

8 Finite-difference methods for BVPs

In this Lecture, we consider methods for solving BVPs whose (methods') idea is to replace the BVPs with a system of approximating algebraic equations. In the first five sections, we will (mostly) deal with linear BVPs, so that the corresponding system of equations will be linear. The last section will show how nonlinear BVPs can be approached.

8.1 Matrix problem for the discretized solution

Let us begin by considering a linear BVP with Dirichlet boundary conditions:

$$\begin{aligned} y'' + P(x)y' + Q(x)y &= R(x), \\ y(a) &= \alpha, \quad y(b) = \beta. \end{aligned} \tag{8.1}$$

As before, we assume that P , Q , and R are twice continuously differentiable, so that y is a four times continuously differentiable function of x . Also, we consider the case where $Q(x) \leq 0$ on $[a, b]$, so that the BVP (8.1) is guaranteed by Theorem 6.2 to have a unique solution.

Let us replace y'' and y' in (8.1) by their second-order accurate discretizations:

$$y'' = \frac{1}{h^2}(y_{n+1} - 2y_n + y_{n-1}) + O(h^2), \tag{8.2}$$

$$y' = \frac{1}{2h}(y_{n+1} - y_{n-1}) + O(h^2). \tag{8.3}$$

Upon omitting the $O(h^2)$ -terms and substituting (8.2) and (8.3) into (8.1), we obtain the following system of linear equations:

$$\begin{aligned} Y_0 &= \alpha; \\ (1 + \frac{h}{2}P_n)Y_{n+1} - (2 - h^2Q_n)Y_n + (1 - \frac{h}{2}P_n)Y_{n-1} &= h^2R_n, \quad 1 \leq n \leq N-1; \\ Y_N &= \beta. \end{aligned} \tag{8.4}$$

In the matrix form, this is

$$A\vec{Y} = \vec{r}, \tag{8.5}$$

where

$$A = \begin{pmatrix} -(2 - h^2Q_1) & (1 + \frac{h}{2}P_1) & 0 & 0 & \cdots & 0 \\ (1 - \frac{h}{2}P_2) & -(2 - h^2Q_2) & (1 + \frac{h}{2}P_2) & 0 & \cdots & 0 \\ 0 & (1 - \frac{h}{2}P_3) & -(2 - h^2Q_3) & (1 + \frac{h}{2}P_3) & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & \cdots & 0 & (1 - \frac{h}{2}P_{N-2}) & -(2 - h^2Q_{N-2}) & (1 + \frac{h}{2}P_{N-2}) \\ 0 & \cdots & 0 & 0 & 0 & (1 - \frac{h}{2}P_{N-1}) & -(2 - h^2Q_{N-1}) \end{pmatrix}, \tag{8.6}$$

$\vec{Y} = [Y_1, Y_2, \dots, Y_{N-1}]^T$, and

$$\vec{r} = \left[h^2R_1 - \left(1 - \frac{h}{2}P_1\right)\alpha, \quad h^2R_2, \quad h^2R_3, \quad \dots, \quad h^2R_{N-2}, \quad h^2R_{N-1} - \left(1 + \frac{h}{2}P_{N-1}\right)\beta \right]^T; \tag{8.7}$$

the superscript 'T' in (8.7) denotes the transpose.

From Linear Algebra, it is known that the linear system (8.5) has a unique solution if (and only if) the matrix A is nonsingular. Therefore, below we list some results that will allow us to guarantee that under certain conditions, a particular A is nonsingular.

Gerschgorin Circles Theorem, 8.1 Let a_{ij} be entries of an $M \times M$ matrix A and let λ_k , $k = 1, \dots, M$ be the eigenvalues of A . Then:

(i) Each eigenvalue lies in the union of “row circles” R_i , where

$$R_i = \left\{ z : |z - a_{ii}| \leq \sum_{j=1, j \neq i}^M |a_{ij}| \right\}. \quad (8.8)$$

(In words, $\sum_{j=1, j \neq i}^M |a_{ij}|$ is the sum of all off-diagonal entries in the i th row.)

(ii) Similarly, since A and A^T have the same eigenvalues,³⁰ then each eigenvalue also lies in the union of “column circles” C_j , where

$$C_j = \left\{ z : |z - a_{jj}| \leq \sum_{i=1, i \neq j}^M |a_{ij}| \right\}. \quad (8.9)$$

(In words, $\sum_{i=1, i \neq j}^M |a_{ij}|$ is the sum of all off-diagonal entries in the j th column.)

(iii) Let $\bigcup_{i=k}^{i=l} R_i$ be a cluster of $(l - k + 1)$ row circles that is disjoint from all the other row circles. Then it contains exactly $(l - k + 1)$ eigenvalues.

A similar statement holds for column circles.

Example Use the Gerschgorin Circles Theorem to estimate eigenvalues of a matrix

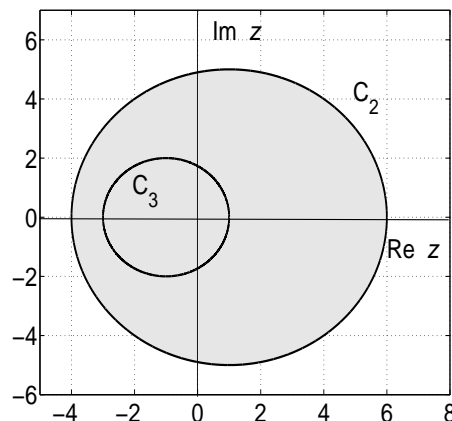
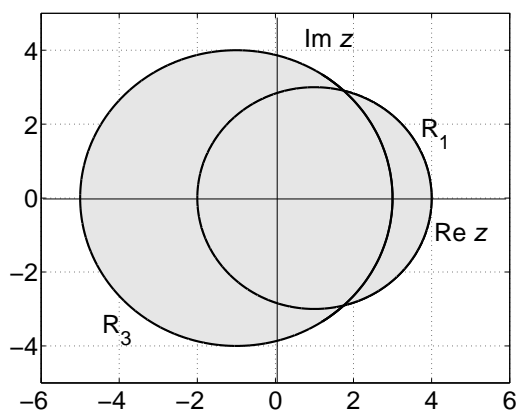
$$E = \begin{pmatrix} 1 & 2 & -1 \\ 1 & 1 & 1 \\ 1 & 3 & -1 \end{pmatrix}. \quad (8.10)$$

Solution The circles are listed and sketched below:

$$R_1 = \{z : |z - 1| \leq |2| + |1| = 3\} \quad C_1 = \{z : |z - 1| \leq |1| + |1| = 2\}$$

$$R_2 = \{z : |z - 1| \leq |1| + |1| = 2\} \quad C_2 = \{z : |z - 1| \leq |2| + |3| = 5\}$$

$$R_3 = \{z : |z + 1| \leq |1| + |3| = 4\} \quad C_3 = \{z : |z + 1| \leq |1| + |-1| = 2\}$$



³⁰This was covered in your Linear Algebra class (without proof).

(Circles R_2 and C_1 are not sketched because they are concentric with, and lie entirely within, circles R_1 and C_2 , respectively.)

Gerschgorin Circles Theorem says that the eigenvalues of E must lie *both* within $\bigcup_{i=1}^3 R_i$ and within $\bigcup_{j=1}^3 C_j$. Therefore, they must lie within the intersection of $\bigcup_{i=1}^3 R_i$ and $\bigcup_{j=1}^3 C_j$. For example, this gives that $|\operatorname{Re} \lambda| \leq 4$ for each and any of the eigenvalues.

Before we can apply the Gerschgorin Circles Theorem, we need to introduce some new terminology.

Definition Matrix A is called diagonally dominant if

$$\text{either } |a_{ii}| \geq \sum_{j=1, j \neq i}^M |a_{ij}| \text{ or } |a_{ii}| \geq \sum_{j=1, j \neq i}^M |a_{ji}|, \quad 1 \leq i \leq M, \quad (8.11)$$

with the *strict inequality holding for at least one i* .

A matrix is called strictly diagonally dominant (SDD) if

$$\text{either } |a_{ii}| > \sum_{j=1, j \neq i}^M |a_{ij}| \text{ or } |a_{ii}| > \sum_{j=1, j \neq i}^M |a_{ji}| \quad \text{for all } i = 1, \dots, M. \quad (8.12)$$

In other words, in a SDD matrix, the sums of the off-diagonal elements along *either* every row *or* every column are less than the corresponding diagonal entries.

Theorem 8.2 If a matrix A is SDD, then it is nonsingular.

