \rho(A_{2D})$ since A_{2D} is symmetric. Then Theorem 7.52 yields

$$\lim_{h \rightarrow 0} \|A_{2D}^{-1}\|_2 = \lim_{h \rightarrow 0} \frac{1}{\min |\lambda_{ij}|} = \lim_{h \rightarrow 0} \frac{h^2}{8 \sin^2 \frac{\pi h}{2}} = \frac{1}{2\pi^2} = O(1).$$

By Definition 7.21, the method is stable. The proof is completed by (7.60), Definition 7.19, Theorem 7.22, and Lemma 7.42. \square

7.6.3 Convergence in the max-norm via a discrete maximum principle

Theorem 7.54. The FD method in Example 7.39 is second-order convergent in the max-norm.

Proof. Let \mathbf{X}_I be the grid obtained by removing from \mathbf{X} in (7.57) those grid points satisfying $i = 0, m+1$ or $j = 0, m+1$. Define a linear map $\hat{A}_{2D} : \{\mathbf{X} \rightarrow \mathbb{R}\} \rightarrow \{\mathbf{X}_I \rightarrow \mathbb{R}\}$,

$$\begin{aligned}\hat{A}_{2D}U_{i,j} &:= (\hat{A}_{2D}\mathbf{U})_{i,j} \\ &= \frac{1}{h^2}(4U_{i,j} - U_{i+1,j} - U_{i-1,j} - U_{i,j+1} - U_{i,j-1}).\end{aligned}\quad (7.74)$$

The matrix of \hat{A}_{2D} is different from A_{2D} in Example 7.39. (How?) Define a *comparison function* $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}$,

$$\phi(x, y) := \left(x - \frac{1}{2}\right)^2 + \left(y - \frac{1}{2}\right)^2 \quad (7.75)$$

and write $\phi_{i,j} := \phi(ih, jh)$. (7.74) and (7.75) yield

$$\hat{A}_{2D}\phi_{i,j} = -4. \quad (7.76)$$

Let $E, \tau : \mathbf{X} \rightarrow \mathbb{R}$ be the solution error and the LTE, respectively, and write $\tau_m := \max_{i,j} |\tau_{i,j}|$. Construct a grid function $\psi : \mathbf{X} \rightarrow \mathbb{R}$,

$$\psi_{i,j} := E_{i,j} + \frac{1}{4}\tau_m\phi_{i,j} \quad (7.77)$$

and we have, from $\hat{A}_{2D}E_{i,j} = -\tau_{i,j}$ (why?) and (7.76),

$$\hat{A}_{2D}\psi_{i,j} = -\tau_{i,j} - \tau_m \leq 0.$$

By (7.74), $\hat{A}_{2D}\psi_{i,j} \leq 0$ dictates that $\psi_{i,j}$ be no greater than at least one of its neighbors $\psi_{i+1,j}$, $\psi_{i-1,j}$, $\psi_{i,j+1}$, and $\psi_{i,j-1}$. Therefore, *the maximum value of ψ must occur at a boundary point*, i.e., a point with $i = 0, m+1$ or $j = 0, m+1$. Consequently, there exists some constant $C > 0$ such that

$$E_{i,j} \leq \psi_{i,j} \leq \frac{1}{8}\tau_m < Ch^2,$$

where the first step follows from (7.77) and $\tau_m\phi_{i,j} \geq 0$, the second step from $E_{i,j}$ being zero at all boundary points and the fact that the maximum of ϕ is $\frac{1}{2}$ at the domain corners, and the last step from (7.60).

By similar arguments, the grid function

$$\chi_{i,j} := -E_{i,j} + \frac{1}{4}\tau_m\phi_{i,j} \quad (7.78)$$

satisfies $\hat{A}_{2D}\chi_{i,j} \leq 0$ for all grid points and thus

$$-E_{i,j} \leq \chi_{i,j} \leq \frac{1}{8}\tau_m < Ch^2.$$

To sum up, we have $|E_{i,j}| = O(h^2)$ for all grid points. \square

Notation 4. Consider discretizing a BVP on domain Ω . Denote by \mathbf{X}_Ω the set of *equation-discretization points* where the BVP is discretized and where values of the unknown function u are sought. Let $\mathbf{X}_{\partial\Omega} \subset \partial\Omega$ be the set of *boundary points* so that each point $Q \in \mathbf{X}_{\partial\Omega}$ satisfies

- the BVP is not discretized at Q ,
- $u(Q)$ is prescribed by a Dirichlet condition,
- this Dirichlet condition at Q is involved in discretizing the BVP at some $P \in \mathbf{X}_\Omega$.

