

Notes from *Spectral Methods: Algorithms, Analysis, and
Applications*

by Jie Shen, Tao Tang, Li-Lian Wang

taken by Samuel T. Wallace

Publisher's Description

Along with finite differences and finite elements, spectral methods are one of the three main methodologies for solving partial differential equations on computers. This book provides a detailed presentation of basic spectral algorithms, as well as a systematical presentation of basic convergence theory and error analysis for spectral methods. Readers of this book will be exposed to a unified framework for designing and analyzing spectral algorithms for a variety of problems, including in particular high-order differential equations and problems in unbounded domains. The book contains a large number of figures which are designed to illustrate various concepts stressed in the book. A set of basic matlab codes has been made available online to help the readers to develop their own spectral codes for their specific applications.

A Note From the Transcriber

These notes were taken over summer 2020 as part of self-study preparing for PhD in applied math and numerical analysis. I am reading this without much interpolation theory knowledge, and some prior exposure to spectral methods through Trefethen's book *Spectral Methods in MATLAB* (sorry, no notes for that book). This book comes from a suggestion by a professor as my grad school, and it looked temptingly challenging.

Contents

0.1	Introduction	3
0.1.1	Weighted Residual Methods	3
0.1.2	Spectral-Collocation Method	5
0.1.3	Spectral Methods of Galerkin Type	7

0.1 Introduction

0.1.1 Weighted Residual Methods

Consider the following general problem:

$$\partial_t u(x, t) - \mathcal{L}u(x, t) = \mathcal{N}(u)(x, t), \quad t > 0, x \in \Omega \quad (1)$$

Where \mathcal{L} is a leading spatial derivative operator, and \mathcal{N} is a lower-order linear or non-linear operator involving only spatial derivatives. Here, Ω denotes a bounded domain of \mathbb{R}^d , $d = 1, 2$, or 3 . This equation is to be supplemented with an initial condition and suitable boundary conditions.

We shall only consider the WRM for the spatial discretization, and assume that the time derivative is discretized with a suitable time-stepping scheme. Among various time-stepping methods, semi-implicit schemes or linearly-implicit schemes, in which the principal linear operators are treated *implicitly* to reduce the associated stability constraint, while the non-linear equations are treated explicitly to avoid the expensive process of solving nonlinear equations at each time step, are most frequently used in the context of spectral methods.

Let τ be the step size, and $u^k(\cdot)$ be an approximation of $u(\cdot, k\tau)$. As an example, we consider the Crank-Nicolson leap-frog scheme for the equation:

$$\frac{u^{n+1} - u^{n-1}}{2\tau} - \mathcal{L} \left(\frac{u^{n+1} + u^{n-1}}{2} \right) = \mathcal{N}(u^n) \quad n \geq 1 \quad (2)$$

We can rewrite this as

$$\mathbf{L}u(x) := \alpha u(x) - \mathcal{L}u(x) = f(x), \quad x \in \Omega \quad (3)$$

where $u = \frac{u^{n+1} + u^{n-1}}{2}$, $\alpha = \tau^{-1}$ and $f = \alpha u^{n-1} + \mathcal{N}(u^n)$. Hence, at each time step, we need to solve a steady-state problem of the form of (3).

At this point, it is important to emphasize that the construction of efficient numerical solvers for some important equations in the form of (3), such as Poisson-type equations and advection-diffusion equations, is an essential step in solving general nonlinear PDEs. With this in mind, a particular emphasis for equations of the form (3) where \mathcal{L} is a *linear elliptic* operator.

The starting point of the WRM is to approximate the solution u is to approximate (3) by a finite sum

$$u(x) \approx u_N(x) = \sum_{k=0}^N a_k \phi_k(x) \quad (4)$$

where $\{\phi_k\}$ are the *trial (or basis) functions*, and the expansion coefficients are to be determined. Substituting u_N for u in (3) leads to the *residual*

$$\mathbf{R}_N(x) = \mathbf{L}u_N(x) - f(x) \neq 0 \quad x \in \Omega \quad (5)$$

The notion of the WRM is to force the residual to zero by requiring

$$(\mathbf{R}_N, \psi_j)_\omega = \int_{\Omega} \mathbf{R}_N(x) \psi_j(x) \omega(x) dx = 0, \quad 0 \leq j \leq N \quad (6)$$

where $\{\psi_j\}$ are the *test functions*, and ω is a positive weight function; or

$$\langle \mathbf{R}_N, \psi_j \rangle_{N,\omega} := \sum_{k=0}^N \mathbf{R}_N(x_k) \psi_j(x_k) \omega_k = 0, \quad 0 \leq j \leq N \quad (7)$$

where $\{x_k\}_{k=0}^N$ are a set of preselected collocation points, and $\{\omega_k\}_{k=0}^N$ are the weights of a numerical quadrature formula.