Proof If A is SDD, then one of the inequalities (8.12) holds. Suppose the inequality for the rows holds. Comparing that inequality with the r.h.s. of (8.8), we conclude, by the Gerschgorin Circles Theorem, that point $\lambda = 0$ is outside of the union $\bigcup_{i=1}^M R_i$ of Gerschgorin circles. Hence it is automatically outside of the intersection of the unions $\bigcup_{i=1}^M R_i$ and $\bigcup_{i=1}^M C_i$. (Make sure you see why this is so.) Therefore, $\lambda = 0$ is not an eigenvalue of A , hence A is nonsingular. **q.e.d.**

Theorem 8.3 Consider the BVP (8.1). If $Q(x) \leq 0$, and if $P(x)$ is bounded on $[a, b]$ (i.e. $|P(x)| \leq \mathcal{P}$ for some \mathcal{P}), then the discrete version (8.4) of the BVP in question has a unique solution, provided that the step size satisfies $h\mathcal{P} \leq 2$.

Proof Case (a): $Q(x) < 0$. In this case, matrix A in (8.6) is SDD, provided that $h\mathcal{P} \leq 2$. Indeed, A is tridiagonal, with the diagonal elements being

$$a_{ii} = -(2 - h^2 Q_i), \quad \Rightarrow \quad |a_{ii}| > 2. \quad (8.13)$$

The sum of the absolute values of the off-diagonal elements is

$$|a_{i,i+1}| + |a_{i,i-1}| = \left| 1 + \frac{h}{2} P_i \right| + \left| 1 - \frac{h}{2} P_i \right| = \left(1 + \frac{h}{2} P_i \right) + \left(1 - \frac{h}{2} P_i \right) = 2. \quad (8.14)$$

In removing the absolute value signs in the above equation, we have used the fact that $h\mathcal{P} \leq 2$. Now, comparing (8.13) with (8.14), we see that

$$|a_{i,i+1}| + |a_{i,i-1}| < |a_{ii}|,$$

which means that A is SDD and hence by Theorem 8.2, (8.4) has a unique solution. **q.e.d.**

Case (b): $Q(x) \leq 0$ requires a more involved proof, which we omit.

Note: We emphasize that Theorem 8.3 gives a bound for the step size,

$$h \cdot \max_{x \in [a, b]} |P(x)| \leq 2, \quad (8.15)$$

which is a sufficient condition for the discretization (8.4) (with $Q(x) \leq 0$) to be solvable.

8.2 Thomas algorithm

In the previous section, we have considered the issue of the *possibility* to obtain the unique solution of the linear system (8.4), which (the system) approximates the BVP (8.1). In this section, we will consider the issue of solving (8.4) *in an efficient manner*. The key fact that will allow us to do so is the tridiagonal form of A ; that is, A has nonzero entries only on the main diagonal and on the subdiagonals directly above and below it. It can be shown that in order to solve³¹ a linear system of the form (8.5) with a *full* $M \times M$ matrix A , one requires $O(M^3)$ operations. However, this would be an inefficient way to solve a linear system with a tridiagonal A . Below we will present an algorithm, which has been discovered independently by several researchers in the 50's, that allows one to solve a linear system with a tridiagonal matrix using only $O(M)$ operations.

One way to numerically solve a linear system of the form

$$A\vec{y} = \vec{r} \quad (8.5)$$

(with *any* matrix A) is via so-called LU decomposition. Namely, one seeks two matrices L (lower triangular) and U (upper triangular) such that³²

$$LU = A. \quad (8.16)$$

Then (8.5) is solved in two steps:

$$\text{Step 1 : } L\vec{z} = \vec{r} \quad \text{and} \quad \text{Step 2 : } U\vec{y} = \vec{z}. \quad (8.17)$$

The linear systems in (8.17) are then solved using forward (for Step 1) and backward (for Step 2) substitutions; details of this will be provided later on.

When A is tridiagonal, both finding the matrices L and U and solving the systems in (8.17) requires only $O(M)$ operations. You will be asked to demonstrate that in a homework problem. Here we will present details of the algorithm itself.

Let

$$A = \begin{pmatrix} b_1 & c_1 & 0 & 0 & \cdot & \cdot & \cdot & 0 \\ a_2 & b_2 & c_2 & 0 & \cdot & \cdot & \cdot & 0 \\ 0 & a_3 & b_3 & c_3 & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & \cdot & 0 & a_{M-1} & b_{M-1} & c_{M-1} \\ 0 & \cdot & \cdot & \cdot & 0 & 0 & a_M & b_M \end{pmatrix}; \quad (8.18)$$

then we seek L and U in the form:

$$LU = \begin{pmatrix} 1 & 0 & 0 & \cdot & 0 \\ \alpha_2 & 1 & 0 & \cdot & 0 \\ 0 & \alpha_3 & 1 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & \alpha_M & 1 \end{pmatrix} \begin{pmatrix} \beta_1 & c_1 & 0 & \cdot & 0 \\ 0 & \beta_2 & c_2 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & \beta_{M-1} & c_{M-1} \\ 0 & \cdot & 0 & 0 & \beta_M \end{pmatrix} \quad (8.19)$$

³¹by Gaussian elimination or by any other direct (i.e., non-iterative) method

³²Disclaimer: The theory of LU decomposition is considerably more involved than the simple excerpt from it given here. We will *not* go into further details of that theory in this course.

Multiplying the matrices in (8.19) and comparing the result with (8.18), we obtain:

$$\begin{aligned} \text{Row 1: } & \beta_1 = b_1; \\ \text{Row 2: } & \alpha_2\beta_1 = a_2, \quad \alpha_2c_1 + \beta_2 = b_2; \\ \text{Row } j: & \alpha_j\beta_{j-1} = a_j, \quad \alpha_jc_{j-1} + \beta_j = b_j. \end{aligned} \tag{8.20}$$

Equations (8.20) are easily solved for the unknown coefficients α_j and β_j :

$$\begin{aligned} \beta_1 &= b_1, \\ \alpha_j &= a_j/\beta_{j-1}, \quad \beta_j = b_j - \alpha_jc_{j-1}, \quad j = 2, \dots, M. \end{aligned} \tag{8.21}$$

Finally, we show how the systems in (8.17) can be solved. Let

$$\vec{z} = [z_1, z_2, \dots, z_M]^T, \quad \text{etc.}$$

Then the forward substitution in $L\vec{z} = \vec{r}$ gives:

$$\begin{aligned} z_1 &= r_1, \\ z_j &= r_j - \alpha_jz_{j-1}, \quad j = 2, \dots, M. \end{aligned} \tag{8.22}$$

The backward substitution in $U\vec{y} = \vec{z}$ gives:

$$\begin{aligned} y_M &= z_M/\beta_M, \\ y_j &= (z_j - c_jy_{j+1})/\beta_j, \quad j = M-1, \dots, 1. \end{aligned} \tag{8.23}$$

Thus, \vec{y} is found in terms of \vec{r} .

The entire procedure (8.21) to (8.23) is fast, as we said earlier, and also requires the storage of only 8 one-dimensional arrays of size $O(M)$ for the coefficients $\{a_j\}$, $\{b_j\}$, $\{c_j\}$, $\{\alpha_j\}$, $\{\beta_j\}$, $\{r_j\}$, $\{z_j\}$, and $\{y_j\}$. Moreover, it is possible to show that when A is SDD, i.e. when

$$|b_j| > |a_j| + |c_j|, \quad j = 1, \dots, M, \tag{8.24}$$

then small round-off (or any other) errors do not get amplified by this algorithm; specifically, the numerical error remains small and independent of the size M of the problem.³³

To conclude this section, we note that similar algorithms exist for other banded (e.g., pentadiagonal) matrices. The details of those algorithms can be found, e.g., in Sec. 3-2 of W. Ames, "Numerical methods for partial differential equations," 3rd ed. (Academic Press, 1992).

8.3 Error estimates, and higher-order discretization

In this section, we will state, without proof, two theorems about the accuracy of the solutions of systems of discretized equations approximating a given BVP.

Theorem 8.4 Let $\{Y_n\}_{n=1}^{N-1}$ be the solution of the discretized problem (8.4) (or, which is the same, (8.5)–(8.7)) and let $y(x)$ be the exact solution of the original BVP (8.1); then $\epsilon_n = y(x_n) - Y_n$ is the error of the numerical solution. In addition, let $\mathcal{P} = \max_{x \in [a, b]} |P(x)|$

³³In line with the earlier disclaimer about the LU decomposition, we also note that condition (8.24) is just the simplest of possible conditions that guarantee boundedness of the error. Other conditions exist and are studied in advanced courses in Numerical Analysis.

and also let $Q(x) \leq \mathcal{Q} < 0$ (recall that $Q(x) \leq 0$ is required for the BVP to have a unique solution). Then the error satisfies the following estimate:

$$\max |\epsilon_n| \leq \frac{1}{h^2 \left(|\mathcal{Q}| + \frac{8}{(b-a)^2} \right)} \left(\frac{1}{12} h^4 (M_4 + \mathcal{P}M_3) + 2\rho \right), \quad (8.25)$$

where $M_3 = \max_{x \in [a, b]} |y'''|$, $M_4 = \max_{x \in [a, b]} |y''''|$, and ρ is the round-off error.

