

**Get Smart! at :**  
**[www.GetPedia.com](http://www.GetPedia.com)**

**\*More than 150,000 articles for  
DUMMIES**

**\*Learn how almost everything  
works**

**\*Get Smart! Get Pedia!**

# Finite Difference Methods for Differential Equations

Randall J. LeVeque

DRAFT VERSION for use in the courses  
AMath 585-6  
University of Washington  
Winter/Spring Quarters, 1998

**WARNING:** These notes are incomplete and may contain errors.  
They are made available primarily for students in my courses.  
Please contact me for other uses.  
`rjl@amath.washington.edu`



# Contents

<b>I</b>	<b>Basic Text</b>	<b>1</b>
<b>1</b>	<b>Finite difference approximations</b>	<b>3</b>
1.1	Truncation errors . . . . .	5
1.2	Deriving finite difference approximations . . . . .	6
1.3	Polynomial interpolation . . . . .	7
1.4	Second order derivatives . . . . .	7
1.5	Higher order derivatives . . . . .	7
<b>2</b>	<b>Boundary Value Problems</b>	<b>9</b>
2.1	The heat equation . . . . .	9
2.2	Boundary conditions . . . . .	10
2.3	The steady-state problem . . . . .	10
2.4	A simple finite difference method . . . . .	11
2.5	Local truncation error . . . . .	12
2.6	Global error . . . . .	13
2.7	Stability . . . . .	14
2.8	Consistency . . . . .	14
2.9	Convergence . . . . .	15
2.10	Stability in the 2-norm . . . . .	15
2.11	Neumann boundary conditions . . . . .	17
2.12	Green's functions and max-norm stability . . . . .	19
2.13	Existence and uniqueness . . . . .	22
2.14	A general linear second order equation . . . . .	23
2.15	Nonlinear Equations . . . . .	25
2.15.1	Discretization of the nonlinear BVP . . . . .	26
2.15.2	Nonconvergence . . . . .	28
2.15.3	Nonuniqueness . . . . .	28
2.15.4	Accuracy . . . . .	29
2.16	Singular perturbations and boundary layers . . . . .	30
2.16.1	Interior layers . . . . .	33
2.17	Nonuniform grids and adaptive refinement . . . . .	33
<b>3</b>	<b>Elliptic Equations</b>	<b>35</b>
3.1	Steady-state heat conduction . . . . .	35
3.2	The five-point stencil for the Laplacian . . . . .	36
3.3	Solving the linear system . . . . .	39
3.3.1	Gaussian elimination . . . . .	39
3.3.2	Fast Poisson solvers . . . . .	40
3.3.3	Iterative methods . . . . .	42
3.4	Accuracy and stability . . . . .	43
3.5	The nine-point Laplacian . . . . .	45

<b>4</b>	<b>Function Space Methods</b>	<b>47</b>
4.1	Collocation . . . . .	47
4.2	Spectral methods . . . . .	48
4.2.1	Matrix interpretation . . . . .	50
4.2.2	Accuracy . . . . .	50
4.2.3	Stability . . . . .	51
4.2.4	Collocation property . . . . .	51
4.3	The finite element method . . . . .	51
4.3.1	Two space dimensions . . . . .	54
<b>5</b>	<b>Iterative Methods for Sparse Linear Systems</b>	<b>57</b>
5.1	Matrix splitting methods . . . . .	57
5.1.1	Rate of convergence . . . . .	59
5.1.2	SOR . . . . .	61
5.2	Conjugate gradient methods . . . . .	62
5.2.1	Preconditioners . . . . .	62
5.3	Multigrid methods . . . . .	62
<b>6</b>	<b>The Initial Value Problem for ODE's</b>	<b>63</b>
6.1	Lipschitz continuity . . . . .	64
6.1.1	Existence and uniqueness of solutions . . . . .	65
6.1.2	Systems of equations . . . . .	66
6.1.3	Significance of the Lipschitz constant . . . . .	66
6.1.4	Limitations . . . . .	67
6.2	Some basic numerical methods . . . . .	68
6.3	Truncation errors . . . . .	69
6.4	One-step errors . . . . .	69
6.5	Linear Multistep Methods . . . . .	70
6.5.1	Local truncation error . . . . .	71
6.5.2	Characteristic polynomials . . . . .	72
6.5.3	Starting values . . . . .	72
6.6	Taylor series methods . . . . .	73
6.7	Runge-Kutta Methods . . . . .	74
<b>7</b>	<b>Zero-Stability and Convergence for Initial Value Problems</b>	<b>77</b>
7.1	Convergence . . . . .	77
7.2	Linear equations and Duhamel's principle . . . . .	78
7.3	One-step methods . . . . .	78
7.3.1	Euler's method on linear problems . . . . .	78
7.3.2	Relation to stability for BVP's . . . . .	80
7.3.3	Euler's method on nonlinear problems . . . . .	81
7.3.4	Realistic error estimates . . . . .	82
7.3.5	General 1-step methods . . . . .	83
7.4	Zero-stability of linear multistep methods . . . . .	83
7.4.1	Solving linear difference equations . . . . .	84
<b>8</b>	<b>Absolute Stability for ODEs</b>	<b>89</b>
8.1	Unstable computations with a zero-stable method . . . . .	89
8.2	Absolute stability . . . . .	91
8.3	Stability regions for LMMs . . . . .	91
8.4	The Boundary Locus Method . . . . .	96
8.5	Systems of equations . . . . .	98
8.5.1	Chemical Kinetics . . . . .	98

8.5.2	Linear systems . . . . .	99
8.5.3	Nonlinear systems . . . . .	100
8.6	Choice of stepsize . . . . .	101
<b>9</b>	<b>Linear Multistep Methods as 1-step Methods</b>	<b>103</b>
9.1	Absolute stability . . . . .	104
9.2	Convergence and zero-stability . . . . .	104
<b>10</b>	<b>Stiff ODEs</b>	<b>107</b>
10.1	Numerical Difficulties . . . . .	107
10.2	Characterizations of stiffness . . . . .	109
10.3	Numerical methods for stiff problems . . . . .	110
10.3.1	A-stability . . . . .	110
10.3.2	L-stability . . . . .	110
10.4	BDF Methods . . . . .	112
<b>11</b>	<b>Some basic PDEs</b>	<b>115</b>
11.1	Classification of differential equations . . . . .	115
11.1.1	Second-order equations . . . . .	115
11.1.2	Elliptic equations . . . . .	115
11.1.3	Parabolic equations . . . . .	116
11.1.4	Hyperbolic equations . . . . .	116
11.2	Derivation of PDEs from conservation principles . . . . .	117
11.3	Advection . . . . .	117
11.4	Diffusion . . . . .	119
11.5	Source terms . . . . .	119
11.5.1	Reaction-diffusion equations . . . . .	120
<b>12</b>	<b>Fourier Analysis of Linear PDEs</b>	<b>121</b>
12.1	Fourier transforms . . . . .	121
12.2	Solution of differential equations . . . . .	121
12.3	The heat equation . . . . .	123
12.4	Dispersive waves . . . . .	123
12.5	Even vs. odd order derivatives . . . . .	124
<b>13</b>	<b>Diffusion Equations</b>	<b>125</b>
13.1	Local truncation errors and order of accuracy . . . . .	127
13.2	Method of Lines discretizations . . . . .	127
13.3	Stability theory . . . . .	129
13.4	Stiffness of the heat equation . . . . .	129
13.5	Convergence . . . . .	130
13.5.1	PDE vs. ODE stability theory . . . . .	132
13.6	von Neumann analysis . . . . .	133
13.7	Multi-dimensional problems . . . . .	137
13.8	The LOD method . . . . .	138
13.8.1	Boundary conditions . . . . .	138
13.8.2	Accuracy and stability . . . . .	139
13.8.3	The ADI method . . . . .	140

<b>14 Advection Equations</b>	<b>141</b>
14.1 MOL discretization . . . . .	142
14.1.1 Forward Euler time discretization . . . . .	143
14.1.2 Leapfrog . . . . .	144
14.1.3 Lax-Friedrichs . . . . .	144
14.2 The Lax-Wendroff method . . . . .	145
14.2.1 Stability analysis . . . . .	147
14.3 Upwind methods . . . . .	148
14.3.1 Stability analysis . . . . .	149
14.3.2 The Beam-Warming method . . . . .	149
14.4 Characteristic tracing . . . . .	149
<b>15 The CFL Condition</b>	<b>151</b>
<b>16 Modified Equations — Numerical Diffusion and Dispersion</b>	<b>155</b>
16.1 Upwind . . . . .	155
16.2 Lax-Wendroff . . . . .	158
16.3 Beam-Warming . . . . .	158
<b>17 Hyperbolic Systems and High-Resolution Methods</b>	<b>159</b>
17.1 Scalar equations . . . . .	159
17.1.1 The linear advection equation . . . . .	159
17.1.2 Burgers' equation . . . . .	160
17.1.3 Shock formation . . . . .	160
17.1.4 Weak solutions . . . . .	161
17.1.5 The Riemann problem . . . . .	162
17.1.6 Shock speed . . . . .	163
17.2 Linear hyperbolic systems . . . . .	163
17.2.1 Characteristic variables . . . . .	164
17.2.2 The Riemann Problem . . . . .	165
17.2.3 The phase plane . . . . .	167
17.3 Finite volume methods . . . . .	168
17.4 Importance of conservation form — incorrect shock speeds . . . . .	169
17.5 Numerical flux functions . . . . .	170
17.6 Godunov's method . . . . .	171
17.6.1 Godunov's method on scalar equations . . . . .	173
17.7 High-resolution methods . . . . .	173
17.7.1 Reconstruct–Solve–Average . . . . .	174
17.7.2 Piecewise linear reconstruction . . . . .	175
17.7.3 Choice of slopes . . . . .	175
17.7.4 Oscillations . . . . .	176
17.7.5 Total variation . . . . .	177
17.7.6 Slope-limiter methods . . . . .	177
17.7.7 Flux formulation with piecewise linears . . . . .	178
17.7.8 Flux limiters . . . . .	179
17.7.9 TVD limiters . . . . .	180
17.7.10 Linear systems . . . . .	182
17.7.11 Implementation and <b>CLAWPACK</b> . . . . .	183
17.8 Boundary conditions . . . . .	183
17.8.1 Periodic boundary conditions . . . . .	184
17.8.2 Outflow boundaries . . . . .	184

<b>18 Mixed Equations and Fractional Step Methods</b>	<b>187</b>
18.1 Advection-reaction equations . . . . .	187
18.1.1 Unsplit methods . . . . .	187
18.1.2 Fractional step methods . . . . .	189
18.2 General formulation of fractional step methods . . . . .	191
18.3 Strang splitting . . . . .	193
<b>II Appendices</b>	<b>A–1</b>
<b>A1 Measuring Errors</b>	<b>A–1</b>
A1.1 Errors in a scalar value . . . . .	A–1
A1.1.1 Absolute error . . . . .	A–1
A1.1.2 Relative error . . . . .	A–2
A1.2 “Big-oh” and “little-oh” notation . . . . .	A–2
A1.3 Errors in vectors . . . . .	A–3
A1.3.1 Norm equivalence . . . . .	A–4
A1.3.2 Matrix norms . . . . .	A–5
A1.4 Errors in functions . . . . .	A–5
A1.5 Errors in grid functions . . . . .	A–6
A1.5.1 Norm equivalence . . . . .	A–7
<b>A2 Estimating errors in numerical solutions</b>	<b>A–9</b>
A2.1 Estimates from the true solution . . . . .	A–10
A2.2 Estimates from a fine-grid solution . . . . .	A–10
A2.3 Estimates from coarser solutions . . . . .	A–11
A2.4 Extrapolation methods . . . . .	A–12
<b>A3 Convergence of iterative methods</b>	<b>A–13</b>
<b>A4 Some Matrix Properties</b>	<b>A–17</b>
A4.1 Adjoints and symmetry . . . . .	A–17
A4.2 Positive definiteness . . . . .	A–19
A4.3 The maximum principle . . . . .	A–20
A4.4 Diagonal dominance . . . . .	A–20
A4.5 Quadratic forms . . . . .	A–21





**Part I**  
**Basic Text**



# Chapter 1

## Finite difference approximations

Our goal is to approximate solutions to differential equations, *i.e.*, to find a function (or some discrete approximation to this function) which satisfies a given relationship between various of its derivatives on some given region of space and/or time, along with some boundary conditions along the edges of this domain. In general this is a difficult problem. A finite difference method proceeds by replacing the derivatives in the differential equations by finite difference approximations, which gives a large algebraic system of equations to be solved in place of the differential equation.

Before tackling this problem, we first consider the more basic question of how we can approximate the derivatives of a known function by finite difference formulas based only on values of the function itself at discrete points. Besides providing a basis for the later development of finite difference methods for solving differential equations, this allows us to investigate several key concepts such as the *order of accuracy* of an approximation in the simplest possible setting.

Let  $u(x)$  represent a function of one variable that, unless otherwise stated, will always be assumed to be smooth, meaning that we can differentiate the function several times and each derivative is a well-defined bounded function over an interval containing a particular point of interest  $\bar{x}$ .

Suppose we want to approximate  $u'(\bar{x})$  by a finite difference approximation based only on values of  $u$  at a finite number of points near  $\bar{x}$ . One obvious choice would be to use

$$D_+u(\bar{x}) \equiv \frac{u(\bar{x} + h) - u(\bar{x})}{h} \quad (1.1)$$

for some small value of  $h$ . This is motivated by the standard definition of the derivative as the limiting value of this expression as  $h \rightarrow 0$ . Note that  $D_+u(\bar{x})$  is the slope of the tangent line interpolating  $u$  at the points  $\bar{x}$  and  $\bar{x} + h$  (see Figure 1.1).

The expression (1.1) is a *one-sided* approximation to  $u'$  since  $u$  is evaluated only at values of  $x \geq \bar{x}$ . Another one-sided approximation would be

$$D_-u(\bar{x}) \equiv \frac{u(\bar{x}) - u(\bar{x} - h)}{h}. \quad (1.2)$$

Each of these formulas gives a *first order accurate* approximation to  $u'(\bar{x})$ , meaning that the size of the error is roughly proportional to  $h$  itself.

Another possibility is to use the *centered approximation*

$$D_0u(\bar{x}) \equiv \frac{u(\bar{x} + h) - u(\bar{x} - h)}{2h} = \frac{1}{2}(D_+u(\bar{x}) + D_-u(\bar{x})). \quad (1.3)$$

This is the slope of the tangent line interpolating  $u$  at  $\bar{x} - h$  and  $\bar{x} + h$ , and is simply the average of the two one-sided approximations defined above. From Figure 1.1 it should be clear that we would expect  $D_0u(\bar{x})$  to give a better approximation than either of the one-sided approximations. In fact this gives a

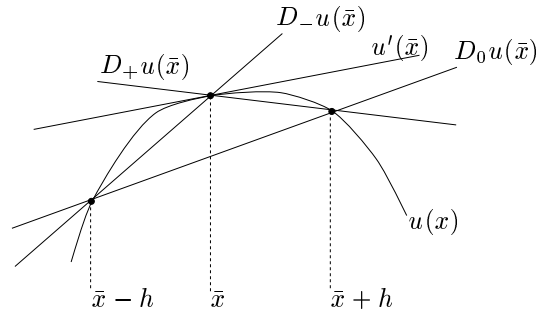


Figure 1.1: Various approximations to  $u'(\bar{x})$  interpreted as the slope of secant lines.

Table 1.1: Errors in various finite difference approximations to  $u'(\bar{x})$ .

h	D+	D-	D0	D3
1.0000e-01	-4.2939e-02	4.1138e-02	-9.0005e-04	6.8207e-05
5.0000e-02	-2.1257e-02	2.0807e-02	-2.2510e-04	8.6491e-06
2.5000e-02	-1.0574e-02	1.0462e-02	-5.6280e-05	1.0885e-06
1.2500e-02	-5.2732e-03	5.2451e-03	-1.4070e-05	1.3651e-07

*second order accurate* approximation — the error is proportional to  $h^2$  and hence is much smaller than the error in a first order approximation when  $h$  is small.

Other approximations are also possible, for example

$$D_3u(\bar{x}) \equiv \frac{1}{6h}[2u(\bar{x} + h) + 3u(\bar{x}) - 6u(\bar{x} - h) + u(\bar{x} - 2h)]. \quad (1.4)$$

It may not be clear where this came from or why it should approximate  $u'$  at all, but in fact it turns out to be a third order accurate approximation — the error is proportional to  $h^3$ .

Our first goal is to develop systematic ways to derive such formulas and to analyze their accuracy and relative worth. First we will look at a typical example of how the errors in these formulas compare.

**Example 1.1.** Let  $u(x) = \sin(x)$  and  $\bar{x} = 1$ , so we are trying to approximate  $u'(1) = \cos(1) = 0.5403023$ . Table 1.1 shows the error  $Du(\bar{x}) - u'(\bar{x})$  for various values of  $h$  for each of the formulas above.

We see that  $D_+u$  and  $D_-u$  behave similarly though one exhibits an error that is roughly the negative of the other. This is reasonable from Figure 1.1 and explains why  $D_0u$ , the average of the two, has an error that is much smaller than either.

We see that

$$\begin{aligned} D_+u(\bar{x}) &\approx -0.42h \\ D_0u(\bar{x}) &\approx -0.09h^2 \\ D_3u(\bar{x}) &\approx 0.007h^3 \end{aligned}$$

confirming that these methods are first order, second order, and third order, respectively.

Figure 1.2 shows these errors plotted against  $h$  on a log-log scale. This is a good way to plot errors when we expect them to behave like some power of  $h$ , since if the error  $E(h)$  behaves like

$$E(h) \approx Ch^p$$

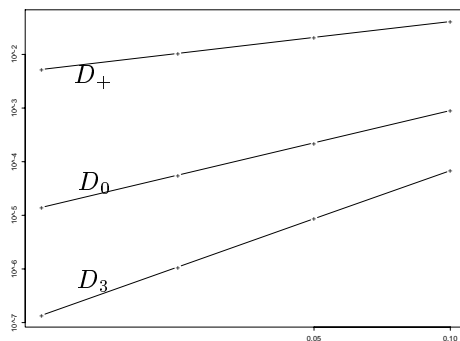


Figure 1.2: The errors in  $Du(\bar{x})$  plotted against  $h$  on a log-log scale.

then

$$\log E(h) \approx \log C + p \log h.$$

So on a log-log scale the error behaves linearly with a slope that is equal to  $p$ , the order of accuracy.

## 1.1 Truncation errors

The standard approach to analyzing the error in a finite difference approximation is to expand each of the function values of  $u$  in a *Taylor series* about the point  $\bar{x}$ , e.g.,

$$u(\bar{x} + h) = u(\bar{x}) + hu'(\bar{x}) + \frac{1}{2}h^2u''(\bar{x}) + \frac{1}{6}h^3u'''(\bar{x}) + O(h^4) \quad (1.5a)$$

$$u(\bar{x} - h) = u(\bar{x}) - hu'(\bar{x}) + \frac{1}{2}h^2u''(\bar{x}) - \frac{1}{6}h^3u'''(\bar{x}) + O(h^4) \quad (1.5b)$$

These expansions are valid provided that  $u$  is sufficiently smooth. Readers unfamiliar with the “big-oh” notation  $O(h^4)$  are advised to read Appendix A1 at this point since this notation will be heavily used and a proper understanding of its use is critical.

Using (1.5a) allows us to compute that

$$D_+u(\bar{x}) = \frac{u(\bar{x} + h) - u(\bar{x})}{h} = u'(\bar{x}) + \frac{1}{2}hu''(\bar{x}) + \frac{1}{6}h^2u'''(\bar{x}) + O(h^3).$$

Recall that  $\bar{x}$  is a fixed point so that  $u''(\bar{x})$ ,  $u'''(\bar{x})$ , etc., are fixed constants independent of  $h$ . They depend on  $u$  of course, but the function is also fixed as we vary  $h$ .

For  $h$  sufficiently small, the error will be dominated by the first term  $\frac{1}{2}hu''(\bar{x})$  and all the other terms will be negligible compared to this term, so we expect the error to behave roughly like a constant times  $h$ , where the constant has the value  $\frac{1}{2}u''(\bar{x})$ .

Note that in Example 1.1, where  $u(x) = \sin x$ , we have  $\frac{1}{2}u''(1) = -0.4207355$  which agrees with the behavior seen in Table 1.1.

Similarly, from (1.5b) we can compute that the error in  $D_-u(\bar{x})$  is

$$D_-u(\bar{x}) - u'(\bar{x}) = -\frac{1}{2}hu''(\bar{x}) + \frac{1}{6}h^2u'''(\bar{x}) + O(h^3)$$

which also agrees with our expectations.

Combining (1.5a) and (1.5b) shows that

$$u(\bar{x} + h) - u(\bar{x} - h) = 2hu'(\bar{x}) + \frac{1}{3}h^3u'''(\bar{x}) + O(h^5)$$

so that

$$D_0u(\bar{x}) - u'(\bar{x}) = \frac{1}{6}h^2u'''(\bar{x}) + O(h^4). \quad (1.6)$$

This confirms the second order accuracy of this approximation and again agrees with what is seen in Table 1.1, since in the context of Example 1.1 we have

$$\frac{1}{6}u'''(\bar{x}) = -\frac{1}{6}\cos(1) = -0.09005038.$$

Note that all of the odd order terms drop out of the Taylor series expansion (1.6) for  $D_0u(\bar{x})$ . This is typical with *centered* approximations and typically leads to a higher order approximation.

In order to analyze  $D_3u$  we need to also expand  $u(\bar{x} - 2h)$  as

$$u(\bar{x} - 2h) = u(\bar{x}) - 2hu'(\bar{x}) + \frac{1}{2}(2h)^2u''(\bar{x}) - \frac{1}{6}(2h)^3u'''(\bar{x}) + O(h^4). \quad (1.7)$$

Combining this with (1.5a) and (1.5b) shows that

$$D_3u(\bar{x}) = u'(\bar{x}) + \frac{1}{12}h^3u'''(\bar{x}) + O(h^4). \quad (1.8)$$

## 1.2 Deriving finite difference approximations

Suppose we want to derive a finite difference approximation to  $u'(\bar{x})$  based on some given set of points. We can use Taylor series to derive an appropriate formula, using the *method of undetermined coefficients*.

**Example 1.2.** Suppose we want a one-sided approximation to  $u'(\bar{x})$  based on  $u(\bar{x})$ ,  $u(\bar{x} - h)$  and  $u(\bar{x} - 2h)$ , of the form

$$D_2u(\bar{x}) = au(\bar{x}) + bu(\bar{x} - h) + cu(\bar{x} - 2h). \quad (1.9)$$

We can determine the coefficients  $a$ ,  $b$ , and  $c$  to give the best possible accuracy by expanding in Taylor series and collecting terms. Using (1.5b) and (1.7) in (1.9) gives

$$\begin{aligned} D_2u(\bar{x}) &= (a + b + c)u(\bar{x}) - (b + 2c)hu'(\bar{x}) + \frac{1}{2}(b + 4c)h^2u''(\bar{x}) \\ &\quad - \frac{1}{6}(b + 8c)h^3u'''(\bar{x}) + \dots \end{aligned}$$

If this is going to agree with  $u'(\bar{x})$  to high order then we need

$$\begin{aligned} a + b + c &= 0 \\ b + 2c &= -1/h \\ b + 4c &= 0 \end{aligned} \quad (1.10)$$

We might like to require that higher order coefficients be zero as well, but since there are only three unknowns  $a$ ,  $b$ , and  $c$  we cannot in general hope to satisfy more than three such conditions. Solving the linear system (1.10) gives

$$a = 3/2h \quad b = -2/h \quad c = 1/2h$$

so that the formula is

$$D_2u(\bar{x}) = \frac{1}{2h}[3u(\bar{x}) - 4u(\bar{x} - h) + u(\bar{x} - 2h)]. \quad (1.11)$$

The error in this approximation is clearly

$$\begin{aligned} D_2u(\bar{x}) - u'(\bar{x}) &= -\frac{1}{6}(b + 8c)h^3u'''(\bar{x}) + \dots \\ &= \frac{1}{12}h^2u'''(\bar{x}) + O(h^3). \end{aligned}$$

### 1.3 Polynomial interpolation

There are other ways to derive the same approximations. One way is to approximate the function  $u(x)$  by some polynomial  $p(x)$  and then use  $p'(\bar{x})$  as an approximation to  $u'(\bar{x})$ . If we determine the polynomial by interpolating  $u$  at an appropriate set of points, then we obtain the same finite difference methods as above.

**Example 1.3.** To derive the method of Example 1.2 in this way, let  $p(x)$  be the quadratic polynomial that interpolates  $u$  at  $\bar{x}$ ,  $\bar{x} - h$  and  $\bar{x} - 2h$  and then compute  $p'(\bar{x})$ . The result is exactly (1.11).

### 1.4 Second order derivatives

Approximations to the second derivative  $u''(x)$  can be obtained in an analogous manner. The standard second order centered approximation is given by

$$\begin{aligned} D^2u(\bar{x}) &= \frac{1}{h^2}[u(\bar{x} - h) - 2u(\bar{x}) + u(\bar{x} + h)] \\ &= u''(\bar{x}) + \frac{1}{2}h^2u''''(\bar{x}) + O(h^4). \end{aligned}$$

Again, since this is a symmetric centered approximation all of the odd order terms drop out. This approximation can also be obtained by the method of undetermined coefficients, or alternatively by computing the second derivative of the quadratic polynomial interpolating  $u(x)$  at  $\bar{x} - h$ ,  $\bar{x}$  and  $\bar{x} + h$ .

Another way to derive approximations to higher order derivatives is by repeatedly applying first order differences. Just as the second derivative is the derivative of  $u'$ , we can view  $D^2u(\bar{x})$  as being a difference of first differences. In fact,

$$D^2u(\bar{x}) = D_+D_-u(\bar{x})$$

since

$$\begin{aligned} D_+(D_-u(\bar{x})) &= \frac{1}{h}[D_-u(\bar{x} + h) - D_-u(\bar{x})] \\ &= \frac{1}{h} \left[ \left( \frac{u(\bar{x} + h) - u(\bar{x})}{h} \right) - \left( \frac{u(\bar{x}) - u(\bar{x} - h)}{h} \right) \right] \\ &= D^2u(\bar{x}). \end{aligned}$$

Alternatively,  $D^2(\bar{x}) = D_-D_+u(\bar{x})$  or we can also view it as a centered difference of centered differences, if we use a step size  $h/2$  in each centered approximation to the first derivative. If we define

$$\hat{D}_0u(x) = \frac{1}{h}(u(x + h/2) - u(x - h/2))$$

then we find that

$$\hat{D}_0(\hat{D}_0u(\bar{x})) = \frac{1}{h} \left( \left( \frac{u(\bar{x} + h) - u(\bar{x})}{h} \right) - \left( \frac{u(\bar{x}) - u(\bar{x} - h)}{h} \right) \right) = D^2u(\bar{x}).$$

### 1.5 Higher order derivatives

Finite difference approximations to higher order derivatives can also be obtained using any of the approaches outlined above. Repeatedly differencing approximations to lower order derivatives is a particularly simple way.



**Example 1.4.** As an example, here are two different approximations to  $u'''(\bar{x})$ . The first one is uncentered and first order accurate:

$$\begin{aligned}D_+ D^2 u(\bar{x}) &= \frac{1}{h^3}(u(\bar{x} + 2h) - 3u(\bar{x} + h) + 3u(\bar{x}) - u(\bar{x} - h)) \\ &= u'''(\bar{x}) + \frac{1}{2}hu''''(\bar{x}) + O(h^2).\end{aligned}$$

The next approximation is centered and second order accurate:

$$\begin{aligned}D_0 D_+ D_- u(\bar{x}) &= \frac{1}{2h^3}(u(\bar{x} + 2h) - 2u(\bar{x} + h) + 2u(\bar{x} - h) - u(\bar{x} - 2h)) \\ &= u'''(\bar{x}) + \frac{1}{4}h^2 u''''(\bar{x}) + O(h^4).\end{aligned}$$

Finite difference approximations of the sort derived above are the basic building blocks of finite difference methods for solving differential equations.

## Chapter 2

# Boundary Value Problems

We will first consider ordinary differential equations that are posed on some interval  $a < x < b$ , together with some boundary conditions at each end of the interval. In the next chapter we will extend this to more than one space dimension, and study *elliptic partial differential equations* that are posed in some region of the plane or three-dimensional space, and are solved subject to some boundary conditions specifying the solution and/or its derivatives around the boundary of the region. The problems considered in these two chapters are generally *steady state* problems in which the solution varies only with the spatial coordinates but not with time. (But see Section 2.15 for a case where  $[a, b]$  is a time interval rather than an interval in space.)

Steady-state problems are generally associated with some time-dependent problem that describes the dynamic behavior, and the 2-point boundary value problem or elliptic equation results from considering the special case where the solution is steady in time, and hence the time-derivative terms are equal to zero, simplifying the equations.

### 2.1 The heat equation

As a specific example, consider the flow of heat in a rod made out of some heat-conducting material, subject to some external heat source along its length and some boundary conditions at each end. If we assume that the material properties, the initial temperature distribution, and the source vary only with  $x$ , the distance along the length, and not across any cross-section, then we expect the temperature distribution at any time to vary only with  $x$  and we can model this with a differential equation in one space dimension. Since the solution might vary with time, we let  $u(x, t)$  denote the temperature at point  $x$  at time  $t$ , where  $a < x < b$  along some finite length of the rod. The solution is then governed by the *heat equation*

$$u_t(x, t) = (\kappa(x)u_x(x, t))_x + \psi(x, t) \quad (2.1)$$

where  $\kappa(x)$  is the coefficient of heat conduction, which may vary with  $x$ , and  $\psi(x, t)$  is the heat source (or sink, if  $\psi < 0$ ). More generally (2.1) is often called the *diffusion equation* since it models diffusion processes more generally, and the diffusion of heat is just one example. It is assumed that the basic theory of this equation is familiar to the reader. See Appendix ?? for a brief derivation, or standard PDE books such as [Kev90] for a better introduction. In general it is extremely valuable to understand where the equation one is attempting to solve comes from, since a good understanding of the physics (or biology, or whatever) is generally essential in understanding the development and behavior of numerical methods for solving the equation.

## 2.2 Boundary conditions

If the material is homogeneous then  $\kappa(x) \equiv \kappa$  is independent of  $x$  and the heat equation (2.1) reduces to

$$u_t(x, t) = \kappa u_{xx}(x, t) + \psi(x, t). \quad (2.2)$$

Along with the equation we need initial conditions,

$$u(x, 0) = u^0(x),$$

and boundary conditions, for example the temperature might be specified at each end,

$$u(a, t) = \alpha(t), \quad u(b, t) = \beta(t). \quad (2.3)$$

Such boundary conditions, where the value of the solution itself is specified, are called *Dirichlet boundary conditions*. Alternatively, one or both ends might be insulated, in which case there is zero heat flux at that end and so  $u_x = 0$  at that point. This boundary condition, which is a condition on the derivative of  $u$  rather than on  $u$  itself, is called a *Neumann boundary condition*. To begin with we will consider the Dirichlet problem for equation (2.2), with boundary conditions (2.3).

## 2.3 The steady-state problem

In general we expect the temperature distribution to change with time. However, if  $\psi(x, t)$ ,  $\alpha(t)$ , and  $\beta(t)$  are all time-independent, then we might expect the solution to eventually reach a *steady-state* solution  $u(x)$  which then remains essentially unchanged at later times. Typically there will be an initial *transient* time, as the initial data  $u^0(x)$  approaches  $u(x)$  (unless  $u^0(x) \equiv u(x)$ ), but if we are only interested in computing the steady state solution itself, then we can set  $u_t = 0$  in (2.2) and obtain an ordinary differential equation in  $x$  to solve for  $u(x)$ :

$$u''(x) = f(x) \quad (2.4)$$

where we introduce  $f(x) = -\psi(x)$  to avoid minus signs below. This is a second order ODE and from basic theory we expect to need two boundary conditions in order to specify a unique solution. In our case we have the boundary conditions

$$u(a) = \alpha, \quad u(b) = \beta. \quad (2.5)$$

**Remark 2.1** *Having two boundary conditions does not necessarily guarantee there exists a unique solution for a general second order equation — see Section 2.13.*

The problem (2.4), (2.5) is called a *two-point boundary value problem* since one condition is specified at each of the two endpoints of the interval where the solution is desired. If instead we had 2 data values specified at the same point, say  $u(a) = \alpha, u'(a) = \sigma$ , then we would have an *initial value problem* instead. These problems are discussed in Chapter 6.

One approach to computing a numerical solution to a steady state problem is to choose some initial data and march forward in time using a numerical method for the time-dependent equation (2.2), as discussed in Section ???. However, this is typically not an efficient way to compute the steady-state solution if this is all we want. Instead we can discretize and solve the two-point boundary value problem given by (2.4) and (2.5) directly. This is the first boundary value problem that we will study in detail, starting in the next section. Note that this problem is sufficiently simple... that we can solve it explicitly (integrate  $f(x)$  twice and choose the two constants of integration so that the boundary conditions are satisfied), but studying finite difference methods for this simple equation will reveal some of the essential features of all such analysis, particularly the relation of the global error to the local truncation error and the use of stability in making this connection. Later in this chapter we will consider some other boundary value problems, including more challenging nonlinear equations.



First we must clarify what we mean by the error in the discrete values  $U_1, \dots, U_m$  relative to the true solution  $u(x)$ , which is a function. Since  $U_j$  is supposed to approximate  $u(x_j)$ , it is natural to use the pointwise errors  $U_j - u(x_j)$ . If we let  $\hat{U}$  be the vector of true values

$$\hat{U} = \begin{bmatrix} u(x_1) \\ u(x_2) \\ \vdots \\ u(x_m) \end{bmatrix} \quad (2.11)$$

then the error vector  $E$  defined by

$$E = U - \hat{U}$$

contains the errors at each grid point.

Our goal is now to obtain a bound on the magnitude of this vector, showing that it is  $O(h^2)$  as  $h \rightarrow 0$ . To measure the magnitude of this vector we must use some *norm*, for example the max-norm

$$\|E\|_\infty = \max_{1 \leq j \leq m} |E_j| = \max_{1 \leq j \leq m} |U_j - u(x_j)|.$$

This is just the largest error over the interval. If we can show that  $\|E\|_\infty = O(h^2)$  then it follows that each pointwise error must be  $O(h^2)$  as well.

Other norms are often used to measure grid functions, either because they are more appropriate for a given problem or simply because they are easier to bound since some mathematical techniques work only with a particular norm. Other norms that are frequently used include the 1-norm

$$\|E\|_1 = h \sum_{j=1}^m |E_j|$$

and the 2-norm

$$\|E\|_2 = \left( h \sum_{j=1}^m |E_j|^2 \right)^{1/2}.$$

Note the factor of  $h$  that appears in these definitions. See Appendix A1 for a more thorough discussion of grid function norms and how they relate to standard vector norms.

Now let's return to the problem of estimating the error in our finite difference solution to the boundary value problem obtained by solving the system (2.9). The technique we will use is absolutely basic to the analysis of finite difference methods in general. It involves two key steps. We first compute the *local truncation error* of the method and then use some form of *stability* to show that the *global error* can be bounded in terms of the local truncation error.

The global error simply refers to the error  $U - \hat{U}$  that we are attempting to bound. The local truncation error (LTE) refers to the error in our finite difference approximation of derivatives, and hence is something that can be easily estimated using Taylor series expansions as we have seen in Chapter 1. Stability is the magic ingredient that allows us to go from these easily computed bounds to the estimates we really want for the global error. Let's look at each of these in turn.

## 2.5 Local truncation error

The LTE is defined by replacing  $U_j$  by the true solution  $u(x_j)$  in the finite difference formula (2.8). Of course the true solution  $u(x_j)$  won't satisfy this equation exactly and the discrepancy is the LTE, which we denote by  $\tau_j$ :

$$\tau_j = \frac{1}{h^2}(u(x_{j-1}) - 2u(x_j) + u(x_{j+1})) - f(x_j) \quad (2.12)$$

for  $j = 1, 2, \dots, m$ . Of course in practice we don't know what the true solution  $u(x)$  is, but if we assume it is smooth then by the Taylor series expansions (1.5a) we know that

$$\tau_j = \left[ u''(x_j) + \frac{1}{12}h^2 u''''(x_j) + O(h^4) \right] - f(x_j). \quad (2.13)$$

Using our original differential equation (2.6) this becomes

$$\tau_j = \frac{1}{12}h^2 u''''(x_j) + O(h^4).$$

Although  $u''''$  is in general unknown, it is some fixed function independent of  $h$  and so  $\tau_j = O(h^2)$  as  $h \rightarrow 0$ .

If we define  $\tau$  to be the vector with components  $\tau_j$ , then

$$\tau = A\hat{U} - F$$

where  $\hat{U}$  is the vector of true solution values (2.11), and so

$$A\hat{U} = F + \tau. \quad (2.14)$$

## 2.6 Global error

To obtain a relation between the local error  $\tau$  and the global error  $E = U - \hat{U}$ , we subtract the equation (2.14) from the equation (2.9) that defines  $U$ , obtaining

$$AE = -\tau. \quad (2.15)$$

This is simply the matrix form of the system of equations

$$\frac{1}{h^2}(E_{j-1} - 2E_j + E_{j+1}) = -\tau(x_j) \text{ for } j = 1, 2, \dots, m.$$

with the boundary conditions

$$E_0 = E_{m+1} = 0$$

since we are using the exact boundary data  $U_0 = \alpha$  and  $U_{m+1} = \beta$ . We see that the global error satisfies a set of finite difference equations that has exactly the same form as our original difference equations for  $U$ , but with the right hand side given by  $-\tau$  rather than  $F$ .

From this it should be clear why we expect the global error to be roughly the same magnitude as the local error  $\tau$ . We can interpret the system (2.15) as a discretization of the ODE

$$e''(x) = -\tau(x) \text{ for } 0 < x < 1 \quad (2.16)$$

with boundary conditions

$$e(0) = 0, \quad e(1) = 0.$$

Since  $\tau(x) \approx \frac{1}{12}h^2 u''''(x)$ , integrating twice shows that the error should be roughly

$$e(x) \approx -\frac{1}{2}h^2 u''(x) + \frac{1}{2}h^2 (u''(0) + x(u''(1) - u''(0)))$$

and hence the error should be  $O(h^2)$ .

## 2.7 Stability

The above argument is not completely convincing because we are relying on the assumption that solving the difference equations gives a decent approximation to the solution of the underlying differential equations. (Actually the converse now, that the solution to the differential equation (2.16) gives a good indication of the solution to the difference equations (2.15).) Since it is exactly this assumption we are trying to prove, the reasoning is rather circular.

Instead, let's look directly at the discrete system (2.15) which we will rewrite in the form

$$A^h E^h = -\tau^h \quad (2.17)$$

where the superscript  $h$  indicates that we are on a grid with mesh spacing  $h$ . This serves as a reminder that these quantities change as we refine the grid. In particular, the matrix  $A^h$  is an  $m \times m$  matrix with  $h = 1/(m + 1)$  so that its dimension is growing as  $h \rightarrow 0$ .

Let  $(A^h)^{-1}$  be the inverse of this matrix. Then solving the system (2.17) gives

$$E^h = -(A^h)^{-1} \tau^h$$

and taking norms gives

$$\begin{aligned} \|E^h\| &= \|(A^h)^{-1} \tau^h\| \\ &\leq \|(A^h)^{-1}\| \|\tau^h\|. \end{aligned}$$

We know that  $\|\tau^h\| = O(h^2)$  and we are hoping the same will be true of  $\|E^h\|$ . It is clear what we need for this to be true: we need  $\|(A^h)^{-1}\|$  to be bounded by some constant independent of  $h$  as  $h \rightarrow 0$ :

$$\|(A^h)^{-1}\| \leq C \text{ for all } h \text{ sufficiently small.}$$

Then we will have

$$\|E^h\| \leq C \|\tau^h\| \quad (2.18)$$

and so  $\|E^h\|$  goes to zero at least as fast as  $\|\tau^h\|$ . This motivates the following definition of *stability* for linear boundary value problems.

**Definition 2.7.1** *Suppose a finite difference method for a linear boundary value problem gives a sequence of matrix equations of the form  $A^h U^h = F^h$  where  $h$  is the mesh width. We say that the method is stable if  $(A^h)^{-1}$  exists for all  $h$  sufficiently small (for  $h < h_0$ , say) and if there is a constant  $C$ , independent of  $h$ , such that*

$$\|(A^h)^{-1}\| \leq C \text{ for all } h < h_0. \quad (2.19)$$

## 2.8 Consistency

We say that a method is *consistent* with the differential equation and boundary conditions if

$$\|\tau^h\| \rightarrow 0 \text{ as } h \rightarrow 0. \quad (2.20)$$

This simply says that we have a sensible discretization of the problem. Typically  $\|\tau^h\| = O(h^p)$  for some integer  $p > 0$ , and then the method is certainly consistent.

## 2.9 Convergence

A method is said to be *convergent* if  $\|E^h\| \rightarrow 0$  as  $h \rightarrow 0$ . Combining the ideas introduced above we arrive at the conclusion that

$$\text{consistency} + \text{stability} \implies \text{convergence}. \quad (2.21)$$

This is easily proved by using (2.19) and (2.20) to obtain the bound

$$\|E^h\| \leq \|(A^h)^{-1}\| \|\tau^h\| \leq C \|\tau^h\| \rightarrow 0 \text{ as } h \rightarrow 0.$$

Although this has been demonstrated only for the linear boundary value problem, in fact most analyses of finite difference methods for differential equations follow this same two-tier approach, and the statement (2.21) is sometimes called the *fundamental theorem of finite difference methods*. In fact, as our above analysis indicates, this can generally be strengthened to say that

$$O(h^p) \text{ local truncation error} + \text{stability} \implies O(h^p) \text{ global error}. \quad (2.22)$$

Consistency (and the order of accuracy) is usually the easy part to check. Verifying stability is the hard part. Even for the linear boundary value problem just discussed it is not at all clear how to check the condition (2.19) since these matrices get larger as  $h \rightarrow 0$ . For other problems it may not even be clear how to define stability in an appropriate way. As we will see, there are many different definitions of “stability” for different types of problems. The challenge in analyzing finite difference methods for new classes of problems is often to find an appropriate definition of “stability” that allows one to prove convergence using (2.21) while at the same time being sufficiently manageable that we can verify it holds for specific finite difference methods. For nonlinear PDEs this frequently must be tuned to each particular class of problems, and relies on existing mathematical theory and techniques of analysis for this class of problems.

Whether or not one has a formal proof of convergence for a given method, it is always good practice to check that the computer program is giving convergent behavior, at the rate expected. Appendix A2 contains a discussion of how the error in computed results can be estimated.

## 2.10 Stability in the 2-norm

Returning to the boundary value problem at the start of the chapter, let’s see how we can verify stability and hence second-order accuracy. The technique used depends on what norm we wish to consider. Here we will consider the 2-norm and see that we can show stability by explicitly computing the eigenvectors and eigenvalues of the matrix  $A$ . In Section 2.12 we show stability in the max-norm by different techniques.

Since the matrix  $A$  from (2.10) is symmetric, the 2-norm of  $A$  is equal to its spectral radius (see Appendix A1):

$$\|A\|_2 = \rho(A) = \max_{1 \leq p \leq m} |\lambda^p|.$$

The matrix  $A^{-1}$  is also symmetric and the eigenvalues of  $A^{-1}$  are simply the inverses of the eigenvalues of  $A$ , so

$$\|A^{-1}\|_2 = \rho(A^{-1}) = \max_{1 \leq p \leq m} |(\lambda^p)^{-1}| = \left( \min_{1 \leq p \leq m} |\lambda^p| \right)^{-1}.$$

So all we need to do is compute the eigenvalues of  $A$  and show that they are bounded away from zero as  $h \rightarrow 0$ . Of course we have an infinite set of matrices  $A^h$  to consider, as  $h$  varies, but since the structure of these matrices is so simple, we can obtain a general expression for the eigenvalues of each  $A^h$ . For more complicated problems we might not be able to do this, but it is worth going through in detail for this problem because one often considers model problems for which such analysis is possible. We will also need to know these eigenvalues for other purposes when we discuss parabolic equations later.



We will now focus on one particular value of  $h = 1/(m + 1)$  and drop the superscript  $h$  to simplify the notation. Then the  $m$  eigenvalues of  $A$  are given by

$$\lambda^p = \frac{2}{h^2}(\cos(p\pi h) - 1), \text{ for } p = 1, 2, \dots, m. \quad (2.23)$$

The eigenvector  $u^p$  corresponding to  $\lambda^p$  has components  $u_j^p$  for  $j = 1, 2, \dots, m$  given by

$$u_j^p = \sin(p\pi j h). \quad (2.24)$$

This can be verified by checking that  $Au^p = \lambda^p u^p$ . The  $j$ th component of the vector  $Au^p$  is

$$\begin{aligned} (Au^p)_j &= \frac{1}{h^2}(u_{j-1}^p - 2u_j^p + u_{j+1}^p) \\ &= \frac{1}{h^2}(\sin(p\pi(j-1)h) - 2\sin(p\pi j h) + \sin(p\pi(j+1)h)) \\ &= \frac{1}{h^2}(\sin(p\pi j h)\cos(p\pi h) - 2\sin(p\pi j h) + \sin(p\pi j h)\cos(p\pi h)) \\ &= \lambda^p u_j^p. \end{aligned}$$

Note that for  $j = 1$  and  $j = m$  the  $j$ th component of  $Au^p$  looks slightly different (the  $u_{j-1}^p$  or  $u_{j+1}^p$  term is missing) but that the above form and trigonometric manipulations are still valid provided that we define

$$u_0^p = u_{m+1}^p = 0,$$

as is consistent with (2.24). From (2.23) we see that the smallest eigenvalue of  $A$  (in magnitude) is

$$\begin{aligned} \lambda^1 &= \frac{2}{h^2}(\cos(\pi h) - 1) \\ &= \frac{2}{h^2}\left(-\frac{1}{2}\pi^2 h^2 + \frac{1}{24}\pi^4 h^4 + O(h^6)\right) \\ &= -\pi^2 + O(h^2). \end{aligned}$$

This is clearly bounded away from zero as  $h \rightarrow 0$ , and so we see that the method is stable in the 2-norm. Moreover we get an error bound from this:

$$\|E^h\|_2 \leq \|(A^h)^{-1}\|_2 \|\tau^h\|_2 \approx \frac{1}{\pi^2} \|\tau^h\|_2.$$

Since  $\tau_j^h \approx \frac{1}{2}h^2 u''''(x_j)$ , we expect  $\|\tau^h\|_2 \approx \frac{1}{2}h^2 \|u''''\|_2 = \frac{1}{2}h^2 \|f''\|_2$ , where now  $\|\cdot\|_2$  refers to the function space norm (A1.11).

Note that the eigenvector (2.24) is closely related to the eigenfunction of the corresponding differential operator  $\frac{\partial^2}{\partial x^2}$ . The functions

$$u^p(x) = \sin(p\pi x), \quad p = 1, 2, 3, \dots$$

satisfy the relation

$$\frac{\partial^2}{\partial x^2} u^p(x) = \mu^p u^p(x)$$

with eigenvalue  $\mu^p = -p^2\pi^2$ . These functions also satisfy  $u^p(0) = u^p(1) = 0$  and hence they are eigenfunctions of  $\frac{\partial^2}{\partial x^2}$  on  $[0, 1]$  with homogeneous boundary conditions. The discrete approximation to this operator given by the matrix  $A$  has only  $m$  eigenvalues instead of an infinite number, and the corresponding eigenvectors (2.24) are simply the first  $m$  eigenfunctions of  $\frac{\partial^2}{\partial x^2}$  evaluated at the grid

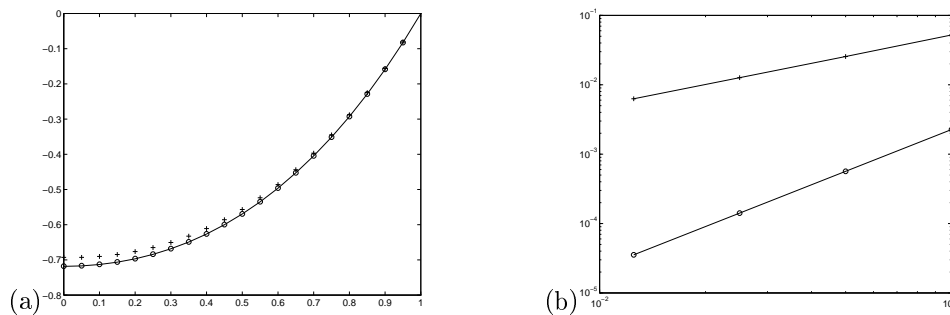


Figure 2.1: (a) Sample solution to the steady-state heat equation with a Neumann boundary condition at the left boundary. Solid line is the true solution. + shows solution on a grid with 20 points using (2.26). o shows the solution on the same grid using (2.28). (b) A log-log plot of the max-norm error as the grid is refined is also shown for each case.

points. The eigenvalue  $\lambda^p$  is not exactly the same as  $\mu^p$ , but at least for small values of  $p$  it is very nearly the same, since by Taylor series expansion of the cosine in (2.23) gives

$$\begin{aligned}\lambda^p &= \frac{2}{h^2} \left( -\frac{1}{2}p^2\pi^2 h^2 + \frac{1}{24}p^4\pi^4 h^4 + \dots \right) \\ &= -p^2\pi^2 + O(h^2) \quad \text{as } h \rightarrow 0 \text{ for } p \text{ fixed.}\end{aligned}$$

This relationship will be illustrated further when we study numerical methods for the heat equation (2.1).

## 2.11 Neumann boundary conditions

Suppose now that we have one or more Neumann boundary conditions, instead of Dirichlet boundary conditions, meaning that a boundary condition on the derivative  $u'$  is given rather than a condition on the value of  $u$  itself. For example, in our heat conduction example we might have one end of the rod insulated so that there is no heat flux through this end and hence  $u' = 0$  there. More generally we might have heat flux at a specified rate giving  $u' = \sigma$  at this boundary.

We first consider the equation (2.4) with boundary conditions

$$u'(0) = \sigma, \quad u(1) = \beta. \quad (2.25)$$

Figure 2.1 shows the solution to this problem with  $f(x) = e^x$ ,  $\sigma = 0$ , and  $\beta = 0$  as one example.

