

From Stepwise Integrations and Low-rank Updates to a Pseudospectral Solution Operator Matrix for the Helmholtz Operator $\frac{d}{dx}a(x)\frac{d}{dx} + c(x)$

Yung-Ta Li, Ping-Hsuan Tsai and Chun-Hao Teng*

Abstract. In this study we propose a construction framework utilizing stepwise integrations and the Sherman-Morrison-Woodbury formula to seek pseudospectral integration preconditioning matrices for differential operators. We illustrate this framework through formulating an inverse matrix for the Helmholtz differential operator of the form $\frac{d}{dx}a(x)\frac{d}{dx} + c(x)$. Numerical experiments were conducted to examine the performance of the derived operator. The results show that the inverse matrix is an effective solution operator to numerically solve general second order differential equations.

1. Introduction

Spectral and pseudospectral methods [21, 22], due to their exponential convergence property of approximating smooth functions, have been applied to accurately solve differential equations for decades. However, to enjoy the advantage of the methods, great care must be exercised to effectively resolve those ill-conditioned systems of equations, resulting from the pseudospectral differentiation matrices. Generally speaking, the condition number of a pseudospectral differentiation matrix is scaled with $\mathcal{O}(N^{2m})$, where m is the order of the approximated differential operator, and N refers to the grid resolution of a mesh. As a result, one often encounters inaccurate computations due to this ill-conditioned issue, especially when high-order differential equations are solved on dense grids.

To overcome this obstacle, many methods have been proposed to construct operating matrices to effectively invert numerical differentiations. In [4–6, 23–25, 27–29], operating matrices are derived based on recombined basis functions satisfying boundary conditions. In [10, 20], operating matrices are formulated based on orthogonal polynomials with boundary conditions enforced through different techniques, for instance, boundary bordering [20] or auxiliary equations deduced from boundary conditions [10]. Based on an integration approach, operating or preconditioning matrices have been developed [15] for basic differentiation matrix operators with boundary conditions enforced weakly in a

Received April 16, 2020; Accepted August 12, 2020.

Communicated by Suh-Yuh Yang.

2010 *Mathematics Subject Classification.* 74S25, 65N12, 65F05.

Key words and phrases. spectral/pseudospectral methods, integration preconditioning, low-rank, Sherman-Morrison-Woodbury.

*Corresponding author.

penalty manner [8, 11, 12, 16, 17]. Arbitrary-grids-based integration preconditioning matrices have been constructed [13]. Beyond the commonly used Lagrange's interpolation formulation, pseudospectral integration matrix operators have been developed based on a suitable Birkhoff interpolation [30]. Indeed, applying these preconditioning operators significantly reduces the condition numbers of the targeted differentiation matrices, and numerical equations are solved effectively.

In the aforementioned studies, most of the developed operating matrices are primarily for pure differential operators, and the derivations rely upon recurrence relationships between orthogonal polynomials and their derivatives. A nice feature shared by these matrices is that they are banded matrices in spectral spaces. Thus, once discretized equations are transformed into spectral spaces, problems can be solved efficiently by exploring the sparsity of those banded matrices. However, for mixed differential operators, variable coefficients are sandwiched by differential operators. The resulting operating matrices constructed based on those recurrence relationships for mixed differential operators are generally full matrices, and the advantage of solving problems in spectral spaces is lost.

Numerical differentiation involving variable coefficients also encounters a similar situation, if variable coefficients and differential operators are integrated into a single matrix through using the orthogonal polynomial recurrence identities. However, as proposed in the seminal study [21, 22], a simple and widely used technique for mixed differentiation is through a differentiation-multiplication-differentiation procedure. This long-established technique has been brought to our attention, and a question is raised. Is it possible to seek inverse matrices for mixed differential operators through reversing the stepwise procedure? The study [19] gives an affirmative answer for the radial direction Laplace differential operators in polar coordinates.

The goal of this study is to realize an integration-division-integration framework to construct integration preconditioning matrices for differential operators. Our construction framework is based on two natural properties of nodal-based pseudospectral formulations: (1) a general pseudospectral differentiation operator, whether it is pure or mixed, is composed of multiple matrices with each one being either the first order differentiation matrix or a simple diagonal one, and (2) boundary operators are by themselves low-rank matrices. These two properties enlighten us to consider combining stepwise integrations and low-rank updates to construct inverse operators in a factored form for differentiation matrices. We demonstrate this approach through deriving an inverse matrix for the Helmholtz differential operator, $\frac{d}{dx}a(x)\frac{d}{dx} + c(x)$, based on the Legendre pseudospectral approximation method with boundary conditions enforced weakly through a penalty approach.

The construction procedures are summarized below. In [15], two boundary condition penalized (BCP) first order differentiation matrices were proposed and their inverses were

developed. We adopt these inverse operators as building blocks. Employing these inverse operators, we express the BCP pseudospectral Laplace operator, resulting from discretizing the differential operator $\frac{d}{dx}a(x)\frac{d}{dx}$, into a factored form, a diagonal matrix perturbed by a low-rank one, inserted in between the two basic BCP differentiation matrices. Since the inverse matrices of the BCP basic differentiation operators are available, seeking the inverse of the BCP pseudospectral Laplace differentiation operator boils down to finding the inverse of the sandwiched matrix. We complete this step by applying the Sherman-Morrison-Woodbury (SMW) formula. As a result, the inverse of the BCP pseudospectral Laplace operator is the matrix product of these three inverse matrices. Once the inverse of the BCP pseudospectral Laplace operator is obtained, we then recursively apply the SMW formula to construct the inverse matrix of the BCP pseudospectral Helmholtz operator.

The derived inverse pseudospectral Helmholtz matrix serves as a solution operator for advection-diffusion-reaction problems. We conduct various numerical experiments to examine the performance of the inverse operator. The convergence results show that the derived operator is effective, even for nonlinear problems.

The potential applications of the proposed framework advertise the significance of the framework itself. In general, numerical differential equations are systems of equations composed of pure/mixed differentiation matrices and low-rank numerical boundary operators. By exploring these properties as we will show in this study for the Helmholtz operator, one may adopt the present divide-and-conquer approach to construct efficient preconditioning matrices for other numerical differential equations. We will address relevant issues at the end of this study.

The rest of the paper is organized as follows. In the next section we demonstrate the proposed framework through seeking an inverse pseudospectral matrix for the targeted mixed differential operator. Section 3 is devoted to the validations of the derived operator. Conclusions are given in the last section.

2. Formulation

2.1. Basic concepts of the Legendre pseudospectral method and the SMW formula

Let $\mathbb{I} = [-1, 1]$ be the pseudospectral reference domain, $N \geq 0$ be an integer, and $P_N(x)$ defined on \mathbb{I} be the N -th degree Legendre polynomial satisfying the differential equation

$$\left((1-x^2)P'_N(x)\right)' = -N(N+1)P_N(x), \quad x \in \mathbb{I},$$

where $'$ denotes the differentiation with respect to the function argument. The roots of the polynomial $(1-x^2)P'_N(x)$, ordered ascendantly as $-1 = x_0 < x_1 < \dots < x_{N-1} < x_N = 1$, are known as the Gauss-Lobatto-Legendre (GLL) grid points. Denoted by $l_j(x)$, the

Lagrange basis polynomials based on the GLL points are given as

$$l_j(x) = -\frac{(1-x^2)P'_N(x)}{N(N+1)(x-x_j)P_N(x_j)}, \quad j = 0, 1, \dots, N$$

satisfying the property $l_j(x_i) = \delta_{ij}$ with δ_{ij} being the Kronecker delta function.

Consider a function $u(x)$ defined on I . Employing the Lagrange basis polynomials, $l_j(x)$, we can construct an N -th degree polynomial to approximate $u(x)$ and its derivative $u'(x)$ as follows:

$$u(x) \approx \mathcal{I}_N[u(x)] = u_N(x) = \sum_{j=0}^N l_j(x)u(x_j), \quad u' \approx u'_N(x) = \sum_{j=0}^N l'_j(x)u(x_j),$$

where \mathcal{I}_N is called the interpolation operator. Introducing the differentiation matrix \mathbf{D} with entries $d_{ij} = l'_j(x_i)$ for $0 \leq i, j \leq N$, we can compute the values of $u'_N(x)$ at the GLL nodes by a matrix-vector multiplication approach:

$$\mathbf{u}'_N = \mathbf{D}\mathbf{u}_N, \quad \mathbf{u}_N = [u(x_0), \dots, u(x_N)]^T, \quad \mathbf{u}'_N = [u'_N(x_0), \dots, u'_N(x_N)]^T,$$

where T denotes the vector transpose.

Notice that the numerical differentiation becomes an exact differentiation, provided that $u(x)$ is a polynomial of degree at most N . As a consequence, the numerical differentiation of the n -th degree Legendre polynomial, $P_n(x)$, is exact for $n \leq N$, i.e.,

$$\mathbf{p}'_n = \mathbf{D}\mathbf{p}_n, \quad \mathbf{p}_n = [P_n(x_0), \dots, P_n(x_N)]^T, \quad \mathbf{p}'_n = [P'_n(x_0), \dots, P'_n(x_N)]^T.$$

Among these derivative grid vectors, \mathbf{p}'_0 , \mathbf{p}'_1 and \mathbf{p}'_N are important to the present study. For the sake of convenience, we introduce the following vector expressions

$$\mathbf{e}_- = [1, 0, \dots, 0]^T, \quad \mathbf{e}_+ = [0, \dots, 0, 1]^T, \quad \mathbf{1} = [1, 1, \dots, 1]^T, \quad \mathbf{0} = [0, 0, \dots, 0]^T.$$

Then, the expressions of \mathbf{p}'_0 , \mathbf{p}'_1 and \mathbf{p}'_N are simply

$$\mathbf{p}'_0 = \mathbf{0}, \quad \mathbf{p}'_1 = \mathbf{p}_0 = \mathbf{1}, \quad \mathbf{p}'_N = \mathbf{D}\mathbf{p}_N = \frac{N(N+1)}{2}(\mathbf{e}_+ - (-1)^N\mathbf{e}_-),$$

where the first and the second expressions are due to the fact that $P_0(x) = 1$ and $P_1(x) = x$, and the last one results from the fact that $P'_N(x_i)$ vanishes at each interior GLL grid point and $P'_N(\pm 1) = (\pm 1)^{N+1}N(N+1)/2$.

Associated with the GLL nodes is the quadrature integration rule:

$$\int_{-1}^1 u(x) dx = \sum_{i=0}^N u(x_i)\omega_i, \quad \omega_i = \frac{2}{N(N+1)(P_N(x_i))^2},$$

where ω_i for $i = 0, 1, \dots, N$ are the quadrature weights, and the exactness of the integration rule is held true provided that $u(x)$ is a polynomial of degree at most $2N - 1$. Employing the integration rule, we define the diagonal mass matrix \mathbf{M} and the stiffness matrix \mathbf{S} with their entries m_{ij} and s_{ij} , respectively, given by

$$\begin{aligned} m_{ij} &= \sum_{k=0}^N l_i(x_k) l_j(x_k) \omega_k = \omega_i \delta_{ij}, & 0 \leq i, j \leq N, \\ s_{ij} &= \sum_{k=0}^N l_i(x_k) l'_j(x_k) \omega_k = \omega_i l'_j(x_i), & 0 \leq i, j \leq N. \end{aligned}$$

Three important properties concerning the matrices \mathbf{M} , \mathbf{S} , and \mathbf{D} are shown in [3]. We directly quote the results here: (1) the mass matrix \mathbf{M} is positive definite and diagonal, and thus invertible, (2) the stiffness matrix \mathbf{S} is almost skew-symmetric,

$$\begin{aligned} \mathbf{S} + \mathbf{S}^T &= \mathbf{I}_+ - \mathbf{I}_-, \\ \mathbf{I}_+ &= \mathbf{e}_+ \mathbf{e}_+^T = \text{diag}(0, \dots, 0, 1), \\ \mathbf{I}_- &= \mathbf{e}_- \mathbf{e}_-^T = \text{diag}(1, 0, \dots, 0), \end{aligned} \tag{2.1}$$

and (3) the matrices, \mathbf{M} , \mathbf{S} , and \mathbf{D} , are related as

$$\mathbf{D} = \mathbf{M}^{-1} \mathbf{S}. \tag{2.2}$$

We now review concepts related to integration preconditioning. In [15] the invertible pseudospectral first order differentiation matrix has the form

$$\mathbf{D}_+ = \mathbf{D} - \eta \mathbf{I}_+, \quad \eta = \frac{N(N+1)}{4} = \frac{1}{2\omega_0} = \frac{1}{2\omega_N},$$

where η is the penalty parameter. In the present study we adopt the formulation of \mathbf{D}_+ and define the \mathbf{D}_- and \mathbf{D}_+ operators, called the left-ended and right-ended BCP differentiation matrix, respectively, with their expressions given as

$$\mathbf{D}_\pm = \mathbf{D} \mp \eta \mathbf{I}_\pm, \quad \eta = \frac{N(N+1)}{2} = \frac{1}{\omega_0} = \frac{1}{\omega_N}.$$