When the round-off error is neglected, estimate (8.25) yields that the discrete approximation (8.4) to the BVP produces the error on the order of $O(h^2)$ (i.e., in other words, is *second-order accurate*). At first sight, this is similar to how the finite-difference approximations (8.2) and (8.3) led to second-order accurate methods for IVPs. In fact, for both the IVPs and BVPs, the expression inside the largest parentheses on the r.h.s. of (8.25) is the local truncation error. However, the interpretation of the $O(1/h^2)$ factor in front of those parentheses differs from the interpretation of the similar factor for IVPs. In the latter case, that factor arose from accumulation of the error over $O(1/h)$ steps. *On the contrary, for the BVPs*, that factor arises due to solving the linear system (8.5) and, more specifically, due to needing to invert matrix A . This is explained in more details in Sec. 8.6 below: see the discussion leading to (8.68). In brief, the multiplication of the vector on the r.h.s. of (8.5) by A^{-1} causes a factor $O(1/h^2)$ in the solution. This fact, along with the local truncation error $\vec{\delta}\mathbf{r}$ being $O(h^4)$, implies estimate (8.25) for the solution error.

The numerical error of the discrete approximation of the BVP (8.1) can be significantly reduced if instead of the simple central-difference approximation to $y'' = f(x, y)$, as in (8.2), one uses the Numerov's formula (5.18). Specifically, if y' does not enter the BVP, i.e. when the BVP is

$$y'' = f(x, y), \quad y(a) = \alpha, \quad y(b) = \beta, \quad (8.26)$$

then Numerov's formula leads to the following system of discrete equations:

$$\begin{aligned} Y_0 &= \alpha; \\ Y_{n+1} - 2Y_n + Y_{n-1} &= \frac{h^2}{12} (f_{n+1} + 10f_n + f_{n-1}) \quad 1 \leq n \leq N-1; \\ Y_N &= \beta. \end{aligned} \quad (8.27)$$

In particular, when $f(x, y) = -Q(x)y + R(x)$, system (8.27) is linear. Since the local truncation error of Numerov's method is $O(h^6)$, the solution of system (8.27) must be a fourth-order accurate approximation to the exact solution of the BVP (8.26). More precisely, the following theorem holds true:

Theorem 8.5 Let $\{Y_n\}_{n=1}^{N-1}$ be the solution of the discretized problem (8.27), where $f(x, y) = -Q(x)y + R(x)$ (so that this system is linear). Then the error $\epsilon_n = y(x_n) - Y_n$ satisfies the following estimate:

$$\max |\epsilon_n| \leq \frac{1}{h^2 \left(|\mathcal{Q}| + \frac{8}{(b-a)^2} \right)} \left(\frac{1}{240} h^6 M_6 + 2\rho \right), \quad (8.28)$$

where the notations are the same as in Theorem 8.4 and, in addition, $M_6 = \max_{x \in [a, b]} |y^{(6)}|$.

8.4 Neumann and mixed boundary conditions

Suppose that at, say, $x = a$, the boundary condition is

$$A_1 y(a) + A_2 y'(a) = \alpha, \quad (8.29)$$

where A_1 and A_2 are some constants. (Recall that for $A_1 = 0$, this is the Neumann boundary condition, while for $A_1 \cdot A_2 \neq 0$, this condition is of the mixed type.) Below we present two methods that allow one to construct generalizations of systems (8.4) (simple central-difference discretization) and (8.27) (Numerov's formula) for the case of the boundary condition (8.29) instead of

$$y(a) = \alpha.$$

Note that the problem we need to handle for such a generalization is obtaining an approximation to $y'(a)$ in (8.29) that would have the accuracy consistent with the accuracy of the discretization scheme. That is, the accuracy needs to be second-order for a generalization of (8.4) and fourth-order for a generalization of (8.27).

Method 1

This method is efficient only for a generalization of the second-order accurate approximation. Thus, we will look for a second-order accurate finite-difference approximation to $y'(a)$.

Introduce a fictitious point $x_{-1} = a - h$ and let the approximate solution at that point be Y_{-1} . Once one knows Y_{-1} , the approximation sought is

$$y'(a) = \frac{Y_1 - Y_{-1}}{2h} + O(h^2), \quad (8.30)$$

so that the equation approximating the *boundary condition* (8.29) is

$$\underline{x = a} : \quad A_1 Y_0 + A_2 \frac{Y_1 - Y_{-1}}{2h} = \alpha. \quad (8.31)$$

The discrete analog of the ODE itself at $x = a$ is:

$$\left(1 + \frac{h}{2}P_0\right) Y_1 - (2 - h^2Q_0)Y_0 + \left(1 - \frac{h}{2}P_0\right) Y_{-1} = h^2R_0. \quad (8.32)$$

Equations (8.31) and (8.32) should replace the first equation,

$$Y_0 = \alpha, \quad (8.33)$$

in (8.4), so that the dimension of the resulting system of equations in (8.4) increases from $(N - 1) \times (N - 1)$ to $(N + 1) \times (N + 1)$. Note that adding two new equations (8.31) and (8.32) to system (8.4) is consistent with introducing two new unknowns Y_{-1} and Y_0 .

Later on we will see that, rather than dealing with two equations (8.31) and (8.32), it is more practical to solve (8.31) for Y_{-1} and substitute the result in (8.32). Then, instead of two equations (8.31) and (8.32), we will have one equation that needs to replace (8.33):

$$2Y_1 - \left[2 - h^2Q_0 - 2h\frac{A_1}{A_2} \left(1 - \frac{h}{2}P_0\right)\right] Y_0 = h^2R_0 + 2h\frac{\alpha}{A_2} \left(1 - \frac{h}{2}P_0\right). \quad (8.34)$$

(We are justified to assume in (8.34) that $A_2 \neq 0$, since otherwise the boundary condition (8.29) becomes Dirichlet.) The resulting system then contains $1 + (N - 1) = N$ equations for the N unknowns: Y_0 through Y_{N-1} , and hence can be solved for a unique solution $\{Y_n\}_{n=0}^{N-1}$, unless the coefficient matrix of this system is singular. We will remark on the latter possibility after we describe the other method.

Method 2

This method does *not* use a fictitious point to approximate the boundary condition with the

required accuracy. We will first show how this method can be used for the second-order accurate approximation (8.4) and then indicate how it (the method) can be modified for the fourth-order accurate approximation (8.27).

By analogy with Eq. (3.10) of Lecture 3, one can obtain:

$$\frac{y_{n+1} - y_n}{h} = y'_n + \frac{h}{2}y''_n + O(h^2). \quad (8.35)$$

Then, using the ODE in (8.1) to express y'' as $R - Py' - Qy$, one finds:

$$\begin{aligned} y'_n &= \frac{y_{n+1} - y_n}{h} - \frac{h}{2}(R_n - P_n y'_n - Q_n y_n) + O(h^2) \\ &= \frac{y_{n+1} - y_n}{h} - \frac{h}{2} \left(R_n - P_n \left[\frac{y_{n+1} - y_n}{h} + O(h) \right] - Q_n y_n \right) + O(h^2). \end{aligned} \quad (8.36)$$

Therefore, the boundary condition (8.29) can be approximated by

$$A_1 Y_0 + A_2 \left[\frac{Y_1 - Y_0}{h} - \frac{h}{2} \left(R_0 - P_0 \frac{Y_1 - Y_0}{h} - Q_0 Y_0 \right) \right] = \alpha, \quad (8.37)$$

with the discretization error being $O(h^2)$. In an extra-credit homework problem, you will be given the chance to show that (8.37) can be put in a form similar to that of (8.34). Such a rearranged Eq. (8.37) should then be used along with the $N - 1$ equations in (8.4) for $\{Y_n\}_{n=1}^{N-1}$. Thus we have a system of N equations for N unknowns $\{Y_n\}_{n=0}^{N-1}$, which can be solved (again, unless its coefficient matrix is singular).