**Exercise 2.1** Determine the function shown in Figure 2.1 by solving the problem exactly.

To solve this problem numerically, we need to introduce one more unknown than we previously had:  $U_0$  at the point  $x_0 = 0$  since this is now an unknown value. We also need to augment the system (2.9) with one more equation that models the boundary condition (2.25).

**First attempt.** As a first try, we might use a one-sided expression for  $u'(0)$ , such as

$$\frac{U_1 - U_0}{h} = \sigma. \quad (2.26)$$

If we append this equation to the system (2.9), we obtain the following system of equations for the

unknowns  $U_0, U_1, \dots, U_m$ :

$$\frac{1}{h^2} \begin{bmatrix} -1 & 1 & & & & & \\ 1 & -2 & 1 & & & & \\ & 1 & -2 & 1 & & & \\ & & 1 & -2 & 1 & & \\ & & & \ddots & \ddots & \ddots & \\ & & & & 1 & -2 & 1 \\ & & & & & 1 & -2 \end{bmatrix} \begin{bmatrix} U_0 \\ U_1 \\ U_2 \\ U_3 \\ \vdots \\ U_{m-1} \\ U_m \end{bmatrix} = \begin{bmatrix} \sigma/h \\ f(x_1) \\ f(x_2) \\ f(x_3) \\ \vdots \\ f(x_{m-1}) \\ f(x_m) - \beta/h^2 \end{bmatrix}. \quad (2.27)$$

Solving this system of equations does give an approximation to the true solution (see Figure 2.1) but checking the errors shows that this is only first order accurate. Figure 2.1 also shows a log-log plot of the max-norm errors as we refine the grid. The problem is that the local truncation error of the approximation (2.26) is  $O(h)$ , since

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

This translates into a global error that is only  $O(h)$  as well.

REMARK: It is often possible to achieve second order accuracy even if the local truncation error is  $O(h)$  at a single point as long as it is  $O(h^2)$  everywhere else. But this is not true in the case we are now discussing. See Chapter ?? for more discussion of this.

**Exercise 2.2** Show that the method we have just discussed is stable in the 2-norm, following the technique of Section 2.10. Hint: The eigenvectors of the matrix in (2.27) are again obtained by simply discretizing the eigenfunctions of the corresponding differential operator  $\frac{\partial^2}{\partial x^2}$ , with boundary conditions  $u'(0) = u(1) = 0$ .

**Second attempt.** To obtain a second-order accurate method, we should use a centered approximation to  $u'(0) = \sigma$  instead of the one-sided approximation (2.26). We can introduce another unknown  $U_{-1}$  and instead of the single equation (2.26) use the following two equations:

$$\begin{aligned} \frac{1}{h^2}(U_{-1} - 2U_0 + U_1) &= f(x_0) \\ \frac{1}{2h}(U_1 - U_{-1}) &= \sigma \end{aligned} \quad (2.28)$$

This results in a system of  $m + 2$  equations. (What is the matrix?)

Introducing the unknown  $U_{-1}$  outside the interval  $[0, 1]$  where the original problem is posed may seem unsatisfactory. We can avoid this by eliminating the unknown  $U_{-1}$  from the two equations (2.28), resulting in a single equation that can be written as

$$\frac{1}{h}(-U_0 + U_1) = \sigma + \frac{h}{2}f(x_0). \quad (2.29)$$

We have now reduced the system to one with only  $m + 1$  equations for the unknowns  $U_0, U_1, \dots, U_m$ . The matrix is exactly the same as the matrix in (2.27), which came from the one-sided approximation. The only difference in the linear system is that the first element in the right hand side of (2.27) is now changed from  $\sigma/h$  to  $\sigma/h + \frac{1}{2}f(x_0)$ . We can interpret this as using the one-sided approximation to  $u'(0)$ , but with a modified value for this Neumann boundary condition that adjusts for the fact that the approximation has an  $O(h)$  error by introducing the same error in the data  $\sigma$ .

**Exercise 2.3** Show that the local truncation error  $\tau_0$  is now  $O(h^2)$ . Hint: Note that  $f(x_0) = u''(x_0)$ .

Combining the results of Exercise 2.2 and Exercise 2.3 shows that this method is second order accurate in the 2-norm.

## 2.12 Green's functions and max-norm stability

In Section 2.10 we demonstrated that  $A$  from (2.10) is stable in the 2-norm, and hence that  $\|E\|_2 = O(h^2)$ . Suppose, however, that we want a bound on the maximum error over the interval, i.e., a bound on  $\|E\|_\infty = \max |E_j|$ . We can obtain one such bound directly from the bound we have for the 2-norm. From (A1.16) we know that

$$\|E\|_\infty \leq \frac{1}{\sqrt{h}} \|E\|_2 = O(h^{3/2}) \text{ as } h \rightarrow 0.$$

However, this does not show the second order accuracy that we hope to have. To show that  $\|E\|_\infty = O(h^2)$  we will explicitly calculate the inverse of  $A$  and then show that  $\|A^{-1}\|_\infty = O(1)$ , and hence

$$\|E\|_\infty \leq \|A^{-1}\|_\infty \|\tau\|_\infty = O(h^2)$$

since  $\|\tau\|_\infty = O(h^2)$ . As in the computation of the eigenvalues in the last section, we can only do this because our model problem (2.6) is so simple. In general it would be impossible to obtain closed form expressions for the inverse of the matrices  $A^h$  as  $h$  varies. But again it is worth working out the details for this case because it gives a great deal of insight into the nature of the inverse matrix and what it represents more generally.

Each column of the inverse matrix can be interpreted as the solution of a particular boundary value problem. The columns are discrete approximations to the *Green's functions* that are commonly introduced in the study of the differential equation. An understanding of this is very valuable in developing an intuition for what happens if we introduce relatively large errors at a few points within the interval. Such difficulties arise frequently in practice, typically at the boundary or at an internal interface where there are discontinuities in the data or solution.

Let  $e_j \in \mathbb{R}^m$  be the  $j$ th coordinate vector or *unit vector* with the value 1 as its  $j$ th element and all other elements equal to 0. If  $B$  is any matrix then the vector  $Be_j$  is simply the  $j$ th column of the matrix  $B$ . So the  $j$ th column of  $A^{-1}$  is  $A^{-1}e_j$ . Let's call this vector  $v$  for the time being. Then  $v$  is the solution of the linear system

$$Av = e_j. \tag{2.30}$$

This can be viewed as an approximation to a boundary value problem of the form (2.6),(2.7) where  $\alpha = \beta = 0$  and  $f(x_i) = 0$  unless  $i = j$ , with  $f(x_j) = 1$ . This may seem like a rather strange function  $f$ , but it corresponds to the delta function that is used in defining Green's functions (or more exactly to a delta function scaled by  $h$ ). We will come back to the problem of determining the  $j$ th column of  $A^{-1}$  after a brief review of delta functions and Green's functions for the differential equation.

The delta function,  $\delta(x)$ , is not an ordinary function but rather the mathematical idealization of a sharply peaked function that is nonzero over an interval  $(-\epsilon, \epsilon)$  near the origin and has the property that

$$\int_{-\infty}^{\infty} \phi_\epsilon(x) dx = \int_{-\epsilon}^{\epsilon} \phi_\epsilon(x) dx = 1. \tag{2.31}$$

The exact shape of  $\phi_\epsilon$  is not important, but note that it must attain a height that is  $O(1/\epsilon)$  in order for the integral to have the value 1. We can think of the delta function as being a sort of limiting case of such functions as  $\epsilon \rightarrow 0$ . These functions naturally arise when we differentiate functions that are discontinuous. For example, consider the *Heaviside function* (or step function)  $H(x)$  that is defined by

$$H(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0. \end{cases} \tag{2.32}$$

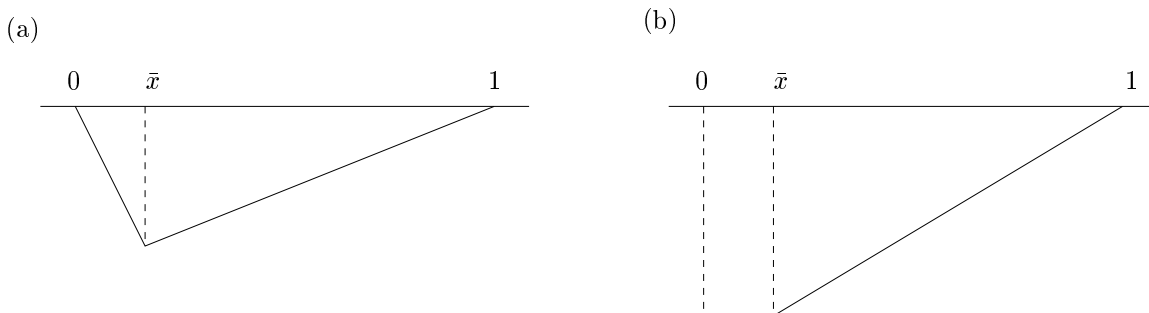


Figure 2.2: (a) The Green's function (2.33) for the Dirichlet problem. (b) The Green's function for the mixed problem with  $u'(0) = u(1) = 0$ .

What is the derivative of this function? For  $x \neq 0$  the function is constant and so  $H'(x) = 0$ . At  $x = 0$  the derivative is not defined in the classical sense. But if we smooth out the function a little bit, making it continuous and differentiable by changing  $H(x)$  only on the interval  $(-\epsilon, \epsilon)$ , then the new function  $H_\epsilon(x)$  is differentiable everywhere and has a derivative  $H'_\epsilon(x)$  that looks something like  $\phi_\epsilon(x)$ . The exact shape of  $H'_\epsilon(x)$  depends on how we choose  $H_\epsilon(x)$ , but note that regardless of its shape, its integral must be 1, since

$$\begin{aligned} \int_{-\infty}^{\infty} H'_\epsilon(x) dx &= \int_{-\epsilon}^{\epsilon} H'_\epsilon(x) dx \\ &= H_\epsilon(\epsilon) - H_\epsilon(-\epsilon) \\ &= 1 - 0 = 1. \end{aligned}$$

This explains the normalization (2.31). By letting  $\epsilon \rightarrow 0$ , we are led to define

$$H'(x) = \delta(x).$$

This expression makes no sense in terms of the classical definition of derivatives, but can be made rigorous mathematically through the use of “distribution theory”. For our purposes it suffices to think of the delta function as being a very sharply peaked function with total integral 1.

Now consider the function  $G(x)$  shown in Figure 2.2(a):

$$G(x) = \begin{cases} (\bar{x} - 1)x & \text{for } x \leq \bar{x} \\ \bar{x}(x - 1) & \text{for } x > \bar{x} \end{cases} \quad (2.33)$$

where  $\bar{x}$  is some point in the interval  $[0, 1]$ . The derivative of  $G(x)$  is piecewise constant:

$$G'(x) = \begin{cases} \bar{x} - 1 & \text{for } x < \bar{x} \\ \bar{x} & \text{for } x > \bar{x}. \end{cases}$$

If we define  $G'(\bar{x}) = \bar{x}$ , then we can write this as

$$G'(x) = \bar{x} - 1 + H(x - \bar{x}).$$

Differentiating again then gives

$$G''(x) = \delta(x - \bar{x}).$$

It follows that the function  $G(x)$  is the solution to the boundary value problem

$$u''(x) = \delta(x - \bar{x}) \text{ for } 0 < x < 1 \quad (2.34)$$

with homogeneous boundary conditions

$$u(0) = 0, \quad u(1) = 0. \quad (2.35)$$

This function is called the *Green's function* for the problem (2.6),(2.7) and is generally written as  $G(x; \bar{x})$  to show the dependence of the function  $G(x)$  on the parameter  $\bar{x}$ , the location of the delta function source term.

The solution to the boundary value problem (2.6),(2.7) for more general  $f(x)$  and boundary conditions can be written as

$$u(x) = \alpha(1-x) + \beta x + \int_0^1 G(x; \xi) f(\xi) d\xi. \quad (2.36)$$

This integral can be viewed as a linear combination of the functions  $G(x; \xi)$  at each point  $\xi$ , with weights given by the strength of the source term at each such point.

Returning now to the question of determining the columns of the matrix  $A^{-1}$  by solving the systems (2.30), we see that the right hand side  $e_j$  can be viewed as a discrete version of the delta function, scaled by  $h$ . So the system (2.30) is a discrete version of the problem

$$v''(x) = h\delta(x - x_j)$$

with homogeneous boundary conditions, whose solution is  $v(x) = hG(x; x_j)$ . We therefore expect the vector  $v$  to approximate this function. In fact it is easy to confirm that we can obtain the vector  $v$  by simply evaluating the function  $v(x)$  at each grid point, so  $v_i = hG(x_i; x_j)$ . This can be easily checked by verifying that multiplication by the matrix  $A$  gives the unit vector  $e_j$ .

If we now let  $G$  be the inverse matrix,  $G = A^{-1}$ , then the  $j$ th column of  $G$  is exactly the vector  $v$  found above, and so the elements of  $G$  are:

$$G_{ij} = \begin{cases} h(x_j - 1)x_i & i = 1, 2, \dots, j \\ h(x_i - 1)x_j & i = j, j + 1, \dots, m. \end{cases} \quad (2.37)$$

Note that each of the elements of  $G$  is bounded by  $h$  in magnitude. From this we obtain an explicit bound on the max-norm of  $G$ :

$$\|A^{-1}\|_{\infty} = \max_{1 \leq i \leq m} \sum_{j=1}^m |G_{ij}| \leq mh < 1.$$

This is uniformly bounded as  $h \rightarrow 0$  and so the method is stable in the max-norm. Since  $\|\tau\|_{\infty} = O(h^2)$ , the method is second order accurate in the max-norm and the pointwise error at each grid point is  $O(h^2)$ .

Note, by the way, that the representation (2.36) of the true solution  $u(x)$  as a superposition of the Green's functions  $G(x; \xi)$ , weighted by the values of the right hand side  $f(\xi)$ , has a direct analog for the solution  $U$  to the difference equation  $AU = F$ . The solution is  $U = A^{-1}F = Gf$ , which can be written as

$$U_i = \sum_{j=1}^m G_{ij} F_j.$$

Using the form of  $F$  from (2.10) and the expression (2.37) shows that

$$\begin{aligned} U_i &= -\frac{\alpha}{h^2} G_{i1} - \frac{\beta}{h^2} G_{im} + \sum_{j=2}^{m-1} G_{ij} f(x_j) \\ &= \alpha(1-x_i) + \beta x_i + h \sum_{j=2}^{m-1} G(x_i; x_j) f(x_j), \end{aligned}$$

which is simply the discrete form of (2.36).

**Exercise 2.4** Write out the  $4 \times 4$  matrix  $A$  and inverse matrix  $A^{-1}$  explicitly for the case  $h = 0.25$ . If  $f(x) = x$ , determine the discrete approximation to the solution of the boundary value problem on this grid and sketch this solution and the 4 Green's functions whose sum gives this solution.



Now suppose that we have a source term  $f(x)$  that is not identically zero, say  $f(x) < 0$  everywhere. Then we are constantly adding heat to the rod (recall that  $f = -\psi$ ). Since no heat can escape through the insulated ends, we expect the temperature to keep rising without bound. In this case we never reach a steady state, and the boundary value problem has no solution.

**Exercise 2.6** *What condition is required on the function  $f(x)$  to guarantee that a steady state exists for the original differential equation? What is the condition if  $\sigma_0$  and/or  $\sigma_1$  are nonzero? How do these conditions relate to the requirement that the right hand side of (2.38) must lie in the range of the matrix in order for the discrete approximation to have solutions?*

**Exercise 2.7** *Consider the following linear boundary value problem with Dirichlet boundary conditions:*

$$\begin{aligned} u''(x) + u(x) &= 0 \quad \text{for } 0 < x < \pi \\ u(0) &= \alpha, \quad u(\pi) = \beta. \end{aligned} \tag{2.40}$$

*For what values of  $\alpha$  and  $\beta$  does this have solutions? Sketch a family of solutions in a case where there are infinitely many solutions.*

## 2.14 A general linear second order equation

We now consider the more general linear equation

$$a(x)u''(x) + b(x)u'(x) + c(x)u(x) = f(x) \tag{2.41}$$

together with two boundary conditions, say the Dirichlet conditions

$$u(a) = \alpha, \quad u(b) = \beta. \tag{2.42}$$

This equation can be discretized to second order by

$$a_i \left( \frac{U_{i-1} - 2U_i + U_{i+1}}{h^2} \right) + b_i \left( \frac{U_{i+1} - U_{i-1}}{2h} \right) + c_i U_i = f_i \tag{2.43}$$

where, for example,  $a_i = a(x_i)$ . This gives the linear system  $AU = F$  where  $A$  is the tridiagonal matrix

$$A = \frac{1}{h^2} \begin{bmatrix} (h^2 c_1 - 2a_1) & (a_1 + hb_1/2) & & & \\ (a_2 - hb_2/2) & (h^2 c_2 - 2a_2) & (a_2 + hb_2/2) & & \\ & \ddots & & \ddots & \\ & & (a_{m-1} - hb_{m-1}/2) & (h^2 c_{m-1} - 2a_{m-1}) & (a_{m-1} + hb_{m-1}/2) \\ & & & (a_m - hb_m/2) & (h^2 c_m - 2a_m) \end{bmatrix} \tag{2.44}$$

and

$$U = \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_{m-1} \\ U_m \end{bmatrix}, \quad F = \begin{bmatrix} f_1 - (a_1/h^2 + b_1/2h)\alpha \\ f_2 \\ \vdots \\ f_{m-1} \\ f_m - (a_m/h^2 - b_m/2h)\beta \end{bmatrix}. \tag{2.45}$$

This linear system can be solved with standard techniques, assuming the matrix is nonsingular. A singular matrix would be a sign that the discrete system does not have a unique solution, which may occur if the original problem, or a nearby problem, is not well posed (see Section 2.13).





modeled by the discrete approximation to the differential equation. One can then show, for example, that the solution to the difference equations satisfies a *maximum principle* of the same type as the solution to the differential equation: for the homogeneous equation with  $f(x) \equiv 0$ , the values of  $u(x)$  lie between the values of the boundary values  $\alpha$  and  $\beta$  everywhere, so the maximum and minimum values of  $u$  arise on the boundaries. See Appendix A4 for more discussion.

When solving the resulting linear system by iterative methods (see Chapter 3 and Chapter 5) it is also often desirable that the matrix have properties such as negative definiteness, since some iterative methods (e.g., the conjugate-gradient method, Section 5.2) depend on such properties.

**Exercise 2.8** Show that the matrix  $A$  appearing in (2.50) is negative definite provided that  $\kappa > 0$  everywhere. Recall that  $A$  is negative definite if  $U^T A U < 0$  for all vectors  $U$  that are not identically zero.

**Exercise 2.9** Show that solutions to the homogeneous problem satisfy the maximum principle when the discretization (2.50) is used.

## 2.15 Nonlinear Equations

We next consider a nonlinear boundary value problem to illustrate the new complications that arise in this case. We will consider a specific example which has a simple physical interpretation that makes it easy to understand and interpret solutions. This example also illustrates that not all 2-point BVP's are steady-state problems,

Consider the motion of a pendulum with mass  $m$  at the end of a rigid (but massless) bar of length  $L$ , and let  $\theta(t)$  be the angle of the pendulum from vertical at time  $t$ , as illustrated in Figure 2.3. Ignoring the mass of the bar and forces of friction and air resistance, the differential equation for the pendulum motion can be well approximated by

$$\theta''(t) = -(g/L) \sin(\theta(t)) \quad (2.51)$$

where  $g$  is the gravitational constant. Taking  $g/L = 1$  for simplicity we have

$$\theta''(t) = -\sin(\theta(t)) \quad (2.52)$$

as our model problem.

For small amplitudes  $\theta$  it is possible to approximate  $\sin(\theta) \approx \theta$  and obtain the approximate *linear* differential equation

$$\theta''(t) = -\theta(t) \quad (2.53)$$

with general solutions of the form  $A \cos(t) + B \sin(t)$ . The motion of a pendulum that is oscillating only a small amount about the equilibrium at  $\theta = 0$  can be well approximated by this sinusoidal motion, which has period  $2\pi$  independent of the amplitude. For larger amplitude motions, however, solving (2.53) does not give good approximations to the true behavior. Figures 2.3(b) and (c) show some sample solutions to the two equations.

To fully describe the problem we also need to specify two auxiliary conditions in addition to the second-order differential equation (2.52). For the pendulum problem the *initial value problem* is most natural — we set the pendulum swinging from some initial position  $\theta(0)$  with some initial velocity  $\theta'(0)$ , which gives two initial conditions that are enough to determine a unique solution at all later times.

To obtain instead a BVP, consider the situation in which we wish to set the pendulum swinging from some initial given location  $\theta(0) = \alpha$  with some unknown velocity  $\theta'(0)$ , in such a way that the pendulum will be at a desired location  $\theta(T) = \beta$  at some specified later time  $T$ . Then we have a 2-point BVP

$$\begin{aligned} \theta''(t) &= -\sin(\theta(t)) \quad \text{for } 0 < t < T, \\ \theta(0) &= \alpha, \quad \theta(T) = \beta. \end{aligned} \quad (2.54)$$

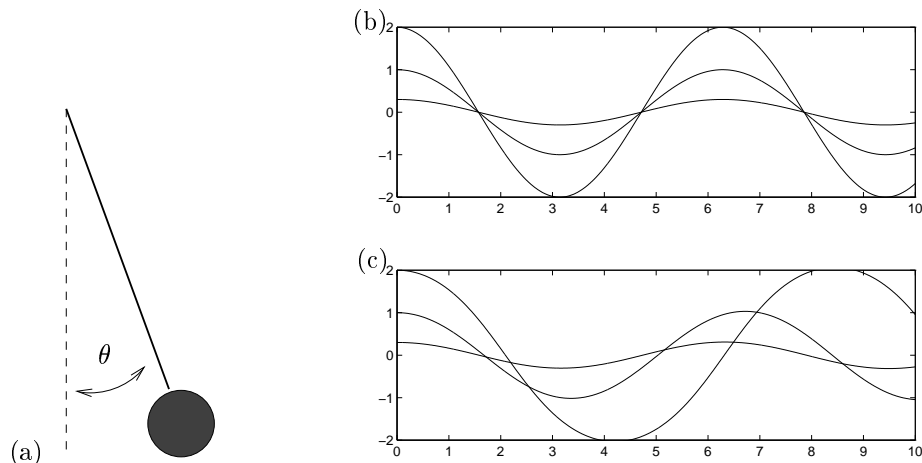


Figure 2.3: (a) Pendulum. (b) Solutions to the linear equation (2.53) for various initial  $\theta$  and zero initial velocity. (c) Solutions to the nonlinear equation (2.52) for various initial  $\theta$  and zero initial velocity.

Similar BVP's do arise in more practical situations, for example trying to shoot a missile in such a way that it hits a desired target. In fact this latter example gives rise to the name *shooting method* for another approach to solving 2-point BVP's that is discussed in [AMR88], [Kel76], for example.

### 2.15.1 Discretization of the nonlinear BVP

We can discretize the nonlinear problem (2.52) in the obvious manner, following our approach for linear problems, to obtain the system of equations

$$\frac{1}{h^2}(\theta_{i-1} - 2\theta_i + \theta_{i+1}) + \sin(\theta_i) = 0 \quad (2.55)$$

for  $i = 1, 2, \dots, m$ , where  $h = T/(m + 1)$  and we set  $\theta_0 = \alpha$  and  $\theta_{m+1} = \beta$ . As in the linear case, we have a system of  $m$  equations for  $m$  unknowns. However, this is now a *nonlinear system* of equations of the form

$$G(\theta) = 0 \quad (2.56)$$

where  $G : \mathbb{R}^m \rightarrow \mathbb{R}^m$ , which cannot be solved as easily as the tridiagonal linear systems encountered so far. Instead of a direct method we must generally use some *iterative method* such as Newton's method. If  $\theta^{[k]}$  is our approximation to  $\theta$  in Step  $k$ , then *Newton's method* is derived via the Taylor series expansion

$$G(\theta^{[k+1]}) = G(\theta^{[k]}) + G'(\theta^{[k]})(\theta^{[k+1]} - \theta^{[k]}) + \dots$$

Setting  $G(\theta^{[k+1]}) = 0$  as desired, and dropping the higher order terms, results in

$$0 = G(\theta^{[k]}) + G'(\theta^{[k]})(\theta^{[k+1]} - \theta^{[k]}).$$

This gives the Newton update

$$\theta^{[k+1]} = \theta^{[k]} + \delta^{[k]} \quad (2.57)$$

where  $\delta^{[k]}$  solves the linear system

$$J^{[k]}\delta^{[k]} = -G(\theta^{[k]}). \quad (2.58)$$



Table 2.1: Change  $\|\delta^{[k]}\|_\infty$  in solution in each iteration of Newton's method.

$k$	Figure 2.4(a)	Figure 2.5
0	3.2841e-01	4.2047e+00
1	1.7518e-01	5.3899e+00
2	3.1045e-02	8.1993e+00
3	2.3739e-04	7.7111e-01
4	1.5287e-08	3.8154e-02
5	5.8197e-15	2.2490e-04
6	1.5856e-15	9.1667e-09
7		1.3395e-15

The program prints out  $\|\delta^{[k]}\|_\infty$  in each iteration, which measures the change in the solution. These are shown in Table 2.1. As expected, Newton's method appears to be converging quadratically (see Appendix A3).

If we start with a different initial guess  $\theta^{[0]}$  (but still close enough to this solution), we would find that the method still converges to this same solution. For example, Figure 2.4(b) shows the iterates  $\theta^{[k]}$  for  $k = 0, 1, 2, \dots$  with a different choice  $\theta^{[0]} \equiv 0.7$ .

### 2.15.2 Nonconvergence

Newton's method can be shown to converge if we start with an initial guess that is sufficiently close to a solution. How close is needed depends on the nature of the problem and is described by the *Newton-Kantorovich Theorem*, see, e.g., [?]. For the problem considered above one need not start very close to the solution to converge, as seen in the examples, but for more sensitive problems one might have to start extremely close. In such cases it may be necessary to use a technique such as *continuation* to find suitable initial data.

### 2.15.3 Nonuniqueness

The nonlinear problem does not have an infinite family of solutions the way the linear equation does on the interval  $[0, 2\pi]$ , and the solution found above is an *isolated solution* in the sense that there are no other solutions very nearby (it is also said to be *locally unique*). However, it does not follow that this is the unique solution to the BVP (2.54). In fact physically we should expect other solutions. The solution we found corresponds to releasing the pendulum with nearly zero initial velocity. It swings through nearly one complete cycle and returns to the initial position at time  $T$ .

Another possibility would be to propel the pendulum upwards so that it rises towards the top (an unstable equilibrium) at  $\theta = \pi$ , before falling back down. By specifying the correct velocity we should be able to arrange it so that the pendulum has fallen back to  $\theta = 0.7$  again at  $T = 2\pi$ . In fact it is possible to find such a solution for any  $T > 0$ .

Physically it seems clear that there is a second solution to the BVP. In order to find it numerically we can use the same iteration as before, but with a different initial guess  $\theta^{[0]}$  that is sufficiently close to this solution. Since we are now looking for a solution where  $\theta$  initially increases and then falls again, let's try a function with this general shape. In Figure 2.5 we see the iterates  $\theta^{[k]}$  generated with data  $\theta_i^{[0]} = 0.7 + \sin(x_i/2)$ . We have gotten lucky here on our first attempt, and we get convergence to a solution of the desired form. (See Table 2.1.) Different guesses with the same general shape might not work. Note that some of the iterates  $\theta^{[k]}$  obtained along the way in Figure 2.5 do not make physical sense (since  $\theta$  goes above  $\pi$  and then back down — what does this mean?), but the method still converges.

**Exercise 2.10** See if you can find yet another solution to this BVP.

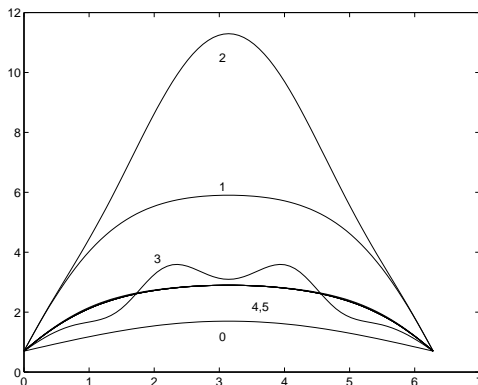


Figure 2.5: Convergence of Newton iterates towards a different solution of the pendulum problem, by starting with initial guess  $\theta_i^{[0]} = 0.7 + \sin(x_i)$ . The iterates  $\theta_i^{[k]}$  for  $nu = 1, 2, \dots$  are denoted by the number  $k$  in the plots.

**Exercise 2.11** Find a numerical solution to this BVP with the same general behavior as seen in Figure 2.5 for the case a longer time interval, say  $T = 20$ , again with  $\alpha = \beta = 0.7$ . Try larger values of  $T$ . What does  $\max_i \theta_i$  approach as  $T$  is increased? Note that for large  $T$  this solution exhibits “boundary layers” (see Section 2.16).

### 2.15.4 Accuracy

The solutions plotted above are not exact solutions to the BVP (2.54). They are only solutions to the discrete system of equations (2.55) with  $h = 1/80$ . How well do they approximate true solutions of the differential equation? Since we have used a second order accurate centered approximation to the second derivative in (2.8), we again hope to obtain second-order accuracy as the grid is refined. In this section we will investigate this.

Note that it is very important to keep clear the distinction between the convergence of Newton’s method to a solution of the finite difference equations, and the convergence of this finite-difference approximation to the solution of the differential equation. Table 2.1 indicates that we have obtained a solution to machine accuracy (roughly  $10^{-15}$ ) of the nonlinear system of equations by using Newton’s method. This does **not** mean that our solution agrees with the true solution of the differential equation to the same degree. This depends on the size of  $h$ , the size of the truncation error in our finite-difference approximation, and the relation between the local truncation error and the resulting global error.

Let’s start by computing the local truncation error of the finite-difference formula. Just as in the linear case, we define this by inserting the true solution of the differential equation into the finite-difference equations. This will not satisfy the equations exactly, and the residual is what we call the *local truncation error*:

$$\begin{aligned} \tau_i &= \frac{1}{h^2}(\theta(t_{i-1}) - 2\theta(t_i) + \theta(t_{i+1})) + \sin(\theta(t_i)) \\ &= (\theta''(t_i) + \sin(\theta(t_i))) + \frac{1}{12}h^2\theta''''(t_i) + O(h^4) \\ &= \frac{1}{12}h^2\theta''''(t_i) + O(h^4). \end{aligned} \tag{2.60}$$

Note that we have used the differential equation to set  $\theta''(t_i) + \sin(\theta(t_i)) = 0$ , which holds exactly since  $\theta(t)$  is the exact solution. The local truncation error is  $O(h^2)$  and has exactly the same form as we found in the linear case. (For a more complicated nonlinear problem it might not work out so simply, see Exercise ??, but similar expressions result.) The vector  $\tau$  with components  $\tau_i$  is simply  $G(\hat{\theta})$ , where

$\hat{\theta}$  is the vector made up of the true solution at each grid point. We now want to obtain an estimate on the global error  $E$  based on this local error. We can attempt to follow the path used in Section 2.6 for linear problems. We have

$$\begin{aligned} G(\theta) &= 0 \\ G(\hat{\theta}) &= \tau_i \end{aligned}$$

and subtracting gives

$$G(\theta) - G(\hat{\theta}) = -\tau. \quad (2.61)$$

If  $G$  were linear ( $G(\theta) = A\theta - F$ ) we would have  $G(\theta) - G(\hat{\theta}) = A\theta - A\hat{\theta} = A(\theta - \hat{\theta}) = AE$ , giving an expression in terms of the global error  $E = \theta - \hat{\theta}$ . This is what we used in Section 2.7.

In the nonlinear case we cannot express  $G(\theta) - G(\hat{\theta})$  directly in terms of  $\theta - \hat{\theta}$ . However, we can use Taylor series expansions to write

$$G(\theta) = G(\hat{\theta}) + J(\hat{\theta})E + O(\|E\|^2)$$

where  $J(\hat{\theta})$  is again the Jacobian matrix of the difference formulas, evaluated now at the exact solution. Combining this with (2.61) gives

$$J(\hat{\theta})E = -\tau + O(\|E\|^2).$$

If we ignore the higher order terms then we again have a linear relation between the local and global errors.

This motivates the following definition of stability. Here we let  $\hat{j}^h$  denote the Jacobian matrix of the difference formulas evaluated at the true solution on a grid with grid-spacing  $h$ .

**Definition 2.15.1** *The nonlinear difference method  $G(\theta) = 0$  is stable in some norm  $\|\cdot\|$  if the matrices  $(\hat{j}^h)^{-1}$  are uniformly bounded in this norm as  $h \rightarrow 0$ , i.e., there exist constants  $C$  and  $h_0$  such that*

$$\|(\hat{j}^h)^{-1}\| \leq C \quad \text{for all } h < h_0. \quad (2.62)$$

It can be shown that if the method is stable in this sense, and consistent in this norm ( $\|\tau^h\| \rightarrow 0$ ), then the method converges and  $\|E^h\| \rightarrow 0$  as  $h \rightarrow 0$ . This is not obvious in the nonlinear case: we obtain a linear system for  $E$  only by dropping the  $O(\|E\|^2)$  nonlinear terms. Since we are trying to show that  $E$  is small, we can't necessarily assume that these terms are negligible in the course of the proof, at least not without some care. See [Kel76] for a proof.

It makes sense that it is uniform boundedness of the inverse Jacobian at the exact solution that is required for stability. After all, it is essentially this Jacobian matrix that is used in solving linear systems in the course of Newton's method, once we get very close to the solution.

WARNING: A final reminder that there is a difference between convergence of the difference method as  $h \rightarrow 0$  and convergence of Newton's method, or some other iterative method, to the solution of the difference equations for some particular  $h$ . Stability of the difference method does not imply that Newton's method will converge from a poor initial guess. (Though it can be shown that for a stable method it will converge from a sufficiently good initial guess — see [Kel76].) Also, the fact that Newton's method has converged to a solution of the nonlinear system of difference equations, with an error of  $10^{-15}$ , say, does not mean that we have a good solution to the original differential equation. The global error of the difference equations determines this.

## 2.16 Singular perturbations and boundary layers

In this section we consider some singular perturbation problems to illustrate the difficulties that can arise in numerically solving problems with boundary layers or other regions where the solution varies

rapidly. See [Kev90], [KC81] for more detailed discussions of singular perturbation problems. In particular the example used here is very similar to one that can be found in [Kev90], where solution by matched asymptotic expansions is discussed.

As a simple example we consider a steady-state advection-diffusion equation. The time-dependent equation is discussed in Appendix ?? and has the form

$$u_t + au_x = \kappa u_{xx} + \psi \quad (2.63)$$

in the simplest case. This models the temperature  $u(x, t)$  of a fluid flowing through a pipe with constant velocity  $a$ , where the fluid has constant heat diffusion coefficient  $\kappa$ , and  $\psi$  is a source term from heating through the walls of the tube.

If  $a > 0$  then we naturally have a boundary condition at the left boundary (say  $x = 0$ ),

$$u(0, t) = \alpha(t)$$

specifying the temperature of the incoming fluid. At the right boundary (say  $x = 1$ ) the fluid is flowing out and so it may seem that the temperature is determined only by what is happening in the pipe and no boundary condition is needed here. This is correct if  $\kappa = 0$  since the first order advection equation needs only one boundary condition and we are allowed to specify  $u$  only at the left boundary. However, if  $\kappa > 0$  then heat can diffuse upstream, and we need to also specify  $u(1, t) = \beta(t)$  in order to determine a unique solution.

If  $\alpha$ ,  $\beta$ , and  $\psi$  are all independent of  $t$  then we expect a steady state solution, which we hope to find by solving the linear 2-point boundary value problem

$$\begin{aligned} au'(x) &= \kappa u''(x) + \psi(x) \\ u(0) &= \alpha, \quad u(1) = \beta. \end{aligned} \quad (2.64)$$

This can be discretized using the approach of Section 2.4. If  $a$  is small relative to  $\kappa$ , then this problem is easy to solve. In fact for  $a = 0$  this is just the steady-state heat equation discussed in Section 2.14 and for small  $a$  the solution looks nearly identical.

But now suppose  $a$  is large relative to  $\kappa$  (i.e., we crank up the velocity, or we decrease the ability of heat to diffuse with the velocity  $a > 0$  fixed). More properly we should work in terms of the nondimensional *Péclet number* which measures the ratio of advection velocity to transport speed due to diffusion. Here we introduce a parameter  $\epsilon$  which is like the inverse of the Péclet number,  $\epsilon = \kappa/a$ , and rewrite the equation (2.64) in the form

$$\epsilon u''(x) - u'(x) = f(x). \quad (2.65)$$

Then taking  $a$  large relative to  $\kappa$  (large Péclet number) corresponds to the case  $\epsilon \ll 1$ .

We should expect difficulties physically in this case where advection overwhelms diffusion. It would be very difficult to maintain a fixed temperature at the outflow end of the tube in this situation. If we had a thermal device that was capable of doing so by instantaneously heating the fluid to the desired temperature as it passes the right boundary, independent of the temperature of the fluid flowing towards this point, then we would expect the temperature distribution to be essentially discontinuous at this boundary.

Mathematically we expect trouble as  $\epsilon \rightarrow 0$  because in the limit  $\epsilon = 0$  the equation (2.65) reduces to a *first order* equation

$$-u'(x) = f(x) \quad (2.66)$$

which allows only one boundary condition, rather than two. For  $\epsilon > 0$ , no matter how small, we have a second order equation that needs two conditions, but we expect to perhaps see strange behavior at the outflow boundary as  $\epsilon \rightarrow 0$ , since in the limit we are overspecifying the problem.



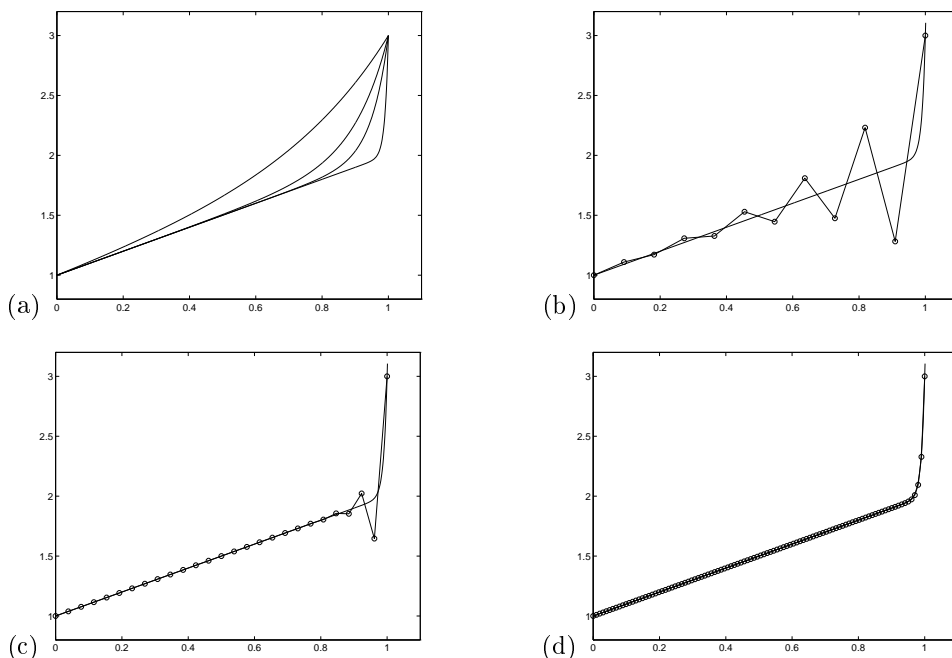


Figure 2.6: (a) Solutions to the steady-state advection-diffusion equation (2.65) for different values of  $\epsilon$ . The four lines correspond to  $\epsilon = 0.3, 0.1, 0.05$  and  $0.01$  from top to bottom. (b) Numerical solution with  $\epsilon = 0.01$  and  $h = 1/10$ . (c)  $h = 1/25$ . (d)  $h = 1/100$ .

Figure 2.6(a) shows how solutions to equation (2.65) look for various values of  $\epsilon$  in the case  $\alpha = 1$ ,  $\beta = 3$ , and  $f(x) = -1$ . In this case the exact solution is

$$u(x) = \alpha + x + (\beta - \alpha - 1) \left( \frac{e^{x/\epsilon} - 1}{e^{1/\epsilon} - 1} \right). \quad (2.67)$$

Note that as  $\epsilon \rightarrow 0$  the solution tends towards a discontinuous function that jumps to the value  $\beta$  at the last possible moment. This region of rapid transition is called the *boundary layer* and it can be shown that for this problem the width of this layer is  $O(\epsilon)$  as  $\epsilon \rightarrow 0$ .

The equation (2.64) with  $0 < \epsilon \ll 1$  is called a *singularly perturbed equation*. It is a small perturbation of the equation (2.66), but this small perturbation changes the character of the equation completely (from a first order equation to a second order equation). Typically any differential equation having a small parameter multiplying the highest order derivative will give a singular perturbation problem.

By contrast, going from the pure diffusion equation  $\kappa u_{xx} = f$  to an advection diffusion equation  $\kappa u_{xx} - au_x = f$  for very small  $a$  is a *regular perturbation*. Both of these equations are second order differential equations requiring the same number of boundary conditions. The solution of the perturbed equation looks nearly identical with the solution of the unperturbed equation for small  $a$ , and the difference in solutions is  $O(a)$  as  $a \rightarrow 0$ .

Singular perturbation problems cause numerical difficulties because the solution changes rapidly over a very small interval in space. In this region derivatives of  $u(x)$  are large, giving rise to large errors in our finite difference approximations. Recall that the error in our approximation to  $u''(x)$  is proportional to  $h^2 u''''(x)$ , for example. If  $h$  is not small enough, then the local truncation error will be very large in the boundary layer. Moreover, even if the truncation error is large only in the boundary layer, the resulting global error may be large everywhere. (Recall that the global error  $E$  is obtained from the truncation error  $\tau$  by solving a linear system  $AE = -\tau$ , which means that each element of  $E$  depends on *all* elements of  $\tau$  since  $A^{-1}$  is a dense matrix.) This is clearly seen in Figure 2.6(b) where the numerical solution with  $h = 1/10$  is plotted. Errors are large even in regions where the exact

solution is nearly linear and  $u'''' \approx 0$ .

On finer grids the solution looks better, see Figure 2.6(c) and (d), and as  $h \rightarrow 0$  the method does exhibit second order accurate convergence. But it is necessary to have a sufficiently fine grid before reasonable results are obtained.

### 2.16.1 Interior layers

The above example has a boundary layer, a region of rapid transition at one boundary. Other problems may have *interior layers* instead. In this case the solution is smooth except for some thin region interior to the interval where a rapid transition occurs. Such problems can be even more difficult to solve since we often don't know to begin with where the interior layer will be. Perturbation theory can often be used to analyse singular perturbation problems and predict where the layers will occur, how wide they will be (as a function of the small parameter  $\epsilon$ ) and how the solution behaves. The use of perturbation theory to obtain good approximations to problems of this type is a central theme of classical applied mathematics. (See the texts mentioned at the beginning of Section 2.16 for examples.)

These analytic techniques can often be used to good advantage along with numerical methods, for example to obtain a good initial guess for Newton's method, or to choose an appropriate nonuniform grid as discussed in the next section. In some cases it is possible to develop special numerical methods that have the correct singular behavior built into the approximation in such a way that far better accuracy is achieved than with a naive numerical method.

## 2.17 Nonuniform grids and adaptive refinement

From Figure 2.6 it is clear that we need to choose our grid fine enough that several points are within the boundary layer in order to obtain a reasonable solution. If we wanted high accuracy within the boundary layer we would have to choose a much finer grid than shown in this Figure. With a uniform grid this means using a very large number of grid points, the vast majority of which are in the region where the solution is very smooth and could be represented well with far fewer points. This waste of effort may be tolerable for simple one-dimensional problems, but can easily be intolerable for more complicated problems, particularly in more than one dimension.

Instead it is preferable to use a nonuniform grid for such calculations, with grid points clustered in regions where they are most needed. This requires developing formulas that are sufficiently accurate on nonuniform grids as in Homework #1. There we saw that a 4-point stencil could be used to obtain second order accuracy for the second derivative operator. Using this for a linear problem would give a banded matrix with 4 nonzero diagonals. A little extra care is needed at the boundaries.

Finite element methods are often easier to define on nonuniform grids than finite difference methods. See Section 4.3.

One way to specify nonuniform grid points is to start with a uniform grid in some artificial coordinate  $\xi$ , which we will denote by  $\xi_i = ih$  for  $i = 0, 1, \dots, m+1$  where  $h = 1/(m+1)$ , and then use some appropriate *grid mapping* function  $X(\xi)$  to define the physical grid points  $x_i = X(\xi_i)$ . This is illustrated in Figure 2.7, where  $\xi$  is plotted on the vertical axis and  $x$  is on the horizontal axis. The curve plotted represents a function  $X(\xi)$ , though with this choice of axes it is more properly the graph of the inverse function  $\xi = X^{-1}(x)$ . The horizontal and vertical lines indicate how the uniform grid points on the  $\xi$  axis are mapped to nonuniform points in  $x$ . If the problem is posed on the interval  $[a, b]$ , then the function  $X(\xi)$  should be monotonically increasing and satisfy  $X(0) = a$  and  $X(1) = b$ .

Note that grid points are clustered in regions where the curve is steepest, which means that  $X(\xi)$  varies slowly with  $\xi$ , and spread apart in regions where  $X(\xi)$  varies rapidly with  $\xi$ .

Readers who wish to use methods on nonuniform grids are strongly encouraged to investigate some of the software which is freely available. There are good packaged routines which will automatically choose an appropriate grid for a given problem (perhaps with some initial guidance) and take care of all the details of discretizing on this grid. The COLSYS collocation package is one such software package,

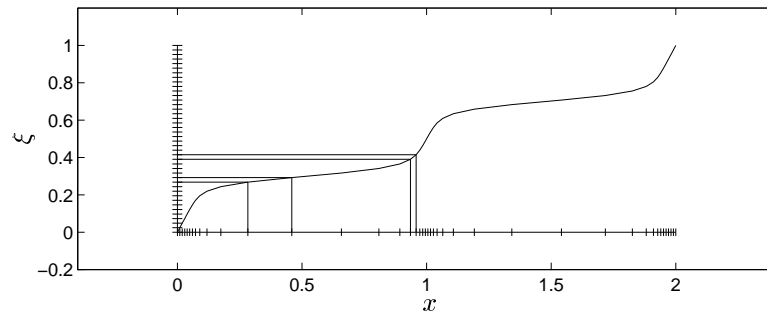


Figure 2.7: Grid mapping from a uniform grid in  $0 \leq \xi \leq 1$  (vertical axis) to the nonuniform grid in physical  $x$ -space shown on the horizontal axis. This particular mapping clusters grid points near the endpoints  $x = 0$  and  $x = 2$  and also near the center.

available from NETLIB. Strategies for *adaptive mesh selection* (i.e., choosing a grid that adapts to the nature of the solution) are discussed, for example, in [AMR88].

## Chapter 3

# Elliptic Equations

In more than one space dimension, the steady-state equations discussed in Chapter 2 generalize naturally to *elliptic* partial differential equations. In two space dimension a constant-coefficient elliptic equation has the form

$$a_1 u_{xx} + a_2 u_{xy} + a_3 u_{yy} + a_4 u_x + a_5 u_y + a_6 u = f \quad (3.1)$$

where the coefficients  $a_1, a_2, a_3$  satisfy

$$a_2^2 - 4a_1 a_3 < 0. \quad (3.2)$$

This equation must be satisfied for all  $(x, y)$  in some region of the plane  $\Omega$ , together with some boundary conditions on  $\partial\Omega$ , the boundary of  $\Omega$ . For example we may have Dirichlet boundary conditions in which case  $u(x, y)$  is given at all points  $(x, y) \in \partial\Omega$ . If the ellipticity condition (3.2) is satisfied then this gives a well-posed problem. If the coefficients vary with  $x$  and  $y$  then this condition must be satisfied at each point in  $\Omega$ .

### 3.1 Steady-state heat conduction

Equations of elliptic character often arise as steady-state equations in some region of space, associated with some time-dependent physical problem. For example, the diffusion or heat conduction equation in two space dimensions takes the form

$$u_t = (\kappa u_x)_x + (\kappa u_y)_y + \psi \quad (3.3)$$

where  $\kappa(x, y) > 0$  is a diffusion or heat conduction coefficient that may vary with  $x$  and  $y$ , and  $\psi(x, y, t)$  is a source term. The solution  $u(x, y, t)$  will generally vary with time as well as space. We also need initial conditions  $u(x, y, 0)$  in  $\Omega$  and boundary conditions at each point in time at every point on the boundary of  $\Omega$ . If the boundary conditions and source terms are independent of time then we expect a steady state to exist, which we can find by solving the elliptic equation

$$(\kappa u_x)_x + (\kappa u_y)_y = f \quad (3.4)$$

where again we set  $f(x, y) = -\psi(x, y)$ , together with the boundary conditions. Note that (3.2) is satisfied at each point provided  $\kappa > 0$  everywhere.

We first consider the simplest case where  $\kappa \equiv 1$ . Variable coefficients will be considered in Section ?? . We then have

$$u_{xx} + u_{yy} = f \quad (3.5)$$

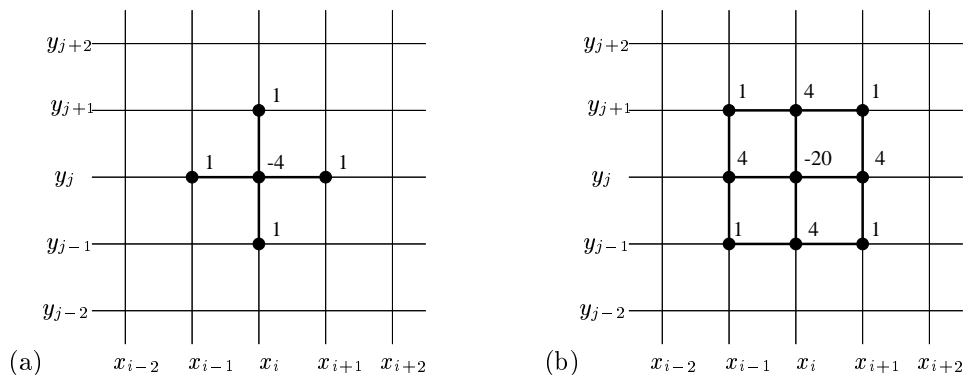


Figure 3.1: Portion of the computational grid for a two-dimensional elliptic equation. (a) The 5-point stencil for the Laplacian about the point  $(i, j)$  is also indicated. (b) The 9-point stencil is indicated, which is discussed in Section 3.5.

which is called the *Poisson problem*. In the special case  $f \equiv 0$  this reduces to *Laplace's equation*,

$$u_{xx} + u_{yy} = 0. \quad (3.6)$$

In one space dimension the corresponding equation  $u''(x) = 0$  is trivial: the solution is a linear function connecting the two boundary values. In two dimensions even this simple equation is nontrivial to solve, since boundary values can now be specified at every point along the curve defining the boundary. Solutions to Laplace's equation are called *harmonic functions*. You may recall from complex analysis that if  $g(z)$  is any complex analytic function of  $z = x + iy$ , then the real and imaginary parts of this function are harmonic. For example,  $g(z) = z^2 = (x^2 - y^2) + 2ixy$  is analytic and the functions  $x^2 - y^2$  and  $2xy$  are both harmonic.