With this slight change in the value of the penalty parameter, we follow the approach shown in [19] and find that the inverse matrices of \mathbf{D}_+ and \mathbf{D}_- have the following factored form

$$\mathbf{D}_\pm^{-1} = \mathbf{P} \mathbf{J}_\pm^{-1} \mathbf{P}^T \mathbf{M}, \quad \mathbf{P} = \begin{bmatrix} \mathbf{p}_0 & \mathbf{p}_1 & \cdots & \mathbf{p}_N \end{bmatrix},$$

$$\mathbf{J}_{\pm}^{-1} = \frac{1}{2} \begin{bmatrix} \mp 1 & -1 & 0 & \cdots & \cdots & 0 & 0 \\ 1 & 0 & -1 & \ddots & & & 0 \\ 0 & 1 & 0 & \ddots & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 & -1 & 0 \\ \vdots & & & \ddots & 1 & 0 & -1 \\ 0 & \cdots & \cdots & \cdots & 0 & 1 & \mp 1 \end{bmatrix}.$$

Let \mathbf{f} be the grid vector with components $f(x_i)$ for $i = 0, 1, \dots, N$. Then $\mathbf{D}_{\pm}^{-1}\mathbf{f}$ is an approximation solution to the differential equation $u'(x) = f(x)$ subject to the left boundary condition $u(-1) = 0$. Likewise, $\mathbf{D}_{\mp}^{-1}\mathbf{f}$ returns an approximation solution to the same differential equation subject to the right boundary condition $u(+1) = 0$. Thus, we call \mathbf{D}_{\pm}^{-1} the first order numerical solution operators.

Due to the simple matrix structure of \mathbf{J}_{\pm}^{-1} , it is worth introducing identities associated with \mathbf{D}_{\pm}^{-1} for later use. Notice that

$$\mathbf{P}^T \mathbf{e}_{\pm} = [P_0(\pm 1), P_1(\pm 1), \dots, P_N(\pm 1)]^T = [(\pm 1)^0, (\pm 1)^1, \dots, (\pm 1)^N]^T.$$

Then we obtain the following identities

$$(2.3) \quad \begin{aligned} \eta \mathbf{D}_{\pm}^{-1} \mathbf{e}_{\pm} &= \mp \mathbf{p}_0, & \eta \mathbf{D}_{\mp}^{-1} \mathbf{e}_{\pm} &= (\pm 1)^{N-1} \mathbf{p}_N, \\ \mathbf{e}_{\pm}^T \mathbf{D}_{\pm}^{-1} &= -(\pm 1)^{N-1} \mathbf{p}_N^T \mathbf{M}, & \mathbf{e}_{\pm}^T \mathbf{D}_{\mp}^{-1} &= \pm \mathbf{p}_0^T \mathbf{M} \end{aligned}$$

by direct computations.

We have summarized the main concepts related to the Legendre pseudospectral method for the present study. For further details of the pseudospectral methods we refer the readers to [1, 2, 18, 26].

In this study, we also need the SMW formula to conduct low rank updates. For the purpose of self-containedness, we summarize some relevant concepts.

The SMW formula relates the inverse of a matrix perturbed by a low rank matrix with the inverse of the original matrix. Assume that a matrix \mathbf{A} of order n -by- n is invertible. Now consider the sum of \mathbf{A} and a low rank matrix of the form \mathbf{UCV} , where \mathbf{U} , \mathbf{V} , and \mathbf{C} are of conformable sizes. Then the inverse of \mathbf{A} after being perturbed by the low rank matrix \mathbf{UCV} can be computed using the following formula

$$(2.4) \quad (\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{C}^{-1} - \mathbf{V} \mathbf{A}^{-1} \mathbf{U})^{-1} \mathbf{V} \mathbf{A}^{-1}$$

provided that $(\mathbf{C}^{-1} - \mathbf{V} \mathbf{A}^{-1} \mathbf{U})$ is invertible. Thus, when \mathbf{A}^{-1} is available, the formula provides a convenient approach to finding the inverse of the \mathbf{A} matrix after being perturbed

by the low rank matrix, especially if the dimension of the matrix \mathbf{C} is much smaller than the dimension of \mathbf{A} . A well known case is the following. If both $\mathbf{U} = \mathbf{u}$ and $\mathbf{V}^T = \mathbf{v}$ are column vectors, and \mathbf{C} is a nonzero scalar, say 1 without losing generality, then (2.4) becomes the Sherman-Morrison formula:

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}^T\mathbf{A}^{-1}}{1 - \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u}}.$$

For further details of the SMW formula, we refer the readers to [14] and the references therein.

2.2. BCP pseudospectral discretization of the $\frac{d}{dx}a(x)\frac{d}{dx} + c(x)$ operator

Let $u(x)$ be a function defined on the interval \mathbf{l} satisfying the problem:

$$(2.5) \quad \begin{aligned} \frac{d}{dx} \left(a(x) \frac{du(x)}{dx} \right) + c(x)u(x) &= f(x), & x \in \mathbf{l}, \quad a(x) > 0, \\ \alpha_{\pm}u(\pm 1) \pm \beta_{\pm}u'(\pm 1) &= g_{\pm}, & \alpha_{\pm}, \beta_{\pm} \geq 0, \quad \alpha_{\pm}^2 + \beta_{\pm}^2 \neq 0, \end{aligned}$$

where $a(x)$, $c(x)$ and $f(x)$ are given smooth functions, the parameters α_+ , α_- , β_+ and β_- are non-negative reals, and g_+ and g_- are real numbers.

To solve the problem we introduce the GLL nodes on \mathbf{l} and denote $a_j = a(x_j)$ and $f_j = f(x_j)$ as the pointwise function values at these points. To approximate $u(x)$, we collocate grid function values v_j for $j = 0, 1, \dots, N$ and seek a numerical solution $v(x)$ of the form $v(x) = \sum_{j=0}^N l_j(x)v_j$ to satisfy the collocation equations

$$(2.6) \quad \begin{aligned} f_i &= (\mathcal{I}_N[av'])' \Big|_i + c(x_i)v(x_i) \\ &- \left(\tau_+\eta^2 a_N \delta_{Ni} - \chi_+ \frac{a_N l'_i(+1)}{\omega_i} \right) (\alpha_+ v_N + \beta_+ v'_N - g_+) \\ &- \left(\tau_-\eta^2 a_0 \delta_{0i} + \chi_- \frac{a_0 l'_i(-1)}{\omega_i} \right) (\alpha_- v_0 - \beta_- v'_0 - g_-), \quad i = 0, 1, \dots, N, \end{aligned}$$

where τ_{\pm} and χ_{\pm} are called penalty parameters related by

$$(2.7) \quad 1 - \tau_- \beta_- \eta = \chi_- \alpha_-, \quad 1 - \tau_+ \beta_+ \eta = \chi_+ \alpha_+.$$

Depending on the enforced boundary conditions, the values of the parameters τ_- and τ_+ are given as follows:

- non-Neumann condition enforced at $x = \pm 1$:

$$\tau_{\pm} > \frac{1}{\alpha_{\pm} + \beta_{\pm} \eta}, \quad \chi_{\pm} = \frac{1 - \tau_{\pm} \beta_{\pm} \eta}{\alpha_{\pm}}.$$

- Neumann condition enforced at $x = \pm 1$:

$$\tau_{\pm} = \frac{1}{\eta}, \quad \chi_{\pm} = 0.$$

Notice that for the Neumann boundary condition case, (2.7) is full-filled for any value of χ_{\pm} . Thus, for simplicity we set $\chi_{\pm} = 0$.

To show the solvability of (2.6), we introduce the following vector and matrix notations

$$\begin{aligned} \mathbf{v} &= [v_0, v_1, \dots, v_N]^T, & \mathbf{f} &= [f_0, f_1, \dots, f_N]^T, \\ \mathbf{A} &= \text{diag}(a_0, a_1, \dots, a_N), & \mathbf{C} &= \text{diag}(c_0, c_1, \dots, c_N). \end{aligned}$$

In these matrix and vector notations, the scheme (2.6) takes the following form

$$(2.8) \quad \begin{aligned} (\mathbf{L} + \mathbf{C})\mathbf{v} &= \mathbf{f} - \mathbf{M}^{-1}(\tau_+ \mathbf{M}^{-1} - \chi_+ \mathbf{D}^T) \mathbf{A} \mathbf{e}_+ g_+ \\ &\quad - \mathbf{M}^{-1}(\tau_- \mathbf{M}^{-1} + \chi_- \mathbf{D}^T) \mathbf{A} \mathbf{e}_- g_- \end{aligned}$$

with

$$(2.9) \quad \begin{aligned} \mathbf{L} &= \mathbf{D} \mathbf{A} \mathbf{D} - \mathbf{M}^{-1}(\tau_+ \mathbf{M}^{-1} - \chi_+ \mathbf{D}^T) \mathbf{A} \mathbf{e}_+ (\alpha_+ \mathbf{e}_+^T + \beta_+ \mathbf{e}_+^T \mathbf{D}) \\ &\quad - \mathbf{M}^{-1}(\tau_- \mathbf{M}^{-1} + \chi_- \mathbf{D}^T) \mathbf{A} \mathbf{e}_- (\alpha_- \mathbf{e}_-^T - \beta_- \mathbf{e}_-^T \mathbf{D}). \end{aligned}$$

The matrix \mathbf{L} is called a BCP pseudospectral Laplace differentiation matrix, which is a discrete analogy of the Laplace differential operator $\frac{d}{dx}a(x)\frac{d}{dx}$ with the boundary conditions taken into account. The resultant matrix of $\mathbf{L} + \mathbf{C}$ is a BCP pseudospectral matrix representation of the Helmholtz operator $\frac{d}{dx}a(x)\frac{d}{dx} + c(x)$.

To seek the inverse matrix of targeted operator $\mathbf{L} + \mathbf{C}$ we first seek the inverse operator of \mathbf{L} . We start by showing that the operator \mathbf{L} is invertible. Multiplying matrix $-\mathbf{M}$ to (2.9) from the left, and employing (2.1) and (2.2) as well as (2.7), we have

$$-\mathbf{M}\mathbf{L} = \mathbf{S}^T \mathbf{A} \mathbf{M}^{-1} \mathbf{S} - \mathbf{R}_+ + \mathbf{R}_-,$$

where

$$(2.10) \quad \mathbf{R}_{\pm} = \chi_{\pm} \alpha_{\pm} (\mathbf{I}_{\pm} \mathbf{A} \mathbf{D}_{\pm} + (\mathbf{I}_{\pm} \mathbf{A} \mathbf{D}_{\pm})^T) \mp \tau_{\pm} \alpha_{\pm} \eta \mathbf{A} \mathbf{I}_{\pm} \mp \beta_{\pm} \chi_{\pm} \mathbf{D}^T (\mathbf{A} \mathbf{I}_{\pm}) \mathbf{D}.$$

Since the matrices, $\mathbf{S}^T \mathbf{A} \mathbf{M}^{-1} \mathbf{S}$, \mathbf{R}_- and \mathbf{R}_+ , are all symmetric, the matrix $-\mathbf{M}\mathbf{L}$ is also symmetric, indicating that $-\mathbf{M}\mathbf{L}$ is diagonalizable and all the eigenvalues of $-\mathbf{M}\mathbf{L}$ are real.

We now show that these eigenvalues are positive by considering the eigenvalue-eigenvector problem:

$$(2.11) \quad \lambda \mathbf{u} = -\mathbf{M}\mathbf{L}\mathbf{u} = (\mathbf{S}^T \mathbf{A} \mathbf{M}^{-1} \mathbf{S} - \mathbf{R}_+ + \mathbf{R}_-) \mathbf{u}, \quad \mathbf{u} = [u_0, u_1, \dots, u_N]^T,$$

where λ is an eigenvalue of $-\mathbf{ML}$ and \mathbf{u} is the associated eigenvector of unit length, that is, $\mathbf{u}^T \mathbf{u} = \sum_{i=0}^N u_i^2 = 1$.