WLOG, we also assume $\mathbf{X} = \mathbf{X}_\Omega \cup \mathbf{X}_{\partial\Omega}$.

Example 7.55. A boundary point in Notation 4 must be on the domain boundary $\partial\Omega$, but an equation-discretization point P might also belong to $\partial\Omega$: if we have a Neumann condition for P , then the *value* of the unknown at P might still be needed. In general, a grid \mathbf{X} that corresponds to a straightforward discretization of the problem domain might contain certain points that belong neither to \mathbf{X}_Ω nor to $\mathbf{X}_{\partial\Omega}$, e.g., the four corners of the square domain $(0, 1)^2$. In Notation 4, we have assumed $\mathbf{X} = \mathbf{X}_\Omega \cup \mathbf{X}_{\partial\Omega}$ to exclude these abnormal points so that hereafter the analysis can be relatively simple.

Lemma 7.56 (Discrete maximum principle). Suppose that an FD discretization of a linear BVP yields

$$\forall P \in \mathbf{X}_\Omega, \quad L_h U_P - f_P + g_P = 0, \quad (7.79)$$

where f_P corresponds to the RHS of (7.3), g_P corresponds to all boundary data other than Dirichlet conditions, and L_h and \mathbf{X}_Ω satisfy

(DMP-1) for each equation-discretization point $P \in \mathbf{X}_\Omega$, L_h is of the form

$$L_h U_P = c_P U_P - \sum_{Q \in Q_P} c_Q U_Q, \quad (7.80)$$

where $Q_P \subset \mathbf{X}_\Omega \cup \mathbf{X}_{\partial\Omega}$, $c_P > 0$ and each $c_Q > 0$. The set $\{P\} \cup Q_P$ in (7.80) is called the *P-stencil* or the *stencil of L_h at P* ;

(DMP-2) $\forall P \in \mathbf{X}_\Omega$, $c_P \geq \sum_{Q \in Q_P} c_Q$;

(DMP-3) \mathbf{X}_Ω is *connected*, i.e.,

$$\begin{aligned}\forall P_0, P_m \in \mathbf{X}_\Omega, \exists P_1, P_2, \dots, P_{m-1} \text{ s.t.} \\ \forall r = 1, 2, \dots, m, \quad P_r \text{ is in the } P_{r-1}\text{-stencil};\end{aligned} \quad (7.81)$$

(DMP-4) at least one equation (7.80) involves a boundary value U_Q given by a Dirichlet condition.

Then for any grid function $\psi : \mathbf{X} \rightarrow \mathbb{R}$ satisfying

$$\begin{cases} \max_{P \in \mathbf{X}} \psi_P \geq 0, \\ \forall P \in \mathbf{X}_\Omega, \quad L_h \psi_P \leq 0, \end{cases} \quad (7.82)$$

we have

$$\max_{P \in \mathbf{X}_\Omega} \psi_P \leq \max_{Q \in \mathbf{X}_{\partial\Omega}} \psi_Q. \quad (7.83)$$

Proof. Suppose

$$M_\Omega := \max_{Q \in \mathbf{X}_\Omega} \psi_Q > M_{\partial\Omega} := \max_{Q \in \mathbf{X}_{\partial\Omega}} \psi_Q$$

and let P be the point where ψ attains M_Ω . Then

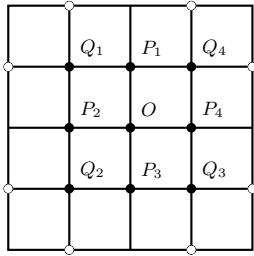
$$(*) : M_\Omega = \psi_P \leq \frac{1}{c_P} \sum_{Q \in Q_P} c_Q \psi_Q \leq \frac{1}{c_P} \sum_{Q \in Q_P} c_Q M_\Omega \leq M_\Omega,$$