The operator  $\nabla^2$  defined by

$$\nabla^2 u = u_{xx} + u_{yy}$$

is called the *Laplacian*. The notation  $\nabla^2$  comes from the fact that, more generally,

$$(\kappa u_x)_x + (\kappa u_y)_y = \nabla \cdot (\kappa \nabla u)$$

where  $\nabla u$  is the gradient of  $u$ ,

$$\nabla u = \begin{bmatrix} u_x \\ u_y \end{bmatrix} \quad (3.7)$$

and  $\nabla \cdot$  is the divergence operator,

$$\nabla \cdot \begin{bmatrix} u \\ v \end{bmatrix} = u_x + v_y. \quad (3.8)$$

The symbol  $\Delta$  is also often used for the Laplacian, but would lead to confusion in numerical work where  $\Delta x$  and  $\Delta y$  will be used for grid spacing.

## 3.2 The five-point stencil for the Laplacian

To discuss discretizations, first consider the Poisson problem (3.5) on the unit square  $0 \leq x \leq 1$ ,  $0 \leq y \leq 1$  and suppose we have Dirichlet boundary conditions. We will use a uniform Cartesian grid consisting of grid points  $(x_i, y_j)$  where  $x_i = i\Delta x$  and  $y_j = j\Delta y$ . A section of such a grid is shown in Figure 3.1.

Let  $u_{ij}$  represent an approximation to  $u(x_i, y_j)$ . In order to discretize (3.5) we replace both the  $x$ - and  $y$ -derivatives with centered finite differences, which gives

$$\frac{1}{(\Delta x)^2}(u_{i-1,j} - 2u_{ij} + u_{i+1,j}) + \frac{1}{(\Delta y)^2}(u_{i,j-1} - 2u_{ij} + u_{i,j+1}) = f_{ij}. \tag{3.9}$$

For simplicity of notation we will consider the special case where  $\Delta x = \Delta y \equiv h$ , though it is easy to handle the general case. We can then rewrite (3.9) as

$$\frac{1}{h^2}(u_{i-1,j} + u_{i+1,j} + u_{i,j-1} + u_{i,j+1} - 4u_{ij}) = f_{ij}. \tag{3.10}$$

This finite difference scheme can be represented by the *5-point stencil* shown in Figure 3.1. We have both an unknown  $u_{ij}$  and an equation of the form (3.10) at each of  $m^2$  grid points for  $i = 1, 2, \dots, m$  and  $j = 1, 2, \dots, m$ , where  $h = 1/(m + 1)$  as in one dimension. We thus have a linear system of  $m^2$  unknowns. The difference equations at points near the boundary will of course involve the known boundary values, just as in the one-dimensional case, that will be moved to the right-hand side.

If we collect all of these equations together into a matrix equation, we will have an  $m^2 \times m^2$  matrix which is very *sparse*, i.e., most of the elements are zero. Since each equation involves at most 5 unknowns (less near the boundary), each row the matrix has at most 5 nonzeros and at least  $m^2 - 5$  elements that are zero. This is analogous to the tridiagonal matrix (2.9) seen in the one-dimensional case, in which each row has at most 3 nonzeros.

Unfortunately, in two space dimensions the structure of the matrix is not as simple as in one dimension, and in particular the nonzeros will not all be as nicely clustered near the main diagonal. The exact structure of the matrix depends on the order in which we order the unknowns and write down the equations, as we will see below, but no ordering is ideal.

Note that in general we are always free to change the order of the equations in a linear system without changing the solution. Modifying the order corresponds to permuting the rows of the matrix and right-hand side. We are also free to change the ordering of the unknowns in the vector of unknowns, which corresponds to permuting the columns of the matrix. As an example, consider the one-dimensional difference equations given by (2.9). Suppose we reordered the unknowns by listing first the unknowns at odd numbered grid points and then the unknowns at even numbered grid points, so that  $\tilde{U} = [U_1, U_3, U_5, \dots, U_2, U_4, \dots]^T$ . If we also reorder the equations in the same way, i.e., we write down first the difference equation centered at  $U_1$ , then at  $U_3, U_5$ , etc., then we would obtain the following system:

$$\frac{1}{h^2} \left[ \begin{array}{cccc|cccc} -2 & & & & 1 & & & \\ & -2 & & & 1 & 1 & & \\ & & -2 & & & 1 & 1 & \\ & & & \ddots & & & \ddots & \ddots \\ & & & & & & & 1 & 1 \\ \hline 1 & 1 & & & -2 & & & \\ & & 1 & 1 & & -2 & & \\ & & & 1 & 1 & & -2 & \\ & & & & \ddots & \ddots & & \\ & & & & & & & 1 & -2 \end{array} \right] \begin{bmatrix} U_1 \\ U_3 \\ U_5 \\ \vdots \\ U_{m-1} \\ U_2 \\ U_4 \\ U_6 \\ \vdots \\ U_m \end{bmatrix} = \begin{bmatrix} f(x_1) - \alpha/h^2 \\ f(x_3) \\ f(x_5) \\ \vdots \\ f(x_{m-1}) \\ f(x_2) \\ f(x_4) \\ f(x_6) \\ \vdots \\ f(x_m) - \beta/h^2 \end{bmatrix}. \tag{3.11}$$

This linear system has the same solution as (2.9) modulo the reordering of unknowns, but looks very different. For this one-dimensional problem there is no point in reordering things this way, and the natural ordering  $[U_1, U_2, U_3, \dots]^T$  clearly gives the optimal matrix structure for the purpose of applying Gaussian elimination. By ordering the unknowns so that those which occur in the same equation are close to one another in the vector, we keep the nonzeros in the matrix clustered near the diagonal.

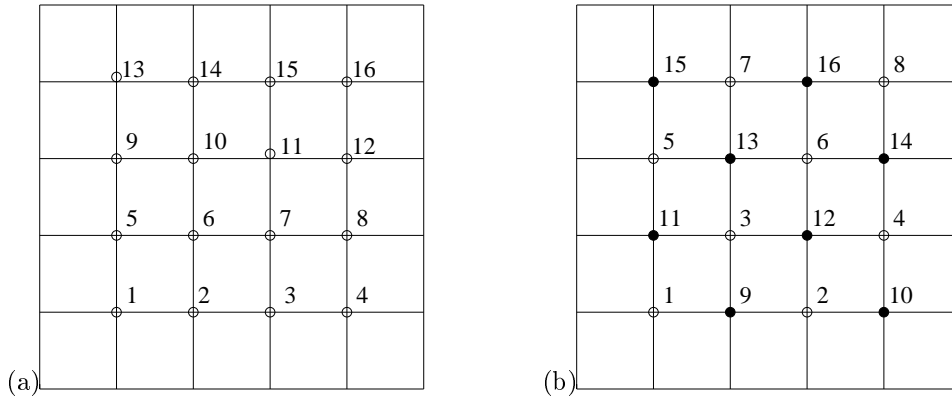


Figure 3.2: (a) The natural rowwise order of unknowns and equations on a  $4 \times 4$  grid. (b) The red-black ordering.

Returning to the two-dimensional problem, it should be clear that there is no way to order the unknowns so that all nonzeros are clustered adjacent to the diagonal. About the best we can do is to use the *natural rowwise ordering*, where we take the unknowns along the bottom row,  $u_{11}, u_{21}, u_{31}, \dots, u_{m1}$ , followed by the unknowns in the second row,  $u_{12}, u_{22}, \dots, u_{m2}$ , and so on, as illustrated in Figure 3.2(a). This gives a matrix equation where  $A$  has the form

$$A = \begin{bmatrix} T & I & & & \\ I & T & I & & \\ & I & T & I & \\ & & \ddots & \ddots & \ddots \\ & & & I & T \end{bmatrix}$$

which is an  $m \times m$  *block tridiagonal matrix* in which each block  $T$  or  $I$  is itself an  $m \times m$  matrix,

$$T = \begin{bmatrix} -4 & 1 & & & \\ 1 & -4 & 1 & & \\ & 1 & -4 & 1 & \\ & & \ddots & \ddots & \ddots \\ & & & 1 & -4 \end{bmatrix}$$

and  $I$  is the  $m \times m$  identity matrix. While this has a nice structure, the 1 values in the  $I$  matrices are separated from the diagonal by  $m - 1$  zeros, since these coefficients correspond to grid points lying above or below the central point in the stencil and hence are in the next or previous row of unknowns.

Another possibility, which has some advantages in the context of certain iterative methods (see Section ??), is to use the *red-black ordering* (or checkerboard ordering) shown in Figure 3.2. This is the two-dimensional analog of the odd-even ordering that leads to the matrix (3.11) in one dimension. This ordering is significant because all 4 neighbors of a red grid point are black points, and vice versa, and leads to a matrix equation with the structure

$$\begin{bmatrix} D & H \\ H^T & D \end{bmatrix} \begin{bmatrix} u_{\text{red}} \\ u_{\text{black}} \end{bmatrix} = \dots \quad (3.12)$$

where  $D = -4I$  is a diagonal matrix of dimension  $m^2/2$ .

**Exercise 3.1** What is the matrix  $H$  in (3.12) for the grid shown in Figure 3.2?

### 3.3 Solving the linear system

In Chapter 5 we will consider various approaches to solving these linear systems in more detail. Here we will only discuss some of the main issues.

There are two fundamentally different approaches that could be used for solving these linear systems. A *direct method* such as Gaussian elimination produces an exact solution (or at least would in exact arithmetic) in a finite number of operations. An *iterative method* starts with an initial guess for the solution and attempts to improve it through some iterative procedure, halting after a sufficiently good approximation has been obtained. For problems with large sparse matrices, iterative methods are often preferable for various reasons described below.

For certain special problems very fast direct methods can be used, which are much better than standard Gaussian elimination. This is discussed briefly in Section 3.3.2.

#### 3.3.1 Gaussian elimination

This is the best known direct method for solving a linear system  $Au = F$ . The matrix  $A$  is reduced to upper triangular form by taking linear combinations of the rows to introduce zeros below the diagonal. This can be viewed as giving a factorization  $A = LU$  of the matrix into a product of an upper and a lower triangular matrix. Forward and back substitution are then used to solve the triangular systems  $Lc = F$  and then  $Uu = c$ . (See, e.g., [GL89] or [?]).

It is easy to compute that for a general  $N \times N$  *dense* matrix (one with few elements equal to zero), performing Gaussian elimination requires  $O(N^3)$  operations. (There are  $N(N-1)/2 = O(N^2)$  elements below the diagonal to eliminate, and eliminating each one requires  $O(N)$  operations to take a linear combination of the rows.)

Applying a general Gaussian elimination program blindly to the matrices we are now dealing with would be disastrous, or at best extremely wasteful of computer resources. Suppose we are solving the Poisson problem on a  $100 \times 100$  grid for example. Then  $N = m^2 = 10^4$  and  $N^3 = 10^{12}$ . On a reasonably fast workstation (ca. 1996) which can do on the order of  $10^7$  floating point operations per second (10 megaflops), this would take on the order of  $10^5$  seconds, which is roughly 28 minutes. If we went up to a  $1000 \times 1000$  grid (a million by million matrix) this would increase by a factor of  $10^6$  to roughly 3000 years. Things are worse in three dimensions. Solving the Poisson problem on a  $100 \times 100 \times 100$  gives the same matrix dimension  $N = 10^6$  as the  $1000 \times 1000$  grid in two dimensions. More sophisticated methods can solve even these larger problems in a reasonable amount of time, and problems with a million unknowns are not unusual in applications.

Moreover, even if speed were not an issue, memory would be. Storing the full matrix  $A$  in order to modify the elements and produce  $L$  and  $U$  would require  $N^2$  memory locations. In 8-byte arithmetic this requires  $8N^2$  bytes. For the larger problems mentioned above, this would be  $8 \times 10^{12}$  bytes, or 8000 gigabytes. One advantage of iterative methods is that they do not store the matrix at all, and at most need to store the nonzero elements.

Of course with Gaussian elimination it would be foolish to store all the elements of a sparse matrix, since the vast majority are zero, or to apply the procedure blindly without taking advantage of the fact that so many elements are already zero and hence do not need to be eliminated.

As an extreme example, consider the one-dimensional case where we have a tridiagonal matrix as in (2.9). Applying Gaussian elimination requires only eliminating the nonzeros along the subdiagonal, only  $N-1$  values instead of  $N(N-1)/2$ . Moreover, when we take linear combinations of rows in the course of eliminating these values, in most columns we will be taking linear combinations of zeros, producing zero again. If we do not do pivoting, then only the diagonal elements are modified. Even with partial pivoting, at most we will introduce one extra superdiagonal of nonzeros in the upper triangular  $U$  that were not present in  $A$ . As a result, it is easy to see that applying Gaussian elimination to an  $m \times m$  tridiagonal system requires only  $O(m)$  operations, not  $O(m^3)$ , and that the storage required is  $O(m)$  rather than  $O(m^2)$ .

Note that this is the best we could hope for in one dimension, at least in terms of the order of



magnitude. There are  $m$  unknowns and even if we had exact formulas for these values, it would require  $O(m)$  work to evaluate them and  $O(m)$  storage to save them.

In two space dimensions we can also take advantage of the sparsity and structure of the matrix to greatly reduce the storage and work required with Gaussian elimination, though not to the minimum that one might hope to attain. On an  $m \times m$  grid there are  $N = m^2$  unknowns, so the best one could hope for is an algorithm that computes the solution in  $O(N) = O(m^2)$  work using  $O(m^2)$  storage. Unfortunately this cannot be achieved with a direct method.

One approach that is better than working with the full matrix is to observe that the  $A$  is a band matrix with bandwidth  $m$  both above and below the diagonal. Since a general  $N \times N$  banded matrix with  $a$  nonzero bands above the diagonal and  $b$  below the diagonal can be factored in  $O(Nab)$  operations, this results in an operation count of  $O(m^4)$ .

A more sophisticated approach that takes more advantage of the special structure (and the fact that there are already many zeros within the bandwidth) is the *nested dissection* algorithm [GL81]. This algorithm requires  $O(m^3)$  operations. It turns out this is the best that can be achieved with a direct method based on Gaussian elimination. George has proved (see [GL81]) that any elimination method for solving this problem requires at least  $O(m^3)$  operations.

### 3.3.2 Fast Poisson solvers

For this very special system of equations there are other techniques that can be used to solve the system quickly. These are generally called *fast Poisson solvers* (recall we are looking at the Poisson problem (3.10) with constant coefficients on a rectangular domain). One standard technique uses fast Fourier transforms to solve the system in  $O(m^2 \log m)$  work, which is nearly optimal since there are  $m^2$  grid values to be determined.

To explain the main idea of this approach, we consider only the one-dimensional case, where the Poisson problem reduces to  $u''(x) = f(x)$  on the interval  $[0, 1]$ , and the centered 3-point discretization gives the matrix equation (2.9) (for the case of Dirichlet boundary conditions). In one dimension this tridiagonal system can be solved in  $O(m)$  operations and there is no need for a “fast solver”, but the idea is easiest to illustrate in this case.

In Section 2.10 we determined the eigenvalues and eigenvectors of the matrix  $A$ . Based on these we can write down the Jordan Canonical Form of  $A$ :

$$A = R\Lambda R^{-1} \quad (3.13)$$

where  $R$  is the matrix of right eigenvectors (the  $p$ 'th column of  $R$  is the vector  $u^p$  from (2.24)) and  $\Lambda$  is a diagonal matrix with the eigenvalues on the diagonal. So we have

$$R = \begin{bmatrix} \sin(\pi x_1) & \sin(2\pi x_1) & \cdots & \sin(m\pi x_1) \\ \sin(\pi x_2) & \sin(2\pi x_2) & \cdots & \sin(m\pi x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \sin(\pi x_m) & \sin(2\pi x_m) & \cdots & \sin(m\pi x_m) \end{bmatrix} \quad (3.14)$$

and

$$\Lambda = \begin{bmatrix} \frac{2}{h^2}(\cos(\pi h) - 1) & & & \\ & \frac{2}{h^2}(\cos(2\pi h) - 1) & & \\ & & \ddots & \\ & & & \frac{2}{h^2}(\cos(m\pi h) - 1) \end{bmatrix}. \quad (3.15)$$

Since  $A^{-1} = R\Lambda^{-1}R^{-1}$ , one approach to solving the linear system  $Au = F$  is to use this to obtain

$$u = A^{-1}F = R\Lambda^{-1}R^{-1}F.$$

There are at least 4 obvious objections to using this approach to solve a general linear system:

- In general it is not easy to determine the eigenvalues and eigenvectors of a given matrix. Doing so computationally is generally much more expensive than solving the original system by other approaches.
- Even if we know  $R$  and  $\Lambda$ , as we do here, we don't in general know  $R^{-1}$ . Computing an inverse matrix is much more work than solving a single linear system. Of course all we really need here is  $R^{-1}F \equiv \hat{F}$ , say, which can be found by solving a single system  $R\hat{F} = F$ , but this system may be no easier to solve than the original system  $Au = F$  (and perhaps much more expensive since  $R$  is typically dense even if  $A$  is sparse).
- Even if we know  $R$  and  $R^{-1}$  explicitly, storing them typically requires  $O(N^2)$  storage for an  $N \times N$  system since these are dense matrices. We probably want to avoid this if  $A$  is sparse.
- It also requires  $O(N^2)$  work to multiply an  $N \times N$  dense matrix by an  $N$ -vector. In two dimensions, where  $N = m^2$ , this is no better than using Gaussian elimination on the banded matrix  $A$ . (And in one dimension the  $O(m^2)$  work required would be much worse than the  $O(m)$  work needed for the tridiagonal system.)

So why would we ever consider such a method? In general we would not, but for the very special case of a Poisson problem on a simple region, these objections all vanish thanks to the Fast Fourier Transform.

We do know  $\Lambda$  and  $R$  for our simple problem, and moreover we know  $R^{-1}$ , since we can compute using trig identities that

$$R^2 = \left( \frac{m+1}{2} \right) I$$

and so

$$R^{-1} = 2hR.$$

( $R$  is symmetric, and, except for the scaling, an orthogonal matrix.)

Although  $R$  is a dense matrix, it contains only  $m$  distinct values (after applying trig identities) and we do not need to compute or store all of the elements. Moreover, and most importantly, multiplying a vector by  $R$  or  $R^{-1}$  does not require  $O(m^2)$  operations, but can be done in  $O(m \log m)$  operations if  $m$  is a product of small prime factors. Ideally we would choose  $m$  to be a factor of 2, say  $m = 2^k$ . The trick is then to observe (It's not obvious!) that  $R$  can be written as a product of  $k$  matrices,

$$R = R_1 R_2 \cdots R_k$$

where each  $R_j$  is very sparse so that multiplying a vector by  $R_j$  requires only  $O(m)$  operations. Applying each matrix in turn to a vector gives the desired product in  $O(km) = O(m \log m)$  operations.

This algorithm is a special case of the *Fast Fourier Transform*, or FFT, since multiplying a vector by  $R$  corresponds to taking a sine transform of the data. See [Loa97] for a brief introduction to the recursive structure of the FFT.

In one space dimension there is still no advantage to this procedure over solving the tridiagonal system directly, which only requires  $O(m)$  operations. However, this approach carries over directly to 2 or 3 space dimensions, with operation counts of  $O(m^2 \log m)$  and  $O(m^3 \log m)$  respectively, which is asymptotically nearly optimal since the grids have  $m^2$  and  $m^3$  points.

To solve the system  $Au = F$  using this technique (back in one dimension, now), we would first use the FFT algorithm to compute

$$\hat{F} = R^{-1}F = 2hRF, \tag{3.16}$$

then divide each element of  $\hat{F}$  by the corresponding eigenvalue from  $\Lambda$ ,

$$\hat{u}_j = \hat{F}_j / \left( \frac{h^2}{2} (\cos(j\pi h) - 1) \right) \text{ for } j = 1, 1, \dots, m, \tag{3.17}$$

so that  $\hat{u} = \Lambda^{-1}R^{-1}F$ . Finally we do another FFT to compute

$$u = R\hat{u}. \quad (3.18)$$

Unfortunately these techniques generally do not extend to other systems of equations, such as those arising in a variable-coefficient problem. The technique depends on knowing the eigen-decomposition of the matrix  $A$  and on having a fast algorithm for multiplying the dense matrix  $R$  times a vector. For a variable-coefficient problem, we won't in general know  $R$  and  $R^{-1}$  (unless we compute them numerically, which takes much more work than solving the original linear system). Even if we did know them, we wouldn't generally have a fast algorithm to apply them to a vector.

The special case of a constant-coefficient Poisson problem in a rectangular region arises sufficiently often in applications that fast Poisson solvers are often useful, however. Software for this problem is available in `FISHPACK` from `netlib`. For a review of fast Poisson solvers, see [Swa84].

### 3.3.3 Iterative methods

Except when the matrix has very special structure and fast direct methods of the type discussed in the previous section apply, iterative methods are usually the method of choice for large sparse linear systems. In this section two classical iterative methods, Jacobi and Gauss-Seidel, are introduced to illustrate the main issues. It should be stressed at the beginning that these are poor methods in general which converge very slowly, but they have the virtue of being simple to explain. Much more efficient methods are discussed in Chapter 5.

We again consider the Poisson problem where we have the system of equations (3.10). We can rewrite this equation as

$$u_{ij} = \frac{1}{4}(u_{i-1,j} + u_{i+1,j} + u_{i,j-1} + u_{i,j+1}) - \frac{h^2}{4}f_{ij}. \quad (3.19)$$

In particular, note that for Laplace's equation (where  $f_{ij} \equiv 0$ ) this simply states that the value of  $u$  at each grid point should be the average of its four neighbors. This is the discrete analog of the well-known fact that a harmonic function has the following property: The value at any point  $(x, y)$  is equal to the average value around a closed curve containing the point, in the limit as the curve shrinks to the point. Physically this also makes sense if we think of the heat equation. Unless the temperature at this point is equal to the average of the temperature at neighboring points, there will be a net flow of heat towards or away from this point.

The equation (3.19) suggests the following iterative method to produce a new estimate  $u^{[k+1]}$  from a current guess  $u^{[k]}$ :

$$u_{ij}^{[k+1]} = \frac{1}{4}(u_{i-1,j}^{[k]} + u_{i+1,j}^{[k]} + u_{i,j-1}^{[k]} + u_{i,j+1}^{[k]}) - \frac{h^2}{4}f_{ij}. \quad (3.20)$$

This is the *Jacobi* iteration for the Poisson problem, and it can be shown that for this particular problem it converges from any initial guess  $u^{[0]}$  (though very slowly).

Here is a short section of `MATLAB` code that implements the main part of this iteration:

```
for iter=0:maxiter
  for j=2:(m+1)
    for i=2:(m+1)
      unew(i,j) = 0.25*(u(i-1,j) + u(i+1,j) + u(i,j-1) + u(i,j+1) - h^2 * f(i,j));
    end
  end
  u = unew;
end
```

where it is assumed that  $u$  initially contains the guess  $u^{[0]}$  and that boundary data is stored into  $u(1,:)$ ,  $u(m+2,:)$ ,  $u(:,1)$ , and  $u(:,m+2)$ . The indexing is off by one from what one might expect since MATLAB begins arrays with index 1, not 0.

Note that one might be tempted to dispense with the variable `unew` and replace the above code by

```
for iter=0:maxiter
  for j=2:(m+1)
    for i=2:(m+1)
      u(i,j) = 0.25*(u(i-1,j) + u(i+1,j) + u(i,j-1) + u(i,j+1) - h^2 * f(i,j));
    end
  end
end
```

This would not give the same results, however. In the correct code for Jacobi we compute new values of  $u$  based entirely on old data from the previous iteration, as required from (3.20). In the second code we have already updated  $u(i-1,j)$  and  $u(i,j-1)$  before we update  $u(i,j)$ , and these new values will be used instead of the old ones. The latter code thus corresponds to the method

$$u_{ij}^{[k+1]} = \frac{1}{4}(u_{i-1,j}^{[k+1]} + u_{i+1,j}^{[k]} + u_{i,j-1}^{[k+1]} + u_{i,j+1}^{[k]}) - \frac{h^2}{4}f_{ij}. \quad (3.21)$$

This is in fact what is known as the *Gauss-Seidel* method, and it would be a lucky coding error since this method generally converges about twice as fast as Jacobi's method (see Chapter 5).

Convergence of these methods will be discussed in Chapter 5. For now we simply note some features of these iterative methods:

- The matrix  $A$  is never stored. In fact, for this simple constant coefficient problem, we don't even store all the  $5m^2$  nonzeros which all have the value  $1/h^2$  or  $-4/h^2$ . The values 0.25 and  $h^2$  in the code are the only values that are "stored". (For a variable coefficient problem where the coefficients are different at each point, we would in general have to store them all.)
- Hence the storage is optimal — essentially only the  $m^2$  solution values are stored in the Gauss-Seidel method. The above code for Jacobi uses  $2m^2$  since `qnew` is stored as well as `q`, but one could eliminate most of this with more careful coding.
- Each iteration requires  $O(m^2)$  work. The total work required will depend on how many iterations are required to reach the desired level of accuracy. In Chapter 5 we will see that with these particular methods we require  $O(m^2 \log m)$  iterations to reach a level of accuracy consistent with the expected global error in the solution (as  $h \rightarrow 0$  we should require more accuracy in the solution to the linear system). Combining this with the work per iteration gives a total operation count of  $O(m^4 \log m)$ . This looks worse than Gaussian elimination with a banded solver, though since  $\log m$  grows so slowly with  $m$  it is not clear which is really more expensive for a realistic size matrix. (And the iterative method definitely saves on storage.)

Other iterative methods also typically require  $O(m^2)$  work per iteration, but may converge much faster and hence result in less overall work. The ideal would be to converge in a number of iterations that is independent of  $h$ . Multigrid methods can achieve this, not only for Poisson's problem but also for many other elliptic equations.

### 3.4 Accuracy and stability

The discretization of the two-dimensional Poisson problem can be analyzed using exactly the same approach as we used for the one-dimensional boundary value problem. The local truncation error  $\tau_{ij}$

at the  $(i, j)$  grid point is defined in the obvious way,

$$\tau_{ij} = \frac{1}{h^2}(u(x_{i-1}, y_j) + u(x_{i+1}, y_j) + u(x_i, y_{j-1}) + u(x_i, y_{j+1}) - 4u(x_i, y_j)) - f(x_i, y_j),$$

and by splitting this into the second order difference in the  $x$ - and  $y$ -directions it is clear from previous results that

$$\tau_{ij} = \frac{1}{12}h^2(u_{xxxx} + u_{yyyy}) + O(h^4).$$

For this linear system of equations the global error  $E_{ij} = u_{ij} - u(x_i, y_j)$  then solves the linear system

$$A^h E^h = -\tau^h$$

just as in one dimension. The method will be globally second order accurate in some norm provided that it is stable, *i.e.*, that  $\|(A^h)^{-1}\|$  is uniformly bounded as  $h \rightarrow 0$ .

In the 2-norm this is again easy to check, for this simple problem, since we can explicitly compute the spectral radius of the matrix, as we did in one dimension in Section 2.10. The eigenvalues and eigenvectors of  $A$  can now be indexed by 2 parameters  $p$  and  $k$  corresponding to wavenumbers in the  $x$  and  $y$  directions, for  $p, k = 1, 2, \dots, m$ . The  $(p, k)$  eigenvector  $u^{p,k}$  has the  $m^2$  elements

$$u_{ij}^{p,k} = \sin(p\pi ih) \sin(k\pi jh).$$

The corresponding eigenvalue is

$$\lambda^{p,k} = \frac{2}{h^2} ((\cos(p\pi h) - 1) + (\cos(k\pi h) - 1)).$$

The eigenvalues are strictly negative ( $A$  is negative definite) and the one closest to the origin is

$$\lambda^{1,1} = -2\pi^2 + O(h^2).$$

The spectral radius of  $(A^h)^{-1}$ , which is also the 2-norm, is thus

$$\rho((A^h)^{-1}) = 1/\lambda^{1,1} \approx -1/2\pi^2$$

The method is hence stable in the 2-norm.

While we are at it, let's also compute the condition number of the matrix  $A^h$ , since it turns out that this is a critical quantity in determining how rapidly certain iterative methods converge. Recall that the condition number is defined by

$$\text{cond}(A) = \|A\|_2 \|A^{-1}\|_2.$$

We've just seen that  $\|(A^h)^{-1}\|_2 \approx -1/2\pi^2$  for small  $h$ , and the norm of  $A$  is given by its spectral radius. The largest eigenvalue of  $A$  (in magnitude) is

$$\lambda^{m/2, m/2} \approx -\frac{4}{h^2}$$

and so

$$\text{cond}(A) \approx \frac{2}{\pi^2 h^2} = O(1/h^2) \quad \text{as } h \rightarrow 0. \quad (3.22)$$

The fact that the matrix becomes more ill-conditioned as we refine the grid is responsible for the slow-down of iterative methods such as Jacobi or Gauss-Seidel.

### 3.5 The nine-point Laplacian

Above we have used the 5-point Laplacian which we will denote by  $\nabla_5^2 u_{ij}$ , where this denotes the left-hand side of equation (3.10). Another possible approximation is the 9-point Laplacian

$$\nabla_9^2 u_{ij} = \frac{1}{6h^2} [4u_{i-1,j} + 4u_{i+1,j} + 4u_{i,j-1} + 4u_{i,j+1} + u_{i-1,j-1} + u_{i-1,j+1} + u_{i+1,j-1} + u_{i+1,j+1} - 20u_{ij}] \quad (3.23)$$

as indicated in Figure 3.1. If we apply this to the true solution and expand in Taylor series we find that

$$\nabla_9^2 u(x_i, y_j) = \nabla^2 u + \frac{1}{2} h^2 (u_{xxxx} + 2u_{xxyy} + u_{yyyy}) + O(h^4).$$

At first glance this discretization looks no better than the 5-point discretization since the error is still  $O(h^2)$ . However, the additional terms lead to a very nice form for the dominant error term, since

$$u_{xxxx} + 2u_{xxyy} + u_{yyyy} = \nabla^2 (\nabla^2 u).$$

This is the Laplacian of the Laplacian of  $u$ , otherwise known as the *biharmonic*. If we are solving  $\nabla^2 u = f$ , then we have

$$u_{xxxx} + 2u_{xxyy} + u_{yyyy} = \nabla^2 f.$$

Hence we can compute the dominant term in the truncation error easily from the known function  $f$  without knowing the true solution  $u$  to the problem.

In particular, if we are solving Laplace's equation, where  $f = 0$ , or more generally if  $f$  is a harmonic function, then this term in the local truncation error vanishes and the 9-point Laplacian would give a fourth order accurate discretization of the differential equation.

More generally, we can obtain a fourth-order accurate method of the form

$$\nabla_9^2 u_{ij} = F_{ij} \quad (3.24)$$

for arbitrary smooth functions  $f(x, y)$  by defining

$$F_{ij} = f(x_i, y_j) - \frac{h^2}{2} \nabla^2 f(x_i, y_j). \quad (3.25)$$

We can view this as deliberately introducing an  $O(h^2)$  error into the right hand side of the equation that is chosen to cancel the  $O(h^2)$  part of the local truncation error. Taylor series expansion easily shows that the local truncation error of the method (3.24) is now  $O(h^4)$ .

If we only have data  $f_{ij} = f(x_i, y_j)$  at the grid points rather than the function  $f$  (but we know that the underlying function is sufficiently smooth), then we can still achieve fourth order accuracy by using

$$F_{ij} = f_{ij} - \frac{h^2}{2} \nabla_5^2 f_{ij}$$

instead of (3.25).

This is a trick that can often be used in developing numerical methods — introducing an “error” into the equations that is carefully chosen to cancel some other error.

Note that the same trick wouldn't work with the 5-point Laplacian, or at least not as directly. The form of the truncation error in this method depends on  $u_{xxxx} + u_{yyyy}$ . There is no way to compute this directly from the original equation, without knowing  $u$ . The extra points in the 9-point stencil convert this into the Laplacian of  $f$ , which can be computed if  $f$  is sufficiently smooth.

On the other hand a two-pass approach could be used for with the 5-point stencil, in which we first estimate  $u$  by solving with the standard 5-point scheme to get a second order accurate estimate of  $u$ . We then use this estimate of  $u$  to approximate  $u_{xxxx} + u_{yyyy}$  and then solve a second time with a right hand side that is modified to eliminate the dominant term of the local truncation error. This would be more complicated for this particular problem, but this idea can be used much more generally than the above trick which depends on the special form of the Laplacian. This general approach is called the method of *deferred corrections*. For more details see, e.g., [Kel76], [AMR88].



## Chapter 4

# Function Space Methods

Finite difference methods determine an approximate solution only at discrete grid points. A “function space method” on the other hand, determines a function  $U(x)$  on the entire domain that approximates the true solution  $u(x)$ . In order to reduce the original differential equation to a system of algebraic equations that can be solved for a finite number of unknowns, we typically search for the approximation  $U(x)$  from a finite-dimensional function space that is spanned by a fixed set of basis functions  $\phi_j(x)$ , for  $j = 1, 2, \dots, N$  that are specified *a priori*. The function  $U(x)$  then has the form

$$U(x) = \sum_{j=1}^N c_j \phi_j(x).$$

If the functions  $\phi_j(x)$  are a basis for the space of all such functions, then they must be *linearly independent*. This means that the only way we can get the zero function from such a linear combination is if  $c_j = 0$  for all  $j$ . It then follows that any function in the span of these functions has a unique set of coefficients  $c_j$ .

In order to determine the coefficients  $c_j$  that yield a good approximation to  $u(x)$ , we need to derive a set of  $N$  equations to be solved for these coefficients. There are many ways in which this can be done. One example was seen in Section 4.2, where  $\phi_j(x) = \sin(j\pi x)$  were used as basis functions.

### 4.1 Collocation

One approach to determining the function  $U(x)$  is to require that this function satisfy the differential equation at some finite set of *collocation points*. Of course we would really like to have  $U(x)$  satisfy the differential equation at all points  $x$ , in which case it would be the exact solution, but this won't be possible generally unless the exact solution happens to lie in the finite-dimensional function space spanned by the  $\phi_j$ . Since there are  $N$  free parameters and also two boundary conditions that need to be satisfied (for our second order model problem), we can hope to satisfy the differential equation at some  $N - 2$  points, however.

**Example 4.2.** Consider our standard model problem  $u''(x) = f(x)$  with Dirichlet boundary conditions. Requiring that the boundary conditions be satisfied gives the two equations

$$\begin{aligned} \sum_{j=1}^N c_j \phi_j(0) &= \alpha, \\ \sum_{j=1}^N c_j \phi_j(1) &= \beta. \end{aligned} \tag{4.1}$$



We can then require in addition that the differential equation be satisfied at some points  $\xi_1, \dots, \xi_{N-2}$ :

$$\sum_{j=1}^N c_j \phi_j''(\xi_i) = f(\xi_i) \quad (4.2)$$

for  $i = 1, 2, \dots, N-2$ . The equations (4.1) and (4.2) together give  $N$  equations for the  $N$  unknowns  $c_j$ . Note that this may be a dense system of equations unless we are careful to choose our basis functions so that many of these coefficients are zero. A local basis such as a B-spline basis for spline functions is often used in practice. (An example of a local basis for piecewise linear functions is given in Section 4.3, but it is not a basis that is suitable for collocation on a second order differential equation.)

Another approach is to choose a set of functions and collocation points in such a way that “fast transform” methods can be used to solve the dense linear system, as in the fast Poisson solvers discussed in Section 3.3.2. This suggests using Fourier series to represent the solution. This turns out to be a very good idea, not only because the fast Fourier transform can be used to solve the resulting system, but also because smooth functions can be represented very accurately with relatively few terms of a Fourier series. Hence the order of accuracy of such a method can be very high, see Section 4.2.2.

## 4.2 Spectral methods

In Section 3.3.2 we looked at a method based on the FFT for solving the linear system  $Au = F$  that arises from the finite difference discretization studied in Chapter 2. Here we will study an entirely different approach to solving the original differential equation  $u''(x) = f(x)$ , approximating the solution by a Fourier series rather than by a discrete set of points. But the techniques are closely related (at least for this simple problem) as noted at the end of this section.

Here we will consider the simplest possible problem to illustrate the main ideas. The problem we consider is again  $u''(x) = f(x)$  on  $[0, 1]$ , and we will also assume that homogeneous Dirichlet boundary conditions  $u(0) = u(1) = 0$  are specified, and also that  $f(x)$  is a smooth function satisfying  $f(0) = f(1) = 0$ . These conditions are not necessary in general, but in this special case we can easily illustrate the so-called *spectral method* using the Fourier sine transform introduced in Section 3.3.2. The name comes from the fact that the set of eigenvalues of a matrix or operator is called its *spectrum*, and these methods are heavily based on the eigenstructure of the problem (at least for this simple problem).

Note that for this problem the functions  $\sin(j\pi x)$  satisfy the boundary conditions and differentiating twice gives a scalar multiple  $-(j\pi)^2$  times the original function. Hence these are indeed *eigenfunctions* of the operator  $\partial^2/\partial x^2$ . This is the reason using a sine-series expansion of the solution is valuable. Recall also that discretized versions of these functions are eigenvectors of the tridiagonal matrix arising from the standard finite-difference formulation, which is why the FFT can be used in the fast Poisson solver described in Section 3.3.2.

Spectral methods are exceedingly powerful tools for solving a wide class of differential equations, particularly when the solution we seek is a smooth function (so that Fourier series give a good representation) and the domain and boundary conditions are reasonably simple. This is not meant as a general overview of spectral methods; only as a brief introduction to the main idea. Better introductions can be found in many sources, e.g., [CHQZ88], [For96], [GO77].

The function  $f(x)$  can be expressed in a Fourier series as

$$f(x) = \sum_{j=1}^{\infty} \hat{f}_j \sin(j\pi x)$$

where

$$\hat{f}_j = 2 \int_0^1 f(x) \sin(j\pi x) dx. \quad (4.3)$$

Similarly, the unknown function  $u(x)$  can be expanded as

$$u(x) = \sum_{j=1}^{\infty} \hat{u}_j \sin(j\pi x). \quad (4.4)$$

Differentiating this series term by term (assuming this is valid, which it is for smooth functions) gives

$$u''(x) = \sum_{j=1}^{\infty} -(j\pi)^2 \hat{u}_j \sin(j\pi x)$$

and equating this with the series for  $f(x)$  shows that

$$\hat{u}_j = -\frac{1}{(j\pi)^2} \hat{f}_j. \quad (4.5)$$

We have just “solved” the differential equation and determined the exact solution, as a Fourier series. Note that this works only because  $\sin(j\pi x)$  is an eigenfunction of the differential operator. For a different differential equation, one would have to use the appropriate eigenfunctions in the series representation in place of the sine functions to achieve the same success.

If  $f(x)$  is sufficiently simple, we may be able to explicitly calculate the integrals in (4.3) and then the resulting infinite series in (4.4). (Though we are more likely to be able to integrate  $f$  twice to compute the exact solution for this trivial equation!) More generally we will not be able to compute the integrals or infinite series exactly, but we can use this approach as a basis for an approximate numerical method. We can approximate the series (4.4) by a truncated finite series,

$$U(x) = \sum_{j=1}^m \hat{U}_j \sin(j\pi x) \quad (4.6)$$

and approximate  $\hat{U}_j$  using (4.5) where  $\hat{f}_j$  is obtained by approximating the integral in (4.3).

For example, we could approximate  $\hat{f}_j$  by

$$\hat{F}_j = 2h \sum_{i=1}^m \sin(j\pi x_i) f(x_i) \quad (4.7)$$

where  $x_i = ih$  with  $h = 1/(m+1)$ . If we then calculate

$$\hat{U}_j = -\frac{1}{(j\pi)^2} \hat{F}_j \quad (4.8)$$

we can use (4.6) to compute an approximate value  $U(x) \approx u(x)$  at any point  $x$ . In particular, we could compute  $U(x_i)$  for each of the grid points  $x_i$  used in (4.7) to obtain an approximation to  $u(x)$  on a discrete grid, similar to what is produced with a finite difference method. Denote the resulting vector of grid values by  $U = [U(x_1), U(x_2), \dots, U(x_m)]^T$ . What has just been described is the *spectral method* for computing the approximate solution  $U$ .

Note that this is very closely related to what was done in the previous section. In fact, the sum in (4.7) can be written in vector form as

$$\hat{F} = 2hRF$$

exactly as in (3.16), where  $F = [f(x_1), f(x_2), \dots, f(x_m)]^T$  (just as it is in (3.16) in the case of homogeneous boundary conditions). Also, computing the vector  $U$  by evaluating (4.6) at each grid point is exactly the operation

$$U = R\hat{U}$$

as in (3.18). In practice each of these operations would be done with the FFT.

The only difference between the spectral method and the “fast Poisson solver” (for this trivial equation!) is the manner in which we compute  $\hat{U}_j$  from  $\hat{F}_j$ . In (3.17) we divided by the eigenvalues of the matrix  $A$ , since we were solving the linear system  $Au = F$ , whereas in (4.8) we divide by  $-(j\pi)^2$ , which is the  $j$ 'th eigenvalue of the differential operator  $-\partial^2/\partial x^2$  from the original differential equation. Intuitively this seems like a better thing to do, and indeed it is. The spectral approximation  $U$  converges to the solution of the differential equation much faster than the finite difference approximation  $u$ . (Note that in spite of using the exact eigenvalue,  $U$  is not the exact solution because it results from a truncated sine series and approximate integrals.)

### 4.2.1 Matrix interpretation

Note that this spectral method can be interpreted as solving a linear system of equations of the form  $BU = F$ . To see this, first note that (4.5) can be written in matrix form as

$$\hat{u} = M^{-1}\hat{F},$$

where  $M$  is the diagonal matrix  $M = \text{diag}(-\pi^2, -(2\pi)^2, \dots, -(m\pi)^2)$ . The  $j$ 'th diagonal element is just  $\mu^j$  from Section 2.10. Combining the various steps above then gives

$$U = R\hat{u} = RM^{-1}\hat{F} = RM^{-1}R^{-1}F$$

and hence  $U$  solves the system  $BU = F$  with

$$B = RMR^{-1}. \tag{4.9}$$

This matrix  $B$  can be interpreted as an approximation to the second derivative operator, analogous to the tridiagonal matrix  $A$  of (2.10). Unlike  $A$ , however, the spectral matrix  $B$  is a dense matrix. With the finite difference method we approximated each  $u''(x_i)$  using only 3 values  $u(x_{i-1})$ ,  $u(x_i)$ , and  $u(x_{i+1})$ . The spectral approximation uses *all* the values  $u(x_1), \dots, u(x_m)$  to approximate each  $u''(x_i)$ . Of course in practice we wouldn't want to form this matrix and solve the system using Gaussian elimination, since this would require  $O(m^3)$  operations rather than the  $O(m \log m)$  required using FFT's.

### 4.2.2 Accuracy

We can, however, use this matrix interpretation to help analyze the accuracy and stability of the method. We can define the local truncation error as usual by replacing the approximate solution with the true solution, and for this simple problem where  $f(x) = u''(x)$  we see that

$$\tau_i = (Bu)_i - u''(x_i)$$

which is just the error in the spectral approximation to the second derivative.

Based on the convergence properties of Fourier series, it can be shown that if the solution  $u(x)$  has  $p$  continuous derivatives, then the error in  $U$  decays like  $1/m^p$  as  $m \rightarrow \infty$ . Since  $h = 1/(m+1)$ , this means that the method is  $p$ 'th order accurate. For smooth solutions the method has a very high order of accuracy. If the solution is  $C^\infty$  (infinitely differentiable) then this is true for any value of  $p$  and it appears to be “infinite-order accurate”! This does not, however, mean that the error is zero. What it means is that it converges to zero faster than any power of  $1/m$ . Typically it is *exponentially fast* (for example the function  $1/2^m$  decays faster than any power of  $1/m$  as  $m \rightarrow \infty$ ). This is often called *spectral accuracy*.

There is a catch here however. In using a sine series representation of the function  $f$  or  $u$  we are obtaining an odd periodic function. In order for this to converge rapidly to the true function on the interval  $[0, 1]$  as we increase the number of terms in the series, we need the odd periodic extension of these functions to be sufficiently smooth, and not just the function specified in this interval. So in

deciding whether  $u(x)$  is  $p$  times continuously differentiable, we need to look at the function defined by setting

$$u(-x) = -u(x)$$

for  $-1 \leq x \leq 0$  and then extended periodically with period 2 from the interval  $[-1, 1]$  to the whole real line. This requires certain properties in  $u$  at the endpoints  $x = 0$  and  $x = 1$ . In particular, the extended  $u(x)$  is  $C^\infty$  only if all even derivatives of  $u$  vanish at these two points along with  $u$  being  $C^\infty$  in the interior.

Such difficulties mean that spectral methods based on Fourier series are most suitable in certain special cases (for example if we are solving a problem with periodic boundary conditions, in which case we expect the solution to be periodic and have the required smoothness). Methods based on similar ideas can be developed using other classes of functions rather than trigonometric functions, and are often used in practice. For example, families of orthogonal polynomials such as Chebyshev or Legendre polynomials can be used, and fast algorithms developed that achieve spectral accuracy.

### 4.2.3 Stability

To see that the results quoted above for the local error carry over to the global error as we refine the grid, we also need to check that the method is stable. Using the matrix interpretation of the method this is easy to do in the 2-norm. The matrix  $B$  in (4.9) is easily seen to be symmetric (recall that  $R^{-1} = 2hR = 2hR^T$  and so the 2-norm of  $B^{-1}$  is equal to its spectral radius, which is clearly  $1/\pi^2$  independent of  $h$ ). Hence the method is stable in the 2-norm.

### 4.2.4 Collocation property

Though it may not be obvious, the approximation we derived above for  $U(x)$  in fact satisfies  $U''(x_i) = f(x_i)$  at each of the points  $x_1$  through  $x_m$ . In other words this spectral method is also a special form of a collocation method, as described in Section 4.1.

## 4.3 The finite element method

The finite element method determines an approximate solution that is a linear combination of some specified basis functions in a very different way from collocation or expansion in eigenfunctions. This method is typically based on some “weak form” of the differential equation, which roughly speaking means that we have integrated the equation.

Consider, for example, the heat conduction problem in one dimension with a variable conductivity  $\kappa(x)$  so the steady-state equation is

$$(\kappa u')' = f. \tag{4.10}$$

Again for simplicity assume that the boundary conditions are  $u(0) = u(1) = 0$ . If we multiply both sides of the equation (4.10) by an arbitrary smooth function  $v(x)$  and integrate the resulting product over the domain  $[0, 1]$ , we obtain

$$\int_0^1 (\kappa(x)u'(x))'v(x) dx = \int_0^1 f(x)v(x) dx. \tag{4.11}$$

On the left-hand side we can integrate by parts. Since  $v$  is arbitrary, let's restrict our attention to  $v$  that satisfy  $v(0) = v(1) = 0$  so that the boundary terms drop out, yielding

$$-\int_0^1 \kappa(x)u'(x)v'(x) dx = \int_0^1 f(x)v(x) dx. \tag{4.12}$$

It can be shown that if  $u(x)$  satisfies this equation for all  $v$  in some suitable class of functions, then  $u(x)$  is in fact the solution to the original differential equation.

Now suppose we replace  $u(x)$  by an approximation  $U(x)$  in this expression, where  $U(x)$  is a linear combination of specified basis functions,

$$U(x) = \sum_{j=1}^m c_j \phi_j(x). \quad (4.13)$$

Let's suppose that our basis functions are chosen to satisfy  $\phi_j(0) = \phi_j(1) = 0$ , so that  $U(x)$  automatically satisfies the boundary conditions regardless of how we choose the  $c_j$ . Then we could try to choose the coefficients  $c_j$  in  $U(x)$  so that the equality (4.12) is satisfied for a large class of functions  $v(x)$ . Since we only have  $m$  free parameters, we can't require that (4.12) be satisfied for all smooth functions  $v(x)$ , but we can require that it be satisfied for all functions in some  $m$ -dimensional function space. Such a space is determined by a set of  $m$  basis functions  $\psi_i(x)$  (which might or might not be the same as the functions  $\phi_j(x)$ ). If we require that (4.12) be satisfied for the special case where  $v$  is chosen to be any one of these functions, then by linearity (4.12) will also be satisfied for any  $v$  that is an arbitrary linear combination of these functions, and hence for all  $v$  in this  $m$ -dimensional linear space.

Hence we are going to require that

$$-\int_0^1 \kappa(x) \left( \sum_{j=1}^m c_j \phi_j'(x) \right) \psi_i'(x) dx = \int_0^1 f(x) \psi_i(x) dx \quad (4.14)$$

for  $i = 1, 2, \dots, m$ . We can rearrange this to give

$$\sum_{j=1}^m K_{ij} c_j = \int_0^1 f(x) \psi_i(x) dx \quad (4.15)$$

where

$$K_{ij} = -\int_0^1 \kappa(x) \phi_j'(x) \psi_i'(x) dx. \quad (4.16)$$

The equations (4.15) for  $i = 1, 2, \dots, m$  give an  $m \times m$  linear system of equations to solve for the  $c_j$ , which we could write as

$$Kc = F$$

with

$$F_i = \int_0^1 f(x) \psi_i(x) dx. \quad (4.17)$$

The functions  $\psi_i$  are generally called "test functions" while the basis functions  $\phi_i$  for our approximate solution are called "trial functions". Frequently the same basis functions are used for both spaces. The resulting method is known as the *Galerkin method*. If the trial space is different from the test space we have a *Petrov-Galerkin method*.