When one solves the BVP (8.26) with its Dirichlet boundary conditions replaced by the mixed boundary conditions (8.29) while requiring *fourth-order* accuracy, one needs to use the Numerov scheme for grid points inside the interval $[a, b]$ and a fourth-order accurate approximation to $y'(a)$. The latter can be obtained by using two more terms in expansion (8.35):

$$\frac{y_{n+1} - y_n}{h} = y'_n + \frac{h}{2}y''_n + \frac{h^2}{6}y'''_n + \frac{h^3}{24}y''''_n + O(h^4). \quad (8.38)$$

Recall that our goal is to find y'_n (for $n = 0$) using quantities that we can compute. Clearly, one can compute the l.h.s. of (8.38). Furthermore, one can obtain y''_n from $y''_n = f(x_n, y_n) \equiv f_n$. It now remains to compute y'''_n and y''''_n , which can be done as follows. Note that

$$f_{n+1} \equiv f(x_{n+1}, y_{n+1}) = f_n + h \frac{d}{dx} f_n + \dots = y''_n + h y'''_n + \frac{h^2}{2} y''''_n + O(h^3), \quad (8.39)$$

and similarly,

$$f_{n+2} = f_n + (2h) \frac{d}{dx} f_n + \dots = y''_n + (2h) y'''_n + \frac{(2h)^2}{2} y''''_n + O(h^3). \quad (8.40)$$

These two equations can be solved for $h y'''_n$ and $h^2 y''''_n$ in terms of the quantities on the l.h.s. and y''_n . Doing so and substituting the result in (8.38) yields the desired fourth-order approximation for y'_n :

$$y'_n = \frac{y_{n+1} - y_n}{h} - \frac{h}{24}(7f_n + 6f_{n+1} - f_{n+2}) + O(h^4). \quad (8.41)$$

Equation (8.41) with $n = 0$ can then be used to obtain a fourth-order accurate counterpart of (8.37).

Remark 1: Recall that the coefficient matrix in the discretized BVP (8.4) with Dirichlet boundary conditions is SDD and hence nonsingular, provided that h satisfies the condition of Theorem 8.3. We will now state a requirement that the coefficients A_1, A_2 in the non-Dirichlet boundary condition (8.29) must satisfy in order to guarantee that the corresponding BVP has a unique solution. To this end, we consider the situations arising in Methods 1 and 2 separately.

When Eq. (8.34) replaces the Dirichlet boundary condition (8.33) in Method 1, the resulting coefficient matrix is SDD provided that

$$A_1 A_2 \leq 0, \quad (8.42)$$

in addition to the earlier requirements of $Q(x) < 0$ and $h \cdot \max |P(x)| \leq 2$. (You will be asked to verify (8.42) in one of the Questions for self-assessment.) On the other hand, if two individual equations (8.31) and (8.32) are used instead of their combined form (8.34), the corresponding coefficient matrix is *no longer* SDD. This is *one advantage* of using (8.34) instead of (8.31) and (8.32).

Thus, condition (8.42) is sufficient, but not necessary, to guarantee that the corresponding discretized BVP has a unique solution. That is, even when (8.42) does not hold, one can still attempt to solve the corresponding linear system, since strict diagonal dominance is only a sufficient, but not necessary, condition for its matrix A to be nonsingular.

In the case of Method 2, by collecting the coefficients of Y_0 and Y_1 in Eq. (8.37), it is straightforward (although, perhaps, a little tedious) to show that the coefficient matrix is SDD if (8.42) and the two conditions stated one line below it, are satisfied. You will verify this in a homework problem. The fourth-order accurate analog of (8.37), based on Eq. (8.41), also yields the same conditions for strict diagonal dominance of the coefficient matrix.

Remark 2: In the case of the BVP with Dirichlet boundary conditions, the coefficient matrix is tridiagonal (see Eq. (8.6)), and hence the corresponding linear system can be solved efficiently by the Thomas algorithm. In the case of non-Dirichlet boundary conditions, one can show (and you will be asked to do so) that Method 1 based on Eq. (8.34) yields a tridiagonal system, but the same method using Eqs. (8.31) and (8.32) does not. This is the *other advantage* of using the single Eq. (8.34) in this method.

Method 2 for a second-order accurate approximation to the BVP gives a tridiagonal matrix. This can be straightforwardly shown by the same rearrangement of terms in Eq. (8.37) that was used above to show that the corresponding matrix is SDD. However, the fourth-order accurate modification of (8.37), based on (8.41), produces a matrix \tilde{A} that is no longer tridiagonal. One can handle this situation in two ways. First, if one is willing to sacrifice one order of accuracy in exchange for the convenience of having a tridiagonal coefficient matrix, then instead of (8.41), one can use

$$y'_n = \frac{y_{n+1} - y_n}{h} - \frac{h}{6}(2f_n + f_{n+1}) + O(h^3). \quad (8.43)$$

Then the modification of (8.37) based on (8.43) does result in a tridiagonal coefficient matrix. Alternatively, one can see that matrix \tilde{A} obtained with the fourth-order accurate formula (8.41) differs from a tridiagonal one only in having its (1, 3)th entry nonzero (verify). Thus, \tilde{A} is in some sense “very close” to a tridiagonal matrix, and we can hope that this fact could be used to find \tilde{A}^{-1} with only $O(M)$ operations. This can indeed be done by using the algorithm described in the next section.

8.5 Periodic boundary condition; Sherman–Morrison algorithm

In this section, we will consider the case when BVP (8.1) has periodic boundary conditions. We will see that the coefficient matrix arising in this case is *not* tridiagonal, but is, in some sense, close to it. We will then present an algorithm that will allow us to find the inverse of that matrix using only $O(M)$ operations, where $M \times M$ is the dimension of the matrix. The same method can also be used in other situations, including the inversion of matrix \tilde{A} defined in the last paragraph of the preceding section.

We first obtain the analogues of Eqs. (8.5)–(8.7) in the case of the BVP having periodic boundary conditions. Consider the corresponding counterpart of the BVP (8.1):

$$\begin{aligned} y'' + P(x)y' + Q(x)y &= R(x), \\ y(a) &= y(b). \end{aligned} \tag{8.44}$$

The corresponding counterpart of system (8.4) is

$$\begin{aligned} Y_0 &= Y_N; \\ (1 + \frac{h}{2}P_n)Y_{n+1} - (2 - h^2Q_n)Y_n + (1 - \frac{h}{2}P_n)Y_{n-1} &= h^2R_n, \quad 0 \leq n \leq N. \end{aligned} \tag{8.45}$$

Note 1: Since our problem has periodic boundary conditions, it is logical to let

$$Y_{-n} = Y_{N-n} \quad \text{and} \quad Y_{N+n} = Y_n, \quad 0 \leq n \leq N.$$

In particular,

$$Y_{-1} = Y_{N-1} \quad \text{and} \quad Y_N = Y_0 \tag{8.46}$$

(the last equation here is just the original periodic boundary condition).

Note 2: The index n of yet unknown values Y_n in (8.45) runs from 0 to $N - 1$, while in (8.4) it runs from 1 to $N - 1$.

In view of Eq. (8.46), the equations in system (8.45) with $n = 0$ and $n = N - 1$ can be written as follows:

$$\begin{aligned} (1 + \frac{h}{2}P_0)Y_1 - (2 - h^2Q_0)Y_0 + (1 - \frac{h}{2}P_0)Y_{N-1} &= h^2R_0, \\ (1 + \frac{h}{2}P_{N-1})Y_0 - (2 - h^2Q_{N-1})Y_{N-1} + (1 - \frac{h}{2}P_{N-1})Y_{N-2} &= h^2R_{N-1}. \end{aligned} \tag{8.47}$$

With this result, system (8.45) can be written in the matrix form (8.5), where now

$$\vec{Y} = [Y_0, Y_1, \dots, Y_{N-1}]^T, \tag{8.48}$$

$$A = \begin{pmatrix} -(2 - h^2Q_0) & (1 + \frac{h}{2}P_0) & 0 & 0 & \dots & \boxed{(1 - \frac{h}{2}P_0)} \\ (1 - \frac{h}{2}P_1) & -(2 - h^2Q_1) & (1 + \frac{h}{2}P_1) & 0 & \dots & 0 \\ 0 & (1 - \frac{h}{2}P_2) & -(2 - h^2Q_2) & (1 + \frac{h}{2}P_2) & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & (1 - \frac{h}{2}P_{N-2}) & -(2 - h^2Q_{N-2}) & (1 + \frac{h}{2}P_{N-2}) \\ \boxed{(1 + \frac{h}{2}P_{N-1})} & \dots & 0 & 0 & (1 - \frac{h}{2}P_{N-1}) & -(2 - h^2Q_{N-1}) \end{pmatrix} \tag{8.49}$$

and

$$\vec{r} = [h^2R_0, h^2R_1, h^2R_2, \dots, h^2R_{N-2}, h^2R_{N-1}]^T. \tag{8.50}$$

Matrix A in Eq. (8.49) differs from its counterpart in Eq. (8.6) in two respects: (i) its dimension is $N \times N$ rather than $(N - 1) \times (N - 1)$ and (ii) it is not tridiagonal due to the terms in its upper-right and lower-left corners. Such matrices are called *circulant*.