Then multiplying \mathbf{u}^T to (2.11) from the left and invoking the relationship $\mathbf{MD} = \mathbf{S}$ (see (2.2)) and the expressions of \mathbf{R}_+ and \mathbf{R}_- in (2.10), we have

$$(2.12) \quad \lambda = \sum_{i=1}^{N-1} (u'_i)^2 a_i \omega_i + a_N \mathbf{r}_+^T \mathbf{W}_+ \mathbf{r}_+ + a_0 \mathbf{r}_-^T \mathbf{W}_- \mathbf{r}_- \geq 0,$$

where $u'_i = \sum_{j=0}^N l'_j(x_i) u_j$, $\mathbf{r}_+ = [u_N, -u'_N]^T$, and $\mathbf{r}_- = [u_0, u'_0]^T$, and the matrices \mathbf{W}_+ and \mathbf{W}_- are

$$\mathbf{W}_\pm = \begin{bmatrix} \tau_\pm \alpha_\pm \eta & \frac{1}{2}(1 - \beta_\pm \tau_\pm \eta + \chi_\pm \alpha_\pm) \\ \frac{1}{2}(1 - \beta_\pm \tau_\pm \eta + \chi_\pm \alpha_\pm) & \eta^{-1} - \beta_\pm \chi_\pm \end{bmatrix}.$$

Notice that \mathbf{W}_- or \mathbf{W}_+ is semi-positive definite if the enforced boundary condition at the associated end point is of the Neumann type, and is positive definite otherwise. However, if both boundary conditions are of the Neumann type, then the Poisson problem does not have a unique solution. Hence, we exclude the pure Neumann case in this part of the analysis and conclude that either \mathbf{W}_+ or \mathbf{W}_- is positive definite. We will discuss the pure Neumann case for the Helmholtz problem later. Thus, shown in (2.12) the non-negativity of λ follows from the facts that (1) the term $\sum_{i=1}^{N-1} a_i (u'_i)^2 \omega_i$ is nonnegative, and (2) both \mathbf{W}_+ and \mathbf{W}_- are semi-positive definite.

We now claim that λ is strictly positive. For an eigenvector \mathbf{u} , the values of u'_i for $1 \leq i \leq N-1$ are either (1) non-zero for some i , or (2) all zeros. For the former case, from (2.12) we have

$$\lambda \geq \sum_{i=1}^{N-1} a_i (u'_i)^2 \omega_i > 0.$$

For the latter one ($u'_i = 0$ for $i = 1, \dots, N-1$), we identify that u' is a polynomial of degree at most $N-1$, and vanishes at all the interior nodes, implying that $u(x) = \sum_{j=0}^N l_j(x) u_j$, the continuous representation of the eigenvector \mathbf{u} , has the form

$$u(x) = c_1 P_N(x) + c_2, \quad c_1^2 + c_2^2 \neq 0,$$

where c_1 and c_2 are constants. Consequently, \mathbf{r}_+ and \mathbf{r}_- are non-zero vectors. Then from (2.12) we have

$$\lambda = a_0 \mathbf{r}_-^T \mathbf{W}_- \mathbf{r}_- + a_N \mathbf{r}_+^T \mathbf{W}_+ \mathbf{r}_+ > 0$$

since either \mathbf{W}_- or \mathbf{W}_+ is positive definite.

We have shown that the matrix $-\mathbf{ML}$ is symmetric positive definite, and thus, invertible. As a result, the existence of \mathbf{L}^{-1} is established through the relationship that $\mathbf{L}^{-1} = (\mathbf{ML})^{-1} \mathbf{M}$.

2.3. Inverse matrix of the BCP pseudospectral Laplace operator

We now focus on explicitly formulating L^{-1} based on matrix factorizations. Applying the operators D_{\pm} , we rewrite L defined in (2.9) as

$$L = D_- A D_+ + Q,$$

where

$$(2.13) \quad \begin{aligned} Q = & (1 - \chi_+(\alpha_+ + \beta_+\eta))\eta a_N D_- I_+ + (\chi_+ - \tau_+)(\alpha_+ + \beta_+\eta)\eta^2 a_N I_+ \\ & + (\chi_+ - \tau_+)\beta_+\eta^2 a_N I_+ D_+ - \chi_+\beta_+\eta a_N D_- I_+ D_+ - \tau_-\alpha_-\eta^2 a_0 I_- \\ & - (1 - \tau_-\beta_-\eta)\eta a_0 I_- D_+ + \chi_-\alpha_-\eta a_0 D_- I_- - \chi_-\beta_-\eta a_0 D_- I_- D_+. \end{aligned}$$

Since D_- and D_+ are invertible, we express L in a factored form as

$$L = D_- (A + D_-^{-1} Q D_+^{-1}) D_+,$$

which is suggested by the equation itself to formulate L^{-1} as

$$L^{-1} = D_+^{-1} (A + D_-^{-1} Q D_+^{-1})^{-1} D_-^{-1}.$$

We now proceed to find the desired operator $(A + D_-^{-1} Q D_+^{-1})^{-1}$. To illustrate the construction of L^{-1} , we consider the problem subject to Dirichlet boundary conditions $(\alpha_{\pm}, \beta_{\pm}) = (1, 0)$, for the sake of clarity.

For the Dirichlet problem, the corresponding penalty parameters are $\chi_{\pm} = 1$ and we set $\tau_- = \tau_+ = \tau > 1$. Then the Q matrix is simplified as

$$Q = (1 - \tau)\eta^2 a_N I_+ - \tau\eta^2 a_0 I_- - \eta a_0 I_- D_+ + \eta a_0 D_- I_-.$$

We factor Q as follows:

$$Q = \widehat{\Phi} \widehat{\Psi}^T, \quad \widehat{\Phi} = \begin{bmatrix} D_- e_- & \eta e_- & \eta e_+ \end{bmatrix}, \quad \widehat{\Psi}^T = \begin{bmatrix} a_0 \eta e_-^T \\ -a_0 \tau \eta e_-^T - a_0 e_-^T D_+ \\ a_N (1 - \tau) \eta e_+^T \end{bmatrix}.$$

Notice that Q is a low-rank matrix and A is a diagonal one. Then, invoking the SMW formula we have

$$(2.14) \quad (A + D_-^{-1} Q D_+^{-1})^{-1} = (I - A^{-1} \Phi K^{-1} \Psi^T) A^{-1},$$

where

$$(2.15) \quad \Phi = D_-^{-1} \widehat{\Phi}, \quad \Psi^T = \widehat{\Psi}^T D_+^{-1}, \quad K = I_3 + \Psi^T A^{-1} \Phi$$

with \mathbf{I}_3 being the identity matrix of order three. As shown in (2.14), inverting the matrix $\mathbf{A} + \mathbf{D}_-^{-1} \mathbf{Q} \mathbf{D}_+^{-1}$ hinges upon finding $\Phi \mathbf{K}^{-1} \Psi^T$, which is treated next.

Employing the \mathbf{D}_-^{-1} and \mathbf{D}_+^{-1} operators and invoking the expressions shown in (2.3), we compute $\mathbf{D}_-^{-1} \hat{\Phi}$ and $\hat{\Psi}^T \mathbf{D}_+^{-1}$ and have

$$(2.16) \quad \Phi = \mathbf{D}_-^{-1} \hat{\Phi} = \begin{bmatrix} \mathbf{e}_- & \mathbf{p}_0 & \mathbf{p}_N \end{bmatrix},$$

and

$$\Psi^T = \hat{\Psi}^T \mathbf{D}_+^{-1} = \begin{bmatrix} -\eta a_0 \mathbf{p}_0^T \\ \tau_- \eta a_0 \mathbf{p}_0^T - \eta a_0 \mathbf{e}_-^T \\ -(1 - \tau) \eta a_N \mathbf{p}_N^T \end{bmatrix} \mathbf{M} = \Lambda \Phi^T \mathbf{M},$$

where Λ and its inverse are given as

$$\Lambda = \eta \begin{bmatrix} 0 & -a_0 & 0 \\ -a_0 & \tau a_0 & 0 \\ 0 & 0 & (\tau - 1) a_N \end{bmatrix}, \quad \Lambda^{-1} = \frac{1}{\eta} \begin{bmatrix} -\frac{\tau}{a_0} & -a_0^{-1} & 0 \\ -a_0^{-1} & 0 & 0 \\ 0 & 0 & \frac{1}{(\tau - 1) a_N} \end{bmatrix}.$$

Employing the expression of Ψ^T we have

$$\Phi \mathbf{K}^{-1} \Psi^T = \Phi \mathbf{K}^{-1} \Lambda \Phi^T \mathbf{M}.$$

Thus, to find the matrix $(\mathbf{A} + \mathbf{D}_-^{-1} \mathbf{Q} \mathbf{D}_+^{-1})^{-1}$ we are led to find the matrix $\mathbf{K}^{-1} \Lambda$. From the expression of \mathbf{K} in (2.15) we have the inverse of $\mathbf{K}^{-1} \Lambda$ as

$$\Lambda^{-1} \mathbf{K} = \Lambda^{-1} + \Phi^T \mathbf{M} \mathbf{A}^{-1} \Phi.$$

Employing (2.16) we compute $\Phi^T \mathbf{M} \mathbf{A}^{-1} \Phi$ and obtain

$$\Phi^T \mathbf{M} \mathbf{A}^{-1} \Phi = \frac{1}{\eta} \begin{bmatrix} a_0^{-1} & a_0^{-1} & (-1)^N a_0^{-1} \\ a_0^{-1} & \kappa_2 & \kappa_1 \\ (-1)^N a_0^{-1} & \kappa_1 & \kappa_0 \end{bmatrix}, \quad \kappa_\nu = \sum_{i=0}^N \frac{1}{a_i (P_N(\xi))^\nu}.$$

Then the matrix $\Lambda^{-1} \mathbf{K}$ is given as

$$(2.17) \quad \Lambda^{-1} \mathbf{K} = \frac{1}{\eta} \begin{bmatrix} (1 - \tau) \frac{1}{a_0} & 0 & \frac{(-1)^N}{a_0} \\ 0 & \kappa_2 & \kappa_1 \\ \frac{(-1)^N}{a_0} & \kappa_1 & \kappa_0 - \frac{1}{(1 - \tau) a_N} \end{bmatrix},$$

and its inverse is found as

$$\mathbf{K}^{-1} \Lambda = \frac{\eta}{T} \begin{bmatrix} t_{11} & t_{12} & t_{13} \\ t_{12} & t_{22} & t_{23} \\ t_{13} & t_{23} & t_{33} \end{bmatrix},$$

where T is given as

$$T = \frac{(1-\tau)^2}{a_0 a_N} (\kappa_1^2 - \kappa_2 \kappa_0) + \frac{(1-\tau)\kappa_2}{a_0 a_N} \left(\frac{1}{a_0} + \frac{1}{a_N} \right),$$

and t_{ij} for $i = 1, 2, 3$ and $i \leq j \leq 3$ are given as

$$\begin{aligned} t_{11} &= \frac{(1-\tau)}{a_N} (\kappa_1^2 - \kappa_0 \kappa_2) + \frac{\kappa_2}{a_N^2}, & t_{12} &= 0, \\ t_{13} &= \frac{(-1)^N (1-\tau)\kappa_2}{a_0 a_N}, & t_{22} &= \frac{(1-\tau)}{a_0 a_N} \left(\frac{a_0 + a_N}{a_0 a_N} - (1-\tau)\kappa_0 \right), \\ t_{23} &= \frac{(1-\tau)^2 \kappa_1}{a_0 a_N}, & t_{33} &= -\frac{(1-\tau)^2 \kappa_2}{a_0 a_N}. \end{aligned}$$

The expression of the operator $(\mathbf{A} + \mathbf{D}_-^{-1} \mathbf{Q} \mathbf{D}_+^{-1})^{-1}$ is given as

$$(\mathbf{A} + \mathbf{D}_-^{-1} \mathbf{Q} \mathbf{D}_+^{-1})^{-1} = (\mathbf{I} - \mathbf{A}^{-1} \mathbf{\Phi} (\mathbf{K}^{-1} \mathbf{\Lambda}) \mathbf{\Phi}^T \mathbf{M}) \mathbf{A}^{-1},$$

and thus, \mathbf{L}^{-1} is given as

$$\mathbf{L}^{-1} = \mathbf{D}_+^{-1} (\mathbf{I} - \mathbf{A}^{-1} \mathbf{\Phi} (\mathbf{K}^{-1} \mathbf{\Lambda}) \mathbf{\Phi}^T \mathbf{M}) \mathbf{A}^{-1} \mathbf{D}_-^{-1}.$$

We have shown the construction of \mathbf{L}^{-1} for problems subject to the Dirichlet boundary conditions. For problems subject to general boundary constraints, the corresponding vectors $\mathbf{\Phi}$ and $\mathbf{\Psi}^T$, and the matrix $\mathbf{K}^{-1} \mathbf{\Lambda}$ to formulate \mathbf{L}^{-1} are given in Appendix A.