**Example 4.3.** As a specific example, consider the Galerkin method for the above problem with basis functions defined as follows on a uniform grid with  $x_i = ih$ , and  $h = 1/(m+1)$ . The  $j$ 'th basis function  $\phi_j(x)$  is

$$\phi_j(x) = \begin{cases} (x - x_{j-1})/h & \text{if } x_{j-1} \leq x \leq x_j \\ (x_{j+1} - x)/h & \text{if } x_j \leq x \leq x_{j+1} \\ 0 & \text{otherwise} \end{cases} \quad (4.18)$$

Each of these functions is continuous and piecewise linear, and  $\phi_j(x)$  takes the value 1 at  $x_j$  and the value 0 at all other nodes  $x_i$  for  $i \neq j$ . (See Figure 4.1(a).) Note that any linear combination (4.13) of these functions will still be continuous and piecewise linear, and will take the value  $x_i$  at the point  $x_i$

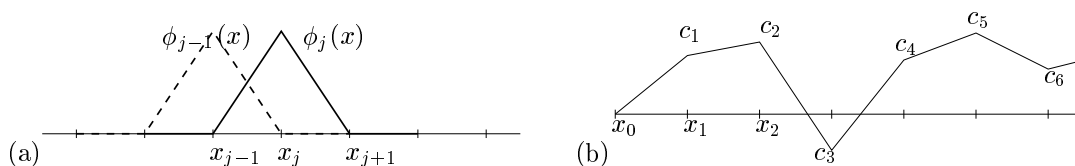


Figure 4.1: (a) Two typical basis functions  $\phi_{j-1}(x)$  and  $\phi_j(x)$  with continuous piecewise linear elements. (b)  $U(x)$ , a typical linear combination such basis functions.

since  $U(x_i) = \sum_j c_j \phi_j(x_i) = c_i$  since all other terms in the sum are zero. Hence the function  $U(x)$  has the form shown in Figure 4.1(b).

The set of functions  $\{\phi_j(x)\}$  form a basis for the space of all continuous piecewise linear functions defined on  $[0, 1]$  with  $u(0) = u(1) = 0$  and with kinks at the points  $x_1, x_2, \dots, x_m$ , which are called the *nodes*. Note that the coefficient  $c_j$  can be interpreted as the value of the approximate solution at the point  $x_j$ .

To use these basis functions in the Galerkin equations (4.14) (with  $\psi_j = \phi_j$ ), we need to compute the derivatives of these basis functions and then the elements of the matrix  $K$  and right-hand side  $F$ . We have

$$\phi_j'(x) = \begin{cases} 1/h & \text{if } x_{j-1} \leq x \leq x_j \\ -1/h & \text{if } x_j \leq x \leq x_{j+1} \\ 0 & \text{otherwise.} \end{cases}$$

For general functions  $\kappa(x)$  we might have to compute an approximation to the integral in (4.16), but as a simple example consider the case  $\kappa(x) \equiv 1$  (so the equation is just  $u''(x) = f(x)$ ). Then we can compute that

$$K_{ij} = - \int_0^1 \phi_j'(x) \phi_i'(x) dx = \begin{cases} 1/h & \text{if } j = i - 1 \text{ or } j = i + 1, \\ -2/h & \text{if } j = i, \\ 0 & \text{otherwise.} \end{cases}$$

The matrix  $K$  is quite familiar (except for the different power of  $h$ ):

$$K = \frac{1}{h} \begin{bmatrix} -2 & 1 & & & & \\ 1 & -2 & 1 & & & \\ & 1 & -2 & 1 & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & -2 & 1 \\ & & & & 1 & -2 \end{bmatrix}. \quad (4.19)$$

In some cases we may be able to evaluate the integral in (4.17) for  $F_i$  explicitly. More generally we might use a discrete approximation. Note that since  $\phi_i(x)$  is nonzero only near  $x_i$ , and  $\int_0^1 \phi_i(x) dx = h$ , this is roughly

$$F_i \approx h f(x_i). \quad (4.20)$$

In fact the trapezoidal method applied to this integral on the same grid would give exactly this result. Using (4.20) in the system  $Kc = F$ , and dividing both sides by  $h$ , gives exactly the same linear system of equations that we obtained in Section 2.4 from the finite difference method (for the case  $\alpha = \beta = 0$  we are considering here).

**Exercise 4.1** *If we have more general Dirichlet boundary conditions  $u(0) = \alpha$  and  $u(1) = \beta$ , we can introduce two additional basis functions  $\phi_0(x)$  and  $\phi_{m+1}(x)$  which are also defined by (4.18). Then we know  $c_0 = \alpha$  and  $c_{m+1} = \beta$  and these terms in the extended sum appearing in (4.12) can be moved to the right hand side. Carry this through to see that we get essentially the system (2.9) in this more general case as well.*

Some comments on this method:

- The matrix  $K$  above is tridiagonal because each  $\phi_j(x)$  is nonzero on only two elements, for  $x_{j-1} < x < x_{j+1}$ . The function  $\phi'_i(x)\phi'_j(x)$  is identically zero unless  $j = i - 1, i$  or  $i + 1$ . More generally, if we choose basis functions that are nonzero only on some region  $x_{j-b} < x < x_{j+a}$ , then the resulting matrix would be banded with  $b$  diagonals of nonzeros below the diagonal and  $a$  bands above. In the finite element method one almost always chooses local basis functions of this sort, that are each nonzero over only a few elements.
- Why did we integrate by parts to obtain equation (4.12), rather than working directly with (4.11)? One could go through the same process based on (4.11), but then we would need an approximate  $U(x)$  with meaningful second derivatives. This would rule out the use of the simple piecewise linear basis functions used above. (Note that the piecewise linear functions don't have meaningful first derivatives at the nodes, but since only integrals of these functions are used in defining the matrix this is not a problem.)
- This is one advantage of the finite element method over collocation, for example. One can often use functions  $U(x)$  for which the original differential equation does not even make sense because  $U$  is not sufficiently smooth.
- There are other good reasons for integrating by parts. The resulting equation (4.12) can also be derived from a variational principle and has physical meaning in terms of minimizing the “energy” in the system. (See, e.g., [SF73].)

### 4.3.1 Two space dimensions

In the last example we saw that the one-dimensional finite element method based on piecewise linear elements is equivalent to the finite difference method derived in Section 2.4. Since it is considerably more complicated to derive via the finite element approach, this may not seem like a useful technique. However, in more than one dimension this method can be extended to irregular grids on complicated regions for which it would not be so easy to derive a finite difference method.

Consider, for example, the Poisson problem with homogeneous Dirichlet boundary conditions on the region shown in Figure 4.2, which also shows a fairly coarse “triangulation” of the region. The points  $(x_j, y_j)$  at the corners of the triangles are called *nodes*. The Galerkin form of the Poisson problem is

$$-\int \int_{\Omega} \nabla q \cdot \nabla v \, dx \, dy = \int \int_{\Omega} f v \, dx \, dy. \quad (4.21)$$

This should hold for all test functions  $v(x, y)$  in some class. Again we can approximate  $u(x, y)$  by some linear combination of specified basis functions:

$$U(x, y) = \sum_{j=1}^N c_j \phi_j(x, y). \quad (4.22)$$

Taking an approach analogous to the one-dimensional case above, we can define a basis function  $\phi_j(x, y)$  associated with each node  $(x_j, y_j)$  to be the unique function that is linear on each triangle, and which takes the value 1 at the node  $(x_j, y_j)$  and 0 at all other nodes. This function is continuous across the boundaries between triangles and nonzero only for the triangles that have Node  $j$  as a corner. For example, Figure 4.2 indicates contour lines for the basis function  $\phi_8(x, y)$  as dashed lines.

Using (4.22) in (4.21) gives an  $N \times N$  linear system of the form  $Kc = F$  where

$$K_{ij} = - \int \int_{\Omega} \nabla \phi_j \cdot \nabla \phi_i \, dx \, dy. \quad (4.23)$$

These gradients are easy to compute and in fact are constant within each triangle since the basis function is linear there. Since  $\nabla \phi_i$  is identically zero in any triangle for which Node  $i$  is not a corner, we see

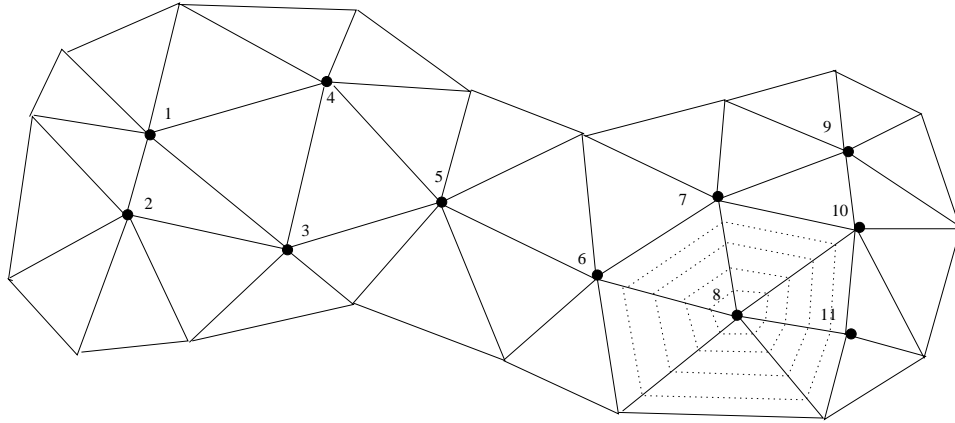


Figure 4.2: Triangulation of a two-dimensional region with 11 nodes. Contourlines for the basis function  $\phi_8(x, y)$  are also shown as dashed lines.

that  $K_{ij} = 0$  unless Nodes  $i$  and  $j$  are two corners of a common triangle. For example, in Figure 4.2 the eighth row of the matrix  $K$  will have nonzeros only in columns 6, 7, 8, 10, and 11.

Note also that  $K$  will be a symmetric matrix, since the expression (4.23) is symmetric in  $i$  and  $j$ . It can also be shown to be positive definite.

For a typical triangulation on a much finer grid, we would have a large but very sparse matrix  $K$ . The structure of the matrix, however, will not generally be as simple as what we would obtain with a finite difference method on a rectangular grid. The pattern of nonzeros will depend greatly on how we order the unknowns and equations. Direct methods for solving such systems rely greatly on algorithms for ordering them to minimize the bandwidth. See [DER86].





## Chapter 5

# Iterative Methods for Sparse Linear Systems

This chapter contains a brief overview of several iterative methods for solving the large sparse linear systems that arise from elliptic equations, either from finite-difference approximations (Chapter 3) or from finite element approximations (Section 4.3). Large sparse linear systems arise from many other practical problems too, of course, and the methods discussed here are important more generally.

The classical Jacobi and Gauss-Seidel methods have already been introduced in Section 3.3.3 for the 5-point Laplacian, and their extension to other problems is straightforward. It is not clear, however, whether these methods should be expected to converge, or how quickly. Here we will analyze these methods based on viewing them in terms of matrix splittings.

The SOR and conjugate gradient methods will also be briefly introduced. The reader can find considerably more theoretical analysis of these methods in the literature. See for example [GO92], [Gre97], [HY81], [Var62], [You71].

### 5.1 Matrix splitting methods

As an example we will consider the one-dimensional analog of the Poisson problem,  $u''(x) = f(x)$  as discussed in Chapter 2. Then we have a tridiagonal system of equations (2.9) to solve. In practice we would never use an iterative method for this system, since it can be solved directly by Gaussian elimination in  $O(m)$  operations, but it is easier to illustrate the iterative methods in the one-dimensional case and all of the analysis done here carries over almost unchanged to the 2-dimensional (or even 3-dimensional) case.

The Jacobi and Gauss-Seidel methods for this problem take the form

$$\begin{aligned} \text{Jacobi:} \quad u_i^{[k+1]} &= \frac{1}{2}(u_{i-1}^{[k]} + u_{i+1}^{[k]} - h^2 f_i), \\ \text{Gauss-Seidel:} \quad u_i^{[k+1]} &= \frac{1}{2}(u_{i-1}^{[k+1]} + u_{i+1}^{[k]} - h^2 f_i), \end{aligned} \tag{5.1}$$

Both of these methods can be analyzed by viewing them as based on a splitting of the matrix  $A$  into

$$A = M - N \tag{5.2}$$

where  $M$  and  $N$  are two  $m \times m$  matrices. Then the system  $Au = f$  can be written as

$$Mu - Nu = f \quad \implies \quad Mu = Nu + f,$$

which suggests the iterative method

$$Mu^{[k+1]} = Nu^{[k]} + f. \tag{5.3}$$



If  $e^{[k]} = u^{[k]} - u^*$  represents the error, then subtracting (5.6) from (5.5) gives

$$e^{[k+1]} = Ge^{[k]},$$

and so after  $k$  steps we have

$$e^{[k]} = G^k e^{[0]}. \quad (5.7)$$

From this we can see that the method will converge from any initial guess  $u^{[0]}$  provided  $G^k \rightarrow 0$  (an  $m \times m$  matrix of zeros) as  $k \rightarrow \infty$ . When is this true?

For simplicity, assume that  $G$  is a diagonalizable matrix, so that we can write

$$G = R\Gamma R^{-1}$$

where  $R$  is the matrix of right eigenvectors of  $G$  and  $\Gamma$  is a diagonal matrix of eigenvalues  $\gamma_1, \gamma_2, \dots, \gamma_m$ . Then

$$G^k = R\Gamma^k R^{-1}, \quad (5.8)$$

where

$$\Gamma^k = \begin{bmatrix} \gamma_1^k & & & \\ & \gamma_2^k & & \\ & & \ddots & \\ & & & \gamma_m^k \end{bmatrix}.$$

Clearly the method converges if  $|\gamma_p| < 1$  for all  $p = 1, 2, \dots, m$ , i.e., if  $\rho(G) < 1$  where  $\rho$  is the spectral radius.

### 5.1.1 Rate of convergence

From (5.7) we can also determine how rapidly the method can be expected to converge in cases where it is convergent. Using (5.8) in (5.7) and using the 2-norm, we obtain

$$\|e^{[k]}\|_2 \leq \|\Gamma^k\|_2 \|R\|_2 \|R^{-1}\|_2 \|e^{[0]}\|_2 = \rho^k \text{cond}(R) \|e^{[0]}\|_2 \quad (5.9)$$

where  $\rho \equiv \rho(G)$  and  $\text{cond}(R) = \|R\|_2 \|R^{-1}\|_2$  is the condition number of the eigenvector matrix.

If the matrix  $G$  is a normal matrix (meaning it commutes with its transpose, in particular if it is symmetric as when Jacobi is applied to the Poisson problem), then the eigenvectors are orthogonal and  $\text{cond}(R) = 1$ . In this case we have

$$\|e^{[k]}\|_2 \leq \rho^k \|e^{[0]}\|_2. \quad (5.10)$$

**Note:** These methods are *linearly* convergent, in the sense that  $\|e^{[k+1]}\| \leq \rho \|e^{[k]}\|$  and it is the first power of  $\|e^{[k]}\|$  that appears on the right. Recall that Newton's method is typically quadratically convergent, and it is the square of the previous error that appears on the right hand side. But Newton's method is for a nonlinear problem, and requires solving a linear system in each iteration. Here we are looking at solving such a linear system.

The *average rate of convergence* for a method after  $k$  iterations is sometimes defined by

$$R_k(G) = -\frac{1}{k} \log \|G^k\|,$$

while the *asymptotic rate of convergence* is

$$R_\infty(G) = -\log(\rho(G)).$$

**Example 5.4.** For the Jacobi method we have

$$G = D^{-1}(D - A) = I - D^{-1}A.$$

If we apply this method to the boundary value problem  $u'' = f$ , then

$$G = I + \frac{h^2}{2}A.$$

The eigenvectors of this matrix are the same as the eigenvectors of  $A$ , and the eigenvalues are hence

$$\gamma_p = 1 + \frac{h^2}{2}\lambda^p$$

where  $\lambda^p$  is given by (2.23). So

$$\gamma_p = \cos(p\pi h), \quad p = 1, 2, \dots, m,$$

where  $h = 1/(m + 1)$ . The spectral radius is

$$\rho(G) = |\gamma_1| = \cos(\pi h) \approx 1 - \frac{1}{2}\pi^2 h^2 + O(h^2). \quad (5.11)$$

The spectral radius is less than one for any  $h > 0$  and the Jacobi method converges, but we see that  $\rho(G) \rightarrow 1$  as  $h \rightarrow 0$  and for small  $h$  is very close to one, resulting in very slow convergence.

How many iterations are required to obtain a good solution? Suppose we want to reduce the error to  $\|e^{[k]}\| \approx \epsilon \|e^{[0]}\|$  (where typically  $\|e^{[0]}\|$  is on the order of 1).<sup>1</sup> Then we want  $\rho^k \approx \epsilon$  and so

$$k \approx \log(\epsilon) / \log(\rho). \quad (5.12)$$

How small should we choose  $\epsilon$ ? To get full machine precision we might choose  $\epsilon$  to be close to the machine roundoff level. However, this would typically be very wasteful. For one thing, we rarely need this many correct digits. More importantly, however, we should keep in mind that even the *exact* solution  $u^*$  of the linear system  $Au = f$  is only an *approximate* solution of the differential equation we are actually solving. If we are using a second order accurate method, as in this example, then  $u_i^*$  differs from  $u(x_i)$  by something on the order of  $h^2$  and so we cannot achieve better accuracy than this no matter how well we solve the linear system. In practice we should thus take  $\epsilon$  to be something related to the expected global error in the solution, e.g.,  $\epsilon = Ch^2$  for some fixed  $C$ .

To estimate the order of work required asymptotically as  $h \rightarrow 0$ , we see that the above choice gives

$$k = (\log(C) + 2\log(h)) / \log(\rho). \quad (5.13)$$

For Jacobi on the boundary value problem we have  $\rho \approx 1 - \frac{1}{2}\pi^2 h^2$  and hence  $\log(\rho) \approx -\frac{1}{2}\pi^2 h^2$ . Since  $h = 1/(m + 1)$ , using this in (5.13) gives

$$k = O(m^2 \log m) \text{ as } m \rightarrow \infty. \quad (5.14)$$

Since each iteration requires  $O(m)$  work in this one-dimensional problem, the total work required to solve the problem goes like

$$\text{total work} = O(m^3 \log m).$$

Of course this tridiagonal problem can be solved exactly in  $O(m)$  work, so we would be very foolish to use an iterative method at all here!

---

<sup>1</sup>Assuming we are using some grid function norm, as discussed in Appendix A1. Note that for the 2-norm in one dimension this requires introducing a factor of  $\sqrt{h}$  in the definition of both  $\|e^{[k]}\|$  and  $\|e^{[0]}\|$ , but these factors cancel out in choosing an appropriate  $\epsilon$ .

For a Poisson problem in 2 or 3 dimensions it can be verified that (5.14) still holds, though now the work required per iteration is  $O(m^2)$  or  $O(m^3)$  respectively if there are  $m$  grid points in each direction. In 2 dimensions we would thus find that

$$\text{total work} = O(m^4 \log m). \quad (5.15)$$

Recall from Section 3.3 that Gaussian elimination on the banded matrix requires  $O(m^4)$  operations while other direct methods can do much better, so Jacobi is still not competitive. Luckily there are much better iterative methods.

For the Gauss-Seidel method applied to the Poisson problem, it can be shown that

$$\rho(G) = 1 - \pi^2 h^2 + O(h^4) \text{ as } h \rightarrow 0. \quad (5.16)$$

This still approaches 1 as  $h \rightarrow 0$ , but is better than (5.11) by a factor of 2 and the number of iterations required to reach a given tolerance will typically be half the number required with Jacobi. The order of magnitude figure (5.15) still holds, however, and this method is also not widely used.

### 5.1.2 SOR

If we look at how iterates  $u^{[k]}$  behave when Gauss-Seidel is applied to a typical problem, we will generally see that  $u_i^{[k+1]}$  is closer to  $u_i^*$  than  $u_i^{[k]}$  was, but only by a little bit. The Gauss-Seidel update moves  $u_i$  in the right direction, but is far too conservative in the amount it allows  $u_i$  to move. This suggests that we use the following two-stage update, illustrated again for the problem  $u'' = f$ :

$$\begin{aligned} u_i^{\text{GS}} &= \frac{1}{2}(u_{i-1}^{[k+1]} + u_{i+1}^{[k]} - h^2 f_i) \\ u_i^{[k+1]} &= u_i^{[k]} + \omega(u_i^{\text{GS}} - u_i^{[k]}) \end{aligned} \quad (5.17)$$

where  $\omega$  is some scalar parameter. If  $\omega = 1$  then  $u_i^{[k+1]} = u_i^{\text{GS}}$  is the Gauss-Seidel update. If  $\omega > 1$  then we move further than Gauss-Seidel suggests. In this case the method is known as *successive overrelaxation* (SOR).

If  $\omega < 1$  then we would be underrelaxing, rather than overrelaxing. This would be even less effective than Gauss-Seidel as a stand-alone iterative method for most problems, though underrelaxation is sometimes used in connection with multigrid methods[?].

The formulas in (5.17) can be combined to yield

$$u_i^{[k+1]} = \frac{\omega}{2}(u_{i-1}^{[k+1]} + u_{i+1}^{[k]} - h^2 f_i) + (1 - \omega)u_i^{[k]}, \quad (5.18)$$

and it can be verified that this corresponds to a matrix splitting method with

$$M = \frac{1}{\omega}(D - \omega L), \quad N = \frac{1}{\omega}((1 - \omega)D + \omega U).$$

Analyzing this method is considerably trickier than the Jacobi or Gauss-Seidel methods because of the form of these matrices. A theorem of Ostrowski states that if  $A$  is symmetric positive definite and  $D - \omega L$  is nonsingular, then the SOR method converges for all  $0 < \omega < 2$ . Young[You50] showed how to find the optimal  $\omega$  to obtain the most rapid convergence for a wide class of problems (including the Poisson problem). This elegant theory can be found in many introductory texts. (For example, see [GO92], [HY81], [Var62], [You71]. See also [LT88] for a different introductory treatment based on Fourier series and modified equations in the sense of Chapter 16, and [ALY88] for applications of this approach to the 9-point Laplacian.)

For the Poisson problem it can be shown that the SOR method converges most rapidly if  $\omega$  is chosen as

$$\omega_{\text{opt}} = \frac{2}{1 + \sin(\pi h)} \approx 2 - 2\pi h.$$

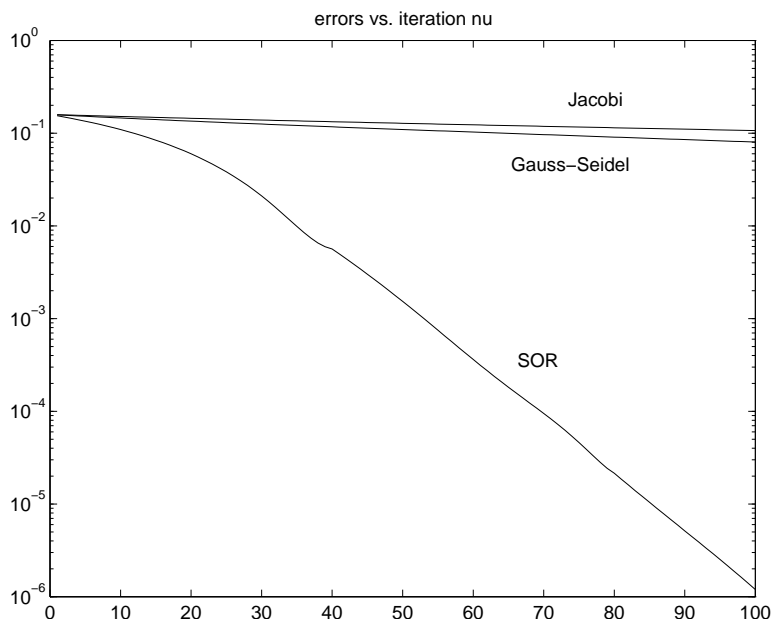


Figure 5.1: Errors vs.  $k$  for three methods.

This is nearly equal to 2 for small  $h$ . One might be tempted to simply set  $\omega = 2$  in general, but this would be a poor choice since SOR does not then converge! In fact the convergence rate is quite sensitive to the value of  $\omega$  chosen. With the optimal  $\omega$  it can be shown that the spectral radius of the corresponding  $G$  matrix is

$$\rho_{\text{opt}} = \omega_{\text{opt}} - 1 \approx 1 - 2\pi h,$$

but if  $\omega$  is changed slightly this can deteriorate substantially.

Even with the optimal  $\omega$  we see that  $\rho_{\text{opt}} \rightarrow 1$  as  $h \rightarrow 0$ , but only linearly in  $h$  rather than quadratically as with Jacobi or Gauss-Seidel. This makes a substantial difference in practice. The expected number of iterations to converge to the required  $O(h^2)$  level, the analogue of (5.14), is now

$$k_{\text{opt}} = O(m \log m).$$

Figure 5.1 shows some computational results for the methods described above on the 2-point boundary value problem.

## 5.2 Conjugate gradient methods

To appear. See [Gre97].

### 5.2.1 Preconditioners

To appear.

## 5.3 Multigrid methods

To appear. See [Bri87], [Jes84].

## Chapter 6

# The Initial Value Problem for ODE's

In this chapter we begin a study of time-dependent differential equations, beginning with the initial value problem (IVP) for a time-dependent ODE. Standard introductory texts are Lambert[Lam73] and Gear[Gea71]. Henrici[Hen62] gives a more complete description of some theoretical issues, although stiff equations are not discussed. Hairer, Norsett, and Wanner[HNW87, HNW93] is a more recent and complete survey of the field.

The initial value problem takes the form

$$u'(t) = f(u(t), t) \text{ for } t > t_0 \quad (6.1)$$

with some initial data

$$u(t_0) = \eta. \quad (6.2)$$

We will often assume  $t_0 = 0$  for simplicity.

In general (6.1) may represent a system of ODEs, i.e.,  $u$  may be a vector with  $s$  components  $u_1, \dots, u_s$  and then  $f(u, t)$  also represents a vector with components  $f_1(u, t), \dots, f_s(u, t)$ , each of which can be a nonlinear function of all the components of  $u$ . The problem is *linear* if

$$f(u, t) = A(t)u + b(t) \quad (6.3)$$

where  $A(t) \in \mathbb{R}^{s \times s}$  and  $b(t) \in \mathbb{R}^s$ .

We will consider only the first order equation (6.1) but in fact this is more general than it appears since we can reduce higher order equations to a system of first order equations.

**Example 6.1.** Consider the initial value problem for the ODE

$$u'''(t) = u'(t)u(t) - 2t(u''(t))^2 \text{ for } t > 0.$$

This third order equation requires three initial conditions, typically specified as

$$\begin{aligned} u(0) &= \eta_1 \\ u'(0) &= \eta_2 \\ u''(0) &= \eta_3 \end{aligned} \quad (6.4)$$

We can rewrite this as a system of the form (6.1) and (6.2) by introducing the variables

$$\begin{aligned} u_1(t) &= u(t) \\ u_2(t) &= u'(t) \\ u_3(t) &= u''(t). \end{aligned}$$



Then the equations take the form

$$\begin{aligned}u_1'(t) &= u_2(t) \\u_2'(t) &= u_3(t) \\u_3'(t) &= u_1(t)u_2(t) - 2tu_3^2(t)\end{aligned}$$

which defines the vector function  $f(u, t)$ . The initial condition is simply (6.2) where the three components of  $\eta$  come from (6.4). More generally, any single equation of order  $m$  can be reduced to  $m$  first order equations, and an  $m$ th order system of  $s$  equations can be reduced to a system of  $ms$  first order equations.

It is also sometimes useful to note that any explicit dependence of  $f$  on  $t$  can be eliminated by introducing a new variable that is simply equal to  $t$ . In the above example we could define

$$u_4(t) = t$$

so that

$$u_4'(t) = 1 \quad \text{and} \quad u_4(0) = 0.$$

The system then takes the form

$$u'(t) = f(u(t)) \tag{6.5}$$

with

$$f(u) = \begin{bmatrix} u_2 \\ u_3 \\ u_1u_2 - 2u_4u_3^2 \\ 1 \end{bmatrix} \quad \text{and} \quad u(0) = \begin{bmatrix} \eta_1 \\ \eta_2 \\ \eta_3 \\ 0 \end{bmatrix}.$$

The equation (6.5) is said to be *autonomous* since it does not depend explicitly on time. It is often convenient to assume  $f$  is of this form since it simplifies notation.

We will always assume that  $f$  is Lipschitz continuous in  $u$  as described in the next section, which implies that the initial value problem has a unique solution over some time interval.

## 6.1 Lipschitz continuity

The standard theory for the existence of a solution to the initial value problem

$$u'(t) = f(u, t), \quad u(0) = \eta \tag{6.6}$$

is discussed in many texts, e.g., [CL55]. To guarantee that there is a unique solution it is necessary to require a certain amount of smoothness in the function  $f(u, t)$  of (6.6). What we require is that the function  $f(u, t)$  be *Lipschitz continuous* in  $u$  over some range of  $t$  and  $u$ , i.e., that there exist some constant  $L > 0$  so that

$$|f(u, t) - f(u^*, t)| \leq L|u - u^*| \tag{6.7}$$

for all  $u$  and  $u^*$  in this range. This is slightly stronger than mere continuity, which only requires that  $|f(u) - f(u^*)| \rightarrow 0$  as  $u \rightarrow u^*$ . Lipschitz continuity requires that  $|f(u) - f(u^*)| = O(|u - u^*|)$  as  $u \rightarrow u^*$ . The function is *uniformly Lipschitz continuous* if there is a single constant  $L$  that works for all  $u$  and  $u^*$ .

If  $f(u, t)$  is differentiable with respect to  $u$  and this derivative  $f_u = \partial f / \partial u$  is bounded then we can take

$$L = \max_u |f_u(u, t)|,$$

since

$$f(u) = f(u^*) + f_u(v)(u - u^*)$$

for some value  $v$  between  $u$  and  $u^*$ .

**Example 6.2.** For the linear problem  $u'(t) = \lambda u(t) + g(t)$ ,  $f'(u) \equiv \lambda$  and we can take  $L = |\lambda|$ . This problem of course has a unique solution for any initial data  $\eta$  given by

$$u(t) = e^{\lambda(t-t_0)}\eta + \int_{t_0}^t e^{\lambda(t-\tau)}g(\tau) d\tau. \quad (6.8)$$

In particular, if  $\lambda = 0$  then  $L = 0$ . In this case  $f(u, t) = g(t)$  is independent of  $u$ . The solution is then obtained by simply integrating the function  $g(t)$ ,

$$u(t) = \eta + \int_{t_0}^t g(\tau) d\tau. \quad (6.9)$$

### 6.1.1 Existence and uniqueness of solutions

The basic existence and uniqueness theorem states that if  $f$  is uniformly Lipschitz continuous over some time period  $0 \leq t \leq T$  then there is a unique solution to the initial value problem (6.6) from any initial value  $\eta$ . If  $f$  is Lipschitz but not uniformly so, then there will be a unique solution through any value  $\eta$  over some finite time interval, but this solution may not exist past some time, as the next example shows.

**Example 6.3.** Consider the initial value problem

$$u'(t) = u(t)^2$$

with initial conditions

$$u(0) = \eta > 0.$$

The function  $f(u) = u^2$  is Lipschitz continuous over any finite interval  $[\eta - a, \eta + a]$  with  $L = 2(\eta + a)$ . From this it can be shown that the initial value problem has a unique solution over *some* time interval  $0 \leq t < \bar{T}$  with  $\bar{T} > 0$ . However, since  $f$  is not uniformly Lipschitz for all  $u$  (i.e., there is not a single value of  $L$  that works for all  $u$ ), we cannot prove that a solution exists for all time, and in fact it does not. The solution to the initial value problem is

$$u(t) = \frac{1}{\eta^{-1} - t}$$

and so  $u(t) \rightarrow \infty$  as  $t \rightarrow 1/\eta$ . There is no solution beyond time  $1/\eta$ .

If the function  $f$  is not Lipschitz continuous at some point then the initial value problem may fail to have a unique solution over any time interval.

**Example 6.4.** Consider the initial value problem

$$u'(t) = \sqrt{u(t)}$$

with initial conditions

$$u(0) = 0.$$

The function  $f(u) = \sqrt{u}$  is not Lipschitz continuous near  $u = 0$  since  $f'(u) = 1/(2\sqrt{u}) \rightarrow \infty$  as  $u \rightarrow 0$ . We cannot find a constant  $L$  so that the bound (6.7) holds for all  $u$  and  $u^*$  near 0.

As a result, this initial value problem does not have a unique solution. In fact it has two distinct solutions:

$$u(t) \equiv 0$$

and

$$u(t) = \frac{1}{4}t^2.$$

### 6.1.2 Systems of equations

For systems of  $m > 1$  ordinary differential equations,  $u(t) \in \mathbb{R}^m$  and  $f(u, t)$  is a function mapping  $\mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}^m$ . We say the function  $f$  is Lipschitz continuous in some norm  $\|\cdot\|$  if there is a constant  $L$  such that

$$\|f(u, t) - f(u^*, t)\| \leq L\|u - u^*\| \quad (6.10)$$

for all  $u$  in a neighborhood of  $u^*$ . By the equivalence of finite-dimensional norms (Appendix A1), if  $f$  is Lipschitz continuous in one norm then it is Lipschitz continuous in any other norm, though the Lipschitz constant may depend on the norm chosen.

The theorems on existence and uniqueness carry over to systems of equations.

**Example 6.5.** Consider the pendulum problem from Section 2.15,

$$\theta''(t) = -\sin(\theta(t)),$$

which can be rewritten as a first order system of two equations by introducing  $v(t) = \theta'(t)$ :

$$u = \begin{bmatrix} \theta \\ v \end{bmatrix}, \quad \frac{d}{dt} \begin{bmatrix} \theta \\ v \end{bmatrix} = \begin{bmatrix} v \\ -\sin(\theta) \end{bmatrix}.$$

Consider the max-norm. We have

$$\|u - u^*\|_\infty = \max(|\theta - \theta^*|, |v - v^*|)$$

and

$$\|f(u) - f(u^*)\|_\infty = \max(|v - v^*|, |\sin(\theta) - \sin(\theta^*)|).$$

To bound  $\|f(u) - f(u^*)\|_\infty$ , first note that  $|v - v^*| \leq \|u - u^*\|_\infty$ . We also have

$$|\sin(\theta) - \sin(\theta^*)| \leq |\theta - \theta^*| \leq \|u - u^*\|_\infty$$

since the derivative of  $\sin(\theta)$  is bounded by 1. So we have Lipschitz continuity with  $L = 1$ :

$$\|f(u) - f(u^*)\|_\infty \leq \|u - u^*\|_\infty.$$

### 6.1.3 Significance of the Lipschitz constant

The Lipschitz constant measures how much  $f(u, t)$  changes if we perturb  $u$  (at some fixed time  $t$ ). Since  $f(u, t) = u'(t)$ , the slope of the line tangent to the solution curve through the value  $u$ , this indicates how the slope of the solution curve will vary if we perturb  $u$ . The significance of this is best seen through some examples.

**Example 6.6.** Consider the trivial equation  $u'(t) = g(t)$ , which has Lipschitz constant  $L = 0$  and solutions given by (6.8). Several solution curves are sketched in Figure 6.1. Note that all of these curves are “parallel”; they are simply shifted depending on the initial data. Several tangent lines at one particular time are indicated. They are all parallel since  $f(u, t) = g(t)$  is independent of  $u$ .

**Example 6.7.** Consider  $u'(t) = \lambda u(t)$  with  $\lambda$  constant and  $L = |\lambda|$ . Then  $u(t) = u(0) \exp(\lambda t)$ . Two situations are shown in Figure 6.2 for negative and positive values of  $\lambda$ . Here the slope of the solution curve does vary depending on  $u$ . The variation in the slope with  $u$  (at fixed  $t$ ) gives an indication of how rapidly the solution curves are converging towards one another (in the case  $\lambda < 0$ ) or diverging away from one another (in the case  $\lambda > 0$ ). If the magnitude of  $\lambda$  is increased, the variation would clearly be more rapid.

The size of the Lipschitz constant is significant if we intend to solve the problem numerically since our numerical approximation will almost certainly produce a value  $U^n$  at time  $t_n$  that is not exactly equal to the true value  $u(t_n)$ . Hence we are on a different solution curve than the true solution. The best we can hope for in the future is that we stay close to the solution curve that we are now on. The size of the Lipschitz constant gives an indication of whether solution curves that start close together can be expected to stay close together or to diverge rapidly.

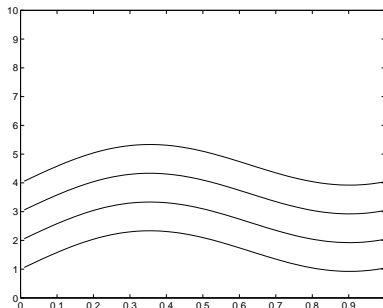


Figure 6.1: Solution curves for Example 6.6, where  $L = 0$ .

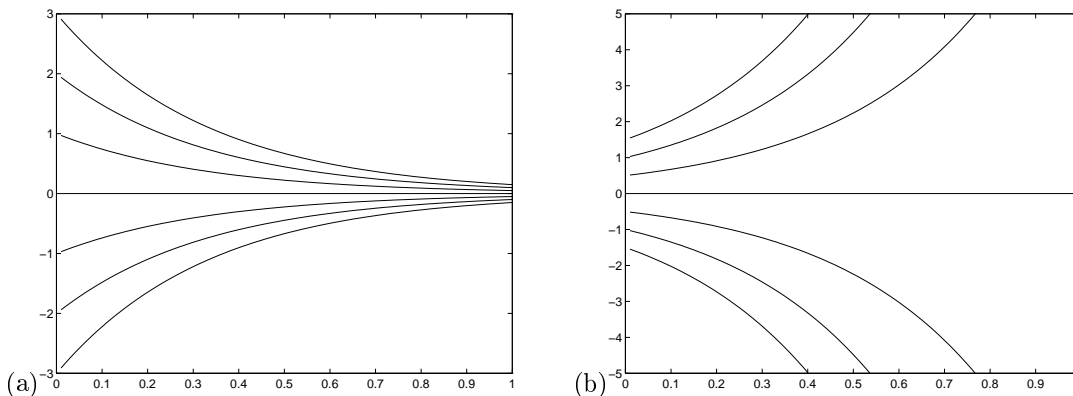


Figure 6.2: Solution curves for Example 6.7 with (a)  $\lambda = -3$  and (b)  $\lambda = 3$ .

### 6.1.4 Limitations

Actually the Lipschitz constant is not the perfect tool for this purpose, since it does not distinguish between rapid divergence and rapid convergence of solution curves. In both Figure 6.2(a) and Figure 6.2(b) the Lipschitz constant has the same value since  $L$  is the absolute value of  $\lambda$ . But we would expect that rapidly convergent solution curves as in Figure 6.2(a) should be easier to handle numerically than rapidly divergent ones. If we make an error at some stage then the effect of this error should decay at later times rather than growing. To some extent this is true and as a result error bounds based on the Lipschitz constant may be orders of magnitude too large in this situation. This is illustrated in Section 7.3.4 where a remedy is proposed for one such error estimate.

However, rapidly converging solution curves can also give serious numerical difficulties, which one might not expect at first glance. This is discussed in detail in Chapter 10 on “stiff equations”.

One should also keep in mind that a small value of the Lipschitz constant does not necessarily mean that two solution curves starting close together will stay close together forever.

**Example 6.8.** Consider two solutions to the pendulum problem from Example 6.5, one with initial data

$$\theta_1(0) = \pi - \epsilon, \quad v_1(0) = 0,$$

and the other with

$$\theta_2(0) = \pi + \epsilon, \quad v_2(0) = 0.$$

The Lipschitz constant is 1 and the data differs by  $2\epsilon$ , which can be arbitrarily small, and yet the solutions soon diverge dramatically, as Solution 1 falls towards  $\theta = 0$  while Solution 2 falls the other way, towards  $\theta = 2\pi$ .

## 6.2 Some basic numerical methods

We begin by listing a few standard approaches to discretizing the equation (6.1). Note that the IVP differs from the BVP considered before in that we are given all the data at the initial time  $t_0 = 0$  and from this we should be able to march forward in time, computing approximations at successive times  $t_1, t_2, \dots$ . We will use  $k$  to denote the time step, so  $t_n = nk$  for  $n \geq 0$ . It is convenient to use the symbol  $k$  that is different from the spatial grid size  $h$  since we will soon study PDEs which involve both spatial and temporal discretizations.

We are given initial data

$$U^0 = \eta \tag{6.11}$$

and want to compute approximations  $U^1, U^2, \dots$  satisfying

$$U^n \approx u(t_n).$$

We will use superscripts to denote the time step index, again anticipating the notation of PDEs where we will use subscripts for spatial indices.

The simplest method is *Euler's method* (also called *Forward Euler*), based on replacing  $u'(t_n)$  by  $D_+U^n = (U^{n+1} - U^n)/k$  from (1.1). This gives the method

$$\frac{U^{n+1} - U^n}{k} = f(U^n), \quad n = 0, 1, \dots \tag{6.12}$$

Rather than viewing this as a system of simultaneous equations as we did for the boundary value problem, it is possible to solve this explicitly for  $U^{n+1}$  in terms of  $U^n$ :

$$U^{n+1} = U^n + kf(U^n). \tag{6.13}$$

From the initial data  $U^0$  we can compute  $U^1$ , then  $U^2$ , and so on. This is called a *time-marching method*.

The *Backward Euler* method is similar, but is based on replacing  $u'(t_{n+1})$  by  $D_-U^{n+1}$ :

$$\frac{U^{n+1} - U^n}{k} = f(U^{n+1}) \tag{6.14}$$

or

$$U^{n+1} = U^n + kf(U^{n+1}). \tag{6.15}$$

Again we can march forward in time since computing  $U^{n+1}$  only requires that we know the previous value  $U^n$ . In the Backward Euler method, however, (6.15) is an equation that must be solved for  $U^{n+1}$  and in general  $f(u)$  is a nonlinear function. We can view this as looking for a zero of the function

$$g(u) = u - kf(u) - U^n$$

which can be approximated using some iterative method such as *Newton's method* discussed in Appendix A3.

Because the Backward Euler method gives an equation that must be solved for  $U^{n+1}$ , it is called an *implicit* method, whereas the Forward Euler method (6.13) is an *explicit* method.

Another implicit method is the *Trapezoidal method*, obtained by averaging the two Euler methods:

$$\frac{U^{n+1} - U^n}{k} = \frac{1}{2}(f(U^n) + f(U^{n+1})). \tag{6.16}$$

As one might expect, this symmetric approximation is second order accurate whereas the Euler methods are only first order accurate.

The above methods are all *one-step methods*, meaning that  $U^{n+1}$  is determined from  $U^n$  alone and previous values of  $U$  are not needed. One way to get higher order accuracy is to use a *multi-step* method that involves other previous values. For example, using the approximation

$$\frac{u(t+k) - u(t-k)}{2k} = u'(t) + \frac{1}{6}k^2 u'''(t) + O(k^3)$$

yields the *Midpoint method* (also called the *Leapfrog method*),

$$\frac{U^{n+1} - U^{n-1}}{2k} = f(U^n) \quad (6.17)$$

or

$$U^{n+1} = U^{n-1} + 2kf(U^n) \quad (6.18)$$

which is a second order accurate explicit 2-step method. The approximation  $D_2u$  from (1.11), rewritten in the form

$$\frac{3u(t+k) - 4u(t) + u(t-k)}{2k} = u'(t+k) + \frac{1}{12}k^2 u'''(t+k) + \dots$$

yields a second order implicit 2-step method

$$3U^{n+1} - 4U^n + U^{n-1} = 2kf(U^{n+1}). \quad (6.19)$$

This is one of the *BDF methods* that will be discussed further in Chapter 10.

### 6.3 Truncation errors

The truncation error for these methods is defined in the same way as in Chapter 2.

**Example 6.9.** The local truncation error of the midpoint method (6.17) is defined by

$$\begin{aligned} \tau^n &= \frac{u(t_{n+1}) - u(t_{n-1})}{2k} - f(u(t_n)) \\ &= \left[ u'(t_n) + \frac{1}{6}k^2 u'''(t_n) + O(k^3) \right] - u'(t_n) \\ &= \frac{1}{6}k^2 u'''(t_n) + O(k^3). \end{aligned}$$

Note that since  $u(t)$  is the true solution of the ODE,  $u'(t_n) = f(u(t_n))$ . The truncation error is  $O(k^2)$  and so we say the method is *second order accurate*, although it is not yet clear that the global error will have this behavior. As always, we need some form of *stability* to guarantee that the global error will exhibit the same rate of convergence as the local truncation error. This will be discussed below.

**Exercise 6.1** Compute the local truncation error of the trapezoidal method and the BDF method discussed above.

### 6.4 One-step errors

In much of the literature concerning numerical methods for IVPs, a slightly different definition of the local truncation error is used that is based on the form (6.18), for example, rather than (6.17). Denoting this value by  $\mathcal{L}^n$ , we have

$$\begin{aligned} \mathcal{L}^n &= u(t_{n+1}) - u(t_{n-1}) - 2kf(u(t_n)) \\ &= \frac{1}{3}k^3 u'''(t_n) + O(k^4). \end{aligned} \quad (6.20)$$

Since  $\mathcal{L}^n = 2k\tau^n$ , this local error is  $O(k^3)$  rather than  $O(k^2)$ , but of course the global error remains the same, and will be only  $O(k^2)$ . Using this alternative definition, many standard results in ODE theory say that a  $p$ th order accurate method should have a local truncation error that is  $O(k^{p+1})$ . With the notation we are using, a  $p$ th order accurate method has a LTE that is  $O(k^p)$ . This notation is consistent with the standard practice for PDEs and leads to a more coherent theory, but one should be aware of this possible source of confusion.

I prefer to call  $\mathcal{L}^n$  the *one-step error*, since this can be viewed as the error that would be introduced in one time step if the past values  $U^n, U^{n-1}, \dots$  were all taken to be the exact values from  $u(t)$ . For example, in the midpoint method (6.18) suppose that

$$U^n = u(t_n) \quad \text{and} \quad U^{n-1} = u(t_{n-1})$$

and we now use these values to compute  $U^{n+1}$ , an approximation to  $u(t_{n+1})$ :

$$\begin{aligned} U^{n+1} &= u(t_{n-1}) + 2kf(u(t_n)) \\ &= u(t_{n-1}) + 2ku'(t_n). \end{aligned}$$

Then the error is

$$u(t_{n+1}) - U^{n+1} = u(t_{n+1}) - u(t_{n-1}) - 2ku'(t_n) = \mathcal{L}^n.$$

From (6.20) we see that in one step the error introduced is  $O(k^3)$ . This is consistent with second order accuracy in the global error if we think of trying to compute an approximation to the true solution  $u(T)$  at some fixed time  $T > 0$ . In order to compute from time  $t = 0$  up to time  $T$ , we need to take  $T/k$  time steps of length  $k$ . A rough estimate of the error at time  $T$  might be obtained by assuming that a new error of size  $\mathcal{L}^n$  is introduced in the  $n$ th time step, and is then simply carried along in later time steps without affecting the size of future local errors and without growing or diminishing itself. Then we would expect the resulting global error at time  $T$  to be simply the sum of all these local errors. Since each local error is  $O(k^3)$  and we are adding up  $T/k$  of them, we end up with a global error that is  $O(k^2)$ .

This viewpoint is in fact exactly right for the simplest ODE

$$u'(t) = f(t)$$

in which  $f$  is independent of  $u$  and the solution is simply the integral of  $f$ , but it is a bit too simplistic for more interesting equations since the error at each time feeds back into the computation at the next step in the case where  $f$  depends on  $u$ . Nonetheless, it is essentially right in terms of the expected order of accuracy, provided the method is stable. In fact it is useful to think of *stability* as exactly what is needed to make this naive analysis essentially correct, by insuring that the old errors from previous time steps do not grow too rapidly in future time steps. This will be investigated in detail in the following chapters.

## 6.5 Linear Multistep Methods

All of the methods introduced so far are members of a class of methods called Linear Multistep Methods (LMMs). In general, an  $r$ -step LMM has the form

$$\sum_{j=0}^r \alpha_j U^{n+j} = k \sum_{j=0}^r \beta_j f(U^{n+j}). \quad (6.21)$$

The value  $U^{n+r}$  is computed from this equation in terms of the previous values  $U^{n+r-1}, U^{n+r-2}, \dots, U^n$ . If  $\beta_r = 0$  then the method (6.21) is explicit, otherwise it is implicit. Note that we can multiply both sides by any constant and have essentially the same method, though the coefficients  $\alpha_j$  and  $\beta_j$  would change. The normalization  $\alpha_r = 1$  is often assumed to fix this scale factor.

There are special classes of methods of this form that are particularly useful and have distinctive names.

**Example 6.10.** The *Adams methods* have the form

$$U^{n+r} = U^{n+r-1} + k \sum_{j=0}^r \beta_j f(U^{n+j}). \quad (6.22)$$

These methods all have

$$\alpha_r = 1, \alpha_{r-1} = -1, \text{ and } \alpha_j = 0 \text{ for } j < r - 1.$$

The  $\beta_j$  coefficients are chosen to maximize the order of accuracy. If we require  $\beta_r = 0$  so the method is explicit then the  $r$  coefficients  $\beta_0, \beta_1, \dots, \beta_{r-1}$  can be chosen so that the method has order  $r$ . This gives the  $r$ -step *Adams-Bashforth* method. Without this restriction we have one more free parameter  $\beta_r$  and so we can eliminate an additional term in the local truncation error, giving an implicit method of order  $r + 1$  called the  $r$ -step *Adams-Moulton* method.

Euler's method is the 1-step Adams-Bashforth method and the trapezoidal method is the 1-step Adams-Moulton method.

**Example 6.11.** The explicit *Nyström methods* have the form

$$U^{n+r} = U^{n+r-2} + k \sum_{j=0}^{r-1} \beta_j f(U^{n+j})$$

with the  $\beta_j$  chosen to give order  $r$ . The midpoint method (6.17) is the 2-step Nyström method.