Thus, to obtain the solution of the BVP (8.44), we will need to solve the linear system (8.5) with the non-tridiagonal matrix A . We will now show how this problem can be reduced to the solution of a system with a tridiagonal matrix. To this end, we first make a preliminary observation. Let \vec{w} be a $N \times 1$ vector whose only nonzero entry is its i th entry and equals w_i . Let \vec{z} be a $N \times 1$ vector whose only nonzero entry is its j th entry and equals z_j . Then $C = \vec{w}\vec{z}^T$ is an $N \times N$ matrix whose only nonzero entry is $C_{ij} = w_i z_j$ (verify). Similarly, if $\vec{w} = [w_1, 0, 0, \dots, w_N]^T$ and $\vec{z} = [z_1, 0, 0, \dots, z_N]^T$, then

$$C = \vec{w}\vec{z}^T = \begin{pmatrix} w_1 z_1 & 0 & \cdots & 0 & w_1 z_N \\ 0 & 0 & \cdots & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 0 & 0 \\ w_N z_1 & 0 & \cdots & 0 & w_N z_N \end{pmatrix}. \quad (8.51)$$

Therefore, the circulant matrix A in Eq. (8.49) can be represented as:

$$A = A_{\text{tridiag}} + \vec{w}\vec{z}^T, \quad (8.52)$$

where A_{tridiag} is some tridiagonal matrix and \vec{w} and \vec{z} are some properly chosen vectors. Note that while the choice of \vec{w} and \vec{z} allows much freedom, the form of A_{tridiag} is unique for a given circulant matrix A , once \vec{w} and \vec{z} have been chosen. In an extra-credit homework problem, you will be given the chance to make a choice for \vec{w} and \vec{z} and consequently come up with the expression for A_{tridiag} , given that A is as in Eq. (8.49).

Linear systems (8.5) with the coefficient matrix A given by (8.52) can be time-efficiently — i.e., in $O(M)$ operations — solved by the so-called Sherman–Morrison algorithm. This algorithm can be found in most textbooks on Numerical Analysis, or online.

8.6 Nonlinear BVPs

The analysis presented in this section is carried out with three restrictions. First, we consider only BVPs with Dirichlet boundary conditions. Generalizations for boundary conditions of the form (8.29) can be done straightforwardly along the lines of Secs. 8.4 and 8.5.

Second, we only consider the BVPs

$$y'' = f(x, y), \quad y(a) = \alpha, \quad y(b) = \beta, \quad (8.26)$$

which does not involve y' . Although the methods described below can also be adopted to a more general BVP

$$y'' = f(x, y, y'), \quad y(a) = \alpha, \quad y(b) = \beta, \quad (7.1)$$

the analysis of *convergence* of those methods to a solution of the BVP (7.1) is significantly more complex than the one presented here. Thus, the methods that we will develop in this section can be applied to (7.1) *without a guarantee* that one will obtain a solution of that BVP.

Third, we will focus our attention on the second-order accurate discretization of the BVPs. The analysis for the fourth-order accurate discretization scheme is essentially the same, and produces similar results.

Consider the BVP (8.26). The counterpart of system (8.4) for this BVP has the same form as (8.5), i.e.:

$$A\vec{Y} = \vec{r}, \quad (8.5)$$

where now

$$A = \begin{pmatrix} -2 & 1 & 0 & 0 & \cdot & \cdot & \cdot & 0 \\ 1 & -2 & 1 & 0 & \cdot & \cdot & \cdot & 0 \\ 0 & 1 & -2 & 1 & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & \cdot & 0 & 1 & -2 & 1 \\ 0 & \cdot & \cdot & \cdot & 0 & 0 & 1 & -2 \end{pmatrix}; \quad \vec{\mathbf{r}} = \begin{pmatrix} h^2 f(x_1, Y_1) - \alpha \\ h^2 f(x_2, Y_2) \\ \dots \\ h^2 f(x_{N-2}, Y_{N-2}) \\ h^2 f(x_{N-1}, Y_{N-1}) - \beta \end{pmatrix}. \quad (8.53)$$

Matrix A above is found from (8.6) upon setting $P(x) \equiv 0$ and $Q(x) \equiv 0$. The r.h.s. vector $\vec{\mathbf{r}}$ is found from (8.7) where we treat the nonlinear function $f(x, y)$ in the same way as we treated the term $R(x)$ in (8.1).

Equations (8.5) and (8.53) constitute a system of nonlinear algebraic equations. Below we consider three methods of *iterative* solution of such a system. Of these methods, Method 1 is an analogue of the fixed-point iteration method for solving a single linear equation³⁴

$$\mathcal{A} \cdot y = r(y), \quad (8.54)$$

Method 2 is a modification of Method 1, and Method 3 is the analog of the Newton–Raphson method for (8.54).

Method 1 (Picard iterations)

The fixed-point iteration scheme, also called the Picard iteration scheme, for the single nonlinear equation (8.54) is simply

$$y^{(k+1)} = \frac{1}{\mathcal{A}} r(y^{(k)}), \quad (8.55)$$

where $y^{(k)}$ denotes the k th iteration of the solution of (8.54). To start the iteration scheme (8.55), one, of course, needs an *initial guess* $y^{(0)}$.

Now, let $\vec{\mathbf{Y}}^{(k)}$ denote the k th iteration of the solution of the matrix nonlinear equation (8.5), (8.53). Then the corresponding Picard iteration scheme is

$$A\vec{\mathbf{Y}}^{(k+1)} = \vec{\mathbf{r}}(\vec{\mathbf{x}}, \vec{\mathbf{Y}}^{(k)}), \quad k = 0, 1, 2, \dots \quad (8.56)$$

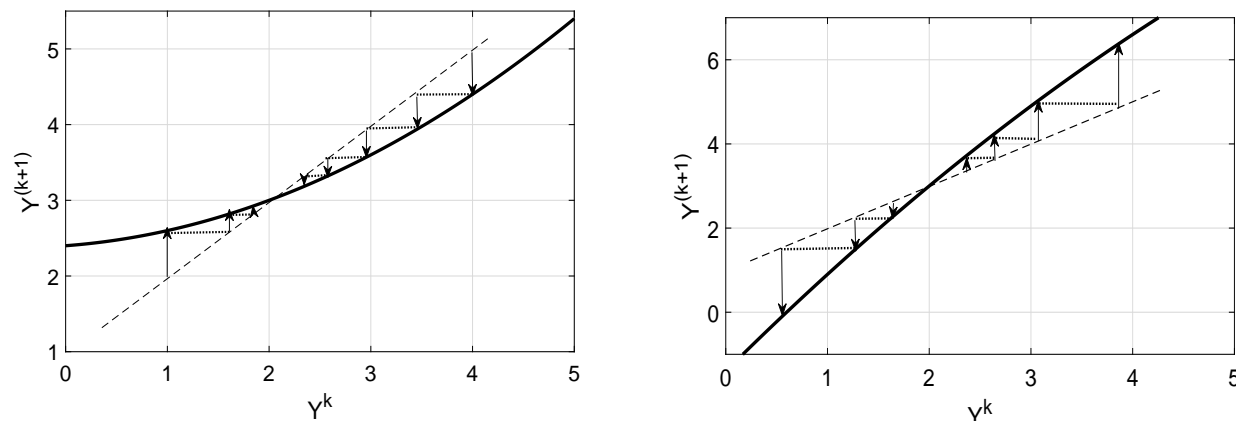
Importantly, unlike the original nonlinear Eqs. (8.5) and (8.53), Eq. (8.56) is a *linear* system. Indeed, at the $(k+1)$ th iteration, the unknown vector $\vec{\mathbf{Y}}^{(k+1)}$ enters linearly, while the nonlinear r.h.s. contains $\vec{\mathbf{Y}}^{(k)}$, which has been determined at the previous iteration.