where the first inequality follows from (7.82) and (7.80), the second from the definition of M_Ω , and the third from (DMP-2) and $M_\Omega \geq 0$. For the second “ \leq ” in $(*)$ to be “ $=$,” we must have $\psi_Q = \psi_P$ for each Q . By (DMP-3), ψ takes the same value M_Ω on all equation-discretization points. Then (DMP-4) implies $M_\Omega = M_{\partial\Omega}$, which contradicts the starting point $M_\Omega > M_{\partial\Omega}$. \square

Example 7.57. Suppose the concept of connectedness in (DMP-3) were defined as

$$\begin{aligned} \forall P_0, P_m \in \mathbf{X}_\Omega, \exists P_1, P_2, \dots, P_{m-1} \text{ s.t. } \forall r = 1, 2, \dots, m, \\ \text{both } U_{P_r} \text{ and } U_{P_{r-1}} \text{ appear in some equation (7.79)}. \end{aligned} \quad (7.84)$$

Then the conclusion (7.83) of Lemma 7.56 would not hold. The following is a counter-example.



Let L_h be

$$\begin{aligned} L_h \psi_{P_1} &= 3\psi_{P_1} - \psi_{P_2} - \psi_{P_3} - \psi_{P_4}, \\ L_h \psi_{P_2} &= 3\psi_{P_2} - \psi_{P_1} - \psi_{P_3} - \psi_{P_4}, \\ L_h \psi_{P_3} &= 3\psi_{P_3} - \psi_{P_1} - \psi_{P_2} - \psi_{P_4}, \\ L_h \psi_{P_4} &= 3\psi_{P_4} - \psi_{P_1} - \psi_{P_2} - \psi_{P_3}, \\ L_h \psi_O &= 4\psi_O - \psi_{P_1} - \psi_{P_2} - \psi_{P_3} - \psi_{P_4}, \end{aligned}$$

and the expressions of L_h at Q_i 's be similar with that of L_h at O such that (DMP-1,2,4) hold. (7.84) also holds. The following distribution of the grid function ψ ,

$$\psi_{\mathbf{x}} = \begin{cases} 10 & \text{if } \mathbf{x} = P_1, P_2, P_3, P_4; \\ 1 & \text{otherwise,} \end{cases}$$

satisfies (7.82), but the conclusion (7.83) does not hold.

The key point of this example is that, in order for two points P_0 and P_m to be called connected in the context of solving elliptic equations, there must simultaneously exist a path from P_0 to P_m and a path from P_m to P_0 . The definition of connectedness in (7.84) is not satisfactory because it fails to capture the direction of the path. In contrast, (7.81) captures the dependence of path connectedness on the direction of the path.

Theorem 7.58. Suppose an FD discretization of a BVP satisfies the conditions (DMP-1,2,3,4) in Lemma 7.56. Then the solution error $E_P := U_P - u(P)$ of the FD method (7.79) is bounded by

$$\forall P \in \mathbf{X}, \quad |E_P| \leq T_{\max} \left(\max_{Q \in \mathbf{X}_{\partial\Omega}} \phi(Q) \right), \quad (7.85)$$

where $T_{\max} = \max_{P \in \mathbf{X}_\Omega} |T_P|$, T_P is the LTE at P and $\phi : \mathbf{X} \rightarrow \mathbb{R}$ is a nonnegative grid function satisfying

$$\forall P \in \mathbf{X}_\Omega, \quad L_h \phi_P \leq -1. \quad (7.86)$$

Proof. Lemma 7.18 implies $L_h E_P = -T_P$. Define

$$\psi_P := E_P + T_{\max} \phi_P$$

and we know from (7.86) that

$$L_h \psi_P \leq -T_P - T_{\max} \leq 0.$$

Furthermore, we have $\max_{P \in \mathbf{X}} \psi_P \geq 0$ because $\phi_P \geq 0$ and

$$(*) : \quad \forall Q \in \mathbf{X}_{\partial\Omega}, \quad E_Q = 0,$$

c.f. Notation 4. Then we have

$$\begin{aligned} E_P &\leq \max_{P \in \mathbf{X}} (E_P + T_{\max} \phi_P) \\ &\leq \max_{Q \in \mathbf{X}_{\partial\Omega}} (E_Q + T_{\max} \phi_Q) = T_{\max} \max_{Q \in \mathbf{X}_{\partial\Omega}} (\phi_Q), \end{aligned}$$

where the first step follows from $T_{\max} \phi_P \geq 0$, the second from Lemma 7.56, and the third from $(*)$.