Before proceeding further we address an issue related to examining the invertibility of the \mathbf{L} operator in the present inverse matrix construction framework. The \mathbf{L} operator is specifically constructed such that $-\mathbf{L}$ becomes symmetric positive definite under the mass matrix weighted discrete 2-norm, and thus invertible. In the construction of \mathbf{L}^{-1} one can also investigate the invertibility of the problem through examining the invertibility of the 3-by-3 matrix $\mathbf{\Lambda}^{-1} \mathbf{K}$ given in (2.17), which does not require knowing properties of the operator \mathbf{L} . Thus, for other pseudospectral differential operators which are not confined to the regular discrete 2-norm as the present \mathbf{L} operator is, one may still establish the invertibility of those operators by examining the corresponding compatible matrices.

2.4. Inverse matrix of the BCP pseudospectral Helmholtz operator

We now use the derived \mathbf{L}^{-1} matrix to construct the $(\mathbf{L} + \mathbf{C})^{-1}$ operator by applying the SMW formula and by exploiting the diagonal matrix structure of \mathbf{C} .

Let \mathbf{e}_i , for $i = 0, 1, \dots, N$, be the i -th column vectors of the identity matrix of order $N + 1$. To find the inverse of $\mathbf{L} + \mathbf{C}$ we express $\mathbf{L} + \mathbf{C}$ as

$$\mathbf{H} = \mathbf{L} + \mathbf{C} = \mathbf{L}_k + \sum_{i=k}^N c_i \mathbf{e}_i \mathbf{e}_i^T, \quad \mathbf{L}_k = \mathbf{L} + \sum_{i=0}^{k-1} c_i \mathbf{e}_i \mathbf{e}_i^T.$$

Then $\mathbf{L}_0 = \mathbf{L}$ and $\mathbf{L}_N = \mathbf{H}$, and furthermore, \mathbf{L}_{k+1} and \mathbf{L}_k are related as

$$\mathbf{L}_{k+1} = \mathbf{L}_k + c_k \mathbf{e}_k \mathbf{e}_k^T.$$

Assume that \mathbf{L}_k^{-1} exists. Then \mathbf{L}_{k+1}^{-1} is given as

$$(2.18) \quad \mathbf{L}_{k+1}^{-1} = \mathbf{L}_k^{-1} - (\mathbf{L}_k^{-1} \mathbf{e}_k) (c_k^{-1} + \mathbf{e}_k^T \mathbf{L}_k^{-1} \mathbf{e}_k)^{-1} (\mathbf{e}_k^T \mathbf{L}_k^{-1})$$

provided that the value $c_k^{-1} + \mathbf{e}_k^T \mathbf{L}_k^{-1} \mathbf{e}_k$ is nonzero. Since $\mathbf{L}_0^{-1} = \mathbf{L}^{-1}$ is available we can recursively apply (2.18) $N + 1$ times to construct \mathbf{H}^{-1} .

For $c(x) < 0$, \mathbf{H}^{-1} definitely exists, because $\mathbf{L} + \mathbf{C}$ is a negative definite operator weighted by the mass matrix \mathbf{M} . However, for a general $c(x)$ we are unable to analytically draw out any result regarding the inverse of $(\mathbf{L} + \mathbf{C})$. In this situation, we may try constructing \mathbf{H}^{-1} by the updating procedure (2.18). If the construction is successfully executed, then one may try to use the resultant \mathbf{H}^{-1} operator to compute numerical solutions. If the updating procedure fails to construct the \mathbf{H}^{-1} operator, for example, if at any k step the value $c_k^{-1} + \mathbf{e}_k^T \mathbf{L}_k^{-1} \mathbf{e}_k$ vanishes, it does not mean that the \mathbf{H} operator is not invertible. For this situation, one would have to use other approaches to seek the \mathbf{H}^{-1} operator, if it exists.

We now discuss the construction of the Helmholtz operator when the boundary conditions at the end points are both of the Neumann type. When we previously discussed the construction of the \mathbf{L} matrix we excluded the pure Neumann case, because the solution to the corresponding Poisson problem is not unique. However, the solution to the Helmholtz equation subject to the Neumann boundary condition at both end points is unique. In what follows we derive the corresponding \mathbf{H}^{-1} operator.

The Helmholtz operator with pure Neumann boundary conditions can be expressed as

$$\mathbf{H} = \mathbf{L} + \mathbf{C} = \mathbf{L} - \varepsilon \mathbf{I}_+ + \mathbf{C} + \varepsilon \mathbf{I}_+,$$

where, deduced from (2.9) with $(\alpha_{\pm}, \beta) = (0, 1)$, the matrix \mathbf{L} given as

$$\mathbf{L} = \mathbf{DAD} - \eta^{-1} \mathbf{M}^{-2} \mathbf{Ae}_+ \mathbf{e}_+^T \mathbf{D} + \eta^{-1} \mathbf{M}^{-2} \mathbf{Ae}_- \mathbf{e}_-^T \mathbf{D},$$

is the pure Neumann type BCP-pseudosectral Laplace operator, and ε is an introduced parameter with its value to be determined such that the resultant matrix $\mathbf{L} - \varepsilon \mathbf{I}_+$ is invertible.

To seek the inverse of the matrix operator $\mathbf{L} - \varepsilon \mathbf{I}_+$ we rewrite the matrix as

$$\mathbf{L} = \mathbf{DAD} - \eta^{-1} \mathbf{M}^{-2} \mathbf{Ae}_+ \mathbf{e}_+^T \mathbf{D} + \eta^{-1} \mathbf{M}^{-2} \mathbf{Ae}_- \mathbf{e}_-^T \mathbf{D} - \varepsilon \mathbf{I}_+ = \mathbf{D}_- \mathbf{AD}_+ + \mathbf{Q}$$

with \mathbf{Q} factored as

$$\mathbf{Q} = \widehat{\mathbf{\Phi}} \widehat{\mathbf{\Psi}}^T, \quad \widehat{\mathbf{\Phi}} = [\eta \mathbf{e}_+ \quad \mathbf{D}_- \mathbf{e}_+], \quad \widehat{\mathbf{\Psi}}^T = \begin{bmatrix} -(\frac{\varepsilon}{\eta} + a_N \eta) \mathbf{e}_+^T - a_N \mathbf{e}_+^T \mathbf{D}_+ \\ \eta a_N \mathbf{e}_+ \end{bmatrix}.$$

Computing $D_-^{-1}\widehat{\Phi}$ and $\widehat{\Psi}^T D_+^{-1}$, we obtain

$$\Phi = D_-^{-1}\widehat{\Phi} = \begin{bmatrix} p_N & e_+ \end{bmatrix}, \quad \Psi^T = \widehat{\Psi}^T D_+^{-1} = \begin{bmatrix} (a_N\eta + \frac{\varepsilon}{\eta})\mathbf{p}_N^T - a_N\eta e_+^T \\ -\eta a_N \mathbf{p}_N^T \end{bmatrix} M,$$

and Ψ^T and Φ are related by

$$\Psi^T = \Lambda \Phi^T M, \quad \Lambda = \begin{bmatrix} \eta a_N + \frac{\varepsilon}{\eta} & -\eta a_N \\ -\eta a_N & 0 \end{bmatrix}.$$

We then compute Λ^{-1} and $\Phi^T M A^{-1} \Phi$ and have

$$\Lambda^{-1} = \frac{-1}{\eta a_N} \begin{bmatrix} 0 & 1 \\ 1 & 1 + \frac{\varepsilon}{\eta^2 a_N} \end{bmatrix}, \quad \Phi^T M A^{-1} \Phi = \frac{1}{\eta a_N} \begin{bmatrix} \kappa_0 a_N & 1 \\ 1 & 1 \end{bmatrix}.$$

Then we have

$$\Lambda^{-1} K = \Lambda^{-1} + \Phi^T M A^{-1} \Phi = \frac{1}{\eta a_N} \begin{bmatrix} \kappa_0 a_N & 0 \\ 0 & -\frac{\varepsilon}{\eta^2 a_N} \end{bmatrix},$$

leading to

$$K^{-1} \Lambda = \text{diag} \left(\frac{\eta}{\kappa_0}, -\frac{\eta^3 a_N^2}{\varepsilon} \right).$$

The inverse of $L - \varepsilon I_+$ is given as

$$(L - \varepsilon I_+)^{-1} = D_+^{-1} (I - A^{-1} \Phi K^{-1} \Lambda \Phi^T M) A^{-1} D_-^{-1}.$$

We have the formula to construct $(L - \varepsilon I_+)^{-1}$, and we can then apply (2.18) to derive the matrix H^{-1} .

2.5. Inhomogeneous boundary conditions

The H^{-1} operator can now be used to compute numerical solutions to the second order problems subject to homogeneous boundary conditions, that is, $g_- = g_+ = 0$, by performing $\mathbf{v} = \mathbf{v}_h = H^{-1} \mathbf{f}$. If inhomogeneous boundary conditions are imposed, then the numerical solution \mathbf{v} to the second order problem is the sum of the homogeneous solution \mathbf{v}_h and the particular solution \mathbf{v}_p satisfying the problem

$$(2.19) \quad \begin{aligned} H \mathbf{v}_p &= -a_0 (\tau_- \eta^2 \mathbf{e}_- + \chi_- M^{-1} D^T \mathbf{e}_-) g_- \\ &\quad - a_N (\tau_+ \eta^2 \mathbf{e}_+ - \chi_+ M^{-1} D^T \mathbf{e}_+) g_+. \end{aligned}$$

To seek the particular solution \mathbf{v}_p we assume that \mathbf{v}_p is of the form

$$\mathbf{v}_p = H^{-1} (-z_1 D \mathbf{a} - z_1 C \mathbf{p}_1 - z_0 C \mathbf{p}_0) + z_1 \mathbf{p}_1 + z_0 \mathbf{p}_0,$$

where z_0 and z_1 are constants to be determined. Performing $\mathbf{H}\mathbf{v}_p$ we obtain

$$\begin{aligned}\mathbf{H}\mathbf{v}_p &= \mathbf{H}\mathbf{H}^{-1}(-z_1\mathbf{D}\mathbf{a} - z_1\mathbf{C}\mathbf{p}_1 - z_0\mathbf{C}\mathbf{p}_0) + \mathbf{L}(z_1\mathbf{p}_1 + z_0\mathbf{p}_0) + \mathbf{C}(z_1\mathbf{p}_1 + z_0\mathbf{p}_0) \\ &= -z_1\mathbf{D}\mathbf{a} - z_1\mathbf{C}\mathbf{p}_1 - z_0\mathbf{C}\mathbf{p}_0 + \mathbf{D}\mathbf{A}\mathbf{D}(z_1\mathbf{p}_1 + z_0\mathbf{p}_0) + z_1\mathbf{C}\mathbf{p}_1 + z_0\mathbf{C}\mathbf{p}_0 \\ &\quad - a_N \left(\frac{\tau_+}{\omega_N^2} \mathbf{e}_+ - \chi_+ \mathbf{M}^{-1} \mathbf{D}^T \mathbf{e}_+ \right) (\alpha_+ \mathbf{e}_+^T + \beta_+ \mathbf{e}_+^T \mathbf{D}) (z_1 \mathbf{p}_1 + z_0 \mathbf{p}_0) \\ &\quad - a_0 \left(\frac{\tau_-}{\omega_0^2} \mathbf{e}_- + \chi_- \mathbf{M}^{-1} \mathbf{D}^T \mathbf{e}_- \right) (\alpha_- \mathbf{e}_-^T - \beta_- \mathbf{e}_-^T \mathbf{D}) (z_1 \mathbf{p}_1 + z_0 \mathbf{p}_0).\end{aligned}$$