### 6.5.1 Local truncation error

For LMMs it is easy to derive a general formula for the local truncation error. We have

$$\tau(t_{n+r}) = \frac{1}{k} \left( \sum_{j=0}^r \alpha_j u(t_{n+j}) - k \sum_{j=0}^r \beta_j u'(t_{n+j}) \right).$$

Assuming  $u$  is smooth and expanding in Taylor series gives

$$\begin{aligned} u(t_{n+j}) &= u(t_n) + jku'(t_n) + \frac{1}{2}(jk)^2 u''(t_n) + \dots \\ u'(t_{n+j}) &= u'(t_n) + jku''(t_n) + \frac{1}{2}(jk)^2 u'''(t_n) + \dots \end{aligned}$$

and so

$$\begin{aligned} \tau(t_{n+r}) &= \frac{1}{k} \left( \sum_{j=0}^r \alpha_j \right) u(t_n) + \left( \sum_{j=0}^r (j\alpha_j - \beta_j) \right) u'(t_n) \\ &\quad + k \left( \sum_{j=0}^r \left( \frac{1}{2} j^2 \alpha_j - j\beta_j \right) \right) u''(t_n) \\ &\quad + \dots + k^{q-1} \left( \sum_{j=0}^r \left( \frac{1}{q!} j^q \alpha_j - \frac{1}{(q-1)!} j^{q-1} \beta_j \right) \right) u^{(q)}(t_n) + \dots \end{aligned}$$

The method is *consistent* if  $\tau \rightarrow 0$  as  $k \rightarrow 0$ , which requires that at least the first two terms in this expansion vanish:

$$\sum_{j=0}^r \alpha_j = 0 \quad \text{and} \quad \sum_{j=0}^r j\alpha_j = \sum_{j=0}^r \beta_j. \quad (6.23)$$



If the first  $p + 1$  terms vanish then the method will be  $p$ th order accurate. Note that these conditions depend only on the coefficients  $\alpha_j$  and  $\beta_j$  of the method and not on the particular differential equation being solved.

### 6.5.2 Characteristic polynomials

It is convenient at this point to introduce the so-called characteristic polynomials  $\rho(\zeta)$  and  $\sigma(\zeta)$  for the LMM:

$$\rho(\zeta) = \sum_{j=0}^r \alpha_j \zeta^j \quad \text{and} \quad \sigma(\zeta) = \sum_{j=0}^r \beta_j \zeta^j. \quad (6.24)$$

The first of these is a polynomial of degree  $r$ . So is  $\sigma(\zeta)$  if the method is implicit, otherwise its degree is less than  $r$ . Note that  $\rho(1) = \sum \alpha_j$  and also that  $\rho'(\zeta) = \sum j\alpha_j \zeta^{j-1}$ , so that the consistency conditions (6.23) can be written quite concisely as conditions on these two polynomials:

$$\rho(1) = 0, \quad \rho'(1) = \sigma(1). \quad (6.25)$$

This, however, is not the main reason for introducing these polynomials. The location of the roots of certain polynomials related to  $\rho$  and  $\sigma$  plays a fundamental role in stability theory as we will see in the next chapter.

### 6.5.3 Starting values

One difficulty with using LMMs if  $r > 1$  is that we need the values  $U^0, U^1, \dots, U^{r-1}$  before we can begin to apply the multistep method. The value  $U^0 = \eta$  is known from the initial data for the problem, but the other values are not and must typically be generated by some other numerical method or methods.

**Example 6.12.** If we want to use the midpoint method (6.17) then we need to generate  $U^1$  by some other method before we begin to apply (6.17) with  $n = 1$ . We can obtain  $U^1$  from  $U^0$  using any 1-step method such as Euler's method or the Trapezoidal method. Since the midpoint method is second order accurate we need to make sure that the value  $U^1$  we generate is sufficiently accurate that this second order accuracy will not be lost. Our first impulse might be to conclude that we need to use a second order accurate method such as the Trapezoidal method rather than the first order accurate Euler method, but in fact this is wrong. The overall method is second order in either case. The reason that we achieve second order accuracy even if Euler is used in the first step is exactly analogous to what was observed earlier for boundary value problems, where we found that we can often get away with one order of accuracy lower in the local error for the boundary conditions than what we have elsewhere.

In the present context this is easiest to explain in terms of the 1-step error. The midpoint method has a 1-step error that is  $O(k^3)$  and because this method is applied in  $O(T/k)$  time steps, the global error is expected to be  $O(k^2)$ . Euler's method has a one-step error that is  $O(k^2)$  but we are applying this method only once.

If  $U^0 = \eta = u(0)$  then the error in  $U^1$  will be  $O(k^2)$ . If the midpoint method is stable then this error will not be magnified unduly in later steps and its contribution to the global error will be only  $O(k^2)$ . The overall second order accuracy will not be affected.

**Example 6.13.** Suppose we solve

$$u'(t) = 3u(t), \quad u(0) = 1$$

with solution  $u(t) = e^{3t}$ . We take  $U^0 = 1$ . Generating  $U^1$  using Euler's method gives  $U^1 = (1 + 3k)U^0 = 1 + 3k$  which agrees with  $u(k) = e^{3k}$  to  $O(k^2)$ . Table 6.1 shows the resulting errors at time  $T = 1$  if we now proceed with the midpoint method. The error with values of  $k = 2^{-j}$ ,  $j = 5, 6, \dots$  along with the ratio of errors for successive values of  $k$ . This ratio approaches 4, confirming second order accuracy.

Table 6.1: Error in the solution of Example 6.13 with  $k = 2^{-j}$  when Euler and Trapezoidal are used to generate  $U^1$ .

j	Error with Euler	Ratio	Error with Trapezoidal	Ratio
5	1.3127e-01	3.90	8.5485e-02	3.80
6	3.3008e-02	3.97	2.1763e-02	3.92
7	8.2641e-03	3.99	5.4816e-03	3.97
8	2.0668e-03	3.99	1.3750e-03	3.98
9	5.1674e-04	3.99	3.4428e-04	3.99
10	1.2919e-04	3.99	8.6134e-05	3.99

Table 6.1 also shows the results obtained if the Trapezoidal method is used to generate  $U^1$ . Although the error is slightly smaller, the order of accuracy is the same.

More generally, with an  $r$ -step method of order  $p$ , we need  $r$  starting values  $U^0, U^1, \dots, U^{r-1}$  and we need to generate these values using a method that has a *one-step error* that is  $O(k^p)$  (corresponding to a local truncation error that is  $O(k^{p-1})$ ). Since the number of times we apply this method ( $r - 1$ ) is independent of  $k$  as  $k \rightarrow 0$ , this is sufficient to give an  $O(h^k)$  global error.

If  $p > 3$  then there is still a difficulty here. To generate  $U^1$  from  $U^0$  we need to use a 1-step method, but 1-step linear multistep methods have order 2 at most and hence the one-step error is at best  $O(k^3)$ , limiting our global accuracy to third order. However, there are other classes of 1-step methods that have greater accuracy and can be used in this context. Two such classes of methods, Taylor series and Runge-Kutta methods, will now be introduced.

## 6.6 Taylor series methods

This approach is not often used in practice because it requires differentiating the differential equation and can entail messy algebraic expressions. However it is such an obvious approach that it is worth mentioning, and in some cases it may be useful.

The idea is to use the first  $p + 1$  terms of the Taylor series expansion

$$u(t_{n+1}) = u(t_n) + ku'(t_n) + \frac{1}{2}k^2u''(t_n) + \dots + \frac{1}{p!}k^pu^{(p)}(t_n)$$

to obtain a  $p$ 'th-order accurate method. The problem is that we are only given

$$u'(t) = f(u(t), t)$$

and we must compute the higher derivatives by repeated differentiation of this function. An example should suffice to illustrate the technique and its limitations.

**Example 6.14.** Suppose we want to solve the equation

$$u'(t) = t^2 \sin(u(t)).$$

Then we can compute

$$\begin{aligned} u''(t) &= 2t \sin(u(t)) + t^2 \cos(u(t)) u'(t) \\ &= 2t \sin(u(t)) + t^4 \cos(u(t)) \sin(u(t)). \end{aligned}$$

A second order method is given by

$$U^{n+1} = U^n + kt_n^2 \sin(U^n) + \frac{1}{2}k^2[2t_n \sin(U^n) + t_n^4 \cos(U^n) \sin(U^n)].$$

Clearly higher order derivatives can be computed and used, but this is cumbersome even for this simple example.

**Exercise 6.2** Derive the third order Taylor series method for the above equation.

## 6.7 Runge-Kutta Methods

Most methods used in practice do not require that the user explicitly calculate higher order derivatives. Instead a higher order finite difference approximation is designed which typically models these terms automatically.

A multistep method can achieve high accuracy by using high order polynomial interpolation through several points. To achieve the same effect with a 1-step method it is typically necessary to use a *multi-stage* method, where intermediate values are generated and used within a single time step.

**Example 6.15.** A 2-stage explicit Runge-Kutta method is given by

$$\begin{aligned}U^* &= U^n + \frac{1}{2}kf(U^n) \\U^{n+1} &= U^n + kf(U^*).\end{aligned}$$

In the first stage an intermediate value is generated which approximates  $u(t_{n+1/2})$  via Euler's method. In the second step the function  $f$  is evaluated at this midpoint to estimate the slope over the full time step. Since this now looks like a centered approximation of the derivative we might hope for second order accuracy, as we'll now verify by computing the local truncation error.

Combining the two steps above, we can rewrite the method as

$$U^{n+1} = U^n + kf\left(U^n + \frac{1}{2}kf(U^n)\right).$$

The truncation error is

$$\tau^n = \frac{1}{k}(u(t_{n+1}) - u(t_n)) - f\left(u(t_n) + \frac{1}{2}kf(u(t_n))\right). \quad (6.26)$$

Note that

$$\begin{aligned}f\left(u(t_n) + \frac{1}{2}kf(u(t_n))\right) &= f\left(u(t_n) + \frac{1}{2}ku'(t_n)\right) \\&= f(u(t_n)) + \frac{1}{2}ku'(t_n)f'(u(t_n)) + \frac{1}{8}k^2(u'(t_n))^2f''(u(t_n)) + \dots\end{aligned}$$

Since  $f(u(t_n)) = u'(t_n)$  and differentiating gives  $f'(u)u' = u''$ , we obtain

$$f\left(u(t_n) + \frac{1}{2}kf(u(t_n))\right) = u'(t_n) + \frac{1}{2}ku''(t_n) + O(k^2).$$

Using this in (6.26) gives

$$\begin{aligned}\tau^n &= \frac{1}{k}\left(ku'(t_n) + \frac{1}{2}k^2u''(t_n) + O(k^3)\right) \\&\quad - \left(u'(t_n) + \frac{1}{2}ku''(t_n) + O(k^2)\right) \\&= O(k^2)\end{aligned}$$

and the method is second order accurate. (Check the  $O(k^2)$  term to see that this does not vanish.)

**Remark.** An easier way to determine the order of accuracy is to apply the method to the special test equation  $u' = \lambda u$ , which has solution  $u(t_{n+1}) = e^{\lambda k}u(t_n)$ , and determine the error on this problem. Here we obtain

$$\begin{aligned}U^{n+1} &= U^n + k\lambda\left(U^n + \frac{1}{2}k\lambda U^n\right) \\&= U^n + (k\lambda)U^n + \frac{1}{2}(k\lambda)^2U^n \\&= e^{k\lambda}U^n + O(k^3).\end{aligned}$$

The one step error is  $O(k^3)$  and hence the local truncation error is  $O(k^2)$ . Of course we have only checked that the local truncation error is  $O(k^2)$  on one particular function  $u(t) = e^{\lambda t}$ , not on all smooth solutions, but this is generally a reliable indication of the order more generally. Applying a method to this special equation is also a fundamental tool in stability analysis — see Chapter 8.

**Example 6.16.** One very popular Runge-Kutta method is the 4'th order 4-stage method given by

$$\begin{aligned}F_0 &= f(U^n, t_n) \\F_1 &= f\left(U^n + \frac{1}{2}kF_0, t_n + \frac{1}{2}k\right) \\F_2 &= f\left(U^n + \frac{1}{2}kF_1, t_n + \frac{1}{2}k\right) \\F_3 &= f(U^n + kF_2, t_{n+1}) \\U^{n+1} &= U^n + \frac{k}{6}(F_0 + 2F_1 + 2F_2 + F_3).\end{aligned}\tag{6.27}$$



## Chapter 7

# Zero-Stability and Convergence for Initial Value Problems

### 7.1 Convergence

In order to discuss the convergence of a numerical method for the Initial Value Problem, we focus on a fixed (but arbitrary) time  $T > 0$  and consider the error in our approximation to  $u(T)$  computed with the method using time step  $k$ . The method converges on this problem if this error goes to zero as  $k \rightarrow 0$ . Note that the number of time steps that we need to take to reach time  $T$  increases as  $k \rightarrow 0$ . If we use  $N$  to denote this value ( $N = T/k$ ), then convergence means that

$$\lim_{\substack{k \rightarrow 0 \\ Nk=T}} U^N = u(T). \quad (7.1)$$

In principle a method might converge on one problem but not on another, or converge with one set of starting values but not with another set. In order to speak of a *method* being *convergent* in general, we require that it converges on *all* problems in a reasonably large class with all reasonable starting values. For an  $r$ -step method we need  $r$  starting values. These values will typically depend on  $k$ , and to make this clear we will write them as  $U^0(k), U^1(k), \dots, U^{r-1}(k)$ . While these will generally approximate  $u(t)$  at the times  $t_0 = 0, t_1 = k, \dots, t_{r-1} = (r-1)k$  respectively, as  $k \rightarrow 0$  each of these times approaches  $t_0 = 0$ . So the weakest condition we might put on our starting values is that they converge to the correct initial value  $\eta$  as  $k \rightarrow 0$ :

$$\lim_{k \rightarrow 0} U^\nu(k) = \eta \text{ for } \nu = 0, 1, \dots, r-1. \quad (7.2)$$

We can now state the definition of convergence.

**Definition 7.1.1** *An  $r$ -step method is said to be convergent if applying the method to any ODE (6.1) with  $f(u, t)$  Lipschitz continuous in  $u$ , and with any set of starting values satisfying (7.2), we obtain convergence in the sense of (7.1) for every fixed time  $T > 0$ .*

In order to be convergent, a method must be *consistent*, meaning as before that the LTE is  $o(1)$  as  $k \rightarrow 0$ , and also *zero-stable*, as described later in this chapter. We will begin to investigate these issues by first proving the convergence of 1-step methods, which turn out to be zero-stable automatically. We start with Euler's method on linear problems, then consider Euler's method on general nonlinear problems and finally extend this to a wide class of 1-step methods.

## 7.2 Linear equations and Duhamel's principle

Much of the theory presented below is based on examining what happens when a method is applied to a simple linear equation of the form

$$u'(t) = \lambda u(t) + g(t) \quad (7.3)$$

with initial data

$$u(t_0) = \eta.$$

Here  $\lambda$  is a constant and  $g(t)$  is a given function. In the special case  $g(t) \equiv 0$  (the *homogeneous equation*), the solution is simply

$$u(t) = e^{\lambda(t-t_0)}\eta.$$

If we let  $\mathcal{S}(t, t_0) = e^{\lambda(t-t_0)}$  be this solution operator, which maps data at time  $t_0$  to the solution at time  $t$ , then we can write the solution of the more general linear problem (7.3) using *Duhamel's principle*, which states that the solution to the nonhomogeneous problem is obtained by adding to the solution of the homogeneous problem a superposition of  $\mathcal{S}(t, \tau)g(\tau)$  over all times  $\tau$  between  $t_0$  and  $t$ . In this sense  $\mathcal{S}(t, \tau)$  acts like a Green's function (compare to (2.36)). So we have

$$\begin{aligned} u(t) &= \mathcal{S}(t, t_0)\eta + \int_{t_0}^t \mathcal{S}(t, \tau)g(\tau) d\tau \\ &= e^{\lambda(t-t_0)}\eta + \int_{t_0}^t e^{\lambda(t-\tau)}g(\tau) d\tau. \end{aligned} \quad (7.4)$$

**Exercise 7.1** Use (7.4) to solve the equation  $u'(t) = 3u(t) + 2t$  with  $u(0) = 1$ .

## 7.3 One-step methods

### 7.3.1 Euler's method on linear problems

If we apply Euler's method to the equation (7.3), we obtain

$$\begin{aligned} U^{n+1} &= U^n + k(\lambda U^n + g(t_n)) \\ &= (1 + k\lambda)U^n + kg(t_n). \end{aligned} \quad (7.5)$$

The local truncation error for Euler's method is given by

$$\begin{aligned} \tau^n &= \left( \frac{u(t_{n+1}) - u(t_n)}{k} \right) - (\lambda u(t_n) + g(t_n)) \\ &= \left( u'(t_n) + \frac{1}{2}ku''(t_n) + O(k^2) \right) - u'(t_n) \\ &= \frac{1}{2}ku''(t_n) + O(k^2). \end{aligned} \quad (7.6)$$

Rewriting this equation as

$$u(t_{n+1}) = (1 + k\lambda)u(t_n) + kg(t_n) + k\tau^n$$

and subtracting this from (7.5) gives a difference equation for the global error  $E^n$ :

$$E^{n+1} = (1 + k\lambda)E^n - k\tau^n. \quad (7.7)$$

Note that this has exactly the same form as (7.5) but with a different nonhomogeneous term:  $\tau^n$  in place of  $g(t_n)$ . This is analogous to equation (2.15) in the boundary value theory and again gives the

relation we need between the local truncation error  $\tau^n$  (which is easy to compute) and the global error  $E^n$  (which we wish to bound). Note again that *linearity* plays a critical role in making this connection. We will consider nonlinear problems below.

Because the equation and method we are now considering are both so simple, we obtain an equation (7.7) that we can explicitly solve for the global error  $E^n$ . Applying the recursion (7.7) repeatedly we see what form the solution should take:

$$\begin{aligned} E^n &= (1 + k\lambda)E^{n-1} - k\tau^{n-1} \\ &= (1 + k\lambda)[(1 + k\lambda)E^{n-2} - k\tau^{n-2}] - k\tau^{n-1} \\ &= \dots \end{aligned}$$

By induction we can easily confirm that in general

$$E^n = (1 + k\lambda)^n E^0 - k \sum_{m=1}^n (1 + k\lambda)^{n-m} \tau^{m-1}. \quad (7.8)$$

(Note that some of the superscripts are powers while others are indices!) This has a form that is very analogous to the solution (7.4) of the corresponding ordinary differential equation, where now  $(1 + k\lambda)^{n-m}$  plays the role of the solution operator of the homogeneous problem — it transforms data at time  $t_m$  to the solution at time  $t_n$ . The expression (7.8) is sometimes called the discrete form of Duhamel's principle.

We are now ready to prove that Euler's method converges on the equation (7.3). We need only observe that

$$|1 + k\lambda| \leq e^{k|\lambda|} \quad (7.9)$$

and so

$$(1 + k\lambda)^{n-m} \leq e^{(n-m)k|\lambda|} \leq e^{nk|\lambda|} \leq e^{|\lambda|T} \quad (7.10)$$

provided that we restrict our attention to the finite time interval  $0 \leq t \leq T$ , so that  $t_n = nk \leq T$ . It then follows from (7.8) that

$$\begin{aligned} |E^n| &\leq e^{|\lambda|T} \left( |E^0| + k \sum_{m=1}^n |\tau^{m-1}| \right) \\ &\leq e^{|\lambda|T} \left( |E^0| + nk \max_{1 \leq m \leq n} |\tau^{m-1}| \right). \end{aligned} \quad (7.11)$$

Let  $N = T/k$  be the number of time steps needed to reach time  $T$  and set

$$\|\tau\|_\infty = \max_{0 \leq n \leq N-1} |\tau^n|.$$

From (7.6) we expect

$$\|\tau\|_\infty \approx \frac{1}{2}k \|u''\|_\infty = O(k)$$

where  $\|u''\|_\infty$  is the maximum value of the function  $u''$  over the interval  $[0, T]$ . Then for  $t = nk \leq T$ , we have from (7.11) that

$$|E^n| \leq e^{|\lambda|T} (|E^0| + T\|\tau\|_\infty).$$

If (7.2) is satisfied then  $E^0 \rightarrow 0$  as  $k \rightarrow 0$ . In fact for this one-step method we would generally take  $U^0 = u(0) = \eta$ , in which case  $E^0$  drops out and we are left with

$$|E^n| \leq e^{|\lambda|T} T \|\tau\|_\infty = O(k) \text{ as } k \rightarrow 0 \quad (7.12)$$





and so

$$\|A^{-1}\|_{\infty} \leq kNe^{|\lambda|T} = Te^{|\lambda|T}.$$

Hence the method is stable and  $\|E\|_{\infty} \leq \|A^{-1}\|_{\infty} \|\tau\|_{\infty} \leq Te^{|\lambda|T} \|\tau\|_{\infty}$  which agrees with the bound (7.12).

### 7.3.3 Euler's method on nonlinear problems

So far we have focused entirely on linear equations. Practical problems are almost always nonlinear, but for the initial value problem it turns out that it is not significantly more difficult to handle this case if we assume that  $f(u)$  is Lipschitz continuous, which is reasonable in light of the discussion in Section 6.1.

Euler's method on  $u' = f(u)$  takes the form

$$U^{n+1} = U^n + kf(U^n) \tag{7.13}$$

and the truncation error is defined by

$$\begin{aligned} \tau^n &= \frac{1}{k}(u(t_{n+1}) - u(t_n)) - f(u(t_n)) \\ &= \frac{1}{2}ku''(t_n) + O(k^3) \end{aligned}$$

just as in the linear case. So the true solution satisfies

$$u(t_{n+1}) = u(t_n) + kf(u(t_n)) + k\tau^n$$

and subtracting this from (7.13) gives

$$E^{n+1} = E^n + k(f(U^n) - f(u(t_n))) - k\tau^n. \tag{7.14}$$

In the linear case  $f(U^n) - f(u(t_n)) = \lambda E^n$  and we get the relation (7.7) for  $E^n$ . In the nonlinear case we cannot express  $f(U^n) - f(u(t_n))$  directly in terms of the error  $E^n$  in general. However, using the Lipschitz continuity of  $f$  we can get a bound on this in terms of  $E^n$ :

$$|f(U^n) - f(u(t_n))| \leq L|U^n - u(t_n)| = L|E^n|.$$

Using this in (7.14) gives

$$|E^{n+1}| \leq |E^n| + kL|E^n| + k|\tau^n| = (1 + kL)|E^n| + k|\tau^n| \tag{7.15}$$

From this inequality we can show by induction that

$$|E^n| \leq (1 + kL)^n |E^0| + k \sum_{m=1}^n (1 + kL)^{n-m} |\tau^{m-1}|$$

and so, using the same steps as in obtaining (7.12), we obtain

$$|E^n| \leq e^{LT} \|\tau\|_{\infty} = O(k) \text{ as } k \rightarrow 0 \tag{7.16}$$

for all  $n$  with  $nk \leq T$ , proving that the method converges. In the linear case  $L = |\lambda|$  and this reduces to exactly (7.12).

### 7.3.4 Realistic error estimates

We have proved the convergence of Euler's method and in the process obtained a bound on the error, (7.16), which goes to zero as  $k \rightarrow 0$ . However, this bound may be totally useless in estimating the actual error for a practical calculation.

**Example 7.1.** Suppose we solve  $u' = -10u$  with  $u(0) = 1$  up to time  $T = 10$  using a time step  $k = 0.01$ . Then the solution is  $u(T) = e^{-100} \approx 3.7 \times 10^{-44}$ . The computed solution is  $U^N = (.9)^{100} \approx 2.65 \times 10^{-5}$  and the error is essentially the same. Since  $L = 10$  for this problem, the error bound (7.16) gives

$$\|E^N\| \leq e^{100} \cdot 10 \cdot \|\tau\|_\infty \approx 2.7 \times 10^{44} \|\tau\|_\infty.$$

Here  $\|\tau\|_\infty = |\tau^0| \approx 50k$ , so this upper bound on the error does go to zero as  $k \rightarrow 0$ , but obviously it is not a realistic estimate of the error, being too large by a factor of  $10^{50}$  or so.

The problem is that the estimate (7.16) is based on the Lipschitz constant, which is always nonnegative even when the solution curves are converging rather than diverging (see Section 6.1.4). If we were instead solving  $u' = +10u$  with  $u(0) = 1$  then again  $L = 10$  and we would get the same error estimate, which in this case would actually be a reasonably accurate estimate. (In this case the true solution is  $e^{100} \approx 2.7 \times 10^{43}$ , so it is not surprising the absolute error is of this magnitude. The relative error is reasonable!)

A more realistic error bound for the case where  $\lambda < 0$  can be obtained by writing (7.13) as

$$U^{n+1} = \Phi(U^n) \tag{7.17}$$

and then determining the Lipschitz constant for the function  $\Phi$ . (Note that if  $f$  is Lipschitz then so is  $\Phi$ .) Call this constant  $M$ . Then we have

$$u(t_{n+1}) = \Phi(u(t_n)) + k\tau^n$$

and subtracting this from (7.17) gives

$$E^{n+1} = \Phi(U^n) - \Phi(u(t_n)) - k\tau^n.$$

Taking norms and using the Lipschitz continuity of  $\Phi$  gives

$$|E^{n+1}| \leq M|E^n| + k|\tau^n|,$$

from which we obtain

$$|E^n| \leq M^n T \|\tau\|_\infty \tag{7.18}$$

for all  $n$  with  $nk \leq T$ .

The advantage of this over (7.16) is seen for the case  $u' = \lambda u$  with  $\lambda < 0$ . We have

$$\Phi(u) = (1 + k\lambda)u,$$

and so

$$M = |1 + k\lambda|.$$

In the above example, where  $\lambda = -10$  and  $k = 0.01$ , we obtain  $M = 0.9$ . Using this in (7.18) gives a much more realistic error bound. The bound (7.16) was obtained by using  $(1 + kL) = (1 + k|\lambda|)$  instead of  $M = 1 + k\lambda$ . When  $\lambda > 0$  the two are the same, but for  $\lambda < 0$  the first is always greater than 1, while  $M < 1$  at least for  $k$  sufficiently small. (If  $k$  is not sufficiently small then the numerical solution may in fact grow exponentially with  $t$  rather than decaying as it should. This form of instability is discussed in Chapter 8.)

### 7.3.5 General 1-step methods

A general explicit one-step method takes the form

$$U^{n+1} = U^n + k\Psi(U^n, t_n, k) \quad (7.19)$$

for some function  $\Psi$ , which depends on  $f$  of course. We will assume that  $\Psi(u, t, k)$  is continuous in  $t$  and  $k$  and Lipschitz continuous in  $u$ , with Lipschitz constant  $L'$  that is generally related to the Lipschitz constant of  $f$ .

The method is *consistent* if

$$\Psi(u, t, 0) = f(u, t)$$

for all  $u, t$ . The local truncation error is

$$\tau^n = \left( \frac{u(t_{n+1}) - u(t_n)}{k} \right) - \Psi(u(t_n), t_n, k).$$

**Exercise 7.2** For the 2-stage Runge-Kutta method of Example 6.15, we have

$$\Psi(u, t, k) = f \left( u + \frac{1}{2}kf(u) \right).$$

Show that if  $f$  is Lipschitz continuous with Lipschitz constant  $L$  then  $\Psi$  has Lipschitz constant  $L' = L + \frac{1}{2}kL^2$ .

Using the same technique as in Section 7.3.4 we can show that any one-step methods satisfying these conditions is convergent. We have

$$u(t_{n+1}) = u(t_n) + k\Psi(u(t_n), t_n, k) + k\tau^n$$

and subtracting this from (7.19) gives

$$E^{n+1} = E^n + k(\Psi(U^n, t_n, k) - \Psi(u(t_n), t_n, k)) - k\tau^n.$$

Using the Lipschitz condition we obtain

$$|E^{n+1}| \leq |E^n| + kL'|E^n| + k|\tau^n|.$$

This has exactly the same form as (7.15) and the proof of convergence proceeds exactly as from there.

**Note:** In general it may be possible to obtain more realistic error bounds, as in Section 7.3.4, by defining  $\Phi(u, t, k) = u + k\Psi(u, t, k)$  and using the Lipschitz constant for  $\Phi$  rather than that for  $\Psi$ .

## 7.4 Zero-stability of linear multistep methods

The convergence proof of the previous section shows that for 1-step methods, each one-step error  $k\tau^{m-1}$  has an effect on the global error that is bounded by  $e^{LT}|k\tau^{m-1}|$ . Although the error is possibly amplified by a factor  $e^{LT}$ , this factor is bounded independent of  $k$  as  $k \rightarrow 0$ . Consequently the method is stable: the global error can be bounded in terms of the sum of all the one-step errors, and hence has the same asymptotic behavior as the local truncation error as  $k \rightarrow 0$ . This form of stability is often called *zero-stability* in ODE theory, to distinguish it from other forms of stability that are of equal importance in practice. The fact that a method is zero-stable (and converges as  $k \rightarrow 0$ ) is no guarantee that it will give reasonable results on the particular grid with  $k > 0$  that we want to use in practice. Other “stability” issues of a different nature will be taken up in the next chapter.

But first we will investigate the issue of zero-stability for general linear multistep methods (LMM’s), where the theory of the previous section does not apply directly. We begin with an example showing

Table 7.1: Solution  $U^N$  to (7.22) with  $U^0 = 0$ ,  $U^1 = k$  and various values of  $k = 1/N$ .

$N$	$U^N$
5	4.2
10	258.4
20	1954408

a consistent LMM that is **not** convergent. Examining what goes wrong will motivate our definition of zero-stability for LMMs.

**Example 7.2.** The LMM

$$U^{n+2} - 3U^{n+1} + 2U^n = -kf(U^n) \quad (7.20)$$

has a local truncation error given by

$$\frac{1}{k}[u(t_{n+2}) - 3u(t_{n+1}) + u(t_n) - ku'(t_n)] = \frac{5}{2}ku''(t_n) + O(k^2)$$

so the method is consistent and “first order accurate”. But in fact the global error will not exhibit first order accuracy, or even convergence, in general. This can be seen even on the trivial initial-value problem

$$u'(t) = 0, \quad u(0) = 0 \quad (7.21)$$

with solution  $u(t) \equiv 0$ . On this equation (7.20) takes the form

$$U^{n+2} - 3U^{n+1} + 2U^n = 0. \quad (7.22)$$

We need two starting values  $U^0$  and  $U^1$ . If we take  $U^0 = U^1 = 0$  then (7.22) generates  $U^n = 0$  for all  $n$  and in this case we certainly converge to correct solution, and in fact we get the exact solution for any  $k$ .

But in general we will not have the exact value  $U^1$  available and will have to approximate this, introducing some error into the computation. Table 7.1 shows results obtained by applying this method with starting data  $U^0 = 0$ ,  $U^1 = k$ . Since  $U^1(k) \rightarrow 0$  as  $k \rightarrow 0$ , this is valid starting data in the context of Definition 7.1.1 of convergence. If the method is convergent we should see that  $U^N$ , the computed solution at time  $T = 1$ , converges to zero as  $k \rightarrow 0$ . Instead it blows up quite dramatically. Similar results would be seen if we applied this method to an arbitrary equation  $u' = f(u)$  and used any one-step method to compute  $U^1$  from  $U^0$ .

The homogeneous linear difference equation (7.22) can be solved explicitly for  $U^n$  in terms of the starting values  $U^0$  and  $U^1$ . We obtain

$$U^n = 2U^0 - U^1 + 2^n(U^1 - U^0). \quad (7.23)$$

It is easy to verify that this satisfies (7.22) and also the starting values. (We’ll see how to solve general linear difference equations in the next section.)

Since  $u(t) = 0$ , the error is  $E^n = U^n$  and we see that any initial errors in  $U^1$  or  $U^0$  are magnified by a factor  $2^n$  in the global error (except in the special case  $U^1 = U^0$ ). This exponential growth of the error is the instability which leads to nonconvergence. In order to rule out this sort of growth of errors, we need to be able to solve a general linear difference equation.

### 7.4.1 Solving linear difference equations

Consider the general homogeneous linear difference equation

$$\sum_{j=0}^r \alpha_j U^{n+j} = 0. \quad (7.24)$$

Eventually we will look for a particular solution satisfying given initial conditions  $U^0, U^1, \dots, U^{r-1}$ , but to begin with we will find the general solution of the difference equation in terms of  $r$  free parameters. We will hypothesize that this equation has a solution of the form

$$U^n = \zeta^n \tag{7.25}$$

for some value of  $\zeta$  (here  $\zeta^n$  is the  $n$ th power!). Plugging this into (7.24) gives

$$\sum_{j=0}^r \alpha_j \zeta^{n+j} = 0$$

and dividing by  $\zeta^n$  yields

$$\sum_{j=0}^r \alpha_j \zeta^j = 0. \tag{7.26}$$

We see that (7.25) is a solution of the difference equation if  $\zeta$  satisfies the equation (7.26), i.e., if  $\zeta$  is a root of the polynomial

$$\rho(\zeta) = \sum_{j=0}^r \alpha_j \zeta^j.$$

Note that this is just the first characteristic polynomial of the LMM introduced in (6.24). In general  $\rho(\zeta)$  has  $r$  roots  $\zeta_1, \zeta_2, \dots, \zeta_r$  and can be factored as

$$\rho(\zeta) = \alpha_r (\zeta - \zeta_1)(\zeta - \zeta_2) \cdots (\zeta - \zeta_r).$$

Since the difference equation is linear, any linear combination of solutions is again a solution. If  $\zeta_1, \zeta_2, \dots, \zeta_r$  are distinct ( $\zeta_i \neq \zeta_j$  for  $i \neq j$ ) then the  $r$  distinct solutions  $\zeta_i^n$  are linearly independent and the general solution of (7.24) has the form

$$U^n = c_1 \zeta_1^n + c_2 \zeta_2^n + \cdots + c_r \zeta_r^n \tag{7.27}$$

where  $c_1, \dots, c_r$  are arbitrary constants. In this case, every solution of the difference equation (7.24) has this form. If initial conditions  $U^0, U^1, \dots, U^{r-1}$  are specified, then the constants  $c_1, \dots, c_r$  can be uniquely determined by solving the  $r \times r$  linear system

$$\begin{aligned} c_1 + c_2 + \cdots + c_r &= U^0 \\ c_1 \zeta_1 + c_2 \zeta_2 + \cdots + c_r \zeta_r &= U^1 \\ &\vdots \\ c_1 \zeta_1^{r-1} + c_2 \zeta_2^{r-1} + \cdots + c_r \zeta_r^{r-1} &= U^{r-1} \end{aligned} \tag{7.28}$$

**Example 7.3.** The characteristic polynomial for the difference equation (7.22) is

$$\rho(\zeta) = 2 - 3\zeta + \zeta^2 = (\zeta - 1)(\zeta - 2) \tag{7.29}$$

with roots  $\zeta_1 = 1, \zeta_2 = 2$ . The general solution has the form

$$U^n = c_1 + c_2 \cdot 2^n$$

and solving for  $c_1$  and  $c_2$  from  $U^0$  and  $U^1$  gives the solution (7.23).

This example indicates that if  $\rho(\zeta)$  has any roots that are greater than one in modulus, the method will not be convergent. It turns out that the converse is nearly true: If all of the roots have modulus

no greater than one, then the method is convergent, with one proviso. There must be no *repeated* roots with modulus equal to one. The next two examples illustrate this.

If the roots are not distinct, say  $\zeta_1 = \zeta_2$  for simplicity, then  $\zeta_1^n$  and  $\zeta_2^n$  are not linearly independent and the  $U^n$  given by (7.27), while still a solution, is not the most general solution. The system (7.28) would be singular in this case. In addition to  $\zeta_1^n$  there is also a solution of the form  $n\zeta_1^n$  and the general solution has the form

$$U^n = c_1\zeta_1^n + c_2n\zeta_1^n + c_3\zeta_3^n + \cdots + c_r\zeta_r^n.$$

If in addition  $\zeta_3 = \zeta_1$ , then the third term would be replaced by  $c_3n^2\zeta_1^n$ . Similar modifications are made for any other repeated roots. Note how similar this theory is to the standard solution technique for an  $r$ 'th order linear ordinary differential equation.

**Example 7.4.** Applying the consistent LMM

$$U^{n+2} - 2U^{n+1} + U^n = \frac{1}{2}k(f(U^{n+2}) - f(U^n)) \quad (7.30)$$

to the differential equation  $u'(t) = 0$  gives the difference equation

$$U^{n+2} - 2U^{n+1} + U^n = 0.$$

The characteristic polynomial is

$$\rho(\zeta) = \zeta^2 - 2\zeta + 1 = (\zeta - 1)^2 \quad (7.31)$$

so  $\zeta_1 = \zeta_2 = 1$ . The general solution is

$$U^n = c_1 + c_2n.$$

For particular starting values  $U^0$  and  $U^1$  the solution is

$$U^n = U^0 + (U^1 - U^0)n.$$

Again we see that the solution grows with  $n$ , though not as dramatically as in Example 7.2 (the growth is linear rather than exponential). But this growth is still enough to destroy convergence. If we take the same starting values as before,  $U^0 = 0$  and  $U^1 = k$ , then  $U^n = kn$  and so

$$\lim_{\substack{k \rightarrow 0 \\ Nk = T}} U^N = kN = T.$$

The method converges to the function  $v(t) = t$  rather than to  $u(t) = 0$ , and hence the LMM (7.30) is not convergent.

This example shows that if  $\rho(\zeta)$  has a *repeated* root of modulus 1, then the method cannot be convergent.

**Example 7.5.** Now consider the consistent LMM

$$U^{n+3} - 2U^{n+2} + \frac{5}{4}U^{n+1} - \frac{1}{4}U^n = \frac{1}{4}hf(U^n). \quad (7.32)$$

Applying this to (7.21) gives

$$U^{n+3} - 2U^{n+2} + \frac{5}{4}U^{n+1} - \frac{1}{4}U^n = 0$$

and the characteristic polynomial is

$$\rho(\zeta) = \zeta^3 - 2\zeta^2 + \frac{5}{4}\zeta - \frac{1}{4} = (\zeta - 1)(\zeta - 0.5)^2. \quad (7.33)$$

So  $\zeta_1 = 1$ ,  $\zeta_2 = \zeta_3 = 1/2$  and the general solution is

$$U^n = c_1 + c_2 \left(\frac{1}{2}\right)^n + c_3n \left(\frac{1}{2}\right)^n.$$

Here there is a repeated root but with modulus less than 1. The linear growth of  $n$  should then be overwhelmed by the decay of  $(1/2)^n$ .

For this 3-step method we need three starting values  $U^0, U^1, U^2$  and we can find  $c_1, c_2, c_3$  in terms of them by solving a linear system similar to (7.28). Each  $c_i$  will be a linear combination of  $U^0, U^1, U^2$  and so if  $U^\nu(k) \rightarrow 0$  as  $k \rightarrow 0$  then  $c_i(k) \rightarrow 0$  as  $k \rightarrow 0$  also. The value  $U^N$  computed at time  $T$  with step size  $k$  (where  $kN = T$ ) has the form

$$U^N = c_1(k) + c_2(k) \left(\frac{1}{2}\right)^N + c_3(k)N \left(\frac{1}{2}\right)^N. \quad (7.34)$$

Now we see that

$$\lim_{\substack{k \rightarrow 0 \\ Nk=T}} U^N = 0$$

and so the method (7.32) converges on  $u' = 0$  with arbitrary starting values  $U^\nu(k)$  satisfying  $U^\nu(k) \rightarrow 0$  as  $k \rightarrow 0$ . (In fact this LMM is convergent in general.)

More generally, if  $\rho(\zeta)$  has a root  $\zeta_j$  that is repeated  $r$  times, then  $U^N$  will involve terms of the form  $N^s \zeta_j^N$  for  $s = 1, 2, \dots, r$ . This converges to zero as  $N \rightarrow \infty$  provided  $|\zeta_j| < 1$ . The algebraic growth of  $N^s$  is overwhelmed by the exponential decay of  $\zeta_j^N$ . This shows that repeated roots are not a problem as long as they have magnitude strictly less than 1.

With the above examples as motivation, we are ready to state the definition of zero-stability.

**Definition 7.4.1** An  $r$ -step Linear Multistep Method is said to be **zero-stable** if the roots of the characteristic polynomial  $\rho(\zeta)$  defined by (6.24) satisfy the following conditions:

$$\begin{aligned} |\zeta_j| &\leq 1 \text{ for } j = 1, 2, \dots, r \\ \text{If } \zeta_j \text{ is a repeated root, then } |\zeta_j| &< 1. \end{aligned} \quad (7.35)$$

If the conditions (7.35) are satisfied for all roots of  $\rho$ , then the polynomial is said to satisfy the *root condition*.

**Example 7.6.** The Adams methods have the form

$$U^{n+r} = U^{n+r-1} + k \sum_{j=1}^r \beta_j f(U^{n+j})$$

and hence

$$\rho(\zeta) = \zeta^r - \zeta^{r-1} = (\zeta - 1)\zeta^{r-1}.$$

The roots are  $\zeta_1 = 1$  and  $\zeta_2 = \dots = \zeta_r = 0$ . The root condition is clearly satisfied and all of the Adams-Bashforth and Adams-Moulton methods are zero stable.

The examples given above certainly do not prove that zero-stability as defined above is a sufficient condition for convergence. We only looked at the simplest possible ODE  $u'(t) = 0$  and saw that things could go wrong if the root condition is *not* satisfied. It turns out, however, that the root condition is all that is needed to prove convergence on the general initial value problem (in the sense of Definition 7.1.1). For the initial value problem we have the general result that

$$\text{consistency} + \text{zero-stability} \implies \text{convergence}. \quad (7.36)$$

This is the analog of the statement (2.21) for the boundary value problem. A proof of this result can be found in [Hen62]

**Note:** A consistent linear multistep method always has one root equal to 1, say  $\zeta_1 = 1$ , called the *principal root*. This follows from (6.25). Hence a consistent 1-step LMM (such as Euler, backward Euler, trapezoidal) is certainly zero-stable. More generally we have proved in Section 7.3.5 that any consistent 1-step method (that is Lipschitz continuous) is convergent. Such methods are automatically “zero-stable” and behave well as  $k \rightarrow 0$ . They may have other stability problems that show up for “large” values of  $k$ , however, which is the subject of the next chapter.





## Chapter 8

# Absolute Stability for ODEs

### 8.1 Unstable computations with a zero-stable method

In the last chapter we investigated zero-stability, the form of stability needed to guarantee convergence of a numerical method as the grid is refined ( $k \rightarrow 0$ ). In practice, however, we are not able to actually compute this limit. Instead we typically perform a single calculation with some particular nonzero time step  $k$  (or some particular sequence of time steps with a variable step size method). Since the expense of the computation increases as  $k$  decreases, we generally want to choose the time step as large as possible consistent with our accuracy requirements. How can we estimate the size of  $k$  required?

Recall that if the method is stable in an appropriate sense, then we expect the global error to be bounded in terms of the local truncation errors at each step, and so we can often use the local truncation error to estimate the time step needed, as illustrated below. But the form of stability now needed is something stronger than zero-stability. We need to know that the error is well behaved for the particular time step we are now using. It is little help to know that things will converge in the limit “for  $k$  sufficiently small”. The potential difficulties are best illustrated with some examples.

**Example 8.7.** Consider the IVP

$$u'(t) = -\sin t, \quad u(0) = 1$$

with solution

$$u(t) = \cos t.$$

Suppose we wish to use Euler’s method to solve this problem up to time  $T = 2$ . The local truncation error is

$$\begin{aligned} \tau(t) &= \frac{1}{2}ku''(t) + O(k^2) \\ &= -\frac{1}{2}k \cos(t) + O(k^2) \end{aligned} \tag{8.1}$$

Since the function  $f(t) = -\sin t$  is independent of  $u$ , it is Lipschitz continuous with Lipschitz constant  $L = 0$ , and so the error estimate (7.12) shows that

$$|E^n| \leq T \|\tau\|_\infty = k \max_{0 \leq t \leq T} |\cos t| = k.$$

Suppose we want to compute a solution with  $|E| \leq 10^{-3}$ . Then we should be able to take  $k = 10^{-3}$  and obtain a suitable solution after  $T/k = 2000$  time steps. Indeed, calculating using  $k = 10^{-3}$  gives a computed value  $U^{2000} = -0.415692$  with an error  $E^{2000} = U^{2000} - \cos(2) = 0.4548 \times 10^{-3}$ .

**Example 8.8.** Now suppose we modify the above equation to

$$u'(t) = \lambda(u - \cos t) - \sin t \tag{8.2}$$

Table 8.1: Errors in the computed solution using Euler’s method in Example 8.9, for different values of the time step  $k$ . Note the dramatic change in behavior of the error for  $k < 0.000952$ .

$k$	Error
0.001000	0.145252E+77
0.000976	0.588105E+36
0.000950	0.321089E-06
0.000800	0.792298E-07
0.000400	0.396033E-07

where  $\lambda$  is some constant. If we take the same initial data as before,  $u(0) = 0$ , then the solution is also the same as before,  $u(t) = \sin t$ . As a concrete example, let’s take  $\lambda = -10$ . Now how small do we need to take  $k$  in order to get an error that is  $10^{-3}$ ? Since the LTE (8.1) depends only on the true solution  $u(t)$ , which is unchanged from Example 8.7, we might hope that we could use the same  $k$  as in that example,  $k = 10^{-3}$ . Solving the problem using Euler’s method with this step size now gives  $U^{2000} = -0.416163$  with an error  $E^{2000} = 0.161 \times 10^{-4}$ . We are again successful. In fact the error is considerably smaller in this case than in the previous example, for reasons that will become clear later.

**Example 8.9.** Now consider the problem (8.2) with  $\lambda = -2100$  and the same data as before. Again the solution is unchanged and so is the LTE. But now if we compute with the same step size as before,  $k = 10^{-3}$ , we obtain  $U^{2000} = -0.2453 \times 10^{77}$  with an error of magnitude  $10^{77}$ . The computation behaves in an “unstable” manner, with an error that grows exponentially in time. Since the method is zero-stable and  $f(u, t)$  is Lipschitz continuous in  $u$  (with Lipschitz constant  $L = 2100$ ), we know that the method is convergent, and indeed with sufficiently small time steps we achieve very good results. Table 8.1 shows the error at time  $T = 2$  when Euler’s method is used with various values of  $k$ . Clearly something dramatic happens between the values  $k = 0.000976$  and  $k = 0.000952$ . For smaller values of  $k$  we get very good results whereas for larger values of  $k$  there is no accuracy whatsoever.

The equation (8.2) is a linear equation of the form (7.3) and so the analysis of Section 7.3.1 applies directly to this problem. From (7.7) we see that the global error  $E^n$  satisfies the recursion relation

$$E^{n+1} = (1 + k\lambda)E^n - k\tau^n \quad (8.3)$$

where the local error  $\tau^n = \tau(t_n)$  from (8.1). The expression (8.3) reveals the source of the exponential growth in the error — in each time step the previous error is multiplied by a factor  $(1 + k\lambda)$ . For the case  $\lambda = -2100$  and  $k = 10^{-3}$ , we have  $1 + k\lambda = -1.1$  and so we expect the local error introduced in step  $m$  to grow by a factor of  $(-1.1)^{n-m}$  by the end of  $n$  steps (recall (7.8)). After 2000 steps we expect the truncation error introduced in the first step to have grown by a factor of roughly  $(-1.1)^{2000} \approx 10^{82}$ , which is consistent with the error actually seen.

Note that in Example 8.8 with  $\lambda = -10$ , we have  $1 + k\lambda = 0.99$ , causing a *decay* in the effect of previous errors in each step. This explains why we got a reasonable result in Example 8.8 and in fact a better result than in Example 8.7, where  $1 + k\lambda = 1$ .

Returning to the case  $\lambda = -2100$ , we expect to observe exponential growth in the error for any value of  $k$  greater than  $2/2100 = 0.00095238$ , since for any  $k$  larger than this we have  $|1 + k\lambda| > 1$ . For smaller time steps  $|1 + k\lambda| < 1$  and the effect of each local error decays exponentially with time rather than growing. This explains the dramatic change in the behavior of the error that we see as we cross the value  $k = 0.00095$  in Table 8.1.

Note that the exponential growth of errors does not contradict zero-stability or convergence of the method in any way. The method does converge as  $k \rightarrow 0$ . In fact the bound (7.12),

$$|E^n| \leq e^{|\lambda|T} T \|\tau\|_\infty = O(k) \text{ as } k \rightarrow 0,$$

that we used to prove convergence, allows the possibility of exponential growth with time. The bound is valid for all values of  $k$ , but since  $T e^{|\lambda|T} = 2e^{4200} = 10^{1825}$  while  $\|\tau\|_\infty = \frac{1}{2}k$ , this bound does not guarantee any accuracy whatsoever in the solution until  $k < 10^{-1825}$ ! This is a good example that a

mathematical convergence proof may be a far cry from what is needed in practice. (Though in this case the convergence proof can be improved to give a much more realistic error bound as in Section 7.3.4.)

## 8.2 Absolute stability

In order to determine whether a numerical method will produce reasonable results with a given value of  $k > 0$ , we need a notion of stability that is different from zero-stability. There are a wide variety of other forms of “stability” that have been studied in various contexts. The one which is most basic and suggests itself from the above examples is *absolute stability*. This notion is based on the linear test equation (7.3) although a study of the absolute stability of a method yields information that is typically directly useful in determining an appropriate time step in nonlinear problems as well.

In fact we can look at the simplest case of the test problem in which  $g(t) = 0$  and we have simply

$$u'(t) = \lambda u(t).$$

Euler’s method applied to this problem gives

$$U^{n+1} = (1 + k\lambda)U^n$$

and we say that this method is *absolutely stable* when  $|1 + k\lambda| \leq 1$ , otherwise it is unstable. Note that there are two parameters  $k$  and  $\lambda$ , but it is only their product  $z \equiv k\lambda$  that matters. The method is stable whenever  $-2 \leq z \leq 0$ , and we say that the *interval of absolute stability* for Euler’s method is  $[-2, 0]$ .

It is more common to speak of the *region of absolute stability* as a region in the complex  $z$  plane, allowing the possibility that  $\lambda$  is complex (of course the time step  $k$  should be real and positive). The region of absolute stability (or simply the *stability region*) for Euler’s method is the disk of radius 1 centered at the point  $-1$ , since within this disk we have  $|1 + k\lambda| \leq 1$  (see Figure 8.1a). Allowing  $\lambda$  to be complex comes from the fact that in practice we are typically solving a nonlinear system of ODEs. After linearization we obtain a linear system of equations and it is the eigenvalues of the resulting Jacobian matrix that are important in determining stability. (This is discussed in Section 8.5.3.) Hence  $\lambda$  represents a typical eigenvalue and these may be complex even if the matrix is real.