Let us investigate the rate of convergence of the iteration scheme (8.56). For its one-variable counterpart, Eq. (8.54), the convergence condition of the fixed-point iterations (8.55) is known to be

$$\left| \mathcal{A}^{-1} \cdot \frac{dr(y)}{dy} \right| < 1 \quad (8.57)$$

for all y sufficiently close to the fixed point. This is illustrated in the figures below, where the thick solid curve represents the r.h.s. of Eq. (8.55) and the thin dashed line represents $Y^{(k+1)} = Y^{(k)}$. The arrows show the evolution of consecutive iterations. In the left figure panel, where condition (8.57) holds, iterations approach the so-called fixed point, which is the solution of (8.55). In the right panel, where the *opposite* of condition (8.57) holds, the iterations are seen to move away from the fixed point.

³⁴The constant \mathcal{A} in (8.54) could have been absorbed into the function $r(y)$, but we kept it so as to mimic the notations of (8.5).



The condition we will establish below will be a direct analog of (8.57) for the vector equation (8.56).

Before we proceed, let us remark that the analysis of convergence of Picard’s iteration scheme (8.56) can proceed in two slightly different ways. In one way, one can transform (8.56) into an iteration scheme for $\vec{\delta}^{(k)} = \vec{Y}^{(k)} - \vec{Y}^{(k-1)}$. Then the conditions that the iterations converge is expressed by

$$\|\vec{\delta}^{(k)}\| < \|\vec{\delta}^{(k-1)}\| \quad \text{for all sufficiently large } k, \tag{8.58}$$

where, as in Lecture 4 (see (4.1)), $\|\dots\|$ denotes the ∞ -norm:

$$\|\vec{\delta}\| \equiv \|\vec{\delta}\|_{\infty} = \max_n |\delta_n|.$$

Indeed, condition (8.58) implies that $\lim_{k \rightarrow \infty} \|\vec{\delta}^{(k)}\| = 0$. This, in its turn, implies that the sequence of iterations $\{\vec{Y}^{(k)}\}$ tends to a limit, which then, according to (8.56), must be a solution of (8.5).

We will, however, proceed in a slightly different way. Namely, we will assume that our starting guess, $\vec{Y}^{(0)}$, is sufficiently close to the exact solution, \vec{Y} , of (8.5). Then, as long as the iterations converge, $\vec{Y}^{(k)}$ will stay close to \vec{Y} for all k , and one can write

$$\vec{Y}^{(k)} = \vec{Y} + \vec{\varepsilon}^{(k)}, \quad \text{where } \|\vec{\varepsilon}^{(k)}\| \ll 1. \tag{8.59}$$

Quantity $\vec{\varepsilon}^{(k)}$ is the “error” in the sense that it shows how different the current iterated solution is close to the exact solution of the *discretized* BVP. Thus, it is different from the error $\vec{\varepsilon}$ due to discretization itself, defined in Sec. 8.3.

The condition that iterations (8.56) converge has the same form as (8.58):

$$\|\vec{\varepsilon}^{(k)}\| < \|\vec{\varepsilon}^{(k-1)}\| \quad \text{for all sufficiently large } k. \tag{8.60}$$

However, its interpretation is slightly different (although equivalent): Now the fact that $\lim_{k \rightarrow \infty} \|\vec{\varepsilon}^{(k)}\| = 0$ implies that $\lim_{k \rightarrow \infty} \vec{Y}^{(k)} = \vec{Y}$, i.e. that the iterative solutions converge to the exact solution of (8.5).

Both methods of convergence analysis described above can be shown to yield the same conclusions. We chose to follow the second method, based on (8.59), because it can be more

naturally related to the *linearization* of the nonlinear equation at hand (see below). Linearization will allow us to replace the analysis of the original nonlinear equation (8.5) by the analysis of a system of linear equations, which can be carried out using well-developed methods of Linear Algebra.

Thus, we begin the convergence analysis of the iteration scheme (8.56) by substituting there expression (8.59) and linearizing the right-hand side using the first two terms of the Taylor expansion near \vec{Y} :

$$A\vec{Y} + A\vec{\varepsilon}^{(k+1)} = \vec{r}(\vec{x}, \vec{Y}) + \frac{\partial \vec{r}}{\partial \vec{Y}} \vec{\varepsilon}^{(k)} + O(\|\vec{\varepsilon}^{(k)}\|^2). \quad (8.61)$$

The first terms on both sides of the above equation cancel out by virtue of the exact equation (8.5). Then, upon discarding the quadratically small terms $O(\|\vec{\varepsilon}^{(k)}\|^2)$ and using the definition of \vec{r} from (8.53), one obtains a linear system

$$A\vec{\varepsilon}^{(k+1)} = \frac{\partial \vec{r}}{\partial \vec{Y}} \vec{\varepsilon}^{(k)}, \quad \text{where } \frac{\partial \vec{r}}{\partial \vec{Y}} = h^2 \text{diag} \left(\frac{\partial f(x_1, Y_1)}{\partial Y_1}, \dots, \frac{\partial f(x_{N-1}, Y_{N-1})}{\partial Y_{N-1}} \right). \quad (8.62)$$

We will use this equation to relate the norms of $\vec{\varepsilon}^{(k+1)}$ and $\vec{\varepsilon}^{(k)}$ and thereby establish a sufficient condition for convergence of iterations (8.62).

Let us assume that

$$\max_{1 \leq n \leq N-1} \left| \frac{\partial f(x_n, Y_n)}{\partial Y_n} \right| = L. \quad (8.63)$$

Multiplying both sides of (8.56) by A^{-1} and taking their norm, one obtains:

$$\|\vec{\varepsilon}^{(k+1)}\| \leq \|A^{-1}\| \cdot h^2 L \cdot \|\vec{\varepsilon}^{(k)}\|, \quad (8.64)$$

where we have also used a known fact from Linear Algebra, stating that for any matrix A and vector \vec{z} ,

$$\|A\vec{z}\| \leq \|A\| \cdot \|\vec{z}\|$$

(actually, the latter inequality follows directly from the definition of a matrix norm). For completeness, let us mention that

$$\|A\|_\infty = \max_{1 \leq i \leq M} \sum_{j=1}^M |a_{ij}|. \quad (8.65)$$

Inequality (8.64) shows that

$$\|\vec{\varepsilon}^{(k)}\| \leq (h^2 L \|A^{-1}\|)^k \|\vec{\varepsilon}^{(0)}\|, \quad (8.66)$$

which implies that Picard iterations converge when

$$h^2 L \|A^{-1}\| < 1. \quad (8.67)$$

As promised earlier, this condition is analogous to (8.57).

It now remains to find $\|A^{-1}\|$. Since matrix A shown in Eq. (8.53) arises in a great many applications, the explicit form of its inverse has been calculated for any size $N = (b - a)/h$. The derivation of A^{-1} can be found on photocopied pages posted on the course website; the corresponding result for $\|A^{-1}\|$, obtained with the use of (8.65), is:

$$\|A^{-1}\| = \frac{(b - a)^2}{8h^2}. \quad (8.68)$$

Substituting (8.68) into (8.67), we finally obtain that *for Picard iterations to converge, it is sufficient (but not necessary) that*

$$\frac{(b-a)^2}{8}L < 1. \tag{8.69}$$

Thus, whether the Picard iterations converge to the discretized solution of the BVP (8.26) depends not only on the function $f(x, y)$ but also on the length of the interval $[a, b]$.

Remark 3: Notice the similarity of the r.h.s. of (8.68) with the denominator of Eq. (8.25), except that here we are dealing with the case $\mathcal{Q} = 0$. This was also briefly discussed in the paragraph after that equation.

Method 2 (modified Picard iterations)

The idea of the modified Picard iterations method can be explained using the example of the single equation (8.54), where we will set $\mathcal{A} = 1$ for convenience and without loss of generality. Suppose the simple fixed-point iteration scheme (8.55) does not converge because at the fixed point \bar{y} , $dr(\bar{y})/dy \approx \kappa > 1$,³⁵ so that the convergence condition (8.57) is violated. Then, instead of iterating (8.54), let us iterate

$$y - \kappa y = r(y) - \kappa y \quad \Rightarrow \quad (1 - \kappa)y^{(k+1)} = r(y^{(k)}) - \kappa y^{(k)} \quad \Rightarrow \quad y^{(k+1)} = \frac{r(y^{(k)}) - \kappa y^{(k)}}{1 - \kappa}. \tag{8.70}$$

Note that the *exact* equation which we start with in (8.70) is equivalent to Eq. (8.54) (with $\mathcal{A} = 1$), but the iteration equation in (8.70) is *different* from (8.55). If our guess κ at the true value of the derivative $dr(\bar{y})/dy$ is “sufficiently close”, then the derivative of the r.h.s. of (8.70) is less than 1, and hence, by (8.57), the iteration scheme (8.70) converges, in contrast to (8.54), which diverges. In other words, by subtracting from both sides of the equation a linear term whose slope closely matches the slope of the nonlinear term at the solution \bar{y} , one drastically reduces the magnitude of the slope of the right-hand side of the iterative scheme, thereby making it converge.