Applying the following relationships

$$\begin{aligned}\mathbf{D}\mathbf{p}_1 &= \mathbf{p}_0, \quad \mathbf{D}\mathbf{p}_0 = \mathbf{0}, \quad \mathbf{A}\mathbf{D}\mathbf{p}_1 = \mathbf{A}\mathbf{p}_0 = \mathbf{a}, \\ (\alpha_\pm \mathbf{e}_\pm^T \pm \beta_\pm \mathbf{e}_\pm^T \mathbf{D})(z_1 \mathbf{p}_1 + z_0 \mathbf{p}_0) &= \alpha_\pm (\pm z_1 + z_0) \pm \beta_\pm z_1,\end{aligned}$$

we arrive at

$$\begin{aligned}\mathbf{H}\mathbf{v}_p &= -a_N (\tau_+ \eta^2 \mathbf{e}_+ - \chi_+ \mathbf{M}^{-1} \mathbf{D}^T \mathbf{e}_+) (\alpha_+ (z_1 + z_0) + \beta_+ z_1) \\ &\quad - a_0 (\tau_- \eta^2 \mathbf{e}_- + \chi_- \mathbf{M}^{-1} \mathbf{D}^T \mathbf{e}_-) (\alpha_- (-z_1 + z_0) - \beta_- z_1).\end{aligned}$$

By matching the coefficients on the right-hand side of (2.19), we have

$$(\alpha_+ + \beta_+) z_1 + \alpha_+ z_0 = g_+, \quad -(\alpha_- + \beta_-) z_1 + \alpha_- z_0 = g_-.$$

Solving the equations we obtain

$$z_0 = \frac{(\alpha_- + \beta_-)g_+ + (\alpha_+ + \beta_+)g_-}{(\alpha_- + \beta_-)\alpha_+ + (\alpha_+ + \beta_+)\alpha_-}, \quad z_1 = \frac{\alpha_- g_+ - \alpha_+ g_-}{(\alpha_- + \beta_-)\alpha_+ + (\alpha_+ + \beta_+)\alpha_-},$$

and \mathbf{v}_p is determined.

If the inhomogeneous boundary conditions at the two end points are both of the Neumann type, then we have the inhomogeneous problem as follows:

$$\mathbf{H}\mathbf{v}_p = -a_0(\eta\mathbf{e}_-)g_- - a_N(\eta\mathbf{e}_+)g_+,$$

implying that

$$\mathbf{v}_p = -\eta a_0 g_- \mathbf{H}^{-1} \mathbf{e}_- - \eta a_N g_+ \mathbf{H}^{-1} \mathbf{e}_+.$$

Notice that the vectors $\mathbf{H}^{-1}\mathbf{e}_\pm$ are available once \mathbf{H}^{-1} is constructed, since they are the first and the last column vectors of \mathbf{H}^{-1} .

Finally, to numerically solve (2.8), we propose the computing steps after initializing the required variables.

Step 1:

$$(2.20) \quad \mathbf{f}^* = \mathbf{f} - \begin{cases} z_1(\mathbf{D}\mathbf{a} + \mathbf{C}\mathbf{p}_1) + z_0\mathbf{C}\mathbf{p}_0 & \text{general case,} \\ \mathbf{0} & \text{pure Neumann case.} \end{cases}$$

Step 2:

$$(2.21) \quad \mathbf{v} = \mathbf{H}^{-1} \mathbf{f}^* + \begin{cases} z_1 \mathbf{p}_1 + z_0 \mathbf{p}_0 & \text{general case,} \\ -\eta a_0 g_- \mathbf{H}^{-1} \mathbf{e}_- - \eta a_N g_+ \mathbf{H}^{-1} \mathbf{e}_+ & \text{pure Neumann case.} \end{cases}$$

2.6. Solving second order differential equations

The constructed \mathbf{H}^{-1} operator can be used as a numerical solution operator for the second order boundary value problem of the form

$$(2.22) \quad a^*(x)u''(x) + b^*(x)u'(x) + c^*(x)u(x) = f^*(x), \quad x \in I, \quad a^*(x) \neq 0, \\ \alpha_{\pm}u(\pm 1) \pm \beta_{\pm}u'(\pm) = g_{\pm}.$$

Notice that (2.22) can be transformed into the form

$$(2.23) \quad (\theta(x)u'(x))' + \theta(x)\frac{c^*(x)}{a^*(x)}u(x) = \theta(x)\frac{f^*(x)}{a^*(x)}, \quad \theta(x) = \exp\left(\int \frac{b^*(x)}{a^*(x)} dx\right),$$

by introducing the integrating factor $\theta(x)$ playing the role of $a(x)$ in (2.5). If θ is not available analytically, we approximate θ numerically as

$$\boldsymbol{\theta} = [\theta_0, \dots, \theta_N] = [e^{\mu_0}, \dots, e^{\mu_N}]^T, \quad [\mu_0, \dots, \mu_N]^T = \mathbf{D}_-^{-1} \left[\frac{b_0^*}{a_0^*}, \dots, \frac{b_N^*}{a_N^*} \right]^T.$$

Employing \mathbf{L} with the diagonal elements of \mathbf{A} by θ_i and denoting

$$\mathbf{C} = \text{diag} \left(\frac{\theta_0 c_0^*}{a_0^*}, \frac{\theta_1 c_1^*}{a_1^*}, \dots, \frac{\theta_N c_N^*}{a_N^*} \right), \quad \mathbf{f} = \left[\theta_0 \frac{f_0^*}{a_0^*}, \theta_1 \frac{f_1^*}{a_1^*}, \dots, \theta_N \frac{f_N^*}{a_N^*} \right]^T,$$

we can apply (2.20) and (2.21) to solve (2.23).

3. Numerical validations and discussions

In this section, we present results obtained by the proposed method. In each numerical experiment, an exact solution $u(x)$, a coefficient function $a(x)$, and a set of parameters $(\alpha_{\pm}, \beta_{\pm})$, are specified. With these pieces of information in hand we compute the corresponding $f(x)$ and g_{\pm} . To examine the performance of the derived inverse operator when solving a problem, we measure the discrete l_2 error, defined as $e_2(N) = (\sum_{i=0}^N |u(x_i) - v(x_i)|^2 \omega_i)^{1/2}$, and the maximum pointwise error, defined as $e_{\infty}(N) = \max_{i=0, \dots, N} |u(x_i) - v(x_i)|$, where $v(x)$ is a numerical solution obtained by the scheme on a grid mesh characterized by N which is the degree of the approximation polynomial $v(x)$.

Our first numerical experiment is solving the boundary value problem

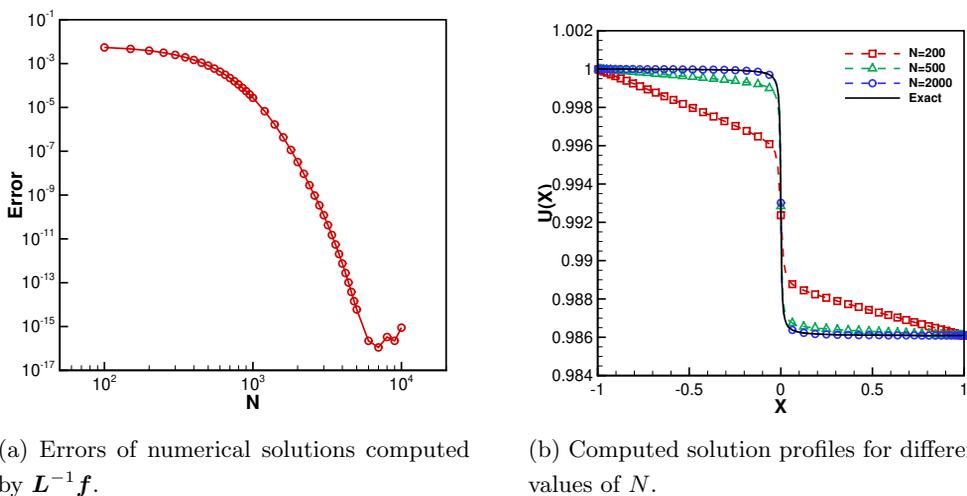
$$(3.1) \quad u''(x) + \frac{1}{ax^2 + 1}u'(x) = f(x), \quad x \in (-1, 1), \quad u(\pm 1) = g_{\pm}$$

with an exact solution $u(x)$ given as

$$u(x) = \exp\left(-\frac{\tan^{-1}(\sqrt{a}x) + \tan^{-1}(\sqrt{a})}{\sqrt{a}}\right), \quad a = 5 \times 10^4.$$

The function $f(x)$ and the values of g_- and g_+ are computed by $u(x)$.

A convergence study for this problem and computed solution profiles are presented in Figure 3.1. We observe that the error decays exponentially as the grid resolution increases. Moreover, the accuracy of the numerical solutions remains at a machine accuracy level during further grid refinements. As pointed out in [20], the function $u(x)$ forms a thin boundary layer near $x = 0$, and a dense grid is needed to resolve $u(x)$ accurately. Similarly, we need to use a dense grid mesh to resolve u to a machine accuracy level, as shown in Figure 3.1(a) which agrees with the result in [20].



(a) Errors of numerical solutions computed by $L^{-1}\mathbf{f}$.

(b) Computed solution profiles for different values of N .

Figure 3.1: Convergence results of (3.1).

The next numerical experiment is solving the problem

$$(3.2) \quad u''(x) - b^2u(x) = -e^{ax}, \quad x \in (-1, 1), \quad u(\pm 1) = 0$$

for the parameter sets $(a, b) = (1, 0)$ and $(a, b) = (2, 1)$ with an exact smooth solution

$$u(x) = \frac{e^{ax} + e^{-bx} \sinh(a - b) \operatorname{csch}(2b) - e^{bx} \sinh(a + b) \operatorname{csch}(2b)}{b^2 - a^2}.$$

Notice that for the case $(a, b) = (1, 0)$ the corresponding exact solution is obtained through taking the limit of the general expression as $b \rightarrow 0$, leading to $u(x) = -e^x$.

N	Present		JG1	JG2
	$e_2(N)$	$e_\infty(N)$	$e_\infty(N)$	$e_\infty(N)$
4	1.42E-04	1.34E-04	2.45E-03	7.96E-04
8	5.81E-10	5.34E-10	3.93E-08	1.57E-08
16	8.66E-16	8.88E-16	7.64E-16	7.13E-16
32	2.74E-16	4.44E-16	-	-
64	6.22E-16	8.88E-16	-	-
128	1.45E-15	1.55E-15	-	-
256	3.05E-16	8.88E-16	-	-
512	6.91E-16	1.11E-15	-	-
1024	8.11E-16	1.33E-15	-	-

Table 3.1: Convergence results for the Poisson problem formulated by (3.2) with $a = 1$ and $b = 0$. Reference results, labeled JG1 and JG2, are collected from the Jacobi polynomial based spectral-Galerkin methods [7]. (see Table 1 in [7] with the parameters $\alpha = \beta = 0$).

N	Present		JG1	JG2
	$e_2(N)$	$e_\infty(N)$	$e_\infty(N)$	$e_\infty(N)$
4	3.05E-03	3.43E-03	2.63E-02	8.45E-03
8	2.05E-07	1.95E-07	7.14E-06	2.81E-06
16	4.06E-16	4.44E-16	6.41E-15	2.93E-15
32	3.83E-16	6.10E-16	-	-
64	6.27E-16	9.43E-16	-	-
128	7.86E-16	1.11E-15	-	-
256	3.66E-16	8.32E-16	-	-
512	4.65E-16	9.99E-16	-	-
1024	6.73E-16	1.55E-15	-	-

Table 3.2: Convergence results for the Helmholtz problem formulated by (3.2) with $a = 2$ and $b = 1$. Reference results, labeled JG1 and JG2, are collected from the Jacobi polynomial based spectral-Galerkin methods [7]. (see Table 1 in [7] with the parameters $\alpha = \beta = 0$).