Repeat the above arguments with $\psi_P = -E_P + T_{\max} \phi_P$ and we have

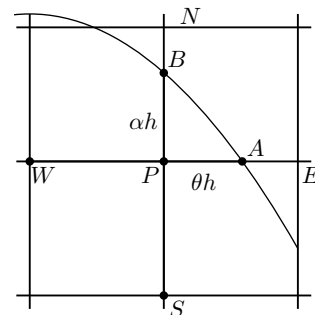
$$-E_P \leq T_{\max} \max_{Q \in \mathbf{X}_{\partial\Omega}} (\phi_Q). \quad \square$$

7.6.4 Convergence on irregular domains

Example 7.59 (An FD method for Poisson's equation in 2D irregular domains). Consider the BVP

$$-\frac{\partial^2}{\partial x^2} u(x, y) - \frac{\partial^2}{\partial y^2} u(x, y) = f(x, y) \quad (7.87)$$

in a 2D irregular domain Ω with Dirichlet conditions.



An equation-discretization point is said to be *regular* if the standard 5-point stencil is applicable; otherwise it is *irregular*. For an irregular point, we modify the FD discretization in Example 7.39 to incorporate the info of local geometry and Dirichlet conditions. For example, in the above plot, the discrete operator becomes

$$L_h U_P := \frac{(1+\theta)U_P - U_A - \theta U_W}{\frac{1}{2}\theta(1+\theta)h^2} + \frac{(1+\alpha)U_P - U_B - \alpha U_S}{\frac{1}{2}\alpha(1+\alpha)h^2}. \quad (7.88)$$

In the resulting linear system $L_h U_P - f_P = 0$, the form of $L_h U_P$ as in (7.88) is different from that in (7.58) at a regular equation-discretization point. Consequently, a global analysis of this linear system is difficult.

Exercise 7.60. Show that, in Example 7.59, the LTE at an irregular equation-discretization point is $O(h)$ while the LTE at a regular equation-discretization point is $O(h^2)$.

Theorem 7.61. Suppose that, in the notation of Theorem 7.58, the set \mathbf{X}_Ω of equation-discretization points can be partitioned as

$$\mathbf{X}_\Omega = \mathbf{X}_1 \cup \mathbf{X}_2, \quad \mathbf{X}_1 \cap \mathbf{X}_2 = \emptyset,$$

the nonnegative function $\phi : \mathbf{X} \rightarrow \mathbb{R}$ satisfies

$$\forall P \in \mathbf{X}_1, \quad L_h \phi_P \leq -C_1 < 0; \quad (7.89a)$$

$$\forall P \in \mathbf{X}_2, \quad L_h \phi_P \leq -C_2 < 0, \quad (7.89b)$$

and the LTE of (7.79) satisfy

$$\forall P \in \mathbf{X}_1, \quad |T_P| < T_1; \quad (7.90a)$$

$$\forall P \in \mathbf{X}_2, \quad |T_P| < T_2. \quad (7.90b)$$

Then the solution error $E_P := U_P - u(P)$ of the FD method (7.79) is bounded by

$$\forall P \in \mathbf{X}, \quad |E_P| \leq \left(\max_{Q \in \mathbf{X}_{\partial\Omega}} \phi(Q) \right) \max \left\{ \frac{T_1}{C_1}, \frac{T_2}{C_2} \right\}.$$

Exercise 7.62. Prove Theorem 7.61 by choosing a function ψ to which Lemma 7.56 applies.

Theorem 7.63. The FD method in Example 7.59 is second-order convergent in the max-norm.