## 8.3 Stability regions for LMMs

For a general LMM of the form (6.21), the region of absolute stability is found by applying the method to  $u' = \lambda u$ , obtaining

$$\sum_{j=0}^r \alpha_j U^{n+j} = k \sum_{j=0}^r \beta_j \lambda U^{n+j}$$

which can be rewritten as

$$\sum_{j=0}^r (\alpha_j - z\beta_j) U^{n+j} = 0. \tag{8.4}$$

Note again that it is only the product  $z = k\lambda$  that is important, not the values of  $k$  or  $\lambda$  separately. Note also that this is a dimensionless quantity since the decay rate  $\lambda$  has dimensions  $\text{time}^{-1}$  while the time step has dimensions of time.

The recurrence (8.4) is a homogeneous linear difference equation of the same form considered in Section 7.4.1. The solution has the general form (7.27) where the  $\zeta_j$  are now the roots of the characteristic polynomial  $\sum_{j=0}^r (\alpha_j - z\beta_j)\zeta^j$ . This polynomial is often called the *stability polynomial* and denoted by  $\pi(\zeta; z)$ . It is a polynomial in  $\zeta$  but its coefficients depend on the value of  $z$ . The stability polynomial can be expressed in terms of the characteristic polynomials for the LMM as

$$\pi(\zeta; z) = \rho(\zeta) - z\sigma(\zeta).$$

The LMM is absolutely stable for a particular value of  $z$  if errors introduced in one time step do not grow in future time steps. According to the theory of Section 7.4.1, this requires that the polynomial  $\pi(\zeta; z)$  satisfy the root condition (7.35).

**Definition 8.3.1** *The region of absolute stability for the LMM (6.21) is the set of points  $z$  in the complex plane for which the polynomial  $\pi(\zeta; z)$  satisfies the root condition (7.35).*

Note that a LMM is zero-stable if and only if the origin  $z = 0$  lies in the stability region. This explains the name.

**Example 8.10.** For Euler's method,

$$\pi(\zeta; z) = \zeta - (1 + z)$$

with the single root  $\zeta_1 = 1 + z$ . We have already seen that the stability region is the circle in Figure 8.1a.

**Example 8.11.** For the Backward Euler method (6.15),

$$\pi(\zeta; z) = (1 - z)\zeta - 1$$

with root  $\zeta_1 = (1 - z)^{-1}$ . We have

$$|(1 - z)^{-1}| \leq 1 \iff |1 - z| \geq 1$$

so the stability region is the *exterior* of the disk of radius 1 centered at  $z = 1$ , as shown in Figure 8.1b.

**Example 8.12.** For the Trapezoidal method (6.16),

$$\pi(\zeta; z) = \left(1 - \frac{1}{2}z\right)\zeta - \left(1 + \frac{1}{2}z\right)$$

with root

$$\zeta_1 = \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z}.$$

This is known as a *linear fractional transformation* in complex analysis and it can be shown that

$$|\zeta_1| \leq 1 \iff \operatorname{Re}(z) \leq 0$$

where  $\operatorname{Re}(z)$  is the real part. So the stability region is the left half plane as shown in Figure 8.1c.

**Example 8.13.** For the Midpoint method (6.17),

$$\pi(\zeta; z) = \zeta^2 - 2z\zeta - 1.$$

The roots are  $\zeta_{1,2} = z \pm \sqrt{z^2 + 1}$ . It can be shown that if  $z$  is a pure imaginary number of the form  $z = i\alpha$  with  $|\alpha| < 1$ , then  $|\zeta_1| = |\zeta_2| = 1$  and  $\zeta_1 \neq \zeta_2$ , and hence the root condition is satisfied. For any other  $z$  the root condition is not satisfied. So the stability region consists only of the open interval from  $-i$  to  $i$  on the imaginary axis, as shown in Figure 8.1d.

Since  $k$  is always real, this means the Midpoint method is only useful on the test problem  $u' = \lambda u$  if  $\lambda$  is pure imaginary. The method is not very useful for scalar problems where  $\lambda$  is typically real, but the method is of great interest in some applications with systems of equations. For example, if the matrix is real but skew symmetric ( $A^T = -A$ ), then the eigenvalues are pure imaginary. This situation arises naturally in the discretization of hyperbolic partial differential equations, as discussed later.

**Example 8.14.** Figures 8.2 and 8.3 show the stability regions for the  $r$ -step Adams-Bashforth and Adams-Moulton methods for various values of  $r$ . For an  $r$ -step method the polynomial  $\pi(\zeta; z)$  has degree  $r$  and there are  $r$  roots. Determining the values of  $z$  for which the root condition is satisfied does not appear simple. However, there is an elegant technique called the *boundary locus method* that makes it simple to determine the regions shown in the figures. This is briefly described in the next section (see also Lambert[Lam73], for example).

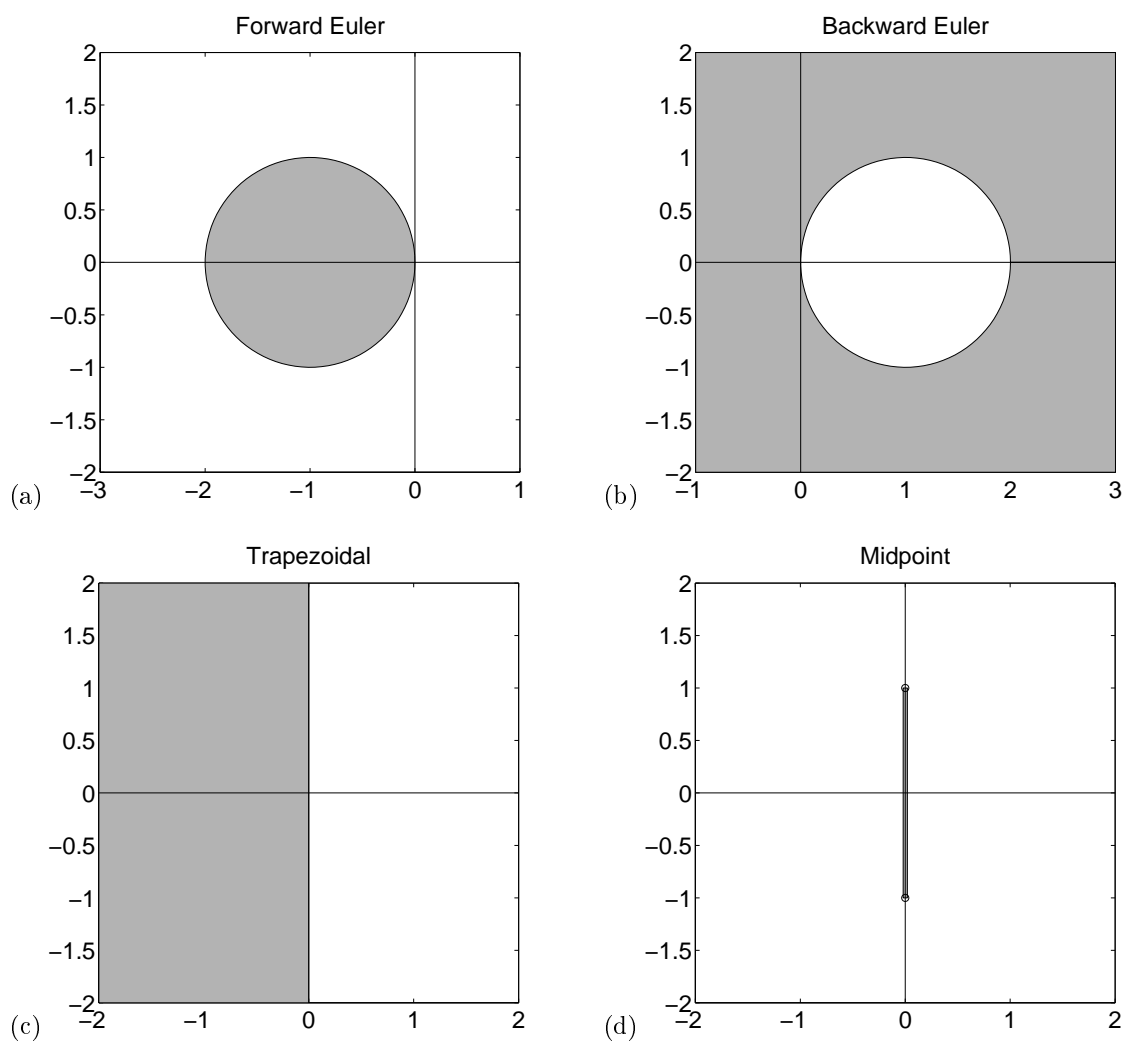


Figure 8.1: Stability regions for (a) Euler — interior of circle, (b) Backward Euler, — exterior of circle (c) Trapezoidal — left half plane, and (d) Midpoint — segment on imaginary axis.

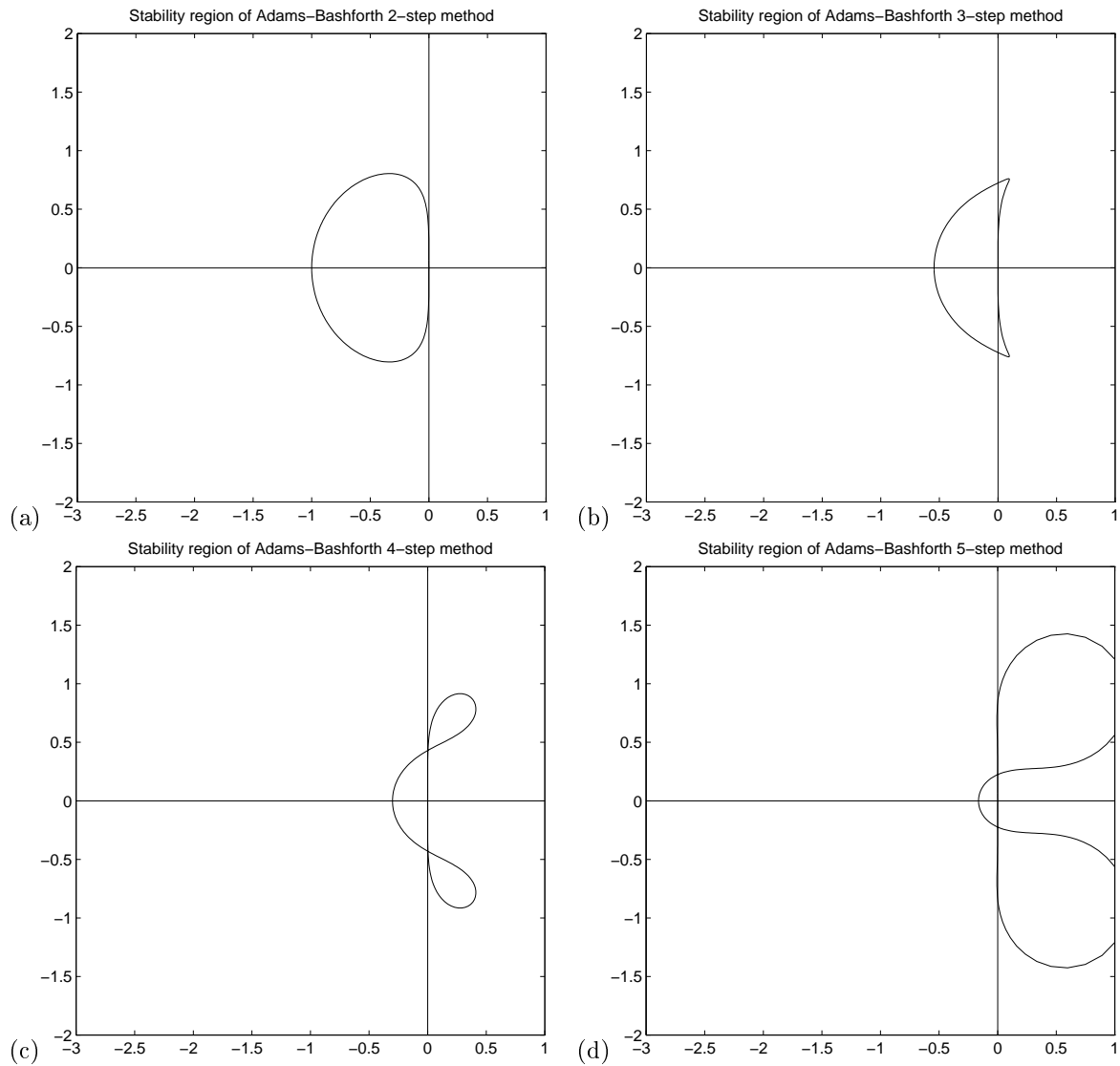


Figure 8.2: Stability regions for some Adams-Bashforth methods. The enclosed region just to the left of the origin is the region of absolute stability. See Section 8.4 for a discussion of the other loops seen in some figures.

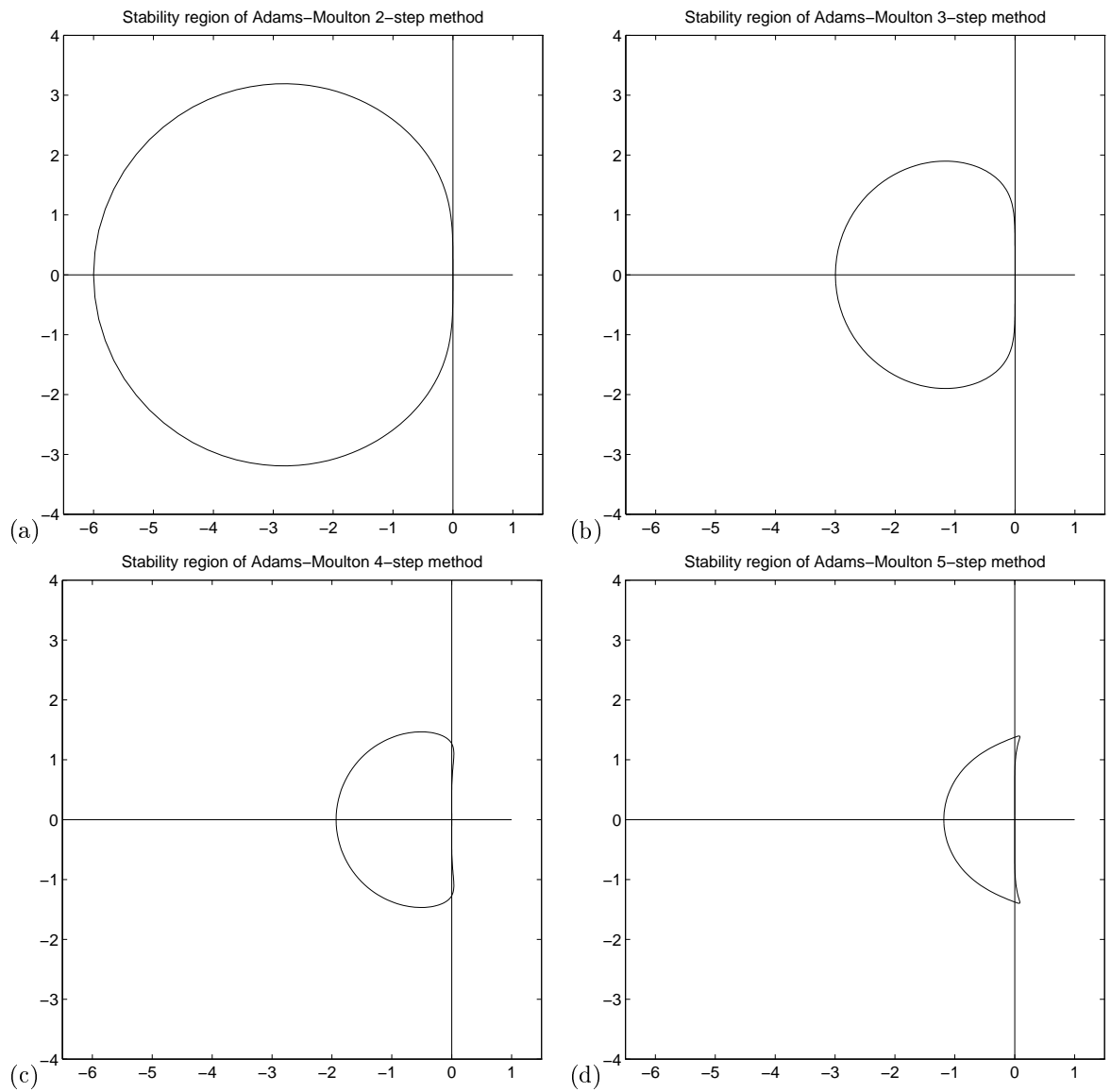


Figure 8.3: Stability regions for some Adams-Moulton methods. The stability region is interior to the curves.



## 8.4 The Boundary Locus Method

A point  $z \in \mathbb{C}$  is in the stability region  $\mathcal{S}$  of a linear multistep method if the stability polynomial  $\pi(\zeta; z)$  satisfies the root condition for this value of  $z$  (see Chapter 8). It follows that if  $z$  is on the boundary of the stability region then  $\pi(\zeta; z)$  must have at least one root  $\zeta_j$  with magnitude exactly equal to 1. Hence  $\zeta_j$  is of the form

$$\zeta_j = e^{i\theta}$$

for some value of  $\theta$  in the interval  $[0, 2\pi]$  (Sorry for the two different uses of  $\pi$ .) Since  $\zeta_j$  is a root of  $\pi$ , we have

$$\pi(e^{i\theta}; z) = 0$$

for this particular combination of  $z$  and  $\theta$ . Recalling the definition of  $\pi$ , this gives

$$\rho(e^{i\theta}) - z\sigma(e^{i\theta}) = 0$$

and hence

$$z = \frac{\rho(e^{i\theta})}{\sigma(e^{i\theta})}.$$

If we know  $\theta$  then we can find  $z$  from this.

Since every point  $z$  on the boundary of  $\mathcal{S}$  must be of this form for some value of  $\theta$  in  $[0, 2\pi]$ , we can simply plot the parametrized curve

$$\tilde{z}(\theta) \equiv \frac{\rho(e^{i\theta})}{\sigma(e^{i\theta})} \tag{8.5}$$

for  $0 \leq \theta \leq 2\pi$  to find the locus of all points which are *potentially* on the boundary of  $\mathcal{S}$ . For simple methods this yields the region  $\mathcal{S}$  directly.

**Example 8.15.** For Euler's method we have  $\rho(\zeta) = \zeta - 1$  and  $\sigma(\zeta) = 1$ , and so

$$\tilde{z}(\theta) = e^{i\theta} - 1.$$

This function maps  $[0, 2\pi]$  to the unit circle centered at  $z = -1$ , which is exactly the boundary of  $\mathcal{S}$ .

To determine which side of this curve is the *interior* of  $\mathcal{S}$ , we need only evaluate the roots of  $\pi(\zeta; z)$  at some random point  $z$  on one side or the other and see if the polynomial satisfies the root condition. For example, at the center of the circle  $z = 1$  we have  $\pi(\zeta; 1) = -\zeta$  with root  $\zeta_1 = 0$ . So  $\pi$  satisfies the root condition here and this point must be in the interior. It follows easily that all other points inside the curve must also be in the interior of  $\mathcal{S}$ , since the roots are continuous functions of  $z$  and so as we vary  $z$  a root can potentially move outside the unit circle only when we cross the boundary locus. Checking the roots at a random point outside the circle shows that these points must be exterior to the stability region.

For some methods the boundary locus may cross itself. In this case we typically find that at most one of the regions cut out of the plane corresponds to the stability region. We can determine which region is  $\mathcal{S}$  by evaluating the roots at some convenient point  $z$  within each region.

**Example 8.16.** The 5-step Adams-Bashforth method gives the boundary locus seen in Figure 8.4. The numbers 0, 1, 2 in the different regions indicate how many roots are outside of the unit circle in each region, as determined by testing the roots at a random point in each region. The stability region is the small semicircular region to the left of the origin where all roots are inside the unit circle. As we cross the boundary of this region one root moves outside. As we cross the boundary locus into one of the regions marked "2", another root moves outside and the method is still unstable in these regions.

**Exercise 8.1** Adapt the idea of the boundary locus method to determine the stability region of the 2-stage Runge-Kutta method from Example 6.15. Note that this is also the stability region of the second order Taylor series method.

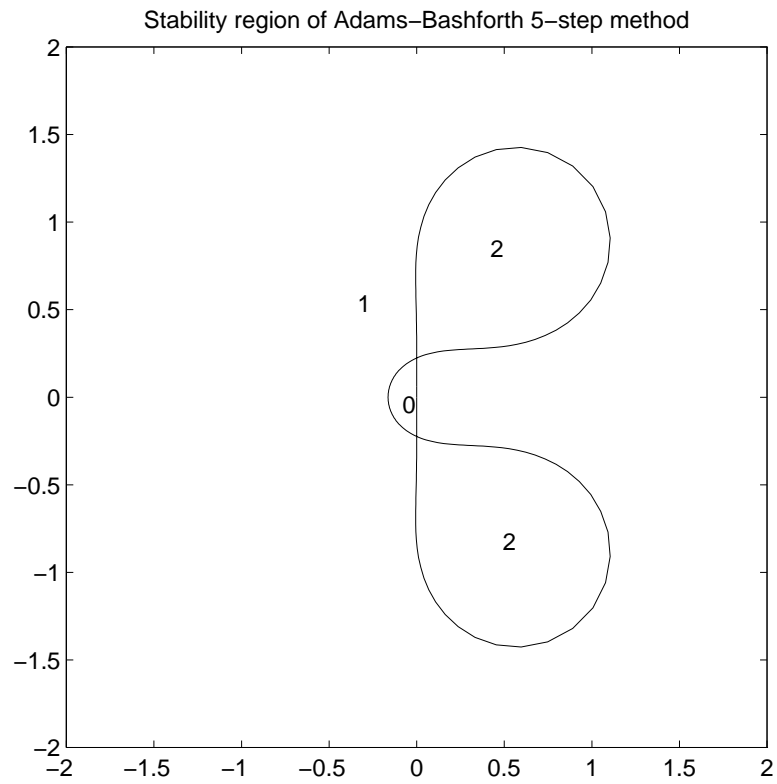


Figure 8.4: Boundary locus for the 5-step Adams-Bashforth method. The numbers indicate how many roots are outside of the unit circle in each region. The region marked “0” is the stability region.

## 8.5 Systems of equations

In this section we will look at how the stability theory carries over to systems of  $m$  differential equations. We will see that for a linear system  $u' = Au + b$ , where  $A$  is an  $m \times m$  matrix, it is the eigenvalues of  $A$  that are important and we need  $k\lambda$  in the stability region for each eigenvalue  $\lambda$  of  $A$ . For general nonlinear systems  $u' = f(u)$ , the theory is more complicated, but a good rule of thumb is that  $k\lambda$  should be in the stability region for each eigenvalue  $\lambda$  of the Jacobian matrix  $f'(u)$ . Before discussing this theory, we will review the theory of chemical kinetics, a field where the solution of systems of ordinary differential equations is very important.

### 8.5.1 Chemical Kinetics

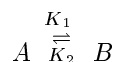
Let  $A, B, C, \dots$  represent chemical compounds and consider a reaction of the form



This represents a reaction in which  $A$  is transformed into  $B$  with rate  $K_1 > 0$ . If we let  $u_1$  represent the concentration of  $A$  and  $u_2$  represent the concentration of  $B$  (often denoted by  $u_1 = [A]$ ,  $u_2 = [B]$ ) then the ODEs for  $u_1$  and  $u_2$  are:

$$\begin{aligned} u_1' &= -K_1 u_1 \\ u_2' &= K_1 u_1. \end{aligned}$$

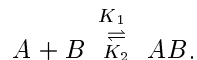
If there is also a reverse reaction at rate  $K_2$ , we write



and the equations then become

$$\begin{aligned} u_1' &= -K_1 u_1 + K_2 u_2 \\ u_2' &= K_1 u_1 - K_2 u_2 \end{aligned} \tag{8.6}$$

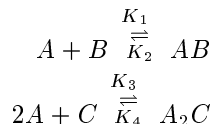
More typically, reactions involve combinations of two or more compounds, e.g.,



Since  $A$  and  $B$  must combine to form  $AB$ , the rate of the forward reaction is proportional to the *product* of  $u_1$  and  $u_2$ , while the backward reaction is proportional to  $u_3 = [AB]$ . The equations become

$$\begin{aligned} u_1' &= -K_1 u_1 u_2 + K_2 u_3 \\ u_2' &= -K_1 u_1 u_2 + K_2 u_3 \\ u_3' &= K_1 u_1 u_2 - K_2 u_3 \end{aligned} \tag{8.7}$$

Often several reactions take place simultaneously, e.g.,



If we now let  $u_4 = [C]$ ,  $u_5 = [A_2C]$ , then the equations are

$$\begin{aligned} u_1' &= -K_1 u_1 u_2 + K_2 u_3 - 2K_3 u_1^2 u_4 + 2K_4 u_5 \\ u_2' &= -K_1 u_1 u_2 + K_2 u_3 \\ u_3' &= K_1 u_1 u_2 - K_2 u_3 \\ u_4' &= -K_3 u_1^2 u_4 + K_4 u_5 \\ u_5' &= K_3 u_1^2 u_4 - K_4 u_5 \end{aligned} \tag{8.8}$$

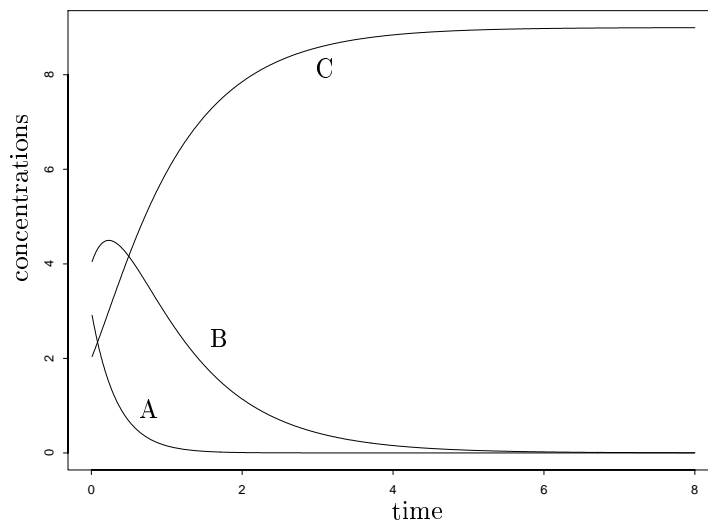
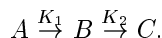


Figure 8.5: Sample solution for the kinetics problem in Example 8.17.

Interesting physical problems can give rise to very large systems of ODEs. Frequently the rate constants  $K_1, K_2, \dots$  are of vastly different orders of magnitude. This leads to **stiff** systems of equations, as discussed in Chapter 10.

**Example 8.17.** One particularly simple system arises from the decay process



Let  $u_1 = [A]$ ,  $u_2 = [B]$ ,  $u_3 = [C]$ . Then the system is linear and has the form  $u' = Au$ , where

$$A = \begin{bmatrix} -K_1 & 0 & 0 \\ K_1 & -K_2 & 0 \\ 0 & K_2 & 0 \end{bmatrix}. \quad (8.9)$$

Note that the eigenvalues are  $-K_1$ ,  $-K_2$  and 0. The general solution thus has the form (assuming  $K_1 \neq K_2$ )

$$u_j(t) = c_{j1}e^{-K_1 t} + c_{j2}e^{-K_2 t} + c_{j3}.$$

In fact, on physical grounds (since  $A$  decays into  $B$  which decays into  $C$ ), we expect that  $u_1$  simply decays to 0 exponentially,

$$u_1(t) = e^{-K_1 t} u_1(0)$$

(which clearly satisfies the first ODE), and also that  $u_2$  decays to 0, while  $u_3$  grows and asymptotically approaches the value  $u_1(0) + u_2(0) + u_3(0)$  as  $t \rightarrow \infty$ . A typical solution for  $K_1 = 3$  and  $K_2 = 1$  with  $u_1(0) = 3$ ,  $u_2(0) = 4$ , and  $u_3(0) = 2$  is shown in Figure 8.5.

## 8.5.2 Linear systems

Consider a linear system  $u' = Au$  where  $A$  is a constant  $m \times m$  matrix, and suppose for simplicity that  $A$  is diagonalizable, which means that it has a complete set of  $m$  linearly independent eigenvectors  $r_p$  satisfying  $Ar_p = \lambda_p r_p$  for  $p = 1, 2, \dots, m$ . Let  $R = [r_1, r_2, \dots, r_m]$  be the matrix of eigenvectors and  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$  be the diagonal matrix of eigenvalues. Then we have

$$A = R\Lambda R^{-1} \text{ and } \Lambda = R^{-1}AR.$$

Now let  $v(t) = R^{-1}u(t)$ . Multiplying  $u' = Au$  by  $R^{-1}$  on both sides and introducing  $I = RR^{-1}$  gives the equivalent equations

$$R^{-1}u'(t) = (R^{-1}AR)(R^{-1}u(t)),$$

i.e.,

$$v'(t) = \Lambda v(t).$$

This is a diagonal system of equations that decouples into  $m$  independent scalar equations, one for each component of  $v$ . The  $p$ th such equation is

$$v_p'(t) = \lambda_p v_p(t).$$

A Linear Multistep Method applied to the linear ODE can also be decoupled in the same way. For example, if we apply Euler's method, we have

$$U^{n+1} = U^n + kAU^n$$

which, by the same transformation, can be rewritten as

$$V^{n+1} = V^n + k\Lambda V^n$$

where  $V^n = R^{-1}U^n$ . This decouples into  $m$  independent numerical methods, one for each component of  $V^n$ . These take the form

$$V_p^{n+1} = (1 + k\lambda_p)V_p^n.$$

We can recover  $U^n$  from  $V^n$  from  $U^n = RV^n$ .

The overall method is stable if each of the scalar problems is stable, and this clearly requires that  $k\lambda_p$  be in the stability region of Euler's method for all values of  $p$ . The same technique can be used more generally to show that a LMM is absolutely stable if  $k\lambda_p$  is in the stability region of the method for each eigenvalue  $\lambda_p$  of the matrix  $A$ .

**Example 8.18.** Consider the linear kinetics problem with  $A$  given by (8.9). Since this matrix is upper triangular, the eigenvalues are the diagonal elements  $\lambda_1 = -K_1$ ,  $\lambda_2 = -K_2$ , and  $\lambda_3 = 0$ . The eigenvalues are all real and the Euler's method is stable provided  $k \max(K_1, K_2) \leq 2$ .

**Example 8.19.** Consider a linearized model for a swinging pendulum, this time with frictional forces added,

$$\theta''(t) = -a\theta(t) - b\theta'(t)$$

which is valid for small values of  $\theta$ . If we introduce  $u_1 = \theta$  and  $u_2 = \theta'$  then we obtain a first order system  $u' = Au$  with

$$A = \begin{bmatrix} 0 & 1 \\ -a & -b \end{bmatrix}. \quad (8.10)$$

The eigenvalues of this matrix are  $\lambda = \frac{1}{2}(-b \pm \sqrt{b^2 - 4a})$ . Note in particular that if  $b = 0$  (no damping) then  $\lambda = \pm\sqrt{-a}$  are pure imaginary. For  $b > 0$  the eigenvalues shift into the left half plane. In the undamped case the Midpoint method would be a reasonable choice, whereas Euler's method might be expected to have difficulties. In the damped case the opposite is true.

### 8.5.3 Nonlinear systems

Now consider a nonlinear system  $u' = f(u)$ . The stability analysis we have developed for the linear problem does not apply directly to this system. However, if the solution is slowly varying relative to the time step, then over a small time interval we would expect a linearized approximation to give a good indication of what is happening. Suppose the solution is near some value  $\bar{u}$ , and let  $v(t) = u(t) - \bar{u}$ . Then

$$v'(t) = u'(t) = f(u(t)) = f(v(t) + \bar{u}).$$

Taylor series expansion about  $\bar{u}$  (assuming  $v$  is small) gives

$$v'(t) = f(\bar{u}) + f'(\bar{u})v(t) + O(v^2).$$

Dropping the  $O(v^2)$  terms gives a linear system

$$v'(t) = Av(t) + b$$

where  $A = f'(\bar{u})$  is the Jacobian matrix evaluated at  $\bar{u}$  and  $b = f(\bar{u})$ . Examining how the numerical method behaves on this linear system (for each relevant value of  $\bar{u}$ ) gives a good indication of how it will behave on the nonlinear system.

**Example 8.20.** Consider the kinetics problem (8.7). The Jacobian matrix is

$$A = \begin{bmatrix} -K_1 u_2 & -K_1 u_1 & K_2 \\ -K_1 u_2 & -K_1 u_1 & K_2 \\ K_1 u_2 & K_1 u_1 & -K_2 \end{bmatrix}$$

with eigenvalues  $\lambda_1 = -K_1(u_1 + u_2) - K_2$  and  $\lambda_2 = \lambda_3 = 0$ . Since  $u_1 + u_2$  is simply the total quantity of species  $A$  and  $B$  present, this can be bounded for all time in terms of the initial data. (For example, we certainly have  $u_1(t) + u_2(t) \leq u_1(0) + u_2(0) + 2u_3(0)$ .) So we can determine the range of  $\lambda_1$  along the negative real axis and hence how small  $k$  must be to stay within the region of absolute stability.

## 8.6 Choice of stepsize

As the examples at the beginning of this chapter illustrated, in order to obtain computed results that are within some error tolerance, we need two conditions satisfied:

1. The time step  $k$  must be small enough that the local truncation error is acceptably small. This gives a constraint of the form  $k \leq k_{\text{acc}}$  that depends on several things:
  - What method is being used, which determines the expansion for the local truncation error.
  - How smooth the solution is, which determines how large the high order derivatives occurring in this expansion are.
  - What accuracy is required.
2. The time step  $k$  must be small enough that the method is absolutely stable on this particular problem. This gives a constraint of the form  $k \leq k_{\text{stab}}$  that depends on the size of  $\lambda$  or, more generally, the eigenvalues of the Jacobian matrix  $f'(u)$ .

Typically we would like to choose our time step based on accuracy considerations. For a given method and problem, we would like to choose  $k$  so that the local error in each step is sufficiently small that the accumulated error will satisfy our error tolerance, assuming some “reasonable” growth of errors. If the errors grow exponentially with time because the method is not absolutely stable, however, then we would have to use a smaller time step in order to get useful results.

If stability considerations force us to use a *much* smaller time step than the local truncation error indicates should be needed, then this particular method is probably not optimal for this problem. We then have a “stiff” problem as discussed in Chapter 10, for which special methods have been developed.



## Chapter 9

# Linear Multistep Methods as 1-step Methods

A linear multistep method can be written in the form of a one-step method applied to a system of equations. Doing so allows us to see how studying absolute stability via the roots of the characteristic polynomial is related to studying the matrix iterations via the eigenvalues of the matrix. Moreover we can see how the condition of zero-stability leads to a convergence proof similar to the proof of convergence of a one-step method seen in Chapter 7.

For simplicity we will only consider *explicit* methods applied to the *linear* problem  $u'(t) = \lambda u(t) + g(t)$ , for which a general  $r$ -step method has the form

$$\sum_{j=0}^r \alpha_j U^{n+j} = k \sum_{j=0}^{r-1} \beta_j [\lambda U^{n+j} + g(t_{n+j})]. \quad (9.1)$$

We can assume without loss of generality that  $\alpha_r = 1$ , since there is one degree of freedom in the coefficients because we can multiply both sides of (9.1) by an arbitrary constant without changing the numerical method. Then we can solve for  $U^{n+r}$  as

$$U^{n+r} = \sum_{j=0}^{r-1} [(-\alpha_j + k\lambda\beta_j)U^{n+j} + k\beta_j g(t_{n+j})].$$

Let  $\gamma_j = -\alpha_j + k\lambda\beta_j$  and define the vector  $V^n$  by

$$V^n = \begin{bmatrix} U^n \\ U^{n+1} \\ \vdots \\ U^{n+r-1} \end{bmatrix}, \quad \text{so that } V^{n+1} = \begin{bmatrix} U^{n+1} \\ U^{n+2} \\ \vdots \\ U^{n+r} \end{bmatrix}.$$

(If  $U^n \in \mathbb{R}^s$  was already a vector for a system of  $s$  equations, then  $V^n \in \mathbb{R}^{rs}$  is a longer vector.) With this notation we see that  $V^n$  is updated by a 1-step method

$$V^{n+1} = AV^n + kb^n$$

where

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & & \ddots & & \vdots \\ 0 & 0 & \cdots & & & 1 \\ \gamma_0 & \gamma_1 & \cdots & & & \gamma_{r-1} \end{bmatrix}, \quad b^n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \sum_{j=0}^{r-1} \beta_j g(t_{n+j}) \end{bmatrix}.$$



Note that we know the initial vector  $V^0$  from the starting values  $U^0, \dots, U^{r-1}$  needed for the  $r$ -step method.

In the usual way, we find that the error vector  $E^n$  satisfies

$$E^{n+1} = AE^n - k\vec{\tau}^n,$$

where

$$\vec{\tau}^n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \tau(t_{n+r}) \end{bmatrix},$$

and hence

$$E^n = A^n E^0 - k \sum_{m=1}^n A^{n-m} \vec{\tau}^{m-1}. \quad (9.2)$$

(Here the superscripts on  $E$  and  $\vec{\tau}$  are time indices while the superscripts on  $A$  are powers.)

## 9.1 Absolute stability

First let's consider the question of absolute stability. How does the error behave as we march forward in time with some fixed  $k$ ? We want the effect of past errors to not be amplified in future steps, so we want a uniform bound on  $\|A^n\|$ . A necessary condition is that the spectral radius of  $A$  must be no larger than 1. If  $A$  is diagonalizable then this is also a sufficient condition.

The eigenvalues  $\mu$  of  $A$  satisfy

$$\det(A - \mu I) = 0.$$

It is easy to verify that this determinant is simply

$$\mu^r - \gamma_{r-1}\mu^{r-1} - \dots - \gamma_1\mu - \gamma_0.$$

Recalling that  $\alpha_r = 1$  and  $\beta_r = 0$ , we see that this is exactly the "stability polynomial"  $\pi(\mu; k\lambda)$  defined in Chapter 8. So for stability we require that the roots of this polynomial (which are the eigenvalues of  $A$ ) be no larger than 1 in modulus. The matrix  $A$  is called the *companion matrix* for this polynomial.

Moreover, it turns out that if  $\pi$  has a repeated root, which is then a multiple eigenvalue of  $A$ , that  $A$  is **not** diagonalizable — it is *defective* and has a nontrivial Jordan block. If this eigenvalue has modulus equal to 1 then  $\|A^n\|$  will grow with  $n$ . Hence for stability we must also require that any multiple root have modulus strictly less than 1.

The upshot is that the condition required for uniform boundedness of  $\|A^n\|$  is exactly the "root condition" on the stability polynomial  $\pi$ , i.e., absolute stability of the method.

## 9.2 Convergence and zero-stability

Now let's examine the question of convergence as  $k \rightarrow 0$ . Looking again at (9.2), and following the proof in Section 7.3, we see that the method converges if  $\|\vec{\tau}^n\|_\infty \rightarrow 0$  (consistency),  $\|E^0\| \rightarrow 0$  (good starting values), and  $\|A^n\|$  is uniformly bounded as  $k \rightarrow 0$  for  $nk \leq T$  (which we will see requires exactly the conditions of zero-stability). In the case of a 1-step method,  $A$  was a scalar bounded by  $1 + L'k$ , and hence  $A^n \leq e^{L'T}$  could be uniformly bounded.

Note that we are looking for something different than the uniform boundedness of the previous section. Rather than looking at how  $A^n$  behaves as  $n \rightarrow \infty$  with  $k$  fixed, we are examining what happens as  $k \rightarrow 0$  and  $n \rightarrow \infty$  simultaneously with  $nk = T$  fixed. To make this clearer, let's now write

the matrix as  $A_k$  to remind ourselves that it depends on  $k$ . As  $k \rightarrow 0$  the value  $\gamma_j$  in the last row of  $A_k$  approaches  $-\alpha_j$ . It can be shown that  $\|A_k^n\|$  is uniformly bounded for all  $n$  and  $k$  with  $nk \leq T$  provided the eigenvalues of the limiting matrix  $A_0$  are no greater than 1 in modulus (and strictly less than 1 for multiple eigenvalues). The determinant of  $A_0 - \mu I$  is just the characteristic polynomial  $\rho(\mu) = \sum_{j=0}^r \alpha_j \mu^j$ , and so this is just the zero-stability condition.

The proof of this statement follows from the fact that the eigenvalues of a matrix are continuous functions of the matrix elements, and moreover that modifying the elements by  $O(k)$  leads to an  $O(k)$  perturbation of the eigenvalues (in the nondefective case). Hence if the characteristic polynomial  $\rho$  satisfies the root condition then the matrix  $A_k$  for  $k > 0$  has spectral radius that can be bounded by  $1 + Ck$  for some constant  $C$ , and hence the  $n$ 'th power can be bounded by some constant times  $(1 + Ck)^n \leq e^{CT}$ , just as in the case of a scalar 1-step method.



## Chapter 10

# Stiff ODEs

The problem of *stiffness* leads to computational difficulty in many practical problems. The classic example is the case of a stiff ordinary differential equation that we will examine in this chapter. In general a problem is *stiff* if, roughly speaking, we are attempting to compute a particular solution that is smooth and slowly varying (relative to the time interval of the computation), but in a context where the nearby solution curves are much more rapidly varying. In other words if we perturb the solution slightly at any time, the resulting solution curve through the perturbed data has rapid variation. Typically this takes the form of a short lived “transient” response that moves the solution back towards a smooth solution.

**Example 10.1.** Consider the ODE (8.2) from the previous chapter,

$$u'(t) = \lambda(\cos t - u) - \sin t. \quad (10.1)$$

One particular solution is the function  $u(t) = \cos t$ , and this is the solution with the initial data  $u(0) = 1$  considered previously. This smooth function is a solution for any value of  $\lambda$ . If we consider initial data of the form  $u(t_0) = \eta$  that does not lie on this curve, then the solution through this point is a different function, of course. However, if  $\lambda < 0$  (or  $\text{Re}(\lambda) < 0$  more generally), this function approaches  $\cos t$  exponentially quickly, with decay rate  $\lambda$ . It is easy to verify that the solution is

$$u(t) = e^{\lambda(t-t_0)}(\eta - \cos(t_0)) + \cos t. \quad (10.2)$$

Figure 10.1 shows a number of different solution curves for this equation with different choices of  $t_0$  and  $\eta$ , with the fairly modest value  $\lambda = -1$ . Figure 10.1b shows the corresponding solution curves when  $\lambda = -10$ .

In this scalar example, when we perturb the solution at some point it quickly relaxes towards the particular solution  $u(t) = \cos t$ . In other stiff problems the solution might move quickly towards some different smooth solution, as seen in the next example.

**Example 10.2.** Consider the kinetics model  $A \rightarrow B \rightarrow C$  developed in Example 8.17. The system of equations is given by (8.9). Suppose that  $K_1 \gg K_2$  so that a typical solution appears as in Figure 10.2(a). (Here  $K_1 = 20$  and  $K_2 = 1$ . Compare this to Figure 8.5.) Now suppose at time  $t = 1$  we perturb the system by adding more of species  $A$ . Then the solution behaves as shown in Figure 10.2(b). The additional  $A$  introduced is rapidly converted into  $B$  (the fast transient response) and then slowly from  $B$  into  $C$ . After the rapid transient the solution is again smooth, though it differs from the original solution since the final asymptotic value of  $C$  must be higher than before by the same magnitude as the amount of  $A$  introduced.

### 10.1 Numerical Difficulties

Stiffness causes numerical difficulties because any finite difference method is constantly introducing errors. The local truncation error acts as a perturbation to the system that moves us away from the

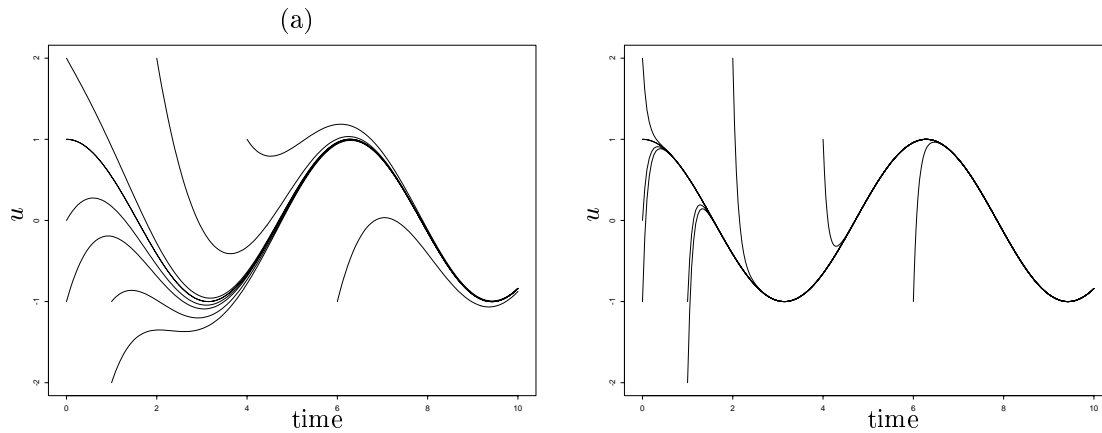


Figure 10.1: Solution curves for the ODE (10.1) for various initial values. (a) With  $\lambda = -1$ . (b) With  $\lambda = -10$  and the same set of initial values.

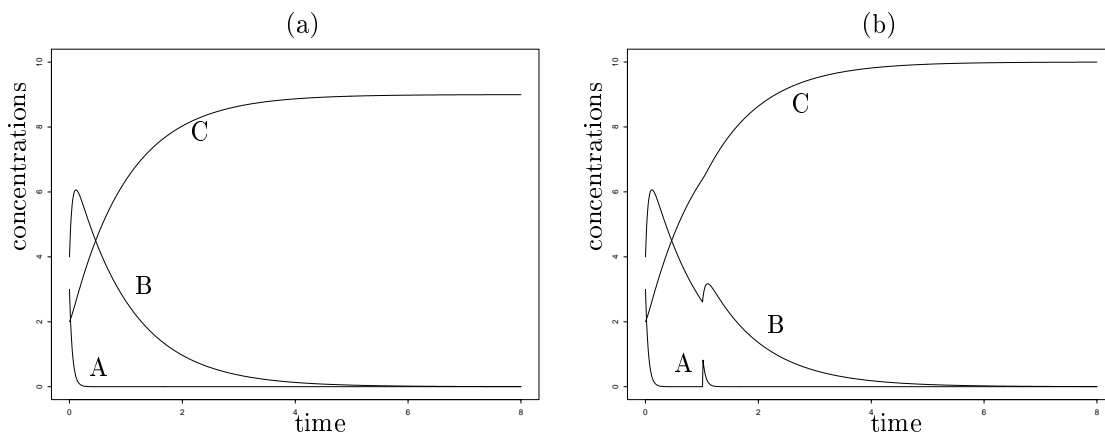


Figure 10.2: Solution curves for the kinetics problem in Example 8.17, with  $K_1 = 20$  and  $K_2 = 1$ . In (b) a perturbation has been made by adding one unit of species  $A$  at time  $t = 1$ .

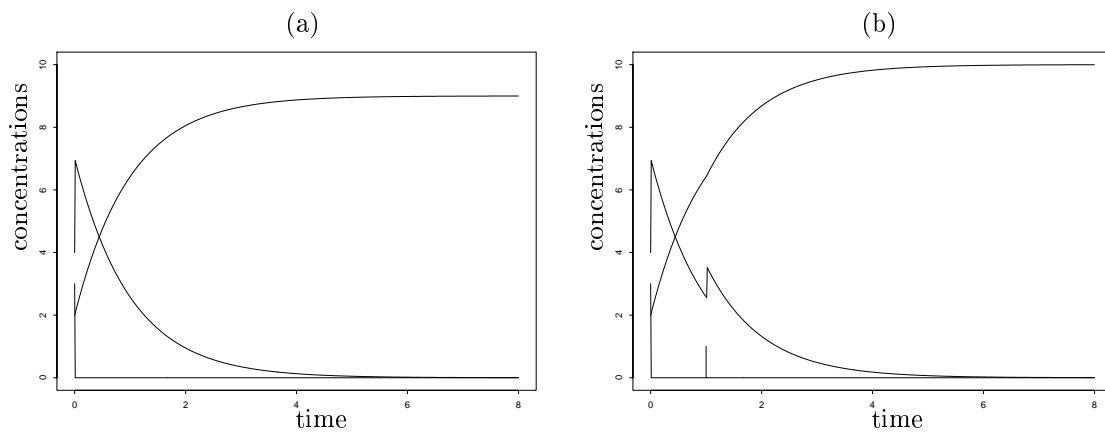


Figure 10.3: Solution curves for the kinetics problem in Example 8.17, with  $K_1 = 10^6$  and  $K_2 = 1$ . In (b) a perturbation has been made by adding one unit of species  $A$  at time  $t = 1$ .

smooth solution we are trying to compute. Why does this cause more difficulty in a stiff system than in other systems? At first glance it seems like the stiffness should work to our advantage. If we are trying to compute the solution  $u(t) = \cos t$  to the ODE (10.1) with initial data  $u(0) = 1$ , for example, then the fact that any errors introduced decay exponentially should help us. The true solution is very robust and the solution is almost completely insensitive to errors made in the past. In fact this *stability* of the true solution does help us, as long as the numerical method is also stable. (Recall that the results in Example 8.8 were much better than in Example 8.7.)

The difficulty arises from the fact that many numerical methods, including all explicit methods, are unstable (in the sense of absolute stability) *unless the time step is small relative to the time scale of the rapid transient*, which in a stiff problem is much smaller than the time scale of the smooth solution we are trying to compute. In the terminology of the previous chapter, this means that  $k_{\text{stab}} \ll k_{\text{acc}}$ . Although the true solution is smooth and it seems that a reasonably large time step would be appropriate, the numerical method must always deal with the rapid transients introduced in every time step and may need a very small time step to do so stably.

## 10.2 Characterizations of stiffness

A stiff ODE can be characterized by the property that  $f'(u)$  is much larger (in absolute value or norm) than  $u'(t)$ . The latter quantity measures the smoothness of the solution being computed, while  $f'(u)$  measures how rapidly  $f$  varies as we move away from this particular solution. Note that stiff problems typically have large Lipschitz constants too.

For *systems* of ODE's, stiffness is sometimes defined in terms of the “stiffness ratio” of the system, which is the ratio

$$\frac{\max |\lambda_p|}{\min |\lambda_p|}$$

over all eigenvalues of the Jacobian matrix  $f'(u)$ . If this is large then there is a large range of time scales present in the problem, a necessary component for stiffness to arise. While this is often a useful quantity, one should not rely entirely on this measure to determine whether a problem is stiff.