Tables 3.1 and 3.2 present convergence studies of our method, with results from the Jacobi polynomial based spectral-Galerkin method [7] for comparison. For the grid resolutions $N = 4$ and $N = 8$, our method is clearly better than the Jacobi polynomial based spectral-Galerkin method [7], for both parameter cases $(a, b) = (1, 0)$ and $(a, b) = (2, 1)$. As the grid resolution increases to $N = 16$, all the numerical solutions have arrived at the machine error accuracy level. As the grid solution N further increases, the accuracy of each numerical solution computed by the present method remains at the machine error level.

The third numerical experiment is solving the problem

$$(3.3) \quad u''(x) + \cos(3x)u'(x) + \sin(3x)u(x) = f(x), \quad x \in (-1, 1), \quad u(\pm 1) = \sin(\pm 6)$$

with an exact solution given as $u(x) = \sin(6x)$. This problem was solved by the arbitrary grid based integration preconditioning matrix method [13]. A convergence study is presented in Table 3.3, with reference results provided in [13] for comparison.

N	Present		GMH-PS	H-PS	PS
	$e_2(N)$	$e_\infty(N)$	$e_\infty(N)$	$e_\infty(N)$	$e_\infty(N)$
8	1.19E-01	1.34E-01	1.70E+01	1.10E-01	2.10E+01
16	6.46E-05	7.87E-05	3.80E-02	1.60E-06	1.50E-01
32	9.70E-12	9.67E-12	5.40E-11	1.60E-05	1.00E-05
64	2.86E-15	4.38E-15	8.90E-12	No Conv.	No Conv.
128	2.65E-15	4.44E-15	-	-	-
256	2.35E-15	5.21E-15	-	-	-
512	7.10E-15	1.04E-14	-	-	-
1024	6.08E-15	1.26E-14	-	-	-

Table 3.3: Convergence results for the Helmholtz problem formulated by (3.3). Reference results, labeled GMH-PS, H-PS, and PS, were obtained by the arbitrary grid based pseudospectral integration precondition matrix [13], the pseudospectral integration preconditioning matrix method [15], and the traditional pseudospectral method, respectively. All these reference results were reported in [13, Table 5]. In [13] the problem was solved by iterative methods, and those computations which failed to converge were labeled as No Conv..

It is clearly seen from the computed errors that the present method is better, except for the case $N = 16$ by the method [15]. Furthermore, the accuracy of the numerical

solutions obtained by the present method is preserved and not ruined by round off errors, even using a dense grid mesh.

The fourth numerical experiment is solving the Helmholtz problem

$$(3.4) \quad \varepsilon u''(x) - xu'(x) - u(x) = 0, \quad x \in (-1, 1), \quad u(\pm 1) = 0$$

for $\varepsilon = 0.1$. This problem was solved by the ultraspherical polynomial based integration preconditioning matrix method [10], with an exact solution given as $u(x) = e^{\frac{x^2-1}{2\varepsilon}}$ which forms a boundary layer of width $\mathcal{O}(1/\varepsilon)$ in the vicinity of each end point.

Table 3.4 presents a convergence study of our method, with reference results provided in [10] for comparison. As shown in the results, for all the grid resolutions used in the computations, our method is as good as the ultraspherical polynomial based integration preconditioning matrix method [10], and both integration type methods are better than the traditional pseudospectral method.

N	Present		PLP	LP
	$e_2(N)$	$e_\infty(N)$	$e_\infty(N)$	$e_\infty(N)$
8	7.59E-03	7.63E-03	-	-
16	8.42E-06	7.81E-06	-	-
32	6.35E-14	5.76E-14	-	-
64	8.03E-16	1.99E-15	9.99E-16	1.59E-13
128	4.71E-16	9.99E-16	1.22E-15	9.82E-14
256	4.43E-16	1.66E-15	1.67E-15	6.61E-13
512	8.39E-16	2.77E-15	2.22E-15	1.11E-11
1024	7.23E-16	2.66E-15	3.71E-15	2.72E-11

Table 3.4: Convergence results for the Helmholtz problem formulated by (3.4) with $\varepsilon = 0.1$. Reference results, labeled PLP and LP, are collected from the ultraspherical polynomial based integration preconditioning matrix methods [10, Table 3].

Our next numerical experiment is solving the boundary value problem

$$(3.5) \quad u''(x) - (1 + \sin x)u'(x) + e^x u(x) = f(x), \quad x \in \mathbb{I}, \quad u(\pm 1) = 1$$

with an exact solution given as $u(x) = e^{(x^2-1)/2}$. The corresponding f is computed by the given u . We use (2.20) and (2.21) to solve the discretized problem.

A convergence study is presented in Table 3.5. Similar to the previous examples, the numerical error vanishes exponentially as the grid resolution increases. Furthermore, the

measured errors are as small as those reported in [30], indicating that the present method is effective and comparable to the approach [30], even during further grid refinements.

N	Present		BCOL	P-LCOL
	$e_2(N)$	$e_\infty(N)$	$e_\infty(N)$	$e_\infty(N)$
4	1.84e-03	1.75e-03	-	-
8	5.49e-07	7.93e-07	-	-
16	1.94e-13	1.98e-13	-	-
32	8.70e-16	8.88e-16	-	-
64	8.10e-16	1.11e-15	5.55e-16	1.67e-15
128	1.02e-15	1.22e-15	1.11e-15	2.44e-15
256	1.79e-15	1.88e-15	1.11e-15	2.55e-15
512	6.95e-15	8.10e-15	1.89e-15	4.77e-15
1024	8.37e-15	6.43e-15	3.44e-15	1.15e-14

Table 3.5: Convergence results for (3.5). Reference results are from [30] by the Birkhoff collocation (BCOL) scheme and the preconditioned Lagrange collocation (P-LCOL) scheme.

In theory, the accuracy of a numerical solution is affected by the truncation errors and the round off errors. In these experiments the results show that the present approach is better than the other methods in most cases, when round off errors are much smaller than the truncation errors of the numerical solutions. For computations based on dense grid meshes and when round off errors and the truncation errors of the numerical solutions are of the same order of magnitude, our method is as good as the other methods.

Encouraged by the performances of our method in the previous experiments for solving linear differential equations, we further test the performance of the present method for solving a nonlinear differential equation.

The following experiment is solving the problem described by the charge conserving Poisson-Boltzmann equation of the form

$$u''(\xi) = \frac{\sinh(u(\xi))}{d^{-1} \int_{-d}^d e^{u(\xi)} d\xi}, \quad \xi \in \left(-\frac{d}{2}, \frac{d}{2}\right), \quad u\left(\pm\frac{d}{2}\right) = \pm V.$$

Introducing the linear coordinate transformation, $\xi(x) = xd/2$, and denoting the transformed variable $u(\xi(x))$ by $v(x)$, we have the transformed problem as follows:

$$(3.6) \quad \frac{d^2 v(x)}{dx^2} = \frac{d^3 \sinh(v(x))}{8 \int_{-1}^1 e^{v(x)} dx}, \quad x \in (-1, 1), \quad v(\pm 1) = \pm V.$$

Let \mathbf{v} be the grid vector of $v(x)$ and \mathbf{f} be the grid vector defined as

$$\mathbf{f} = \begin{bmatrix} f_0 & f_1 & \cdots & f_N \end{bmatrix}, \quad f_i = \frac{\sinh(v(x_i))}{\sum_{k=0}^N e^{v(x_i)} \omega_i}.$$

We solve (3.6) by the following scheme

$$\mathbf{v} = \mathbf{L}^{-1} \left[\frac{d^3}{8} \mathbf{f} - (\tau_- \eta^2 \mathbf{e}_- + \chi_- \mathbf{M}^{-1} \mathbf{D}^T \mathbf{e}_-) V_- - (\tau_+ \eta^2 \mathbf{e}_+ - \chi_+ \mathbf{M}^{-1} \mathbf{D}^T \mathbf{e}_+) V_+ \right]$$

by a nonlinear equation solver.

Figure 3.2 presents computed solutions for various values of N . For certain particular parameter ranges a good approximation of the Dirichlet problem has been derived [9], and we use the derived formula to compute an approximation solution for comparison. In Figure 3.2 it is shown that the difference between the numerical solution and the approximation solution becomes smaller as the grid resolution N increases.

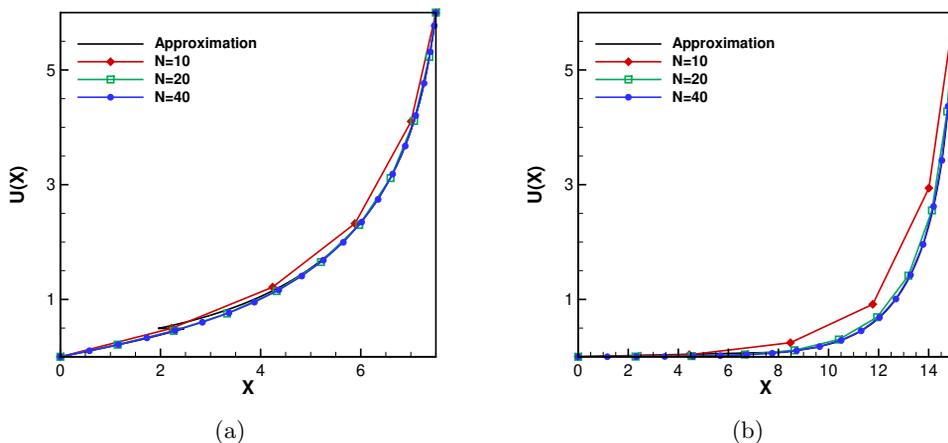


Figure 3.2: Numerical solutions of (3.6) with problem parameters given by $V = 6$ and $d = 15$ (a) and $d = 30$ (b). Approximation solutions are computed by the formula given in [9].

Note that the plotted approximation solution has a zigzag for $2 \leq x \leq 2.5$ in Figure 3.2(a). It is caused by the approximation formula provided in [9] (Eq. (3.15b) in the reference), which is for approximating $x(u)$ (the inverse of u) instead of $u(x)$. The formula was derived by singular perturbation theory, with range splitting and asymptotic matching. Thus, depending on the value of u , different approximation formula were applied. As a result, the plotted curve $x(u)$ (black line) has a jump discontinuity at the matching point in u (approximately at $u = 0.5$), forming a zigzag shape in the plot in this case. For further information about the property of the approximation solution $x(u)$ we refer the readers to [9].

We close this section by demonstrating the performance of the derived \mathbf{H}^{-1} operator when solving a time-dependent partial differential equation. Consider $u(x, t)$ satisfying the model problem

$$(3.7) \quad \begin{cases} \frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + f(x, t), & x \in (-1, 1), t > 0, \\ u(x, 0) = u_0(x), & x \in [-1, 1], \\ \pm \frac{\partial u(\pm 1, t)}{\partial x} = g_{\pm}(t), & t \geq 0 \end{cases}$$

with an exact solution given as

$$u(x, t) = e^{-\pi^2 t} \sin(\pi x) + e^{-4\left(x - \frac{\pi}{12}\right)^2} \cos(6\pi x).$$

The functions $f(x, t)$, $u_0(x)$, and $g_{\pm}(t)$ are evaluated by the given exact solution.

We discretize the problem by the present pseudospectral method in space and the Crank-Nicholson difference method in time, and have the scheme as

$$(3.8) \quad \frac{\mathbf{v}^{n+1} - \mathbf{v}^n}{\Delta t} = \mathbf{L} \left(\frac{\mathbf{v}^{n+1} + \mathbf{v}^n}{2} \right) + \frac{\mathbf{f}^{n+1} + \mathbf{f}^n}{2} + \eta \mathbf{e}_- \frac{g_-^{n+1} + g_-^n}{2} + \eta \mathbf{e}_+ \frac{g_+^{n+1} + g_+^n}{2},$$

$$\mathbf{v}^0 = [u_0(x_0), u_0(x_1), \dots, u_0(x_N)]^T,$$

where Δt is the time step, the superscript n denotes the integer time level, $\mathbf{v}^n = [v_0^n, \dots, v_N^n]^T$ is the solution grid vector at discrete time $t_n = n\Delta t$ with v_i^n being the value approximating $u(x_j, t_n)$, $\mathbf{f}^n = [f_0^n, \dots, f_N^n]^T$ is the grid vector of f at time level t_n , and \mathbf{L} is the BCP pseudospectral Laplace operator defined in (2.9).