Let us return to the iterative solution of Eq. (8.26), where now, instead of iterating, as in Picard’s method, the equation

$$(y^{(k+1)})'' = f(x, y^{(k)}), \tag{8.71}$$

we will iterate the equation

$$(y^{(k+1)})'' - c y^{(k+1)} = f(x, y^{(k)}) - c y^{(k)} \tag{8.72}$$

with some constant c . The corresponding linearized equation in vector form is obtained similarly to (8.62):

$$(A - h^2 c I)\vec{\varepsilon}^{(k+1)} = h^2 \text{diag} \left(\frac{\partial f(x_1, Y_1)}{\partial Y_1} - c, \dots, \frac{\partial f(x_{N-1}, Y_{N-1})}{\partial Y_{N-1}} - c \right) \vec{\varepsilon}^{(k)}, \tag{8.73}$$

where I is the identity matrix of the same size as A . The c -terms in (8.73) are the counterparts of such terms in (8.72), where on the l.h.s. we used the identity $I\vec{\varepsilon} = \vec{\varepsilon}$ that holds for any vector $\vec{\varepsilon}$.

We will now address the question of how one should choose the constant c so as to ensure convergence of (8.73) and hence of the modified scheme (8.72).

³⁵Here ‘ \approx ’ is used instead of ‘ $=$ ’ because $dr(\bar{y})/dy$ is usually not known exactly.

The main difference between the multi-component equation (8.73) and the single-component equation (8.70) is that no single value of c could simultaneously match all of the values $\partial f(x_n, Y_n)/\partial Y_n$, which, in general, are distributed in some interval

$$L^- \leq \frac{\partial f(x_n, Y_n)}{\partial Y_n} \leq L^+, \quad n = 1, \dots, N-1. \quad (8.74)$$

It may be intuitive to suppose that the optimal choice for c may be at the midpoint of that interval, i.e.,

$$c_{\text{opt}} = L^{\text{av}} = \frac{1}{2}(L^- + L^+). \quad (8.75)$$

Below we will show that this is indeed the case.

Specifically, let us only consider the case where $\partial f/\partial y > 0$, when a unique solution of BVP (8.26) is guaranteed to exist by Theorem 6.1 of Lecture 6. Then in (8.74), both

$$L^\pm > 0, \quad (8.76)$$

and hence $L^{\text{av}} > 0$. Next, by following the steps of the derivation of (8.64) one obtains from (8.73):

$$\|\bar{\varepsilon}^{(k+1)}\| \leq \left(\|(A - h^2 c I)^{-1}\| \cdot h^2 \max_{L^- \leq \ell \leq L^+} |\ell - c| \right) \cdot \|\bar{\varepsilon}^{(k)}\|. \quad (8.77)$$

Here ℓ stands for any of $\partial f(x_n, Y_n)/\partial Y_n$, which satisfy (8.74). Note that the maximum above is taken with respect to ℓ while c is assumed to be fixed. Then:

$$L^- - c \leq \ell - c \leq L^+ - c \quad \Rightarrow \quad \max_{L^- \leq \ell \leq L^+} |\ell - c| = \max\{|c - L^-|, |L^+ - c|\}. \quad (8.78)$$

Our immediate goal is to determine for what c the coefficient multiplying $\|\bar{\varepsilon}^{(k)}\|$ in (8.77) is the smallest: this will yield the fastest convergence of the modified Picard iterations (8.72). The entire analysis of this question is a little tedious, since one will have to consider separately the cases where $c \in [L^-, L^+]$ and $c \notin [L^-, L^+]$. Since we have announced that the answer, (8.75), corresponds to the former case, we will present the details only for it. Details for the case $c \notin [L^-, L^+]$ are similar, but will not yield an optimal value of c , and so we will omit them.

Thus, we are looking to determine

$$K \equiv \min_{L^- \leq c \leq L^+} \left(h^2 \|(A - h^2 c I)^{-1}\| \max\{|c - L^-|, |L^+ - c|\} \right), \quad (8.79)$$

for which we will first need to find the norm $\|(A - h^2 c I)^{-1}\|$. Since A is a symmetric matrix (see (8.53)), so are matrices $(A - h^2 c I)$ and $(A - h^2 c I)^{-1}$. (Theorems to that effect were presented in your undergraduate Linear Algebra course.) In graduate-level courses on Linear Algebra it is shown that the norm of a real symmetric matrix B equals $|\lambda_B|_{\text{max}}$, where λ_B are the eigenvalues of B . Since (an eigenvalue of B^{-1}) = 1/(an eigenvalue of B), then

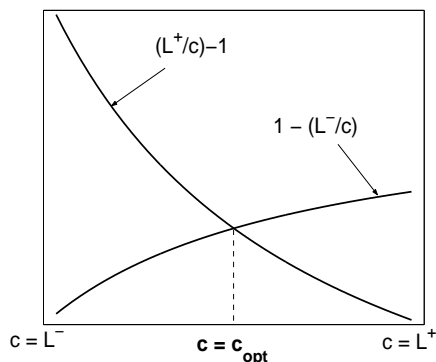
$$\|(A - h^2 c I)^{-1}\| = \frac{1}{\min |\text{eigenvalue of } (A - h^2 c I)|}. \quad (8.80a)$$

To find the lower bound for the latter eigenvalue, we use the result of Problem 2 of HW 8 (which is based on the Gerschgorin Circles Theorem of Section 8.1). Using this result, one finds:

$$\|(A - h^2 c I)^{-1}\| = \frac{1}{h^2 c} \quad (8.80b)$$

(recall that $c > 0$ because we assumed in (8.76) that $L^\pm > 0$ and also since $c \in [L^-, L^+]$). Combining (8.79) and (8.80b), we will now determine

$$\begin{aligned} K &= \min_{L^- \leq c \leq L^+} \left(\max \left\{ \left| \frac{c - L^-}{c} \right|, \left| \frac{L^+ - c}{c} \right| \right\} \right) \\ &= \min_{L^- \leq c \leq L^+} \left(\max \left\{ 1 - \frac{L^-}{c}, \frac{L^+}{c} - 1 \right\} \right). \end{aligned} \tag{8.81}$$



The corresponding optimal c is found as shown in the figure on the left:

$$\begin{aligned} 1 - \frac{L^-}{c_{\text{opt}}} &= \frac{L^+}{c_{\text{opt}}} - 1 \\ \Rightarrow c_{\text{opt}} &= \frac{1}{2}(L^- + L^+), \end{aligned}$$

which is (8.75).

Substituting (8.75) into (8.81) and then using the result in (8.77), one finally obtains:

$$\|\bar{\varepsilon}^{(k+1)}\| \leq \left(\frac{L^+ - L^-}{L^+ + L^-} \right) \|\bar{\varepsilon}^{(k)}\|. \tag{8.82}$$

Since according to our assumption (8.76) $L^\pm > 0$, then the factor $(L^+ - L^-)/(L^+ + L^-) < 1$. Thus, the absolute value of the error of the iterated solution $\bar{\mathbf{Y}}^{(k)}$ decreases with each iterations, and the modified Picard scheme (8.72), (8.75) converges.

The issues of using the modified Picard iterations on BVPs where conditions (8.76) are not met, or using scheme (8.72) with a non-optimal constant c , are considered in homework problems. Let us only emphasize that the modified Picard iterations *can* sometimes converge even when (8.76) and/or (8.75) do not hold.