For one thing, it is possible even for a scalar problem to be stiff (as we have seen in Example 10.1), even though for a scalar problem the stiffness ratio is always 1 since there is only one eigenvalue. There can still be more than one time scale present. In (10.1) the fast time scale is determined by  $\lambda$ , the eigenvalue, and the slow time scale is determined by the inhomogeneous term  $\sin(t)$ . For systems of equations there may also be additional time scales arising from inhomogeneous forcing terms or other time-dependent coefficients that are distinct from the scales imposed by the eigenvalues.

It is also important to note that a system of ODEs which has a large “stiffness ratio” is not necessarily stiff! If the eigenvalue with large amplitude lies close to the imaginary axis, then it leads to highly oscillatory behavior in the solution rather than rapid damping. If the solution is rapidly oscillating then it will probably be necessary to take small time steps for accuracy reasons and  $k_{\text{acc}}$  may be roughly the same magnitude as  $k_{\text{stab}}$  even for explicit methods.

Finally, note that a particular problem may be stiff over some time intervals and nonstiff elsewhere. In particular, if we are computing a solution that has a rapid transient, such as the kinetics problem shown in Figure 10.3(a), then the problem is not stiff over the initial transient period where the true solution is as rapidly varying as nearby solution curves. Only for times greater than  $10^{-6}$  or so does the problem become stiff, once the desired solution curve is much smoother than nearby curves.

For the problem shown in Figure 10.3(b), there is another time interval just after  $t = 1$  over which the problem is again not stiff since the solution again exhibits rapid transient behavior and a small time step would be needed on the basis of accuracy considerations.

### 10.3 Numerical methods for stiff problems

Over time intervals where a problem is stiff, we would like to use a numerical method that has a large region of absolute stability, extending far into the left-half plane. The problem with a method like Euler's method, with a stability region that only extends out to  $\text{Re}(\lambda) = -2$ , is that the time step  $k$  is severely limited by the eigenvalue with largest magnitude, and we need to take  $k \approx 2/|\lambda_{\max}|$ . Over time intervals where this fastest time scale does not appear in the solution, we would like to be able to take much larger time steps. For example, in the problems shown in Figure 10.3, where  $K_1 = 10^6$ , we would need to take  $k \approx 2 \times 10^{-6}$  with Euler's method, requiring 4 million time steps to compute over the time interval shown in the Figure, even though the solution is very smooth over most of this time.

An analysis of stability regions shows that there are basically two different classes of LMM's: those for which the stability region is bounded and extends distance  $O(1)$  from the origin, such as Euler's method, the Midpoint method, or any of the Adams methods (see Figure 8.1 and Figure 10.5), and those for which the stability region is unbounded, such as Backward Euler or trapezoidal. Clearly the first class of methods are inappropriate for stiff problems.

Unfortunately, all explicit methods have bounded stability regions and hence are inefficient on stiff problems. Some implicit methods also have bounded stability regions, such as the Adams-Moulton methods.

#### 10.3.1 A-stability

It seems like it would be optimal to have a method whose stability region contains the entire left half plane. Then any time step would be allowed, provided that all the eigenvalues have negative real parts as is often the case in practice. The Backward Euler and Trapezoidal methods have this property, for example. A method with this property is said to be *A-stable*. Unfortunately, a theorem of Dahlquist states that any A-stable LMM is at most second order accurate, and in fact the Trapezoidal method is the A-stable method with smallest truncation error[Hen62]. Higher order A-stable implicit Runge-Kutta methods do exist, but are more difficult to apply[But87].

#### 10.3.2 L-stability

Notice a major difference between the stability regions for Trapezoidal and Backward Euler: the Trapezoidal method is stable only in the left half plane, whereas Backward Euler is also stable over much of the right half plane. Recall that the root of the stability polynomial for Backward Euler is  $\zeta_1 = (1-z)^{-1}$ . This approaches 0 as  $|z| \rightarrow \infty$  in any direction in the complex plane. The point at infinity (on the Riemann sphere, in the sense of complex analysis) is in the *interior* of the stability region for the Backward Euler method.

For the Trapezoidal method, on the other hand,

$$\zeta_1 = \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z},$$

which approaches 1 in modulus as  $|z| \rightarrow \infty$ . It is less than 1 in modulus in the left half plane but greater than 1 in the right half plane. The point at infinity is on the boundary of the stability region for this method.

For a general linear multistep method, the roots of the stability polynomial

$$\pi(\zeta; z) = \rho(\zeta) - z\sigma(\zeta)$$

are the same as the roots of

$$\frac{1}{z}\pi(\zeta; z) = \frac{1}{z}\rho(\zeta) - \sigma(\zeta)$$

and should behave like the roots of  $\sigma(\zeta)$  as  $|z| \rightarrow \infty$ . (This is actually true only for implicit methods. What happens for an explicit method?) For Backward Euler the root of  $\sigma$  is at 0, while for Trapezoidal the root is at 1.

A method is called *L-stable* if the roots of  $\sigma$  are strictly inside the unit circle,  $|\zeta_j| < 1$ . For an L-stable method the point at infinity is in the interior of the stability region.

If we are solving a stiff equation with initial data such that the solution is smooth from the beginning (no rapid transients), or if we plan to compute rapid transients accurately by taking suitably small time steps in these regions, then it may be fine to use a method such as the Trapezoidal Rule that is not L-stable.

However, in some situations there are rapid transients in the solution that we are not interested in resolving accurately with very small time steps. For these transients we want more than just stability — we want them to be effectively damped in a single time step since we are planning to use at time step that is much larger than the true decay time of the transient. For this purpose an L-stable method is crucial. This is best illustrated with an example.

**Example 10.3.** Again consider the problem  $u'(t) = \lambda(u(t) - \cos(t)) - \sin(t)$  with  $\lambda = -10^6$  and let's see how Trapezoidal and Backward Euler behave in two different situations.

**Case 1:** Take data  $u(0) = 1$ , so that  $u(t) = \cos(t)$  and there is no initial transient. Then both Trapezoidal and Backward Euler behave reasonably and the trapezoidal method gives smaller errors since it is second order accurate. The following table shows the errors at  $T = 3$  with various values of  $k$ .

k	Backw. Euler	Trapezoidal
4.0000e-01	4.7770e-02	4.7770e-02
2.0000e-01	9.7731e-08	4.7229e-10
1.0000e-01	4.9223e-08	1.1772e-10
5.0000e-02	2.4686e-08	2.9409e-11
2.5000e-02	1.2360e-08	7.3508e-12
1.2500e-02	6.1837e-09	1.8373e-12
6.2500e-03	3.0928e-09	4.6030e-13
3.1250e-03	1.5466e-09	1.1757e-13

**Case 2:** Now take data  $u(0) = 1.5$  so there is an initial rapid transient towards  $u = \cos(t)$  on a time scale of about  $10^{-6}$ . Both methods are still absolutely stable, but the results in the next table show that Backward Euler works much better in this case.

k	Backw. Euler	Trapezoidal
4.0000e-01	4.7770e-02	4.5219e-01
2.0000e-01	9.7731e-08	4.9985e-01
1.0000e-01	4.9223e-08	4.9940e-01
5.0000e-02	2.4686e-08	4.9761e-01
2.5000e-02	1.2360e-08	4.9049e-01
1.2500e-02	6.1837e-09	4.6304e-01
6.2500e-03	3.0928e-09	3.6775e-01
3.1250e-03	1.5466e-09	1.4632e-01

To understand what is happening, see Figure 10.4, which shows the true and computed solutions with each method if we use  $k = 0.1$ . The trapezoidal method is stable and the results stay bounded, but since  $k\lambda = -10^5$  we have  $\frac{1 - \frac{1}{2}k\lambda}{1 + \frac{1}{2}k\lambda} = -0.99996 \approx -1$  and the initial deviation from the smooth curve  $\cos(t)$  is essentially negated in each time step.



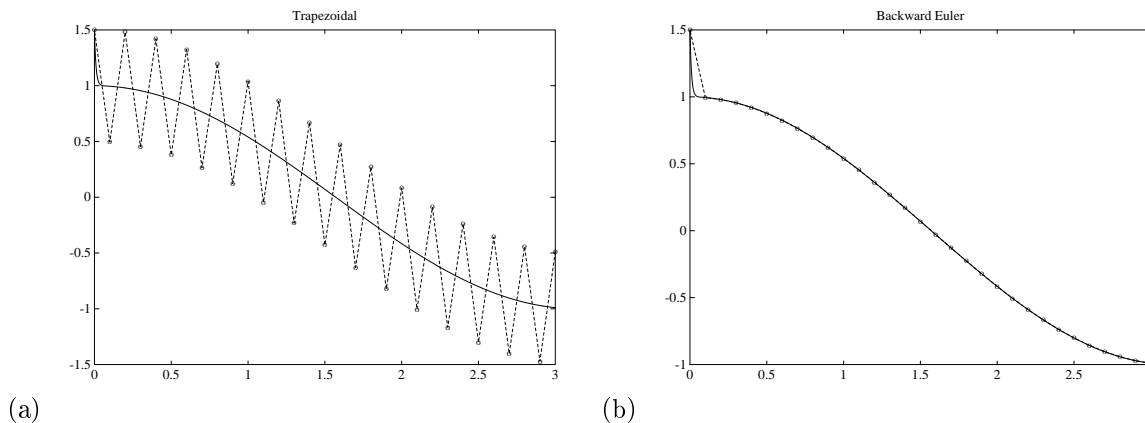


Figure 10.4: Comparison of (a) Trapezoidal method and (b) Backward Euler on a stiff problem with an initial transient (Case 2 of Example 10.3).

Backward Euler, on the other hand, damps the deviation very effectively in the first time step, since  $\frac{1}{1+k\lambda} \approx -10^{-6}$ . This is the proper behavior since the true rapid transient decays in a time period much shorter than a single time step.

## 10.4 BDF Methods

One class of very effective methods for stiff problems are the BDF methods (Backward Differentiation Formulas). These were first used by Curtis and Hirschfelder[CH52] but are often referred to as Gear's methods since he wrote one of the first software packages for stiff problems based them.

These methods result from taking  $\sigma(\zeta) = \beta_r \zeta^r$ , which has all its roots at the origin, resulting in an L-stable method. The method thus has the form

$$\alpha_0 U^n + \alpha_1 U^{n+1} + \cdots + \alpha_r U^{n+r} = k\beta_r f(U^{n+r}), \quad (10.3)$$

with  $\beta_0 = \beta_1 = \cdots = \beta_{r-1} = 0$ . Since  $f(u) = u'$ , this form of method can be derived by approximating  $u'(t_{n+r})$  by a backward difference approximation based on  $u(t_{n+r})$  and  $r$  additional points going backwards in time.

It is possible to derive an  $r$ -step method that is  $r$ 'th order accurate. The 1-step BDF method is simply the Backward Euler method,  $U^{n+1} = U^n + kf(U^{n+1})$ , which is first order accurate. The next few are:

$$\begin{aligned} r = 2 : \quad & 3U^{n+2} - 4U^{n+1} + U^n = 2kf(U^{n+2}) \\ r = 3 : \quad & 11U^{n+3} - 18U^{n+2} + 9U^{n+1} - 2U^n = 6kf(U^{n+3}) \end{aligned}$$

others can be found in Gear[Gea71] or Lambert[Lam73].

These methods have the proper behavior on eigenvalues for which  $\text{Re}(\lambda)$  is very negative, but of course we also have other eigenvalues for which  $z = k\lambda$  is closer to the origin, corresponding to the active time scales in the problem. So deciding its suitability for a particular problem requires looking at the full stability region. These are shown in Figure 10.5 for several values of  $r$ .

In particular, we need to make sure that the method is zero-stable. Otherwise it would not be convergent. This is not guaranteed from our derivation of the methods, since zero-stability depends only on the polynomial  $\rho(\zeta)$ , whose coefficients  $\alpha_j$  are determined by considering the local truncation error and not stability considerations. It turns out that the BDF methods are zero-stable only for  $r \leq 6$ . Higher order BDF methods cannot be used in practice.

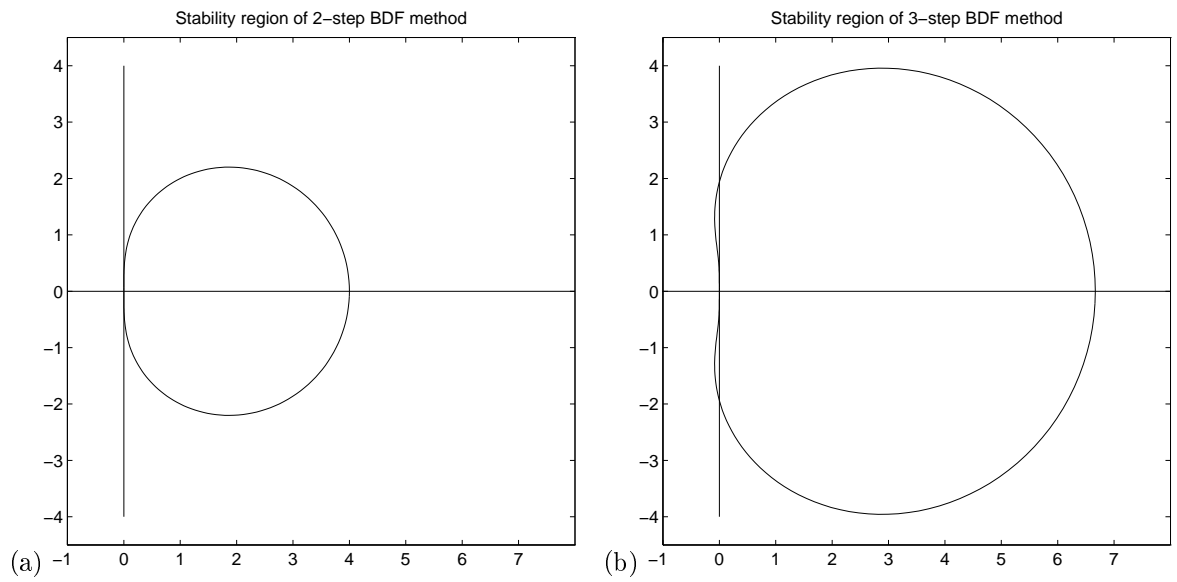


Figure 10.5: Stability regions for some BDF methods. The stability region is the exterior of the curves.



# Chapter 11

## Some basic PDEs

In this chapter we will briefly develop some of the basic PDEs that will be used to illustrate the development of numerical methods. In solving a partial differential equation, we are looking for a function of more than one variable that satisfies some relations between different partial derivatives.

### 11.1 Classification of differential equations

First we review the classification of differential equations into elliptic, parabolic, and hyperbolic equations. Not all PDE's fall into one of these classes, by any means, but many important equations that arise in practice do. These classes of equations model different sorts of phenomena, display different behavior, and require different numerical techniques for their solution. Standard texts on partial differential equations such as Kevorkian[Kev90] give further discussion.

#### 11.1.1 Second-order equations

In most elementary texts the classification is given for a linear second-order differential equation in two independent variables of the form

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

The classification depends on the sign of the discriminant,

$$b^2 - 4ac \begin{cases} < 0 & \implies \text{elliptic} \\ = 0 & \implies \text{parabolic} \\ > 0 & \implies \text{hyperbolic} \end{cases}$$

and the names arise by analogy with conic sections. The canonical example are the Poisson problem  $u_{xx} + u_{yy} = g$  for an elliptic problem, the heat equation  $u_t = \beta u_{xx}$  (with  $\beta > 0$ ) for a parabolic problem, and the wave equation  $u_{tt} = c^2 u_{xx}$  for a hyperbolic problem. In the parabolic and hyperbolic case  $t$  is used instead of  $y$  since these are typically time-dependent problems. These can all be extended to more space dimensions. These equations describe different types of phenomena and require different techniques for their solution (both analytically and numerically), and so it is convenient to have names for classes of equations exhibiting the same general features. There are other equations that have some of the same features and the classification scheme can be extended beyond the second-order linear form given above. Some hint of this is given in the next few sections.

#### 11.1.2 Elliptic equations

The classic example of an elliptic equation is the **Poisson problem**

$$\nabla^2 u = f, \tag{11.1}$$

where  $\nabla^2$  is the Laplacian operator and  $f$  is a given function of  $\vec{x}$  in some spatial domain  $\Omega$ . We seek a function  $u(\vec{x})$  in  $\Omega$  satisfying (11.1) together with some **boundary conditions** all along the boundary of  $\Omega$ . Elliptic equations typically model steady-state or equilibrium phenomena, and so there is no temporal dependence. Elliptic equations may also arise in solving time-dependent problems if we are modeling some phenomena that are always in local equilibrium and equilibrate on time scales that are much faster than the time scale being modeled. For example, in “incompressible” flow the fast acoustic waves are not modeled and instead the pressure is computed by solving a Poisson problem at each time step which models the global effect of these waves.

Elliptic equations give boundary value problems (BVP’s) where the solution at all points must be simultaneously determined based on the boundary conditions all around the domain. This typically leads to a very large sparse system of linear equations to be solved for the values of  $U$  at each grid point. If an elliptic equation must be solved in every time step of a time-dependent calculation, as in the examples above, then it is crucial that these systems be solved as efficiently as possible.

More generally, a linear elliptic equation has the form

$$Lu = f, \tag{11.2}$$

where  $L$  is some **elliptic operator**. This notion will not be discussed further here, but the idea is that mathematical conditions are required on the differential operator  $L$  which insure that the boundary value problem has a unique solution.

### 11.1.3 Parabolic equations

If  $L$  is an elliptic operator then the time-dependent equation

$$u_t = Lu - f \tag{11.3}$$

is called **parabolic**. If  $L = \nabla^2$  is the Laplacian, then (11.3) is known as the **heat equation** or **diffusion equation** and models the diffusion of heat in a material, for example.

Now  $u(\vec{x}, t)$  varies with time and we require **initial data**  $u(\vec{x}, 0)$  for every  $\vec{x} \in \Omega$  as well as boundary conditions around the boundary at each time  $t > 0$ . If the boundary conditions are independent of time, then we might expect the heat distribution to reach a steady state in which  $u$  is independent of  $t$ . We could then solve for the steady state directly by setting  $u_t = 0$  in (11.3), which results in the elliptic equation (11.2).

Marching to steady state by solving the time-dependent equation (11.3) numerically would be one approach to solving the elliptic equation (11.2), but this is typically not the fastest method if all we require is the steady state.

### 11.1.4 Hyperbolic equations

Rather than discretizing second order hyperbolic equations such as the wave equation  $u_{tt} = c^2 u_{xx}$ , we will consider a related form of hyperbolic equations known as *first-order hyperbolic systems*. The linear problem in one space dimension has the form

$$u_t + Au_x = 0 \tag{11.4}$$

where  $u(x, t) \in \mathbb{R}^m$  and  $A$  is an  $m \times m$  matrix. The problem is called **hyperbolic** if  $A$  has *real* eigenvalues and is *diagonalizable*, i.e., has a complete set of linearly independent eigenvectors. These conditions allow us to view the solution in terms of propagating waves, and indeed hyperbolic systems typically arise from physical processes that give wave motion or advective transport.

The simplest example of a hyperbolic equation is the constant-coefficient **advection equation**

$$u_t + au_x = 0, \tag{11.5}$$

where  $u$  is the advection velocity. The solution is simply  $u(x, t) = u(x - at, 0)$ , so any  $u$  profile simply advects with the flow at velocity  $a$ .

As a simple example of a linear system, the equations of linearized acoustics arising from elasticity or gas dynamics can be written as a first-order system of two equations in one space dimension as

$$\begin{bmatrix} p \\ u \end{bmatrix}_t + \begin{bmatrix} 0 & \kappa_0 \\ 1/\rho_0 & 0 \end{bmatrix} \begin{bmatrix} p \\ u \end{bmatrix}_x = 0 \quad (11.6)$$

in terms of pressure and velocity perturbations, where  $\rho_0$  is the background density and  $\kappa_0$  is the “bulk modulus” of the material. Note that if we differentiate the first equation with respect to  $t$ , the second with respect to  $x$ , and then eliminate  $u_{xt} = u_{tx}$  we obtain the second-order wave equation for the pressure:

$$p_{tt} = c^2 p_{xx},$$

where

$$c = \sqrt{\kappa_0/\rho}$$

is the speed of sound in the material.

## 11.2 Derivation of PDEs from conservation principles

Many physically relevant partial differential equations can be derived based on the principle of conservation. We can view  $u(x, t)$  as a *concentration* or *density* function for some substance or chemical that is in dilute suspension in a liquid, for example. Basic equations of the same form arise in many other applications, however. The material presented here is meant to be a brief review, and much more complete discussions are available in many sources. See, for example, [Kev90],[Whi74].

A reasonable model to consider in one space dimension is the concentration or density of a contaminant in a stream or pipe, where the variable  $x$  represents distance along the pipe. The concentration is assumed to be constant across any cross-section, so that its value varies only with  $x$ . The density function  $u(x, t)$  is defined in such a way that integrating the function  $u(x, t)$  between any two points  $x_1$  and  $x_2$  gives the total mass of the substance in this section of the pipe at time  $t$ :

$$\text{Total mass between } x_1 \text{ and } x_2 \text{ at time } t = \int_{x_1}^{x_2} u(x, t) dx.$$

The density function is measured in units such as grams/meter. (Note that this  $u$  really represents the integral over the cross section of the pipe of a density function that is more properly measured in grams/meter<sup>3</sup>.)

The basic form of differential equation that models many physical processes can be derived in the following way. Consider a section  $x_1 < x < x_2$  and the manner in which  $\int_{x_1}^{x_2} u(x, t) dx$  changes with time. This integral represents the total mass of the substance in this section, so if we are studying a substance which is neither created nor destroyed within this section, then the total mass within this section can change only due to the *flux* or flow of particles through the endpoints of the section at  $x_1$  and  $x_2$ . This flux is given by some function  $f$  which, in the simplest case, depends only on the value of  $u$  at the corresponding point.

## 11.3 Advection

If the substance is simply carried along (advected) in a flow at some constant velocity  $a$ , then the flux function is

$$f(u) = au. \quad (11.7)$$

The local density  $u(x, t)$  (in grams/meter, say) multiplied by the velocity (in meters/sec, say) gives the flux of material past the point  $x$  (in grams/sec).

Since the total mass in  $[x_1, x_2]$  changes only due to the flux at the endpoints, we have

$$\frac{d}{dt} \int_{x_1}^{x_2} u(x, t) dx = f(u(x_1, t)) - f(u(x_2, t)).$$

The minus sign on the last term comes from the fact that  $f$  is, by definition, the flux to the right.

If we assume that  $u$  and  $f$  are smooth functions, then this equation can be rewritten as

$$\frac{d}{dt} \int_{x_1}^{x_2} u(x, t) dx = \int_{x_1}^{x_2} \frac{\partial}{\partial x} f(u(x, t)) dx,$$

or, with some further modification, as

$$\int_{x_1}^{x_2} \left[ \frac{\partial}{\partial t} u(x, t) + \frac{\partial}{\partial x} f(u(x, t)) \right] dx = 0.$$

Since this integral must be zero for all values of  $x_1$  and  $x_2$ , it follows that the integrand must be identically zero. This gives, finally, the differential equation

$$\frac{\partial}{\partial t} u(x, t) + \frac{\partial}{\partial x} f(u(x, t)) = 0. \quad (11.8)$$

This form of equation is called a *conservation law*. (For further discussion see [Lax72], [LeV90], [Whi74].)

For the case considered in Section 11.3,  $f(u) = au$  with  $a$  constant and this equation becomes

$$u_t + au_x = 0. \quad (11.9)$$

This is called the *advection equation*.

This equation requires initial conditions and possibly also boundary conditions in order to determine a unique solution. The simplest case is the *Cauchy problem* on  $-\infty < x < \infty$  (with no boundary), also called the pure initial value problem. Then we only need to specify initial data

$$u(x, 0) = \eta(x). \quad (11.10)$$

Physically, we would expect the initial profile of  $\eta$  to simply be carried along with the flow at speed  $a$ , so we should find

$$u(x, t) = \eta(x - at). \quad (11.11)$$

It is easy to verify that this function satisfies the advection equation (11.9) and is the solution of the PDE.

The curves

$$x = x_0 + at$$

through each point  $x_0$  at time 0 are called the *characteristics* of the equation. If we set

$$U(t) = u(x_0 + at, t)$$

then

$$\begin{aligned} U'(t) &= au_x(x_0 + at, t) + u_t(x_0 + at, t) \\ &= 0 \end{aligned}$$

using (11.9). Along these curves the PDE reduces to a simple ODE  $U' = 0$  and the solution must be constant along each such curve, as is also seen from the solution (11.11).

## 11.4 Diffusion

Now suppose that the fluid in the pipe is not flowing, and has zero velocity. Then according to the above equation,  $u_t = 0$  and the initial profile  $\eta(x)$  does not change with time. However, if  $\eta$  is not constant in space then in fact it will tend to slowly change due to molecular diffusion. The velocity  $a$  should really be thought of as a *mean velocity*, the average velocity that the roughly  $10^{23}$  molecules in a given drop of water have. But individual molecules are bouncing around in different directions and so molecules of the substance we are tracking will tend to get spread around in the water, as a drop of ink spreads. There will tend to be a net motion from regions where the density is large to regions where it is smaller, and in fact it can be shown that the flux (in 1D) is proportional to  $-u_x$ . The flux at a point  $x$  now depends on the value of  $u_x$  at this point, rather than on the value of  $u$ , so we write

$$f(u_x) = -\beta u_x,$$

where  $\beta$  is the *diffusion coefficient*. Using this flux in (11.8) gives

$$u_t = \beta u_{xx} \tag{11.12}$$

which is known as the *diffusion equation*. It is also called the *heat equation* since heat diffuses in much the same way. In this case we can think of the one-dimensional equation as modeling the conduction of heat in a rod. The heat conduction coefficient  $\beta$  depends on the material and how well it conducts heat. The variable  $u$  is then the temperature.

In some problems the diffusion coefficient may vary with  $x$ , for example in a rod made of a composite of different materials. Then  $f = -\beta(x)u_x$  and the equation becomes

$$u_t = (\beta(x)u_x)_x.$$

Returning to the example of fluid flow, more generally there would be both advection and diffusion occurring simultaneously. Then the flux is  $f(u, u_x) = au - \beta u_x$ , giving the *advection-diffusion* equation

$$u_t + au_x = \beta u_{xx}. \tag{11.13}$$

The diffusion and advection-diffusion equations are examples of the general class of PDEs called *parabolic*.

## 11.5 Source terms

In some situations  $\int_{x_1}^{x_2} u(x, t) dx$  changes due to effects other than flux through the endpoints of the section, if there is some source or sink of the substance within the section. Denote the density function for such a source by  $\psi(x, t)$ . (Negative values of  $\psi$  correspond to a sink rather than a source.) Then the equation becomes

$$\frac{d}{dt} \int_{x_1}^{x_2} u(x, t) dx = \int_{x_1}^{x_2} \frac{\partial}{\partial x} f(u(x, t)) dx + \int_{x_1}^{x_2} \psi(x, t) dx.$$

This leads to the PDE

$$u_t(x, t) + f(u(x, t))_x = \psi(x, t). \tag{11.14}$$

For example, if we have heat conduction in a rod together with an external source of heat energy distributed along the rod with density  $\psi$ , then we have

$$u_t = \beta u_{xx} + \psi.$$

In some cases the strength of the source may depend on the value of  $u$ . For example, if the rod is immersed in a liquid that is held at constant temperature  $u_0$ , then the flux of heat into the rod at the point  $(x, t)$  is proportional to  $u_0 - u(x, t)$  and the equation becomes

$$u_t(x, t) = \beta u_{xx}(x, t) + \alpha(u_0 - u(x, t)).$$



### 11.5.1 Reaction-diffusion equations

One common form of source terms arises from chemical kinetics. If the components of  $u \in \mathbb{R}^m$  represent concentrations of  $m$  different species reacting with one another, then the kinetics equations have the form  $u_t = \psi(u)$ , as described in Section 8.5.1. This assumes the different species are well-mixed at all times and so the concentrations vary only with time. If there are spacial variations in concentrations, then these equations may be combined with diffusion of each species. This would lead to a system of reaction-diffusion equations of the form

$$u_t = \beta u_{xx} + \psi(u). \quad (11.15)$$

The diffusion coefficient could be different for each species, in which case  $\beta$  would be a diagonal matrix instead of a scalar. Advection terms might also be present if the reactions are taking place in a flowing fluid.

## Chapter 12

# Fourier Analysis of Linear PDEs

For *linear* PDEs, Fourier analysis is often used to obtain solutions or perform theoretical analysis. This is because the functions  $e^{i\xi x} = \cos(\xi x) + i \sin(\xi x)$  are eigenfunctions of the differentiation operator  $\partial_x = \frac{\partial}{\partial x}$ . Differentiating this function gives a scalar multiple of the function, and hence differential equations are simplified. Fourier analysis is equally important in the study of finite difference methods for *linear* PDEs for the same reason: these same functions are eigenfunctions of finite difference operators. This will be exploited in Section 13.6 where von Neumann stability analysis of finite difference methods is discussed. An understanding of Fourier analysis of PDEs is also required in Chapter 16 where finite difference methods are discussed by derived and analyzed by studying “modified equations”.

### 12.1 Fourier transforms

Recall that a function  $v(x)$  is in the space  $L^2$  if it has a finite 2-norm, defined by

$$\|v\|_2 = \int_{-\infty}^{\infty} |v(x)|^2 dx.$$

If  $v \in L^2$ , then we can define its Fourier transform  $\hat{v}(\xi)$  by

$$\hat{v}(\xi) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} v(x) e^{-i\xi x} dx. \quad (12.1)$$

The function  $\hat{v}(\xi)$  is also in  $L^2$  and in fact it has exactly the same 2-norm as  $v$ ,

$$\|\hat{v}\|_2 = \|v\|_2. \quad (12.2)$$

This is known as *Parseval's relation*.

We can express the original function  $v(x)$  as a linear combination of the set of functions  $e^{i\xi x}$  for different values of  $\xi$ , which together form a basis for the infinite dimensional function space  $L^2$ . The Fourier transform  $\hat{v}(\xi)$  gives the coefficients in the expression

$$v(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{v}(\xi) e^{i\xi x} d\xi, \quad (12.3)$$

which is known as the *inverse Fourier transform*. This is analogous to writing a vector as a linear combination of basis vectors.

### 12.2 Solution of differential equations

The Fourier transform plays a fundamental role in the study of linear PDEs because of the fact that the functions  $w(x) = e^{i\xi x}$  (for each fixed  $\xi$ ) are *eigenfunctions* of the differential operator  $\partial_x$ . In general

applying  $\partial_x$  to a function gives a new function that is not simply a scalar multiple of the original function. For example,  $\partial_x(x^3) = 3x^2$  and these functions are shaped quite differently. However, applying  $\partial_x$  to  $w(x) = e^{i\xi x}$  gives  $i\xi e^{i\xi x}$ , which is simply the constant  $i\xi$  times the original function  $w(x)$ . We say that  $e^{i\xi x}$  is an *eigenfunction* of the operator  $\partial_x$  with *eigenvalue*  $i\xi$ .

**Example 12.1.** To see the importance of this fact, let's solve the advection equation  $u_t + au_x = 0$  using Fourier transforms. We will transform in  $x$  only, and denote the transform of  $u(x, t)$  (a function of  $x$  at each fixed  $t$ ) by  $\hat{u}(\xi, t)$ :

$$\hat{u}(\xi, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} u(x, t) e^{-i\xi x} dx. \quad (12.4)$$

Then

$$u(x, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{u}(\xi, t) e^{i\xi x} d\xi \quad (12.5)$$

and differentiating this with respect to  $t$  and  $x$  gives

$$\begin{aligned} u_t(x, t) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{u}_t(\xi, t) e^{i\xi x} dx \\ u_x(x, t) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{u}(\xi, t) i\xi e^{i\xi x} dx. \end{aligned}$$

From this we see that the Fourier transform of  $u_t$  is  $\hat{u}_t$  and the Fourier transform of  $u_x$  is  $i\xi\hat{u}$ . Fourier transforming the advection equation by computing

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (u_t + au_x) e^{-i\xi x} dx = 0$$

thus gives

$$\hat{u}_t(\xi, t) + ai\xi\hat{u}(\xi, t) = 0$$

or

$$\hat{u}_t = -i\xi a\hat{u}.$$

This is a time-dependent ODE for the evolution of  $\hat{u}(\xi, t)$  in time. There are two important points to notice:

- Since differentiation with respect to  $x$  has become multiplication by  $i\xi$  after Fourier transforming, the original PDE involving derivatives with respect to  $x$  and  $t$  has become an ODE in  $t$  alone.
- The ODEs for different values of  $\xi$  are decoupled from one another. We have to solve an infinite number of ODEs, one for each value of  $\xi$ , but they are decoupled scalar equations rather than a coupled system.

In fact it is easy to solve these ODEs. We need initial data  $\hat{u}(\xi, 0)$  at time  $t = 0$  for each value of  $\xi$ , but this comes from Fourier transforming the initial data  $u(x, 0) = \eta(x)$ ,

$$\hat{u}(\xi, 0) = \hat{\eta}(\xi) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \eta(x) e^{-i\xi x} dx.$$

Solving the ODEs then gives

$$\hat{u}(\xi, t) = e^{-i\xi at} \hat{\eta}(\xi). \quad (12.6)$$

We can now Fourier transform back using (12.5) to get the desired solution  $u(x, t)$ :

$$\begin{aligned} u(x, t) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-i\xi at} \hat{\eta}(\xi) e^{i\xi x} d\xi \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{\eta}(\xi) e^{i\xi(x-at)} d\xi \\ &= \eta(x - at). \end{aligned}$$

This last equality comes from noting that we are simply evaluating the inverse Fourier transform of  $\hat{\eta}$  at the point  $x - at$ . We see that we have recovered the standard solution (11.11) of the advection equation in this manner.

### 12.3 The heat equation

Now consider the heat equation,

$$u_t = \beta u_{xx}. \quad (12.7)$$

Since the Fourier transform of  $u_{xx}(x, t)$  is  $(i\xi)^2 \hat{u}(\xi, t) = -\xi^2 \hat{u}(\xi, t)$ , Fourier transforming the equation (12.7) gives the ODE

$$\hat{u}_t(\xi, t) = -\beta \xi^2 \hat{u}(\xi, t). \quad (12.8)$$

Again we have initial data  $\hat{u}(\xi, 0) = \hat{\eta}(\xi)$  from the given initial data on  $u$ . Now solving the ODE gives

$$\hat{u}(\xi, t) = e^{-\beta \xi^2 t} \hat{\eta}(\xi).$$

Note that this has a very different character than (12.6), the Fourier transform obtained from the advection equation. For the advection equation,  $\hat{u}(\xi, t) = e^{ia\xi t} \hat{\eta}(\xi)$  and  $|\hat{u}(\xi, t)| = |\hat{\eta}(\xi)|$  for all  $t$ . Each Fourier component maintains its original amplitude and is modified only in phase, leading to a traveling wave behavior in the solution.

For the heat equation,  $|\hat{u}(\xi, t)|$  decays in time exponentially fast. The decay rate depends on  $\beta$ , the diffusion coefficient, and also on  $\xi$ , the wavenumber. High frequency waves (with  $\xi^2$  large) decay much faster than low frequencies. This results in a *smoothing* of the solution as time evolves. (See Figure 13.3.)

The fact that the solution contains components which decay at very different rates leads us to expect numerical difficulties with stiffness, similar to those discussed for ODE's in Chapter 10. In Section 13.4 we will see that this is indeed the case, and that implicit methods must generally be used in order to efficiently solve the heat equation.

### 12.4 Dispersive waves

Now consider the equation

$$u_t = u_{xxx}. \quad (12.9)$$

Fourier transforming now leads to the ODE

$$\hat{u}_t(\xi, t) = -i\xi^3 \hat{u}(\xi, t),$$

so

$$\hat{u}(\xi, t) = e^{-i\xi^3 t} \hat{\eta}(\xi).$$

This has a character similar to advection problems in that  $|\hat{u}(\xi, t)| = |\hat{\eta}(\xi)|$  for all time and each Fourier component maintains its original amplitude. However, when we recombine with the inverse Fourier transform we obtain

$$u(x, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{\eta}(\xi) e^{i\xi(x - \xi^2 t)} d\xi, \quad (12.10)$$

which shows that the Fourier component with wave number  $\xi$  is propagating with velocity  $\xi^2$ . In the advection equation all Fourier components propagate with the same speed  $a$ , and hence the shape of the initial data is preserved with time. The solution is the initial data shifted over a distance  $at$ .

With the equation (12.9), the shape of the initial data will in general not be preserved, unless the data is simply a single Fourier mode. This behavior is called *dispersive* since the Fourier components disperse relative to one another. Typically smooth data leads to oscillatory solutions since the cancellation of high wave number modes that smoothness depends on will be lost as these modes shift relative to one another. See, for example, Whitham[Whi74] for an extensive discussion of dispersive waves.

## 12.5 Even vs. odd order derivatives

Note that odd order derivatives  $\partial_x, \partial_x^3, \dots$  (as in the advection equation or the dispersive equation (12.9)) have pure imaginary eigenvalues  $i\xi, -i\xi^3, \dots$ , which results in Fourier components that propagate with their magnitude preserved. Even order derivatives, such as the  $\partial_x^2$  in the heat equation, have real eigenvalues ( $-\xi^2$  for the heat equation) which results in exponential decay of the eigencomponents. Another such equation is

$$u_t = -u_{xxxx},$$

in which case  $\hat{u}(\xi, t) = e^{-\xi^4 t} \hat{\eta}(\xi)$ . Solutions to this equation behave much like solutions to the heat equation, but with even more rapid damping of oscillatory data.

## Chapter 13

# Diffusion Equations

We now begin to study numerical methods for time-dependent partial differential equations, where variations in space are related to variations in time. We begin with the heat equation (or diffusion equation) introduced in Chapter 11,

$$u_t = u_{xx}. \tag{13.1}$$

This is the classical example of a *parabolic* equation, and many of the general properties seen here carry over to the design of numerical methods for other parabolic equations.

Along with this equation we need initial conditions,

$$u(x, 0) = \eta(x) \tag{13.2}$$

and also boundary conditions if we are working on a bounded domain, e.g., the Dirichlet conditions

$$\begin{aligned} u(0, t) &= g_0(t) \text{ for } t > 0 \\ u(1, t) &= g_1(t) \text{ for } t > 0 \end{aligned} \tag{13.3}$$

if  $0 \leq x \leq 1$ .

We have already studied the steady state version of this equation and spatial discretizations of  $u_{xx}$  (Chapter 2). We have also studied discretizations of the time derivatives and some of the stability issues that arise with these discretizations in Chapters 6 through 10. Next we will put these two types of discretizations together.

In practice we generally apply a set of finite difference equations on a discrete grid with grid points  $(x_i, t_n)$  where

$$x_i = ih, \quad t_n = nk.$$

Here  $h = \Delta x$  is the mesh spacing on the  $x$ -axis and  $k = \Delta t$  is the time step. Let  $U_i^n \approx u(x_i, t_n)$  represent the numerical approximation at grid point  $(x_i, t_n)$ .

Since the heat equation is an evolution equation that can be solved forward in time, we set up our difference equations in a form where we can march forward in time, determining the values  $U_i^{n+1}$  for all  $i$  from the values  $U_i^n$  at the previous time level, or perhaps using also values at earlier time levels with a multistep formula. Compare this to the case of elliptic equations (boundary value problems) considered in Chapter 3, where we had to solve for all the grid values simultaneously.

As an example, one natural discretization of (13.1) would be

$$\frac{U_i^{n+1} - U_i^n}{k} = \frac{1}{h^2}(U_{i-1}^n - 2U_i^n + U_{i+1}^n). \tag{13.4}$$

This uses our standard centered difference in space and a forward difference in time. This is an *explicit* method since we can compute each  $U_i^{n+1}$  explicitly in terms of the previous data:

$$U_i^{n+1} = U_i^n + \frac{k}{h^2}(U_{i-1}^n - 2U_i^n + U_{i+1}^n). \tag{13.5}$$

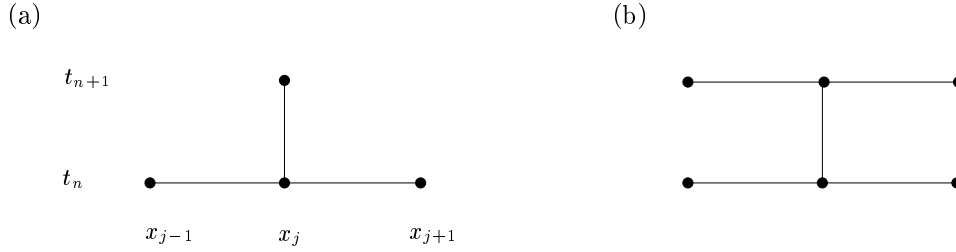


Figure 13.1: Stencils for the methods (13.5) and (13.7).

Figure 13.1(a) shows the *stencil* of this method. This is a one-step method in time, which is also called a two-level method in the context of PDE's since it involves the solution at two different time levels.

Another one-step method, which is much more useful in practice as we will see below, is the *Crank-Nicolson* method,

$$\begin{aligned} \frac{U_i^{n+1} - U_i^n}{k} &= \frac{1}{2}(D^2 U_i^n + D^2 U_i^{n+1}) \\ &= \frac{1}{2h^2}(U_{i-1}^n - 2U_i^n + U_{i+1}^n + U_{i-1}^{n+1} - 2U_i^{n+1} + U_{i+1}^{n+1}), \end{aligned} \quad (13.6)$$

which can be rewritten as

$$U_i^{n+1} = U_i^n + \frac{k}{2h^2}(U_{i-1}^n - 2U_i^n + U_{i+1}^n + U_{i-1}^{n+1} - 2U_i^{n+1} + U_{i+1}^{n+1}) \quad (13.7)$$

or

$$-rU_{i-1}^{n+1} + (1+2r)U_i^{n+1} - rU_{i+1}^{n+1} = rU_{i-1}^n + (1-2r)U_i^n + rU_{i+1}^n \quad (13.8)$$

where  $r = k/2h^2$ . This is an *implicit* method and gives a tridiagonal system of equations to solve for all the values  $U_i^{n+1}$  simultaneously. In matrix form this is

$$\begin{aligned} &\begin{bmatrix} (1+2r) & -r & & & & & \\ -r & (1+2r) & -r & & & & \\ & -r & (1+2r) & -r & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & -r & (1+2r) & -r & \\ & & & & -r & (1+2r) & \end{bmatrix} \begin{bmatrix} U_1^{n+1} \\ U_2^{n+1} \\ U_3^{n+1} \\ \vdots \\ U_{m-1}^{n+1} \\ U_m^{n+1} \end{bmatrix} \\ &= \begin{bmatrix} r(g_0(t_n) + g_0(t_{n+1})) + (1-2r)U_1^n + rU_2^n \\ rU_1^n + (1-2r)U_2^n + rU_3^n \\ rU_2^n + (1-2r)U_3^n + rU_4^n \\ \vdots \\ rU_{m-2}^n + (1-2r)U_{m-1}^n + rU_m^n \\ rU_{m-1}^n + (1-2r)U_m^n + r(g_1(t_n) + g_1(t_{n+1})) \end{bmatrix}. \end{aligned} \quad (13.9)$$

Note how the boundary conditions  $u(0, t) = g_0(t)$  and  $u(1, t) = g_1(t)$  come into these equations.

Since a tridiagonal system of  $m$  equations can be solved with  $O(m)$  work, this method is essentially as efficient per time step as an explicit method. We will see in Section 13.4 that the heat equation is “stiff”, and hence this implicit method, which allows much larger time steps to be taken than an explicit method, is a very efficient method for the heat equation.

### 13.1 Local truncation errors and order of accuracy

We can define the local truncation error as usual — we insert the exact solution  $u(x, t)$  of the PDE into the finite difference equation and determine by how much it fails to satisfy the discrete equation. It is important to use the form of the difference equation that directly models the PDE in order to get meaningful results in terms of the powers of  $k$  and  $h$  which appear.

**Example 13.1.** The local truncation error of the method (13.5) is based on the form (13.4):

$$\tau(x, t) = \frac{u(x, t+k) - u(x, t)}{k} - \frac{1}{h^2}(u(x-h, t) - 2u(x, t) + u(x+h, t)).$$

Again we should be careful to use the form that directly models the differential equation in order to get powers of  $k$  and  $h$  that agree with what we hope to see in the global error. Although we don't know  $u(x, t)$  in general, if we assume it is smooth and use Taylor series expansions about  $u(x, t)$ , we find that

$$\tau(x, t) = \left( u_t + \frac{1}{2}ku_{tt} + \frac{1}{6}k^2u_{ttt} + \cdots \right) - \left( u_{xx} + \frac{1}{12}h^2u_{xxxx} + \cdots \right).$$

Since  $u_t = u_{xx}$ , the  $O(1)$  terms drop out. By differentiating  $u_t = u_{xx}$  we find that  $u_{tt} = u_{txx} = u_{xxxx}$  and so

$$\tau(x, t) = \left( \frac{1}{2}k - \frac{1}{12}h^2 \right) u_{xxxx} + O(k^2 + h^4).$$

This method is said to be *second order accurate in space* and *first order accurate in time* since the truncation error is  $O(h^2 + k)$ .

The Crank-Nicolson method is centered in both space and time, and an analysis of its local truncation error shows that it is second order accurate in both space and time,

$$\tau(x, t) = O(k^2 + h^2).$$

**Exercise 13.1** Compute the dominant term in the truncation error of Crank-Nicolson.

A method is said to be *consistent* if  $\tau(x, t) \rightarrow 0$  as  $k, h \rightarrow 0$ . Just as in the other cases we have studied (boundary value problems and initial value problems for ODE's), we expect that consistency, plus some form of stability, will be enough to prove that the method converges at each fixed point  $(X, T)$  as we refine the grid in both space and time. Moreover we expect that for a stable method the global order of accuracy will agree with the order of the local truncation error, e.g., for Crank-Nicolson we expect that

$$U_i^n - u(X, T) = O(k^2 + h^2)$$

as  $k, h \rightarrow 0$  when  $ih \equiv X$  and  $nk \equiv T$  are fixed.

For linear PDE's, the fact that consistency plus stability is equivalent to stability is known as the *Lax Equivalence Theorem*, and is discussed in Section 13.5 after introducing the proper concept of stability. As usual, it is the definition and study of stability that is the hard (and interesting) part of this theory.

### 13.2 Method of Lines discretizations

To understand how stability theory for time-dependent PDE's relates to the stability theory we have already developed for time-dependent ODE's, it is easiest to first consider the so-called Method of Lines (MOL) discretization of the PDE. In this approach we first discretize in space alone, which gives a large system of ODE's with each component of the system corresponding to the solution at some grid point, as a function of time. The system of ODE's can then be solved using one of the methods for ODE's that we have previously studied.





### 13.3 Stability theory

We can now investigate the stability of schemes like (13.5) or (13.7) since these can be interpreted as standard ODE methods applied to the linear system (13.11). We expect the method to be stable if  $k\lambda \in \mathcal{S}$ , i.e., if the time step  $k$  multiplied by any eigenvalue  $\lambda$  of  $A$  lies in the stability region of the ODE method, as discussed in Chapter 8.

We have determined the eigenvalues of  $A$  in (2.23),

$$\lambda_p = \frac{2}{h^2}(\cos(p\pi h) - 1), \text{ for } p = 1, 2, \dots, m, \quad (13.13)$$

where again  $m$  and  $h$  are related by  $h = 1/(m+1)$ . Note that there is a new wrinkle here relative to the ODE's we considered in Chapter 8: the eigenvalues  $\lambda_p$  depend on the mesh width  $h$ . As we refine the grid and  $h \rightarrow 0$ , the dimension of  $A$  increases, the number of eigenvalues we must consider increases, and the values of the eigenvalues change.

This is something we must bear in mind when we attempt to prove convergence as  $k, h \rightarrow 0$ . To begin with, however, let's consider the simpler question of how the method behaves for some fixed  $k$  and  $h$ , i.e., the question of absolute stability in the ODE sense. Then it is clear that the method is absolutely stable (i.e., the effect of past errors will not grow exponentially in future time steps) provided that  $k\lambda_p \in \mathcal{S}$  for each  $p$ .

For the matrix (13.12) coming from the heat equation, the eigenvalues lie on the negative real axis and the one farthest from the origin is  $\lambda_m \approx -4/h^2$ . Hence we require that  $-4k/h^2 \in \mathcal{S}$  (assuming the stability region is connected along the negative real axis up to the origin, as is generally the case).

**Example 13.2.** If we use Euler's method to obtain the discretization (13.5), then we must require  $|1 + k\lambda| \leq 1$  for each eigenvalue (see Chapter 8) and hence we require  $-2 \leq -4k/h^2 \leq 0$ . This limits the time step allowed to

$$\frac{k}{h^2} \leq \frac{1}{2}. \quad (13.14)$$

This is a severe restriction: the time step must go like  $h^2$  as we refine the grid, which is much smaller than the spatial width  $h$ . The Crank-Nicolson method, on the other hand, is based on the trapezoidal method for the ODE, which is absolutely stable in the whole left half plane. Hence the Crank-Nicolson method is stable for *any* time step  $k > 0$ .

### 13.4 Stiffness of the heat equation

Note that the system of ODE's we are solving is quite stiff, particularly for small  $h$ . The eigenvalues of  $A$  lie on the negative real axis with one fairly close to the origin,  $\lambda_1 \approx -\pi^2$  for all  $h$ , while the largest in magnitude is  $\lambda_m \approx -4/h^2$ . The "stiffness ratio" of the system is  $4/\pi^2 h^2$ , which grows rapidly as  $h \rightarrow 0$ . As a result the explicit Euler method is stable only for very small time steps  $k \leq \frac{1}{2}h^2$ . This is typically much smaller than what we would like to use over physically meaningful times, and an implicit method designed for stiff problems will be more efficient.

The stiffness is a reflection of the very different time scales present in solutions to the physical problem modelled by the heat equation. High frequency spatial oscillations in the initial data will decay very rapidly due to rapid diffusion over very short distances, while smooth data decays much more slowly since diffusion over long distances takes much longer. This is most easily seen by writing down the exact solution to the heat equation on  $0 \leq x \leq 1$  with  $g_0(t) = g_1(t) \equiv 0$  as a Fourier sine series:

$$u(x, t) = \sum_{j=1}^{\infty} \hat{u}_j(t) \sin(j\pi x).$$

Inserting this in the heat equation gives the ODE's

$$\hat{u}'_j(t) = -j^2\pi^2 \hat{u}_j(t), \quad \text{for } j = 1, 2, \dots \quad (13.15)$$

and so

$$\hat{u}_j(t) = e^{-j^2\pi^2 t} \hat{u}_j(0),$$

with the  $\hat{u}_j(0)$  determined as the Fourier coefficients of the initial data  $\eta(x)$ .