The fully discrete scheme (3.8) is stable independent of the size of the time step. To show this fact it is sufficient to consider the homogeneous version of the scheme. Multiplying $(\mathbf{v}^{n+1} + \mathbf{v}^n)\mathbf{M}$ to the scheme from the left and invoking that \mathbf{M} is symmetric positive definite and $\mathbf{M}\mathbf{L}$ is semi-negative definite, we have

$$(\mathbf{v}^{n+1})^T \mathbf{M} \mathbf{v}^{n+1} = (\mathbf{v}^n)^T \mathbf{M} \mathbf{v}^n + \frac{\Delta t}{2} (\mathbf{v}^{n+1} + \mathbf{v}^n)^T \mathbf{M} \mathbf{L} (\mathbf{v}^{n+1} + \mathbf{v}^n) \leq (\mathbf{v}^n)^T \mathbf{M} \mathbf{v}^n.$$

As a result, the discrete energy norm of the numerical solution at time level t_n is bounded as follows:

$$(\mathbf{v}^n)^T \mathbf{M} \mathbf{v}^n \leq (\mathbf{v}^{n-1})^T \mathbf{M} \mathbf{v}^{n-1} \leq \dots \leq (\mathbf{v}^0)^T \mathbf{M} \mathbf{v}^0.$$

Thus, for a fixed terminal time $T = n\Delta t$ the numerical solution \mathbf{v}^n is bounded by the initial data and is independent of the time step, during grid refinements.

For the sake of computation accuracy, we still need to use a proper time step size to advance the solution in time, although the scheme is unconditional stable. For this numerical experiment we use $\Delta t = 1/N$ for computations. As we will see soon, this choice of time step gives satisfactory results.

To compute \mathbf{v}^{n+1} we rewrite the system of equations as

$$(3.9) \quad \begin{aligned} \mathbf{v}^{n+1} = & -\mathbf{v}^n - \mathbf{H}^{-1} \left(\frac{4}{\Delta t} \mathbf{v}^n + \mathbf{f}^{n+1} + \mathbf{f}^n \right) \\ & - \mathbf{H}^{-1} \mathbf{e}_{-\eta} (g_-^{n+1} + g_-^n) - \mathbf{H}^{-1} \mathbf{e}_{+\eta} (g_+^{n+1} + g_+^n), \end{aligned}$$

where \mathbf{H}^{-1} is given as

$$\mathbf{H}^{-1} = \left(\mathbf{L} - \frac{2}{\Delta t} \mathbf{I} \right)^{-1},$$

and is formulated by the method described in Section 2.4. To examine the performance of the scheme, we measure the approximation error defined by $e(N, t_n) = \left(\sum_{i=0}^N |u(x_i, t_n) - v_i^n|^2 \omega_i \right)^{1/2}$.

Figure 3.3(a) presents the error history curves resulting from computing numerical solutions by different grid resolutions. For a given grid resolution N the corresponding curve indicates that the computed solution gradually evolves to the steady state solution within a certain level of accuracy, where the higher the grid resolution, the better the solution accuracy. Figure 3.3(b) shows the convergence rate of the numerical solution at times $t = 2$ and $t = 50$. At $t = 2$ the solution profile has not yet arrived at the steady state, and hence, the convergence rate is of second order due to the Crank-Nicholson difference in time. On the other hand at $t = 50$ the solution profile has reached the steady state, and thus, as the grid resolution increases, the error decays exponentially due to the pseudospectral discretization in space.

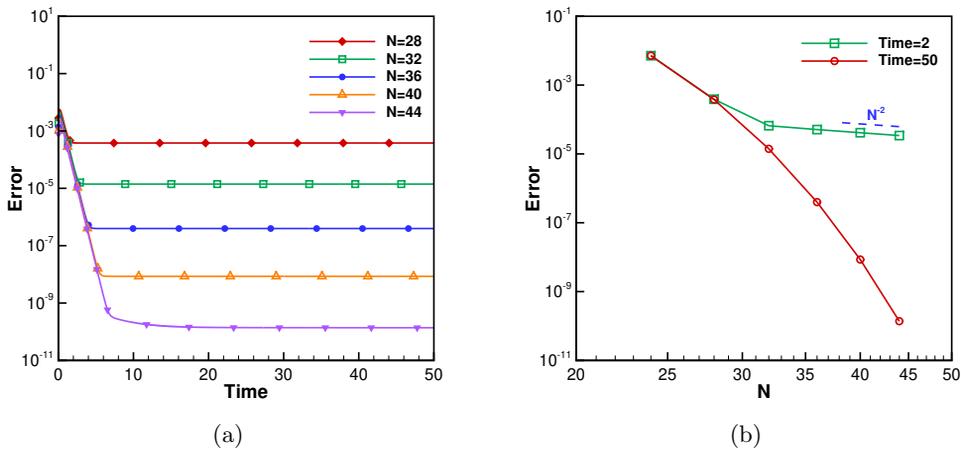


Figure 3.3: (a) Error history curves of the model heat problem (3.7) solved by (3.9) with different grid resolutions. (b) Measured errors at time $t = 2$ and $t = 50$ for different values of N . We take time step as $\Delta = 1/N$ in each computation.

4. Conclusions and future plans

We have presented a framework based on stepwise integrations and low-rank updates to seek an inverse pseudospectral matrix for the Helmholtz differential operator. The resultant inverse matrix was used to solve general second order differential equations, and the computation results clearly illustrated the performance of the operators as expected. Before closing, we would like to address potential applications of the present framework for seeking integration preconditioning matrices.

The present framework for seeking inverse matrices is applied in nodal spaces. Thus, for the basic mixed differential operator, $\frac{d}{dx}a(x)\frac{d}{dx}$, the linkage between differential operators and sandwiched variable coefficients is easily decoupled, and the inverse matrix for the mixed differential operator is formed by a matrix product of the inverses of the decoupled ones. Consequently, existing inverse matrices, developed by other orthogonal polynomial methods for the first order differential operator, can be adopted into the present framework as basic elements to formulate different orthogonal polynomials based operating matrices for the mixed operator as well. Furthermore, using existing inverse spectral/pseudospectral matrices as building blocks in the present divide-and-conquer approach, we can further seek integration preconditioning matrices for other high-order pure, mixed, or even nonlinear differential operators.

We may further apply the present framework to devise preconditioning techniques for multidomain schemes based on pseudospectral, discontinuous Galerkin and summation-by-parts finite difference methods. Similar to the model single domain scheme present in this study, multidomain schemes are composed of element-wise differentiation operators which are commonly products of basic differentiation matrices and element-wise low-rank boundary operators which relate field values between adjacent elements. By exploring these properties within these schemes as shown in this study, we strongly believe that efficient and accurate numerical procedures may be devised. We hope to report these works in the near future.

A. General \mathbf{L}^{-1}

Following a parallel procedure as shown in Section 2.3, we now proceed to finding the matrix $\mathbf{K}^{-1}\mathbf{A}$ to formulate the \mathbf{L}^{-1} matrix for the general case. We factor \mathbf{Q} described in (2.13) as

$$\mathbf{Q} = \widehat{\Phi}\widehat{\Psi}^T,$$

where $\widehat{\Phi}$ and $\widehat{\Psi}^T$ are, respectively, given as

$$\widehat{\Phi} = \begin{bmatrix} \mathbf{D}_- \mathbf{e}_- & \eta \mathbf{e}_- & \eta \mathbf{e}_+ & \mathbf{D}_- \mathbf{e}_+ \end{bmatrix},$$

and

$$\widehat{\Psi}^T = \begin{bmatrix} \eta a_0 \chi_- (\alpha_- \mathbf{e}_-^T - \beta_- \mathbf{e}_-^T \mathbf{D}_+) \\ -\tau_- \alpha_- \eta a_0 \mathbf{e}_-^T - (1 - \tau_- \beta_- \eta) a_0 \mathbf{e}_-^T \mathbf{D}_+ \\ (\chi_+ - \tau_+) (\alpha_+ + \beta_+ \eta) \eta a_N \mathbf{e}_+^T + (\chi_+ - \tau_+) \beta_+ \eta a_N \mathbf{e}_+^T \mathbf{D}_+ \\ (1 - \chi_+ (\alpha_+ + \beta_+ \eta)) \eta a_N \mathbf{e}_+^T - \chi_+ \beta_+ \eta a_N \mathbf{e}_+^T \mathbf{D}_+ \end{bmatrix}.$$

Employing the inverse operators, \mathbf{D}_-^{-1} and \mathbf{D}_+ , we compute $\mathbf{D}_-^{-1} \widehat{\Phi}$ and $\widehat{\Psi}^T \mathbf{D}_+^{-1}$ and have

$$\Phi = \mathbf{D}_-^{-1} \widehat{\Phi} = \begin{bmatrix} \mathbf{e}_- & \mathbf{p}_0 & \mathbf{p}_n & \mathbf{e}_+ \end{bmatrix},$$

and

$$\Psi^T = \widehat{\Psi}^T \mathbf{D}_+^{-1} = \begin{bmatrix} -\chi_- \alpha_- \eta a_0 \mathbf{p}_0^T - \chi_- \beta_- \eta^2 a_0 \mathbf{e}_-^T \\ +\tau_- \alpha_- \eta a_0 \mathbf{p}_0^T - (1 - \tau_- \beta_- \eta) \eta a_0 \mathbf{e}_-^T \\ -(\chi_+ - \tau_+) (\alpha_+ + \beta_+ \eta) \eta a_N \mathbf{p}_N^T + (\chi_+ - \tau_+) \beta_+ \eta^2 a_N \mathbf{e}_+^T \\ -(1 - \chi_+ (\alpha_+ + \beta_+ \eta)) \eta a_N \mathbf{p}_N^T - \chi_+ \beta_+ \eta^2 a_N \mathbf{e}_+ \end{bmatrix} \mathbf{M} = \mathbf{\Lambda} \Phi^T \mathbf{M},$$

where with the use of the relationships in (2.7) the matrix $\mathbf{\Lambda}$ is given as

$$\mathbf{\Lambda} = \eta \begin{bmatrix} -\chi_- \beta_- \eta a_0 & -\chi_- \alpha_- a_0 & 0 & 0 \\ -\chi_- \alpha_- a_0 & \tau_- \alpha_- a_0 & 0 & 0 \\ 0 & 0 & -(\chi_+ - \tau_+) (\alpha_+ + \beta_+ \eta) a_N & (\chi_+ - \tau_+) \beta_+ \eta a_N \\ 0 & 0 & -(\chi_+ - \tau_+) \beta_+ \eta a_N & -\chi_+ \beta_+ \eta a_N \end{bmatrix}.$$

To find $\mathbf{K}^{-1} \mathbf{\Lambda}$ we compute $\Phi^T \mathbf{M} \mathbf{A}^{-1} \Phi$ and $\mathbf{\Lambda}^{-1}$, and the results are

$$\Phi^T \mathbf{M} \mathbf{A}^{-1} \Phi = \frac{1}{\eta} \begin{bmatrix} \frac{1}{a_0} & \frac{1}{a_0} & \frac{(-1)^N}{a_0} & 0 \\ \frac{1}{a_0} & \kappa_2 & \kappa_1 & \frac{1}{a_N} \\ \frac{(-1)^N}{a_0} & \kappa_1 & \kappa_0 & \frac{1}{a_N} \\ 0 & \frac{1}{a_N} & \frac{1}{a_N} & \frac{1}{a_N} \end{bmatrix}, \quad \kappa_\nu = \sum_{i=0}^N \frac{1}{a_i (P_N(x_i))^\nu},$$

and

$$\mathbf{\Lambda}^{-1} = \frac{1}{\eta} \begin{bmatrix} -\frac{\tau_-}{\chi_- a_0} & -\frac{1}{a_0} & 0 & 0 \\ -\frac{1}{a_0} & \frac{\beta_- \eta}{\alpha_- a_0} & 0 & 0 \\ 0 & 0 & -\frac{\chi_+}{(\chi_+ - \tau_+) a_N} & -\frac{1}{a_N} \\ 0 & 0 & -\frac{1}{a_N} & -\frac{1}{a_N} \left(\frac{\alpha_+}{\beta_+ \eta} + 1 \right) \end{bmatrix}.$$