Method 3 (Newton–Raphson method)

Although this method can be described for a general BVP of the form (7.1), we will only do so for a BVP of the more restricted form (8.26):

$$y'' = f(x, y), \quad y(a) = \alpha, \quad y(b) = \beta. \tag{8.26}$$

Let us begin by writing down the system of second-order accurate discrete equations for this BVP for internal (i.e., non-boundary) points $0 < n < N$:

$$Y_{n+1} - 2Y_n + Y_{n-1} = f(x_n, Y_n) \tag{8.83a}$$

Soon, we will need its equivalent form:

$$Y_{n+1} - 2Y_n + Y_{n-1} - f(x_n, Y_n) = 0. \tag{8.83b}$$

Let $Y_n^{(0)}$ be the initial guess for the solution of (8.83) at x_n . Similarly to (8.59), we relate it to the *exact* solution $\{Y_n\}$ of (8.83):

$$Y_n^{(0)} = Y_n + \varepsilon_n^{(0)}, \quad |\varepsilon_n^{(0)}| \ll 1. \tag{8.84a}$$

Equivalently:

$$Y_n = Y_n^{(0)} - \varepsilon_n^{(0)}. \quad (8.84b)$$

Substituting (8.84b) into (8.83b) and using Taylor expansion, we obtain:

$$\left\{ Y_{n+1}^{(0)} - 2Y_n^{(0)} + Y_{n-1}^{(0)} - f(x_n, Y_n^{(0)}) \right\} - \left\{ \varepsilon_{n+1}^{(0)} - 2\varepsilon_n^{(0)} + \varepsilon_{n-1}^{(0)} - \varepsilon_n^{(0)} \frac{\partial f(x_n, Y_n^{(0)})}{\partial Y_n^{(0)}} \right\} + O(\|\bar{\varepsilon}^{(0)}\|^2) = 0. \quad (8.85a)$$

For example, for $f(x, y) = y^2/(2+x)$, used in one of the homework problems in HW 7 and also in HW 8, $\partial f(x, Y)/\partial Y = 2Y/(2+x)$. Now, recall that we know the initial guess $Y_n^{(0)}$ and seek the correction term $(-\varepsilon_n^{(0)})$ that will bring it closer to the exact solution; see (8.84b). Therefore, following the convention for linear systems, we keep the unknown variables on the l.h.s. and move all known quantities on the r.h.s. Also, we will neglect the quadratically small terms. The result is:

$$\left\{ \varepsilon_{n+1}^{(0)} - 2\varepsilon_n^{(0)} + \varepsilon_{n-1}^{(0)} - \varepsilon_n^{(0)} \frac{\partial f(x_n, Y_n^{(0)})}{\partial Y_n^{(0)}} \right\} = \left\{ Y_{n+1}^{(0)} - 2Y_n^{(0)} + Y_{n-1}^{(0)} - f(x_n, Y_n^{(0)}) \right\} \quad (8.85b)$$

If our initial guess satisfies the boundary conditions of the BVP (8.26), then we also have

$$\varepsilon_0^{(0)} = 0, \quad \varepsilon_N^{(0)} = 0. \quad (8.86)$$

System (8.85b), (8.86) is linear and tridiagonal, and so it can be solved time-efficiently by the Thomas algorithm. Thus we obtain $\varepsilon_n^{(0)}$. According to (8.84), this gives the next iteration for our solution $\{Y_n\}$:

$$Y_n \approx Y_n^{(1)} = Y_n^{(0)} - \varepsilon_n^{(0)}. \quad (8.87)$$

Remark 4: The reason that the r.h.s. of (8.87) does not equal the exact solution Y_n *exactly* is that we had to neglect the $O(\bar{\varepsilon}^{(0)})$ term in (8.85b). We had no choice but do that because, while one knows how to solve *linear* systems of equations, there is no formula or algorithm for solving nonlinear systems of equations.

We now substitute

$$Y_n^{(1)} = Y_n + \varepsilon_n^{(1)} \quad (8.88)$$

into (8.83), obtain a system analogous to (8.85), and then solve it for $\varepsilon_n^{(1)}$ in the same way we have solved that system for $\varepsilon_n^{(0)}$. Repeating these steps, we stop when the process converges, i.e. $\|\vec{Y}^{(k+1)} - \vec{Y}^{(k)}\|$ becomes less than a given tolerance.

Remark 5 (brief comparison of Methods 1 – 3): The goal of all these Methods is to solve the *nonlinear* system of equations (8.56). In each Method, this is done by solving a *linear* system. The Methods differ in which linear system they solve.

Using iterative methods in more complicated situations

First, let us note that iterative methods, which we have described above for ODEs, are equally applicable to partial differential equations (PDEs), where the function $f(\vec{x}, y)$ is defined not on an interval $x \in [a, b]$ as above, but over a two- or three-dimensional region in the space where \vec{x} “lives”. This distinguishes iterative methods from the shooting method considered in Lecture 7. Indeed, as we stressed at the end of that Lecture, the shooting method *cannot* be

extended to solving BVPs for PDEs. Thus, iterative methods remain *the only* group of method for solving nonlinear BVPs for PDEs. The ideas of these methods are the same as we have described above for ODEs.

Second, the fixed-point iteration methods (i.e., counterparts of the Picard and modified Picard described above) are not used (or are rarely used) in commercial software. The main reason is that they are considerably slower than the Newton–Raphson method and its variations and also than so-called *Krylov subspace methods*, studied in advanced courses on Numerical Linear Algebra. We will only mention that the most famous of those Krylov subspace methods is the Conjugate Gradient method (CGM) for solving symmetric positive (or negative) definite linear systems. An introduction to this method can be found in many textbooks and also in an online paper by J.R. Shewchuk, “An Introduction to the Conjugate Gradient Method Without the Agonizing Pain”.³⁶ An extension of the CGM to nonlinear BVPs whose linearization yields symmetric matrices is known as the Newton–CGM or, more generally, as a class of Newton–Krylov methods.

Another reason why Krylov subspace methods are used much more widely than fixed-point iterative methods is that convergence conditions of the latter methods are typically restricted by a condition analogous to (8.76). Some of Krylov subspace methods either do not have restrictions, or are less sensitive to them. Also, the Newton–Raphson method does not have any restrictions similar to (8.76). Yet, this method also has its own issues, and a number of books are written about application of the Newton–Raphson method to systems of nonlinear equations.

8.7 Questions for self-assessment

1. Verify that (8.4) follows from (8.1)–(8.3).
2. Explain why the form of the first and last terms of $\vec{\mathbf{r}}$ in (8.7) is different from the form of the other terms.
3. What does the Gerschgorin Circles Theorem allow one to do?
4. What is the difference between a diagonally dominant and strictly diagonally dominant matrices?
5. Make sure you can follow the derivations in the Example about the Gerschgorin Circles Theorem.
6. Make sure you can follow the proof of Theorem 8.2.
7. What can you say about a solution of a linear system (8.5) where A is SDD?
8. Under what condition(s) on h is the discretized BVP (8.4) guaranteed to have a unique solution?
9. Verify (8.20).
10. Verify (8.21) through (8.23).
11. What causes the local truncation error of the discretized BVP to be multiplied by a factor $O(\frac{1}{h^2})$ to give the error of the solution?

³⁶Some pain, however, is to be expected.

12. Describe the idea of Method 1 for handling non-Dirichlet boundary conditions.
13. Derive (8.34).
14. Describe the idea of Method 2 for handling non-Dirichlet boundary conditions.
15. Suppose that you need to solve a BVP with a mixed-type boundary condition and such that (8.42) does *not* hold. Will you even attempt to solve such a BVP? What should your expectations be?
16. Suppose that the BVP has a non-Dirichlet boundary condition of the form (8.29) at $x = b$ (the right end point of the interval) rather than at $x = a$. What will the analog of condition (8.42) be in this case?
Hint 1: Remember that this is a QSA and not an exercise requiring calculations.
Hint 2: To better visualize the situation, suppose $[a, b] = [-1, 1]$, so that $-1 \leq x \leq 1$. What simple mathematical operation, when applied to x , will swap the left and right end points of this interval? How will this operation transform each of the terms on the l.h.s. of (8.29)?
17. Convince yourself (and be prepared to convince the instructor) that the statements in Remark 2 in Sec. 8.4 about the coefficient matrices for the second-order methods being tridiagonal, are correct.
18. Describe the idea of the Picard iterations.
19. Derive (8.64) as explained in the text.
20. What is the condition for the Picard iterations to converge? Is this a necessary or sufficient condition? In other words, if that condition does not hold, should one still attempt to solve the BVP?
21. Describe the idea behind the modified Picard iterations.
22. Make sure you can obtain (8.73).
23. How is the finite-difference *implementation* of (8.72) different from that of (8.56)? That is, to solve (8.72), will you need to define a matrix and a vector differently than when solving (8.56)? You will need to answer this question in order to do a homework problem.
24. Make sure you can obtain (8.77) and (8.78).
25. Explain how the result of Problem 1 of HW 8 leads to (8.80b).
26. Derive (8.75) from the explanation found after (8.81).
27. Obtain (8.82) as explained in the text.
28. Are (8.76) and (8.75) necessary or sufficient conditions for convergence of the modified Picard iterations (8.72)?
29. Make sure you can follow the derivation of (8.85).
30. Write down a linear system satisfied by $\epsilon_n^{(1)}$ defined in (8.88).
31. Describe the idea behind the Newton–Raphson method.

32. Give equation numbers of the three linear systems referred to in Remark 5.
33. Can one use iterative methods when solving nonlinear BVPs for PDEs?