The choice of trial/test functions is one of the main features that distinguishes spectral methods from finite-elements and finite-difference methods. In the latter two methods, the trial/test functions are local in character with finite regularities. In contrast, spectral methods employ globally smooth functions as trial/test functions. The most commonly used trial/test functions are trigonometric functions of orthogonal polynomials (typically, the eigenfunctions of singular Sturm-Liouville problems), which include

- $\phi_k(x) = e^{ikx}$ (Fourier spectral method)
- $\phi_k(x) = T_k(x)$ (Chebyshev spectral method)
- $\phi_k = L_k(x)$ (Legendre spectral method)
- $\phi_k = \mathcal{L}_k(x)$ (Laguerre spectral method)
- $\phi_k(x) = H_k(x)$ (Hermite spectral method)

Here, T_k, L_k, \mathcal{L}_k , and H_k are the Chebyshev, Legendre, Laguerre and Hermite polynomials of degree k respectively.

The choice of test functions distinguishes the following formulations:

- *Galerkin.* The test functions are the same as the trial ones (i.e., $\phi_k = \psi_k$ in (6) or (7)), assuming the boundary conditions are periodic or homogeneous.
- *Petrov-Galerkin.* The test functions are different from the trial ones.
- *Collocation.* The test functions $\{\psi_k\}$ in (7) are the Lagrange basis polynomials such that $\psi_k(x_j) = \delta_{jk}$, where $\{x_j\}$ are preassigned collocation points. Hence, the residual is forced to zero at each x_j , i.e. $\mathbf{R}_N(x_j) = 0$.

Remark 0.1.1. *In the literature, the term of pseudo-spectral methods is often used to describe any spectral method where some operations involve a collocation approach or a numerical quadrature which produces aliasing errors. In this sense, almost all practical spectral methods are pseudo-spectral. In this book, we shall not classify a method as pseudo-spectral or spectral. Instead, it will be classified as Galerkin type or collocation type.*

Remark 0.1.2. *The so-called tau method is a particular class of Petrov-Galerkin method. While the tau method offers some advantages in certain situations, for most problems, it is usually better to use a well-designed Galerkin or Petrov-Galerkin method. So in this book, we shall not touch on this topic, and refer to the references therein for a thorough discussion of this approach.*

In the forthcoming sections, we shall demonstrate how to construct spectral methods for solving differential equations by examining several spectral schemes based on Galerkin, Petrov-Galerkin, and collocation formulas in a general manner. We shall revisit these illustrative examples in a more rigorous fashion in the main body of the book.

0.1.2 Spectral-Collocation Method

To fix the idea, let us consider the following linear problem:

$$\mathbf{L}u(x) = -u''(x) + p(x)u'(x) + q(x)u(x) = f(x), \quad x \in (-1, 1) \quad (8)$$

$$B_{\pm}u(\pm 1) = g_{\pm} \quad (9)$$

Where B_{\pm} are linear operators corresponding to Dirichlet, Neumann, or Robin boundary conditions, and the data p, q, f and g_{\pm} are given such that the above problem is well-posed.

As mentioned above, the collocation method forces the residual to vanish pointwisely at a set of preassigned points. More precisely, let $\{x_j\}_{j=0}^N$ (with $x_0 = -1$ and $x_N = 1$) be a set of Gauss-Lobatto points (see Chap. 3), and let P_N be the set of all real algebraic polynomials of degree $\leq N$. The spectral-collocation method for (8) amounts to finding $u_N \in P_N$ such that

1. the residual $\mathbf{R}_N(x_k) = \mathbf{L}u_N(x_k) - f(x_k) = 0$, $1 \leq k \leq N-1$