Proof. Define a comparison function ϕ as

$$\phi(x, y) := \begin{cases} F_1 \left[(x-p)^2 + (y-q)^2 \right] & \text{if } (x, y) \in \mathbf{X}_\Omega; \\ F_1 \left[(x-p)^2 + (y-q)^2 \right] + F_2 & \text{if } (x, y) \in \mathbf{X}_{\partial\Omega}, \end{cases}$$

where (p, q) is the geometric center of Ω and $F_1, F_2 > 0$ are constants to be chosen later. Both regular points and irregular points belong to \mathbf{X}_Ω . Their difference is that, for a regular point Q , we have

$$L_h \phi_Q = -4F_1$$

while for an irregular point P shown in Example 7.59 the coefficient of U_A is

$$-\frac{2}{\theta(1+\theta)h^2} < -\frac{1}{h^2}$$

because $\theta \in (0, 1)$. Therefore,

$$L_h \phi_P < -4F_1 - \frac{1}{h^2} F_2 < -\frac{1}{h^2} F_2. \quad (7.91)$$

Note that the last upper bound $-\frac{1}{h^2} F_2$ can not be sharpened to $-\frac{2}{h^2} F_2$ or $-\frac{3}{h^2} F_2$ because the stencil of the irregular point p might as well only contains one point outside the domain.

By Exercise 7.60, we write the maximum LTEs on regular and irregular points as $T_1 = K_1 h^2$ and $T_2 = K_2 h$, respectively. Then Theorem 7.61 implies

$$|E_P| \leq (F_1 R^2 + F_2) \max \left\{ \frac{K_1 h^2}{4F_1}, \frac{K_2 h^3}{F_2} \right\}, \quad (7.92)$$

where R is the maximum distance of a point in Ω to the geometric center of Ω . The RHS of the above equation is minimized when we choose $\frac{F_1}{F_2} = \frac{K_1}{4K_2 h}$ so that the two terms in $\max\{\}$ equal. It follows that

$$\forall P \in \mathbf{X}, \quad |E_P| \leq \frac{1}{4} K_1 R^2 h^2 + K_2 h^3. \quad \square$$

7.7 Programming assignments

Write a C++ package to solve the two-dimensional Poisson equation in Definition 7.4 by the FD methods in Examples 7.39 and 7.59.

I. Your package must give the user the following options:

- (a) the problem domain: either $\Omega = (0, 1)^2$ or $\Omega \setminus \mathbf{D}$ where \mathbf{D} is a closed circular disk, of which the radius and the center, specified by the user, must keep $\Omega \setminus \mathbf{D}$ connected and \mathbf{D} must cover at least four equation-discretization points in Notation 4 (you must check the validity of input parameters);
- (b) boundary conditions: Dirichlet, Neumann, or mixed (partly Dirichlet and partly Neumann).

Use a direct method such as the LU factorization to solve the linear system, with your favorite implementation of a BLAS or LAPACK.

II. For the function

$$u(x, y) = \exp(y + \sin(x)), \quad (7.93)$$

derive the corresponding $f(x, y)$ and the boundary conditions. Test your solver for all combinations of (a,b) in I on grids with $n = 8, 16, 32, 64$ along each dimension, report the 1-, 2-, and ∞ -norms of the errors and the corresponding convergence rates on the four grids. You should also design at least two of your own test functions and carry out the same process.

III. The user-specified parameters must be clearly listed in an input file with a key-value syntax; I recommend json (<https://www.json.org/json-en.html>), but you may use or write your own parser. The main program is supposed to read the input file, create objects from the input, call your BVP solver, generate test results, and report errors and convergence rates.

- | | |
|--|---|
| <p>IV. The generation of pictures on solution and errors can be performed outside your program using another tool of your choice.</p> <p>V. Write a GNU <code>makefile</code> under your root directory so that the command “<code>make run</code>” would trigger the compilation of your source code, the production of the ex-</p> | <p>ecutable, and the running of your tests.</p> <p>VI. Write a report to summarize the main points of your numerical experiments, which should be designed to verify the analytic results in the notes. To some extent, your grade depends on the number of key conclusions in the notes confirmed by your numerical experiments.</p> |
|--|---|