Thus, we have $\mathbf{\Lambda}^{-1}\mathbf{K}$ as

$$\mathbf{\Lambda}^{-1}\mathbf{K} = \mathbf{\Lambda}^{-1} + \mathbf{\Phi}^T \mathbf{M} \mathbf{A}^{-1} \mathbf{\Phi} = \frac{1}{\eta} \begin{bmatrix} \left(1 - \frac{\tau_-}{\chi_-}\right) \frac{1}{a_0} & 0 & \frac{(-1)^N}{a_0} & 0 \\ 0 & \kappa_2 + \frac{\beta - \eta}{\alpha - a_0} & \kappa_1 & \frac{1}{a_N} \\ \frac{(-1)^N}{a_0} & \kappa_1 & \kappa_0 - \frac{\chi_+}{(\chi_+ - \tau_+)a_N} & 0 \\ 0 & \frac{1}{a_N} & 0 & -\frac{\alpha_+}{\beta_+ \eta a_N} \end{bmatrix}.$$

Then the inverse of $\mathbf{\Lambda}^{-1}\mathbf{K}$ is found as

$$\mathbf{K}^{-1}\mathbf{\Lambda} = \frac{\eta}{T} \begin{bmatrix} t_{11} & t_{12} & t_{13} & t_{14} \\ t_{12} & t_{22} & t_{23} & t_{24} \\ t_{13} & t_{23} & t_{33} & t_{34} \\ t_{14} & t_{24} & t_{34} & t_{44} \end{bmatrix},$$

where T is given as

$$\begin{aligned} T = & \frac{(\chi_- - \tau_-)(\chi_+ - \tau_+)}{a_0 a_N} \left((\kappa_1^2 - \kappa_0 \kappa_2) \alpha_- \alpha_+ - \frac{\alpha_+ \beta_- \eta \kappa_0}{a_0} \right) \\ & + \frac{1}{a_0 a_N} \left(\kappa_2 \alpha_- \alpha_+ + \frac{\alpha_+ \beta_- \eta}{a_0} \right) \left(\frac{(\chi_- - \tau_-) \chi_+}{a_N} + \frac{(\chi_+ - \tau_+) \chi_-}{a_0} \right) \\ & - \frac{\alpha_- \beta_+ \eta}{a_0 a_N^2} \left((\chi_- - \tau_-)(\chi_+ - \tau_+) \kappa_0 - \frac{(\chi_- - \tau_-) \chi_+}{a_N} - \frac{(\chi_+ - \tau_+) \chi_-}{a_0} \right), \end{aligned}$$

and t_{ij} for $1 \leq i \leq 4$ and $i \leq j \leq 4$ are found as follows:

$$\begin{aligned} t_{11} = & \frac{\chi_-}{a_N} (\chi_+ - \tau_+) (\kappa_1^2 - \kappa_0 \kappa_2) \alpha_- \alpha_+ + \frac{\chi_- \chi_+}{a_N^2} (\kappa_2 \alpha_- \alpha_+) \\ & - \frac{\chi_-}{a_N} \left(\frac{\alpha_+ \beta_- \eta}{a_0} + \frac{\alpha_- \beta_+ \eta}{a_N} \right) \left(\kappa_0 (\chi_+ - \tau_+) - \frac{\chi_+}{a_N} \right), \\ t_{12} = & \frac{(-1)^N \chi_- \alpha_- \beta_+ \eta \kappa_1}{a_0} \left(\kappa_0 (\chi_+ - \tau_+) - \frac{\chi_+}{a_N} \right), \\ t_{13} = & \frac{(-1)^N \chi_- (\chi_+ - \tau_+)}{a_0 a_N} \left(\frac{\alpha_- \beta_+ \eta}{a_N} + \frac{\alpha_+ \beta_- \eta}{a_0} + \kappa_2 \alpha_- \alpha_+ \right), \\ t_{14} = & \frac{(-1)^{N+1} \chi_- (\chi_+ - \tau_+) \kappa_1 \alpha_- \beta_+ \eta}{a_0 a_N}, \\ t_{22} = & \left(\frac{\chi_+ (\chi_- - \tau_-)}{a_N} + \frac{\chi_- (\chi_+ - \tau_+)}{a_0} - (\chi_- - \tau_-)(\chi_+ - \tau_+) \kappa_0 \right) \frac{\alpha_- \alpha_+}{a_0 a_N}, \\ t_{23} = & \frac{(\chi_+ - \tau_+) (\chi_- - \tau_-)}{a_0 a_N} \kappa_1 \alpha_- \alpha_+, \\ t_{24} = & \frac{\alpha_- \beta_+ \eta}{a_0 a_N} \left(\frac{\chi_- (\chi_+ - \tau_+)}{a_0} + \frac{\chi_+ (\chi_- - \tau_-)}{a_N} - \kappa_0 (\chi_- - \tau_-)(\chi_+ - \tau_+) \right), \\ t_{33} = & -\frac{(\chi_- - \tau_-)(\chi_+ - \tau_+)}{a_0 a_N} \left(\frac{\alpha_+ \beta_- \eta}{a_0} + \frac{\alpha_- \beta_+ \eta}{a_N} + \kappa_2 \alpha_- \alpha_+ \right), \end{aligned}$$

$$\begin{aligned}
t_{34} &= (\chi_+ - \tau_+)(\chi_- - \tau_-) \frac{\alpha_- \beta_+ \eta \kappa_1}{a_0 a_N}, \\
t_{44} &= -\frac{\alpha_- \beta_+ \eta}{a_0} \left((\kappa_1^2 - \kappa_0 \kappa_2)(\chi_- - \tau_-)(\chi_+ - \tau_+) + \kappa_2 \left(\frac{\chi_+(\chi_- - \tau_-)}{a_N} + \frac{\chi_-(\chi_+ - \tau_+)}{a_0} \right) \right) \\
&\quad + \frac{\beta_- \beta_+ \eta^2}{a_0^2} \left(\kappa_0(\chi_- - \tau_-)(\chi_+ - \tau_+) - \frac{\chi_+(\chi_- - \tau_-)}{a_N} - \frac{\chi_-(\chi_+ - \tau_+)}{a_0} \right).
\end{aligned}$$

Acknowledgments

We thank the anonymous reviewers for their many insightful comments and suggestions. Li was supported in part by the Ministry of Science and Technology of Taiwan under grant 106-2115-M-030-001. Teng was supported in part by the Ministry of Science and Technology of Taiwan under grant 105-2115-M-005-007.

References

- [1] C. Canuto, M. Y. Hussaini, A. Quarteroni and T. A. Zang, *Spectral Methods: Fundamentals in Single Domains*, Scientific Computation, Springer, Berlin, 2006.
- [2] ———, *Spectral Methods: Evolution to Complex Geometries and Applications to Fluid Dynamics*, Scientific Computation, Springer, Berlin, 2007.
- [3] M. H. Carpenter and D. Gottlieb, *Spectral methods on arbitrary grids*, J. Comput. Phys. **129** (1996), no. 1, 74–86.
- [4] E. H. Doha, *On the construction of recurrence relations for the expansion and connection coefficients in series of Jacobi polynomials*, J. Phys. A **37** (2004), no. 3, 657–675.
- [5] E. H. Doha and W. M. Abd-Elhameed, *Efficient spectral-Galerkin algorithms for direct solution of second-order equations using ultraspherical polynomials*, SIAM J. Sci. Comput. **24** (2002), no. 2, 548–571.
- [6] ———, *Efficient spectral ultraspherical-dual-Petrov-Galerkin algorithms for the direct solution of $(2n + 1)$ th-order linear differential equations*, Math. Comput. Simulation **79** (2009), no. 11, 3221–3242.
- [7] E. H. Doha and A. H. Bhrawy, *Efficient spectral-Galerkin algorithms for direct solution for second-order differential equations using Jacobi polynomials*, Numer. Algorithms **42** (2006), no. 2, 137–164.
- [8] W. S. Don and D. Gottlieb, *The Chebyshev-Legendre method: Implementing Legendre methods on Chebyshev points*, SIAM J. Numer. Anal. **31** (1994), no. 6, 1519–1534.

- [9] D. Elad and N. Gavish, *Finite domain effects in steady state solutions of Poisson-Nernst-Planck equations*, SIAM J. Appl. Math. **79** (2019), no. 3, 1030–1050.
- [10] E. M. E. Elbarbary, *Integration preconditioning matrix for ultraspherical pseudospectral operators*, SIAM J. Sci. Comput. **28** (2006), no. 3, 1186–1201.
- [11] D. Funaro and D. Gottlieb, *A new method of imposing boundary conditions in pseudospectral approximations of hyperbolic equations*, Math. Comp. **51** (1988), no. 184, 599–613.
- [12] ———, *Convergence results for pseudospectral approximations of hyperbolic systems by a penalty-type boundary treatment*, Math. Comp. **57** (1991), no. 196, 585–596.
- [13] F. Ghoreishi and S. M. Hosseini, *A preconditioned implementation of pseudospectral methods on arbitrary grids*, Appl. Math. Comput. **148** (2004), no. 1, 15–34.
- [14] W. W. Hager, *Updating the inverse of a matrix*, SIAM Rev. **31** (1989), no. 2, 221–239.
- [15] J. S. Hesthaven, *Integration preconditioning of pseudospectral operators I: Basic linear operators*, SIAM J. Numer. Anal. **35** (1998), no. 4, 1571–1593.
- [16] ———, *Spectral penalty methods*, Appl. Numer. Math. **33** (2000), no. 1-4, 23–41.
- [17] J. S. Hesthaven and D. Gottlieb, *A stable penalty method for the compressible Navier-Stokes equations I: Open boundary conditions*, SIAM J. Sci. Comput. **17** (1996), no. 3, 579–612.
- [18] J. S. Hesthaven, S. Gottlieb and D. Gottlieb, *Spectral Methods for Time-dependent Problems*, Cambridge Monographs on Applied and Computational Mathematics **21**, Cambridge University Press, Cambridge, 2007.
- [19] Y.-T. Li, P.-Y. Lin and C.-H. Teng, *Inverse matrices for pseudospectral differentiation operators in polar coordinates by stepwise integrations and low-rank updates*, Appl. Numer. Math. **150** (2020), 519–535.
- [20] S. Olver and A. Townsend, *A fast and well-conditioned spectral method*, SIAM Rev. **55** (2013), no. 3, 462–489.
- [21] S. A. Orszag, *Numerical simulation of incompressible flows within simple boundaries: Accuracy*, J. Fluid Mech. **49** (1971), 75–112.
- [22] ———, *Numerical simulation of incompressible flows within simple boundaries I: Galerkin (spectral) representations*, Studies in Appl. Math. **50** (1971), no. 4, 293–327.

- [23] J. Shen, *Efficient spectral-Galerkin method I: Direct solvers of second- and fourth-order equations using Legendre polynomials*, SIAM J. Sci. Comput. **15** (1994), no. 6, 1489–1505.
- [24] ———, *Efficient spectral-Galerkin method II: Direct solvers of second- and fourth-order equations using Chebyshev polynomials*, SIAM J. Sci. Comput. **16** (1995), no. 1, 74–87.
- [25] ———, *A new dual-Petrov-Galerkin method for third and higher odd-order differential equations: Application to the KdV equation*, SIAM J. Numer. Anal. **41** (2003), no. 5, 1595–1619.
- [26] J. Shen, T. Tang and L.-L. Wang, *Spectral Methods: Algorithms, analysis and applications*, Springer, Heidelberg, 2011.
- [27] J. Shen and L.-L. Wang, *Legendre and Chebyshev dual-Petrov-Galerkin methods for hyperbolic equations*, Comput. Methods Appl. Mech. Engrg. **196** (2007), no. 37-40, 3785–3797.
- [28] ———, *Some recent advances on spectral methods for unbounded domains*, Commun. Comput. Phys. **5** (2009), no. 2-4, 195–241.
- [29] J. Shen, Y. Wang and J. Xia, *Fast structured direct spectral methods for differential equations with variable coefficients I: The one-dimensional case*, SIAM J. Sci. Comput. **38** (2016), no. 1, A28–A54.
- [30] L.-L. Wang, M. D. Samson and X. Zhao, *A well-conditioned collocation method using a pseudospectral integration matrix*, SIAM J. Sci. Comput. **36** (2014), no. 3, A907–A929.

Yung-Ta Li

Department of Mathematics, Fu Jen Catholic University, New Taipei City 242, Taiwan

E-mail address: ytli@math.fju.edu.tw

Ping-Hsuan Tsai

Department of Computer Science, University of Illinois at Urbana-Campaign, Urbana, IL 61801, USA

E-mail address: pht2@illinois.edu

Chun-Hao Teng

Department of Applied Mathematics, National Chung Hsing University, Taichung 402, Taiwan

E-mail address: tengch@email.nchu.edu.tw