2. u_N satisfies exactly the boundary conditions, i.e.,

$$B_- u_N(x_0) = g_-, \quad B_+ u_N(x_N) = g_+ \quad (10)$$

The spectral-collocation method is usually implemented in the physical space by seeking approximate solution in the form

$$u_N(x) = \sum_{j=0}^N u_N(x_j) h_j(x) \quad (11)$$

where $\{h_j\}$ are the Lagrange basis polynomials (also referred to as *nodal* basis functions), i.e., $h_j \in P_N$ and $h_j(x_k) = \delta_{jk}$. Hence, in inserting (11) into (9)-(10) leads to the linear system

$$\sum_{j=0}^N [\mathbf{L}h_j(x_k)] u_N(x_j) = f(x_k), \quad (12)$$

$$\sum_{j=0}^N [\mathbf{L}h_j(x_k)] u_N(x_j) = f(x_k), \quad 1 \leq k \leq N-1 \quad (13)$$

$$\sum_{j=0}^N [B_- h_j(x_0)] u_N(x_j) = g_-, \quad \sum_{j=0}^N [B_+ h_j(x_N)] u_N(x_j) = g_+ \quad (14)$$

The above system contains $N+1$ unknowns, so we can rewrite it in a matrix form. To fix the idea, we consider (8) with Dirichlet boundary conditions: $u(\pm 1) = g_{\pm}$. In this case, setting $u_N(x_0) = g_-$ and $u_N(x_N) = g_+$ in the first equation of (12) reduces to

$$\sum_{j=1}^{N-1} [\mathbf{L}h_j(x_k)] u_N(x_j) = f(x_k) - \{[\mathbf{L}h_0(x_k)] g_- + [\mathbf{L}h_N(x_k)] g_+\} \quad (15)$$

for $1 \leq k \leq N-1$. Differentiating (11) m times leads to

$$u_N^{(m)}(x_k) = \sum_{j=0}^N d_{kj}^{(m)} u_N(x_j) \quad (16)$$

where

$$d_{kj}^{(m)} = h_j^{(m)}(x_k) \quad (17)$$

The matrix $D^{(m)} \left(d_{kj}^{(m)} \right)_{k,j=0 \dots N}$ is called the differentiation matrix of order m relative to the $\{s_j\}_{j=0}^N$. If we denote by $\mathbf{u}^{(m)}$ the vector whose components are the values of $u_N^{(m)}$ at the collocation points, it follows from (14) that

$$\mathbf{u}^{(m)} = D^{(m)} \mathbf{u}^{(0)}, \quad m \geq 1 \quad (18)$$

Hence, we have

$$\mathbf{L}h_j(x_k) = -d_{kj}^{(2)} + p(x_k)d_{kj}^{(1)} + q(x_k)\delta_{kj} \quad (19)$$

Denote by \mathbf{f} the vector with $N - 1$ components given by the right-handside of (13). Setting

$$\tilde{D}_m = \left(d_{kj}^{(m)} \right)_{kj=1,\dots,N}, \quad m = 1, 2 \quad (20)$$

$$P = \text{diag}(p(x_1), \dots, p(x_{N-1})), \quad Q = \text{diag}(q(x_1), \dots, q(x_{N-1}))$$

the system (13) reduces to

$$\left(-\tilde{D}_2 + P\tilde{D}_1 + Q \right) \mathbf{u}^{(0)} = \mathbf{f} \quad (21)$$

Observe that the collocation method is easy to implement, once the differentiation matrices are precomputed. Moreover, it is very convenient for solving problems with variable coefficients and/or nonlinear problems, since we work in the physical space and derivatives can be valuated by (14) directly. As a result, the collocation method has been extensively used in practice. However, three important issues should be considered in the implementation and analysis of a collocation method:

- The coefficient matrix of the collocation system is always full with a condition number behaving like $O(N^{2m})$ (m is the order of the differential equation).
- The choice of collocation points is crucial in terms of stability, accuracy, and ease of dealing with boudnary conditions. In general, they are chosen as nodes (typically, zeros of orthogonal polynomials) of Gauss-type quadrature formulas.
- The aforementioned collocation scheme is formulated in a *strong* form. In terms of error analysis, it is more convenient to reformulate it as a (but not always equivalent) *weak* form, see Sect. 1.3.3 and Chap. 4.

0.1.3 Spectral Methods of Galerkin Type