We can view the equations (13.15) as an infinite system of ODE's, but which are decoupled so that the coefficient matrix is diagonal, with eigenvalues  $-j^2\pi^2$  for  $j = 1, 2, \dots$ . By choosing data with sufficiently rapid oscillation (large  $j$ ), we can obtain arbitrarily rapid decay. For general initial data there may be some transient period when any high frequencies are rapidly damped, but then the long-time behavior is dominated by the slower decay rates. See Figure 13.3 for some examples of the time evolution with different sets of data.

If we are solving the problem over the long time periods needed to track this slow diffusion, then we would ultimately (after any physical transients have decayed) like to use rather large time steps. Typically the variation in time is then on roughly the same scale as variations in space, and so we would like to take  $k \approx h$  so that we have roughly the same resolution in time as we do in space. A method that requires  $k \approx h^2$  forces us to take a much finer temporal discretization that we should need to represent smooth solutions. If  $h = 0.001$ , for example, then we would need to take 1000 time steps to cover each time interval that should be well modelled by a single time step. This is the same difficulty we encountered with stiff ODE's in Chapter 10.

**Note:** The remark above that we want  $k \approx h$  is reasonable assuming the method we are using has the same order of accuracy in both space and time. The method (13.5) does not have this property. Since the error is  $O(k + h^2)$  we might want to take  $k = O(h^2)$  just to get the same level of accuracy in both space and time. In this sense the stability restriction  $k = O(h^2)$  may not seem unreasonable, but this is simply another reason for not wanting to use this particular method in practice.

**Note:** The general diffusion equation is  $u_t = \beta u_{xx}$  and in practice the diffusion coefficient  $\beta$  may be different from 1 by many orders of magnitude. How does this affect our conclusions above? We would expect by scaling considerations that we should take  $k \approx h/\beta$  in order to achieve comparable resolution in space and time, *i.e.*, we would like to take  $\beta k/h \approx 1$ . (Note that  $\hat{u}_j(t) = \exp(-j^2\pi^2\beta t)\hat{u}_j(0)$  in this case.) With the MOL discretization we obtain the system (13.11) but  $A$  now has a factor  $\beta$  in front. For stability we thus require  $-4\beta k/h^2 \in \mathcal{S}$ , which requires  $\beta k/h^2$  to be order 1 for any explicit method. This is smaller than what we wish to use by a factor of  $h$ , regardless of the magnitude of  $\beta$ . So our conclusions on stiffness are unchanged by  $\beta$ . In particular, even when the diffusion coefficient is very small it is best to use an implicit method. (At least for the case of pure diffusion. If we are solving an advection-diffusion or reaction-diffusion equation where there are other time scales determined by other phenomena, then if the diffusive term has a very small coefficient we may be able to use an explicit method efficiently because of other restrictions on the time step.)

**Note:** The physical problem of diffusion is “infinitely stiff” in the sense that there are eigenvalues  $-j^2\pi^2$  with arbitrarily large magnitude. Luckily the discrete problem is not this stiff. The reason it is not is that, once we discretize in space, only a finite number of spatial frequencies can be represented. As we refine the grid we can represent higher and higher frequencies, leading to the increasing stiffness ratio as  $h \rightarrow 0$ .

## 13.5 Convergence

So far we have only discussed absolute stability, and determined the relation between  $k$  and  $h$  that must be satisfied to ensure that errors do not grow exponentially as we march forward in time on this fixed grid. We now address the question of convergence at a fixed point  $(X, T)$  as the grid is refined. It turns out that in general exactly the same relation between  $k$  and  $h$  must now be required to hold as we vary  $k$  and  $h$ , letting both go to zero.

In other words, we cannot let  $k$  and  $h$  go to zero at arbitrary independent rates and necessarily expect the resulting approximations to converge to the solution of the PDE. For a particular sequence of grids  $(k_1, h_1), (k_2, h_2), \dots$ , with  $k_j \rightarrow 0$  and  $h_j \rightarrow 0$ , we will expect convergence only if the proper relation

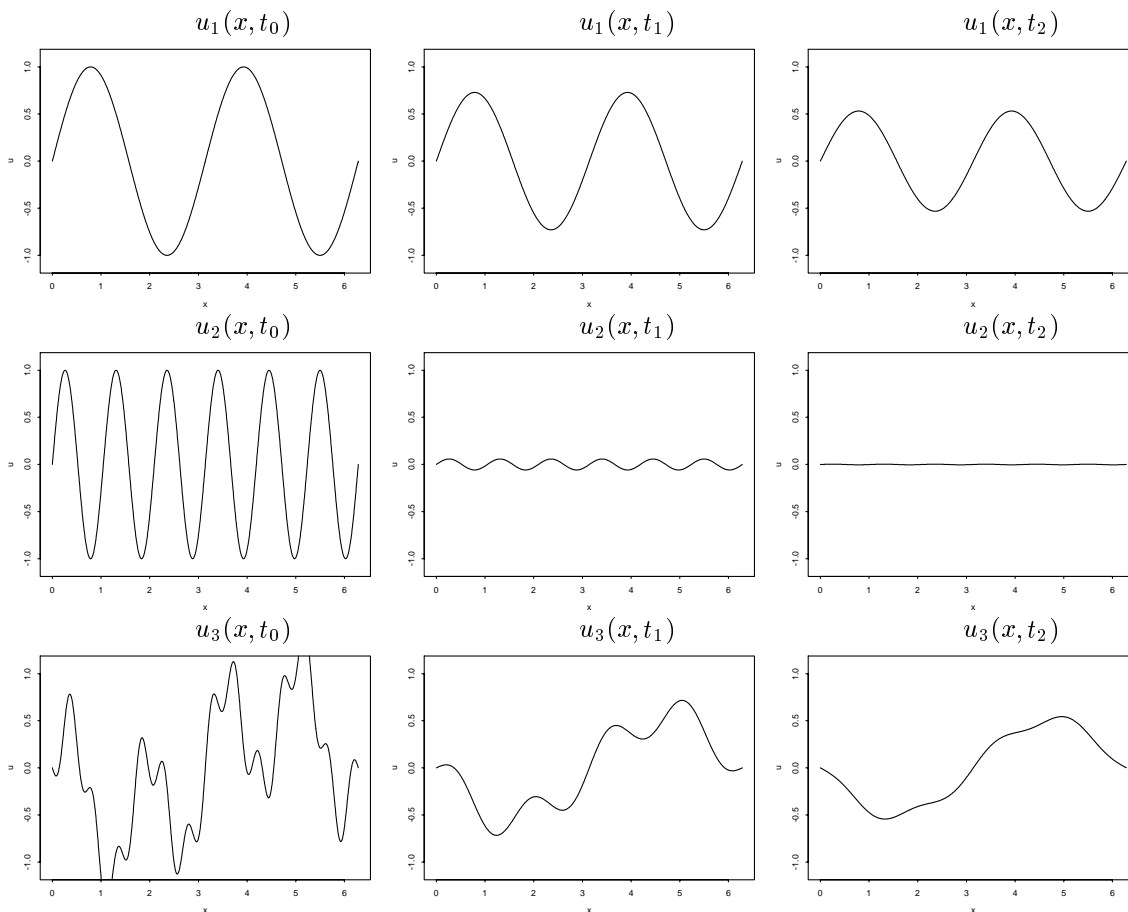


Figure 13.3: Solutions to the heat equation at three different times (columns) shown for three different sets of initial conditions (rows). In the first example  $u_1$  there is only a low frequency which decays slowly. In  $u_2$  there is only a higher frequency which decays more quickly. In  $u_3$  there is a mixture of frequencies, and the high frequencies are most rapidly damped (the initial rapid transient).

ultimately holds for each pair. For the method (13.5), for example, the sequence of approximations will converge only if  $k_j/h_j^2 \leq 1/2$  for all  $j$  sufficiently large.

It is sometimes easiest to think of  $k$  and  $h$  as being related by some fixed rule (e.g., we might choose  $k = 0.4h^2$  for the method (13.5)), so that we can speak of convergence as  $k \rightarrow 0$  with the understanding that this rule applies on each grid.

The methods we have studied so far can be written in the form

$$U^{n+1} = BU^n + b^n \quad (13.16)$$

where  $B \in \mathbb{R}^{m \times m}$  on a grid with  $h = 1/(m+1)$  and  $b^n \in \mathbb{R}^m$ . In the usual way, we can apply the difference equation to the exact solution and obtain

$$u^{n+1} = Bu^n + b^n + k\tau^n \quad (13.17)$$

where

$$u^n = \begin{bmatrix} u(x_1, t_n) \\ u(x_2, t_n) \\ \vdots \\ u(x_m, t_n) \end{bmatrix}, \quad \tau^n = \begin{bmatrix} \tau(x_1, t_n) \\ \tau(x_2, t_n) \\ \vdots \\ \tau(x_m, t_n) \end{bmatrix}$$

Subtracting (13.17) from (13.16) gives the difference equation for the global error  $E^n = U^n - u^n$ :

$$E^{n+1} = BE^n - k\tau^n,$$

and hence, as usual,

$$E^n = B^n E^0 - k \sum_{m=1}^n B^{n-m} \tau^{m-1},$$

from which we obtain

$$\|E^n\| \leq \|B^n\| \|E^0\| + k \sum_{m=1}^n \|B^{n-m}\| \|\tau^{m-1}\|. \quad (13.18)$$

The method *converges* provided it is *consistent*, so that  $\tau^{m-1} \rightarrow 0$  in each step, and *stable*, which now requires that  $\|B^n\|$  be uniformly bounded for all  $k$  and  $n$  with  $nk \leq T$ . In the context of linear PDE's, this is known as the *Lax Equivalence Theorem*. A complete proof can be found in [RM67], but the essential inequality is (13.18). The form of stability required here, the uniform bound on  $\|B^n\|$ , is often called *Lax-Richtmyer stability* in the present context.

Recall that  $B$  depends on both  $k$  and  $h$ , but we are assuming some fixed relationship between these. For the methods we analyzed earlier in this Chapter, we found relations that would guarantee  $\|B\|_2 \leq 1$ , from which Lax-Richtmyer stability follows directly (in the 2-norm at least).

### 13.5.1 PDE vs. ODE stability theory

It may bother you that the stability we need for convergence now seems to depend on absolute stability, and on the shape of the stability region for the time-discretization, which determines the required relationship between  $k$  and  $h$ . Recall that in the case of ODE's all we needed for convergence was "zero-stability", which does not depend on the shape of the stability region except for the requirement that the point  $z = 0$  must lie in this region.

Here is the difference: With ODE's we were studying a fixed system of ODE's and the fixed set of eigenvalues  $\lambda$  were independent of  $k$ . For convergence we needed  $k\lambda$  in the stability region as  $k \rightarrow 0$ , but since these values all converge to 0 it is only the origin that is important, at least in order to prove convergence as  $k \rightarrow 0$ . Hence the need for zero-stability. With PDE's, on the other hand, in our MOL discretization the system of ODE's grows as we refine the grid, and the eigenvalues  $\lambda$  grow as  $k$  (and

hence  $h$ ) go to zero. So it is not clear that  $k\lambda$  will go to zero, and zero-stability is not sufficient. For the heat equation with  $k/h^2$  fixed, these values do not go to zero as  $k \rightarrow 0$ . For convergence we must now require that these values at least lie in the region of absolute stability as  $k \rightarrow 0$ , and this gives the stability restriction relating  $k$  and  $h$ .

Although for the methods considered so far we have obtained  $\|B\| \leq 1$ , this is not really necessary in order to have Lax-Richtmyer stability. If there is a constant  $\alpha$  so that a bound of the form

$$\|B\| \leq 1 + \alpha k \quad (13.19)$$

holds in some norm, then we will have Lax-Richtmyer stability in this norm, since

$$\|B^n\| \leq (1 + \alpha k)^n \leq e^{\alpha T}$$

for  $nk \leq T$ . Since the matrix  $B$  depends on  $k$  and grows in size as  $k \rightarrow 0$ , the general theory of stability in the sense of uniform power boundedness of such families of matrices is often nontrivial. The *Kreiss Matrix Theorem* is one important tool in many practical problems. This is discussed in [RM67] along with some other techniques. See also [Str89] for a good discussion of stability. The recent review paper [LT98] gives an overview of how ODE and PDE stability theory are related, with a discussion of stability theory based on the “energy method”, another important approach.

## 13.6 von Neumann analysis

Although it is useful to go through the MOL formulation in order to understand how stability theory for PDE’s is related to the theory for ODE’s, in practice there is another approach that will typically give the proper stability restrictions more easily.

The von Neumann approach to stability analysis is based on Fourier analysis and hence is generally limited to constant coefficient linear PDE’s. For simplicity it is usually applied to the *Cauchy problem*, which is the PDE on all space with no boundaries,  $-\infty < x < \infty$  in the one-dimensional case. Von Neumann analysis can also be used to study the stability of problems with *periodic boundary conditions*, e.g., in  $0 \leq x \leq 1$  with  $u(0, t) = u(1, t)$  imposed. This is generally equivalent to a Cauchy problem with periodic initial data.

Stability theory for PDE’s with more general boundary conditions can often be quite difficult, as the coupling between the discretization of the boundary conditions and the discretization of the PDE can be very subtle. Von Neumann analysis addresses the issue of stability of the PDE discretization alone. Some discussion of stability theory for initial boundary value problems can be found in [Str89], [RM67].

The Cauchy problem for linear PDE’s can be solved using Fourier transforms — see Chapter 12 for a review. The basic reason this works is that the functions  $e^{i\xi x}$  with wave number  $\xi = \text{constant}$  are eigenfunctions of the differential operator  $\partial_x$ ,

$$\partial_x e^{i\xi x} = i\xi e^{i\xi x},$$

and hence of any constant coefficient linear differential operator. Von Neumann analysis is based on the fact that the related grid function  $W_j = e^{ijh\xi}$  is an eigenfunction of any standard finite difference operator<sup>1</sup>. For example, if we approximate  $v'(x_j)$  by  $D_0 V_j = \frac{1}{2h}(V_{j+1} - V_{j-1})$ , then in general the grid function  $D_0 V$  is not just a scalar multiple of  $V$ . But for the special case of  $W$ , we obtain

$$\begin{aligned} D_0 W_j &= \frac{1}{2h} \left( e^{i(j+1)h\xi} - e^{i(j-1)h\xi} \right) \\ &= \frac{1}{2h} \left( e^{ih\xi} - e^{-ih\xi} \right) e^{ijh\xi} \\ &= \frac{i}{h} \sin(h\xi) e^{ijh\xi} \\ &= \frac{i}{h} \sin(h\xi) W_j. \end{aligned} \quad (13.20)$$

<sup>1</sup>Note: in this section  $i = \sqrt{-1}$  and the index  $j$  is used on the grid functions.

So  $W$  is an “eigengridfunction” of the operator  $D_0$ , with eigenvalue  $\frac{i}{h} \sin(h\xi)$ .

Note the relation between these and the eigenfunctions and eigenvalues of the operator  $\partial_x$  found earlier:  $W_j$  is simply the eigenfunction  $w(x)$  of  $\partial_x$  evaluated at the point  $x_j$ , and for small  $h\xi$  we can approximate the eigenvalue of  $D_0$  by

$$\begin{aligned} \frac{i}{h} \sin(h\xi) &= \frac{i}{h} \left( h\xi - \frac{1}{6} h^3 \xi^3 + O(h^5 \xi^5) \right) \\ &= i\xi - \frac{i}{6} h^2 \xi^3 + \dots \end{aligned}$$

This agrees with the eigenvalue  $i\xi$  of  $\partial_x$  to  $O(h^2 \xi^3)$ .

Suppose we have a grid function  $V_j$  defined at grid points  $x_j = jh$  for  $j = 0, \pm 1, \pm 2, \dots$ , which is an  $l_2$  function in the sense that the 2-norm

$$\|U\|_2 = \left( h \sum_{j=-\infty}^{\infty} |U_j|^2 \right)^{1/2}$$

is finite. Then we can express  $V_j$  as a linear combination of the grid functions  $e^{ijh\xi}$  for all  $\xi$  in the range  $-\pi/h \leq \xi \leq \pi/h$ . Functions with larger wave number  $\xi$  cannot be resolved on this grid. We can write

$$V_j = \frac{1}{\sqrt{2\pi}} \int_{-\pi/h}^{\pi/h} \hat{V}(\xi) e^{ijh\xi} d\xi$$

where

$$\hat{V}(\xi) = \frac{h}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} V_j e^{-ijh\xi}.$$

These are direct analogs of the formulas for a function  $v(x)$  in the discrete case.

Again we have *Parseval's relation*,  $\|\hat{V}\|_2 = \|V\|_2$ , although the 2-norms used for the grid function  $V_j$  and the function  $\hat{V}(\xi)$  defined on  $[-\pi/h, \pi/h]$  are different:

$$\|V\|_2 = \left( h \sum_{j=-\infty}^{\infty} |V_j|^2 \right)^{1/2}, \quad \|\hat{V}\|_2 = \left( \int_{-\pi/h}^{\pi/h} |\hat{V}(\xi)|^2 d\xi \right)^{1/2}.$$

In order to show that a finite difference method is stable in the 2-norm by the techniques discussed earlier in this chapter, we would have to show that  $\|B\|_2 \leq 1 + \alpha k$  in the notation of (13.19). This amounts to showing that there is a constant  $\alpha$  such that

$$\|U^{n+1}\|_2 \leq (1 + \alpha k) \|U^n\|_2$$

for all  $U^n$ . This can be difficult to attack directly because of the fact that computing  $\|U\|_2$  requires summing over all grid points, and each  $U_j^{n+1}$  depends on values of  $U^n$  at neighboring grid points so that all grid points are coupled together. In some cases one can work with these infinite sums directly, but it is rare that this can be done. Alternatively one can work with the matrix  $B$  itself, as we did above, but this matrix is growing as we refine the grid.

Using Parseval's relation, we see that it is sufficient to instead show that

$$\|\hat{U}^{n+1}\|_2 \leq (1 + \alpha k) \|\hat{U}^n\|_2$$

where  $\hat{U}^n$  is the Fourier transform of the grid function  $U^n$ . The utility of Fourier analysis now stems from the fact that after Fourier transforming the finite difference method, we obtain a recurrence relation for each  $\hat{U}^n(\xi)$  that is decoupled from all other wave numbers. For a 2-level method this has the form

$$\hat{U}^{n+1}(\xi) = g(\xi) \hat{U}^n(\xi). \quad (13.21)$$

The factor  $g(\xi)$ , which depends on the method, is called the *amplification factor* for the method at wave number  $\xi$ . If we can show that

$$|g(\xi)| \leq 1 + \alpha k$$

where  $\alpha$  is independent of  $\xi$ , then it follows that the method is stable, since then

$$|\hat{U}^{n+1}(\xi)| \leq (1 + \alpha k)|\hat{U}^n(\xi)| \text{ for all } \xi$$

and so

$$\|\hat{U}^{n+1}\|_2 \leq (1 + \alpha k)\|\hat{U}^n\|_2.$$

Fourier analysis allows us to obtain simple scalar recursions of the form (13.21) for each wave number separately, rather than dealing with a system of equations for  $U_j^n$  that couples together all values of  $j$ .

**Note:** Here we are assuming that  $u(x, t)$  is a scalar, so that  $g(\xi)$  is a scalar. For an system of  $s$  equations we would find that  $g(\xi)$  is an  $s \times s$  matrix for each value of  $\xi$ , so some analysis of matrix eigenvalues is still required to investigate stability. But the dimension of the matrices is  $s$ , independent of the grid spacing, unlike the MOL analysis where the matrix dimension increases as  $h \rightarrow 0$ .

**Example 13.3.** Consider the method (13.5). To apply von Neumann analysis we consider how this method works on a single wavenumber  $\xi$ , *i.e.*, we set

$$U_j^n = e^{ijh\xi}. \quad (13.22)$$

Then we expect that

$$U_j^{n+1} = g(\xi)e^{ijh\xi}, \quad (13.23)$$

where  $g(\xi)$  is the amplification factor for this wavenumber. Inserting these expressions into (13.5) gives

$$\begin{aligned} g(\xi)e^{ijh\xi} &= e^{ijh\xi} + \frac{k}{h^2} \left( e^{i\xi(j-1)h} - 2e^{ijh\xi} + e^{i\xi(j+1)h} \right) \\ &= \left( 1 + \frac{k}{h^2} (e^{-i\xi h} - 2 + e^{i\xi h}) \right) e^{ijh\xi}, \end{aligned}$$

and hence

$$g(\xi) = 1 + 2\frac{k}{h^2}(\cos(\xi h) - 1).$$

Since  $-1 \leq \cos(\xi h) \leq 1$  for any value of  $\xi$ , we see that

$$1 - 4\frac{k}{h^2} \leq g(\xi) \leq 1$$

for all  $\xi$ . We can guarantee that  $|g(\xi)| \leq 1$  for all  $\xi$  if we require

$$4\frac{k}{h^2} \leq 2.$$

This is exactly the stability restriction (13.14) we found earlier for this method. If this restriction is violated, then the Fourier components with some wave number  $\xi$  will be amplified (and, as expected, it is the largest wavenumbers that go unstable first as  $k$  is increased).

**Example 13.4.** The fact that the Crank-Nicolson method is stable for all  $k$  and  $h$  can also be shown using von Neumann analysis. Substituting (13.22) and (13.23) into the difference equations (13.7) and cancelling the common factor of  $e^{ijh\xi}$  gives the following relation for  $g \equiv g(\xi)$ :

$$g = 1 + \frac{k}{2h^2} (e^{-i\xi h} - 2 + e^{i\xi h}) (1 + g)$$

and hence

$$g = \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z} \quad (13.24)$$



where

$$\begin{aligned} z &= \frac{k}{h^2}(e^{-i\xi h} - 2 + e^{i\xi h}) \\ &= \frac{2k}{h^2}(\cos(\xi h) - 1). \end{aligned} \tag{13.25}$$

Since  $z \leq 0$  for all  $\xi$ , we see that  $|g| \leq 1$  and the method is stable for any choice of  $k$  and  $h$ .

Note that (13.24) agrees with the root  $\zeta_1$  found for the Trapezoidal method in Example 8.12, while the  $z$  determined in (13.25), for certain values of  $\xi$ , is simply  $k$  times an eigenvalue  $\lambda_p$  from (13.13), the eigenvalues of the Method of Lines matrix. So there is a close connection between the von Neumann approach and the MOL reduction to a system of ODE's.

## 13.7 Multi-dimensional problems

In two space dimensions the heat equation takes the form

$$u_t = u_{xx} + u_{yy} \quad (13.26)$$

with initial conditions  $u(x, y, 0) = \eta(x, y)$  and boundary conditions all along the boundary of our spatial domain  $\Omega$ . We can discretize in space using a discrete Laplacian of the form considered in Chapter 3, say the five-point Laplacian from Section 3.2:

$$\nabla_h^2 U_{ij} = \frac{1}{h^2}(U_{i-1,j} + U_{i+1,j} + U_{i,j-1} + U_{i,j+1} - 4U_{ij}). \quad (13.27)$$

If we then discretize in time using the trapezoidal method, we will obtain the two-dimensional version of the Crank-Nicolson method,

$$U_{ij}^{n+1} = U_{ij}^n + \frac{k}{2}[\nabla_h^2 U_{ij}^n + \nabla_h^2 U_{ij}^{n+1}]. \quad (13.28)$$

Since this method is implicit, we must solve a system of equations for all the  $U_{ij}$  where the matrix has the same nonzero structure as for the elliptic systems considered in Chapters 3 and 5. This matrix is large and sparse, and we generally do not want to solve the system by a direct method such as Gaussian Elimination. In fact this is even more true for the systems we are now considering than for the elliptic equation, because of the slightly different nature of this system, which makes other approaches even more efficient relative to direct methods. It is also extremely important now that we use the most efficient method possible, because we must now solve a linear system of this form *in every time step*, and we may need to take thousands of time steps to solve the time-dependent problem.

We can rewrite the equations (13.28) as

$$\left(I - \frac{k}{2}\nabla_h^2\right) U_{ij}^{n+1} = \left(I + \frac{k}{2}\nabla_h^2\right) U_{ij}^n. \quad (13.29)$$

The matrix for this linear system has the same pattern of nonzeros as the matrix for  $\nabla_h^2$  (see Chapter 3), but the values are scaled by  $k/2$  and then subtracted from the identity matrix, so that the diagonal elements are fundamentally different. If we call this matrix  $A$ ,

$$A = I - \frac{k}{2}\nabla_h^2,$$

then we find that the eigenvalues of  $A$  are

$$\lambda_{p_1, p_2} = 1 - \frac{k}{h^2}[(\cos(p_1\pi h) - 1) + (\cos(p_2\pi h) - 1)]$$

for  $p_1, p_2 = 1, 2, \dots, m$ , where we have used the expression for the eigenvalues of  $\nabla_h^2$  from Section 3.4. Now the largest eigenvalue of the matrix  $A$  thus has magnitude  $O(k/h^2)$  while the ones closest to the origin are at  $1 + O(k)$ . As a result the condition number of  $A$  is  $O(k/h^2)$ . By contrast, the discrete Laplacian  $\nabla_h^2$  alone has condition number  $O(1/h^2)$  as we found in Section 3.4. The smaller condition number in the present case can be expected to lead to faster convergence of iterative methods.

Moreover, we have an excellent starting guess for the solution  $U^{n+1}$  to (13.28), a fact that we can use to good advantage with iterative methods but not with direct methods. Since  $U_{ij}^{n+1} = U_{ij}^n + O(k)$ , we can use  $U_{ij}^n$ , the values from the previous time step, as initial values  $U_{ij}^{[0]}$  for an iterative method. We might do even better by extrapolating forward in time, using say  $U_{ij}^{[0]} = 2U_{ij}^n - U_{ij}^{n-1}$ , or by using an explicit method, say

$$U_{ij}^{[0]} = (I + k\nabla_h^2)U_{ij}^n.$$

This explicit method (forward Euler) would probably be unstable as a time-marching procedure alone with the value of  $k$  we have in mind, but it can be used in this context successfully.

Because of the combination of a reasonably well-conditioned system and very good initial guess, we can often get away with taking only one or two iterations in each time step, and still get global second order accuracy.

## 13.8 The LOD method

Rather than solving the coupled sparse matrix equation for all the unknowns on the grid simultaneously as in (13.29), an alternative approach is to replace this fully-coupled single time step by a sequence of steps, each of which is coupled in only one space direction, resulting in a set of tridiagonal systems which can be solved much more easily. One example is the *Locally One-Dimensional (LOD)* method:

$$U_{ij}^* = U_{ij}^n + \frac{k}{2}(D_x^2 U_{ij}^n + D_x^2 U_{ij}^*) \quad (13.30)$$

$$U_{ij}^{n+1} = U_{ij}^* + \frac{k}{2}(D_y^2 U_{ij}^* + D_y^2 U_{ij}^{n+1}). \quad (13.31)$$

or, in matrix form,

$$\left(I - \frac{k}{2}D_x^2\right)U^* = \left(I + \frac{k}{2}D_x^2\right)U^n \quad (13.32)$$

$$\left(I - \frac{k}{2}D_y^2\right)U^{n+1} = \left(I + \frac{k}{2}D_y^2\right)U^*. \quad (13.33)$$

In (13.30) we apply Crank-Nicolson in the  $x$ -direction only, solving  $u_t = u_{xx}$  alone over time  $k$ , and we call the result  $U^*$ . Then in (13.31) we take this result and apply Crank-Nicolson in the  $y$ -direction to it, solving  $u_t = u_{yy}$  alone, again over time  $k$ . Physically this corresponds to modeling diffusion in the  $x$ - and  $y$ -directions over time  $k$  as a decoupled process in which we first allow  $u$  to diffuse only in the  $x$ -direction and then only in the  $y$ -direction. If the time steps are very short then this might be expected to give similar physical behavior and hence convergence to the correct behavior as  $k \rightarrow 0$ . In fact, for the constant coefficient diffusion problem, it can even be shown that (in the absence of boundaries at least) this alternating diffusion approach gives *exactly* the same behavior as the original two-dimensional diffusion. Diffusing first in  $x$  alone over time  $k$  and then in  $y$  alone over time  $k$  gives the same result as if the diffusion occurs simultaneously in both directions. This will be shown in Chapter 18 after *fractional step methods* have been introduced in a more general context.

Numerically there is a great advantage in using (13.32) and (13.33) rather than the fully coupled (13.29). In (13.32) the unknowns  $U_{ij}^*$  are coupled together only across each row of the grid. For any fixed value of  $j$  we have a *tridiagonal system* of equations to solve for  $U_{ij}^*$  ( $i = 1, 2, \dots, m$ ). The system obtained for each value of  $j$  is completely decoupled from the system obtained for other values of  $j$ . Hence we have a set of  $m + 2$  tridiagonal systems to solve (for  $j = 0, 1, \dots, m + 1$ ), each of dimension  $m$ , rather than a single coupled system with  $m^2$  unknowns as in (13.29). Since each of these systems is tridiagonal, it is easily solved in  $O(m)$  operations by Gaussian elimination and there is no need for iterative methods. (In the next section we will see why we need to solve these for  $j = 0$  and  $j = m + 1$  as well as at the interior grid points.)

Similarly, (13.31) decouples into a set of  $m$  tridiagonal systems in the  $y$ -direction for  $i = 1, 2, \dots, m$ . Hence taking a single time step requires solving  $2m + 2$  tridiagonal systems of size  $m$ , and thus  $O(m^2)$  work. Since there are  $m^2$  grid points, this is the optimal order and no worse than an explicit method, except for a constant factor.

### 13.8.1 Boundary conditions

In solving the second set of systems (13.31), we need boundary values  $U_{i0}^*$  and  $U_{i0}^{n+1}$  along the bottom boundary and  $U_{i,m+1}^*$  and  $U_{i,m+1}^{n+1}$  along the top boundary, for terms that go on the right-hand side of each tridiagonal system. The values at level  $n + 1$  are available from the given boundary data for the heat equation, by evaluating the boundary conditions at time  $t_{n+1}$  (assuming Dirichlet boundary conditions are given). To obtain the values  $U_{i0}^*$  we solve equation (13.30) for  $j = 0$  and  $j = m + 1$  (along the boundaries) in addition to the systems along each row interior to the grid.

In order to solve the first set of systems (13.30), we need boundary values  $U_{0j}^n$  and  $U_{0j}^*$  along the left boundary and values  $U_{m+1,j}^n$  and  $U_{m+1,j}^*$  along the right boundary. The values at level  $n$  come

from the given boundary conditions, but we must determine the intermediate boundary conditions at level  $*$  along these boundaries. It is not immediately clear what values should be used. One might be tempted to think of level  $*$  as being half way between  $t_n$  and  $t_{n+1}$ , since  $U^*$  is generated in the middle of the two-step procedure used to obtain  $U^{n+1}$  from  $U^n$ . If this were valid, then evaluating the given boundary data at time  $t_{n+1/2} = t_n + k/2$  might provide values for  $U^*$  on the boundary. This is not a good idea, however, and would lead to a degradation of accuracy. The problem is that in the first step, equation (13.30) does not model the full heat equation over time  $k/2$ , but rather models part of the equation (diffusion in  $x$  alone) over the full time step  $k$ . The values along the boundary will in general evolve quite differently in the two different cases.

To determine proper values for  $U_{0j}^*$  and  $U_{m+1,j}^*$ , we can use the equations (13.31) along the left and right boundaries. At  $i = 0$ , for example, this equation gives a system of equations along the left boundary that can be viewed as a tridiagonal linear system or the unknowns  $U_{0j}^*$  in terms of the values  $U_{0j}^{n+1}$ , which are already known from the boundary conditions at time  $t_{n+1}$ . Note that we are solving this equation backwards from the way it will be used in the second step of the LOD process on the interior of the grid, and this works only because we already know  $U_{0j}^{n+1}$  from boundary data.

Since we are solving this equation backwards, we can view this as solving the diffusion equation  $u_t = u_{yy}$  over a time step of length  $-k$ , backwards in time. This makes sense physically — the intermediate solution  $U^*$  represents what is obtained from  $U^n$  by doing diffusion in  $x$  alone, with no diffusion yet in  $y$ . There are in principle two ways to get this, either by starting with  $U^n$  and diffusing in  $x$ , or by starting with  $U^{n+1}$  and “undiffusing” in  $y$ . We are using the latter approach along the boundaries to generate data for  $U^*$ .

Equivalently we can view this as solving the *backward heat equation*  $u_t = -u_{yy}$  over time  $k$ . This may be cause for concern, since the backward heat equation is ill-posed. However, since we are only doing this over one time step starting with given values  $U_{0j}^{n+1}$  in each time step, this turns out to be a stable procedure.

There is still a difficulty at the corners. In order to solve (13.31) for  $U_{0j}^*$ ,  $j = 1, 2, \dots, m$ , we need to know the values of  $U_{00}^*$  and  $U_{0,m+1}^*$  that are the boundary values for this system. These can be approximated using some sort of explicit and uncentered approximation to either  $u_t = u_{xx}$  starting with  $U^n$ , or to  $u_t = -u_{yy}$  starting with  $U^{n+1}$ . For example we might use

$$U_{00}^* = U_{00}^{n+1} - \frac{k}{h^2}(U_{00}^{n+1} - 2U_{01}^{n+1} + U_{02}^{n+1}),$$

which uses the approximation to  $u_{yy}$  centered at  $(x_0, y_1)$ .

Alternatively, rather than solving the tridiagonal systems obtained from (13.31) for  $U_{0j}^*$ , we could simply use an explicit approximation to the backwards heat equation along this boundary,

$$U_{0j}^* = U_{0j}^{n+1} - \frac{k}{h^2}(U_{0,j-1}^{n+1} - 2U_{0j}^{n+1} + U_{0,j+1}^{n+1}), \quad (13.34)$$

for  $j = 1, 2, \dots, m$ . This eliminates the need for values of  $U^*$  in the corners. Again, since this is not iterated but only done starting with given (and presumably smooth) boundary data  $U^{n+1}$  in each time step, this yields a stable procedure.

### 13.8.2 Accuracy and stability

With proper treatment of the boundary conditions, it can be shown that the LOD method is second order accurate. This will be discussed further in Chapter 18 after fractional step methods are introduced more generally.

It can also be shown that this method, like full Crank-Nicolson, is unconditionally stable for any time step.

### 13.8.3 The ADI method

A modification of the LOD method is also often used, in which the two steps each involve discretization in only one spatial direction at the advanced time level (giving decoupled tridiagonal systems again), but coupled with discretization in the *opposite* direction at the old time level. The classical method of this form is:

$$U_{ij}^* = U_{ij}^n + \frac{k}{2}(D_y^2 U_{ij}^n + D_x^2 U_{ij}^*) \quad (13.35)$$

$$U_{ij}^{n+1} = U_{ij}^* + \frac{k}{2}(D_x^2 U_{ij}^* + D_y^2 U_{ij}^{n+1}). \quad (13.36)$$

This is called the *Alternating Direction Implicit (ADI)* method and was first introduced by Douglas and Rachford [?]. This again gives decoupled tridiagonal systems to solve in each step:

$$\left(I - \frac{k}{2}D_x^2\right)U^* = \left(I + \frac{k}{2}D_y^2\right)U^n \quad (13.37)$$

$$\left(I - \frac{k}{2}D_y^2\right)U^{n+1} = \left(I + \frac{k}{2}D_x^2\right)U^*. \quad (13.38)$$

With this method, each of the two steps involves diffusion in both the  $x$ - and  $y$ -directions. In the first step the diffusion in  $x$  is modelled implicitly while diffusion in  $y$  is modelled explicitly, with the roles reversed in the second step. In this case each of the two steps can be shown to give a *first-order accurate* approximation to the full heat equation over time  $k/2$ , so that  $U^*$  represents a first-order accurate approximation to the solution at time  $t_{n+1/2}$ . Because of the symmetry of the two steps, however, the local error introduced in the second step almost exactly cancels the local error introduced in the first step, so that the combined method is in fact *second-order accurate* over the full time step.

Because  $U^*$  does approximate the solution at time  $t_{n+1/2}$  in this case, it is possible to simply evaluate the given boundary conditions at time  $t_{n+1/2}$  to generate the necessary boundary values for  $U^*$ . This will maintain second-order accuracy. A better error constant can be achieved by using slightly modified boundary data which introduces the expected error in  $U^*$  into the boundary data that should be cancelled out by the second step.

## Chapter 14

# Advection Equations

In this chapter we consider numerical methods for the scalar advection equation

$$u_t + au_x = 0$$

where  $a$  is a constant. See Section 11.3 for a discussion of this equation. For the Cauchy problem we also need initial data

$$u(x, 0) = \eta(x).$$

This is the simplest example of a *hyperbolic* equation, and is so simple that we can write down the exact solution,

$$u(x, t) = \eta(x - at). \quad (14.1)$$

One can verify directly that this is the solution (see also Chapter 12). However, many of the issues that arise more generally in discretizing hyperbolic equations can be most easily seen with this equation. Other hyperbolic systems will be discussed later.

The first approach we might consider is the analog of the method (13.4) for the heat equation. Using the centered difference in space and the forward difference in time results in

$$\frac{U_i^{n+1} - U_i^n}{k} = -\frac{a}{2h}(U_{i+1}^n - U_{i-1}^n), \quad (14.2)$$

which can be rewritten as

$$U_i^{n+1} = U_i^n - \frac{ak}{2h}(U_{i+1}^n - U_{i-1}^n). \quad (14.3)$$

This again has the stencil shown in Figure 13.1(a). In practice this method is not useful because of stability considerations, as we will see in the next section.

A minor modification gives a more useful method. If we replace  $U_i^n$  on the right-hand side of (14.3) by the average  $\frac{1}{2}(U_{i-1}^n + U_{i+1}^n)$  then we obtain the *Lax-Friedrichs method*,

$$U_i^{n+1} = \frac{1}{2}(U_{i-1}^n + U_{i+1}^n) - \frac{ak}{2h}(U_{i+1}^n - U_{i-1}^n). \quad (14.4)$$

**Exercise 14.1** *Compute the local truncation error and show that Lax-Friedrichs is first order accurate in both space and time.*

Because of the low accuracy, this method is not commonly used in practice, but it serves to illustrate some stability issues and so we will study this method along with (14.3) before describing higher order methods such as the well-known Lax-Wendroff method.

We will see in the next section that Lax-Friedrichs is Lax-Richtmyer stable (see Section 13.5) and convergent provided

$$\left| \frac{ak}{h} \right| \leq 1. \quad (14.5)$$

Note that this stability restriction allows us to use a time step  $k = O(h)$  even though the method is explicit, unlike the case of the heat equation. The basic reason is that the advection equation involves only the first order derivative  $u_x$  rather than  $u_{xx}$  and so the difference equation involves  $1/h$  rather than  $1/h^2$ .

The time step restriction (14.5) is consistent with what we would choose anyway based on accuracy considerations, and in this sense the advection equation is **not stiff**, unlike the heat equation. This is a fundamental difference between hyperbolic equations and parabolic equations more generally, and accounts for the fact that hyperbolic equations are typically solved with explicit methods while the efficient solution of parabolic equations generally requires implicit methods.

To see that (14.5) gives a reasonable time step, note that

$$u_x(x, t) = \eta'(x - at)$$

while

$$u_t(x, t) = -a\eta'(x - at).$$

The time derivative  $u_t$  is larger in magnitude by a factor  $a$ , and so we would expect the time step required to achieve temporal resolution consistent with the spatial resolution  $h$  to be smaller by a factor of  $a$ . This suggests that the relation  $k = h/a$  would be reasonable in practice. This is completely consistent with (14.5).

## 14.1 MOL discretization

To investigate stability further we will again introduce the method of lines (MOL) discretization as we did in Section 13.2 for the heat equation. To obtain a system of equations with finite dimension we must solve the equation on some bounded domain rather than solving the Cauchy problem. However, in a bounded domain, say  $0 \leq x \leq 1$ , the advection equation can have a boundary condition specified only on one of the two boundaries. If  $a > 0$  then we need a boundary condition at  $x = 0$ , say

$$u(0, t) = g_0(t), \quad (14.6)$$

which is the *inflow* boundary in this case. The boundary at  $x = 1$  is the *outflow* boundary and the solution there is completely determined by what is advecting to the right from the interior. If  $a < 0$  we instead need a boundary condition at  $x = 1$ , which is the inflow boundary in this case.

The symmetric 3-point methods defined above can still be used near the inflow boundary, but not at the outflow boundary. Instead the discretization will have to be coupled with some “numerical boundary condition” at the outflow boundary, say a one-sided discretization of the equation. This issue complicates the stability analysis and will be discussed later.

We can obtain a nice MOL discretization if we consider the special case of *periodic boundary conditions*,

$$u(0, t) = u(1, t) \text{ for } t \geq 0.$$

Physically, whatever flows out at the outflow boundary flows back in at the inflow boundary. This also models the Cauchy problem in the case where the initial data is periodic with period 1, in which case the solution remains periodic and we only need to model a single period  $0 \leq x \leq 1$ .

In this case the value  $U_0(t) = U_{m+1}(t)$  along the boundaries is another unknown, and we must introduce one of these into the vector  $U(t)$ . If we introduce  $U_{m+1}(t)$  then we have the vector of grid

values

$$U(t) = \begin{bmatrix} U_1(t) \\ U_2(t) \\ \vdots \\ U_{m+1}(t) \end{bmatrix}.$$

For  $2 \leq j \leq m$  we have the ODE

$$U'_j(t) = -\frac{a}{2h}(U_{j+1}(t) - U_{j-1}(t)),$$

while the first and last equations are modified using the periodicity:

$$\begin{aligned} U'_1(t) &= -\frac{a}{2h}(U_2(t) - U_{m+1}(t)), \\ U'_{m+1}(t) &= -\frac{a}{2h}(U_1(t) - U_m(t)). \end{aligned}$$

This system can be written as

$$U'(t) = AU(t) \tag{14.7}$$

with

$$A = -\frac{a}{2h} \begin{bmatrix} 0 & 1 & & & -1 \\ -1 & 0 & 1 & & \\ & -1 & 0 & 1 & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 0 & 1 \\ 1 & & & & -1 & 0 \end{bmatrix} \in \mathbb{R}^{(m+1) \times (m+1)}. \tag{14.8}$$

Note that this matrix is skew-symmetric ( $A^T = -A$ ) and so its eigenvalues must be pure imaginary. In fact, the eigenvalues are

$$\lambda_p = \frac{ia}{h} \sin(2\pi ph), \text{ for } p = 1, 2, \dots, m+1. \tag{14.9}$$

The corresponding eigenvector  $u^p$  has components

$$u_j^p = e^{2\pi i p j h}, \text{ for } j = 1, 2, \dots, m+1. \tag{14.10}$$

The eigenvalues lie on the imaginary axis between  $-ia/h$  and  $ia/h$ .

For absolute stability of a time discretization we need the stability region  $\mathcal{S}$  to include this interval. Any method that includes some interval  $iy$ ,  $|y| < b$  of the imaginary axis will be stable provided  $|ak/h| \leq b$ .

### 14.1.1 Forward Euler time discretization

The method (14.3) can be viewed as the forward Euler time discretization of the MOL system of ODE's (14.7). We found in Section 8.3 that this method is only stable if  $|1 + k\lambda| \leq 1$  and the stability region  $\mathcal{S}$  is the unit circle centered at  $-1$ . No matter how small the ratio  $k/h$  is, since the eigenvalues  $\lambda_p$  from (14.9) are imaginary, the values  $k\lambda_p$  will not lie in  $\mathcal{S}$ . Hence the method is *unstable* for any fixed mesh ratio  $k/h$ . See Figure 14.1(a).

The method will be convergent if we let  $k \rightarrow 0$  faster than  $h$ , since then  $k\lambda_p \rightarrow 0$  for all  $p$  and the zero-stability of Euler's method is enough to guarantee convergence. Taking  $k$  much smaller than  $h$  is generally not desirable and the method is not used in practice. However, it is interesting to analyze this situation also in terms of Lax-Richtmyer stability, since it shows an example where the Lax-Richtmyer



stability uses a weaker bound of the form (13.19),  $\|B\| \leq 1 + \alpha k$ , rather than  $\|B\| \leq 1$ . Here  $B = I + kA$ . Suppose we take  $k = h^2$ , for example. Then we have

$$|1 + k\lambda_p|^2 \leq 1 + (ka/h)^2$$

for each  $p$  (using the fact that  $\lambda_p$  is pure imaginary) and so

$$|1 + k\lambda_p|^2 \leq 1 + a^2 h^2 = 1 + a^2 k.$$

Hence  $\|I + kA\|_2^2 \leq 1 + a^2 k$  and if  $nk \leq T$  we have

$$\|(I + kA)^n\|_2 \leq (1 + a^2 k)^{n/2} \leq e^{a^2 T/2},$$

showing the uniform boundedness of  $\|B^n\|$  (in the 2-norm) needed for Lax-Richtmyer stability.

### 14.1.2 Leapfrog

A better time discretization is to use the midpoint method (6.17),

$$U^{n+1} = U^{n-1} + 2kAU^n,$$

which gives the *Leapfrog method* for the advection equation,

$$U_i^{n+1} = U_i^{n-1} + \frac{ak}{h}(U_{i+1}^n - U_{i-1}^n). \quad (14.11)$$

This is a 3-level explicit method, and is second order accurate in both space and time.

Recall from Section 8.3 that the stability region of the midpoint method is the interval  $i\alpha$  for  $-1 < \alpha < 1$  of the imaginary axis. This method is hence stable on the advection equation provided  $|ak/h| < 1$  is satisfied.

### 14.1.3 Lax-Friedrichs

Now consider the Lax-Friedrichs method (14.4). At first glance this does not look like a time discretization of an ODE system of the form  $U'(t) = AU(t)$ . However, we can rewrite (14.4) using the fact that

$$\frac{1}{2}(U_{i-1}^n + U_{i+1}^n) = U_i^n + \frac{1}{2}(U_{i-1}^n - 2U_i^n + U_{i+1}^n),$$

to obtain

$$U_i^{n+1} = U_i^n - \frac{ak}{2h}(U_{i+1}^n - U_{i-1}^n) + \frac{1}{2}(U_{i-1}^n - 2U_i^n + U_{i+1}^n). \quad (14.12)$$

This can be rearranged to give

$$\frac{U_i^{n+1} - U_i^n}{k} + a \left( \frac{U_{i+1}^n - U_{i-1}^n}{2h} \right) = \frac{h^2}{2k} \left( \frac{U_{i-1}^n - 2U_i^n + U_{i+1}^n}{h^2} \right).$$

If we compute the local truncation error from this form we see, of course, that it is consistent with the advection equation  $u_t + au_x = 0$ , since the term on the right hand side vanishes as  $k, h \rightarrow 0$  (assuming  $k/h$  is fixed). However, it looks more like a discretization of the advection-diffusion equation

$$u_t + au_x = \epsilon u_{xx}$$

where  $\epsilon = h^2/2k$ . Actually, we will see later that it is in fact a *second order* accurate method on a slightly different advection-diffusion equation. Viewing Lax-Friedrichs in this way allows us to investigate the diffusive nature of the method quite precisely.

For our present purposes, however, the crucial part is that we can now view (14.12) as resulting from a Forward Euler discretization of the system of ODE's

$$U'(t) = BU(t)$$

with

$$B = -\frac{a}{2h} \begin{bmatrix} 0 & 1 & & & -1 \\ -1 & 0 & 1 & & \\ & -1 & 0 & 1 & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 0 & 1 \\ 1 & & & & -1 & 0 \end{bmatrix} + \frac{\epsilon}{h^2} \begin{bmatrix} -2 & 1 & & & & 1 \\ 1 & -2 & 1 & & & \\ & 1 & -2 & 1 & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & -2 & 1 \\ 1 & & & & 1 & -2 \end{bmatrix} \quad (14.13)$$

where  $\epsilon = h^2/2k$ . The matrix  $B$  differs from the matrix  $A$  of (14.8) by the addition of a small multiple of the second difference operator, which is symmetric rather than skew-symmetric. As a result the eigenvalues of  $B$  are shifted off the imaginary axis and now lie in the left half plane. There is now a hope that each  $k\lambda$  will lie in the stability region of Euler's method if  $k$  is small enough relative to  $h$ .

It can be verified that the eigenvectors (14.10) are also eigenvectors of the second difference operator (with periodic boundary conditions) that appears in (14.13), and hence these are also the eigenvectors of the full matrix  $B$ . We can easily compute that the eigenvalues of  $B$  are

$$\mu_p = \frac{ia}{h} \sin(2\pi ph) - \frac{2\epsilon}{h^2} (1 - \cos(2\pi ph)). \quad (14.14)$$

The values  $k\lambda$  are plotted in the complex plane for various different values of  $\epsilon$  in Figure 14.1. They lie on an ellipse centered at  $-2k\epsilon/h^2$  with semi-axes of length  $2k\epsilon/h^2$  in the  $x$ -direction and  $ak/h$  in the  $y$ -direction. For the special case  $\epsilon = h^2/2k$  used in Lax-Friedrichs, we have  $-2k\epsilon/h^2 = -1$ , provided that  $|ak/h| \leq 1$ . (If  $|ak/h| > 1$  then the top and bottom of the ellipse would extend outside the circle.) The forward Euler method is stable as a time-discretization, and hence the Lax-Friedrichs method is Lax-Richtmyer stable, provided  $|ak/h| \leq 1$ .

**Exercise 14.2** Verify the expression (14.14) for the eigenvalues of  $B$ .

## 14.2 The Lax-Wendroff method

One way to achieve second order accuracy on the advection equation is to use a second order temporal discretization of the system of ODE's (14.7) since this system is based on a second order spatial discretization. This can be done with the midpoint method, for example, which gives rise to the Leapfrog scheme (14.11) already discussed. However, this is a 3-level method and for various reasons it is often much more convenient to use 2-level methods for PDE's whenever possible — in more than one dimension the need to store several levels of data may be restrictive, boundary conditions can be harder to impose, and combining methods using fractional step procedures (as discussed in Chapter 18) may require 2-level methods for each step, to name a few reasons. Moreover, the Leapfrog method is “non-dissipative” in a sense that will be discussed in Chapter 16, leading to potential stability problems if the method is extended to variable coefficient or nonlinear problems.

Another way to achieve second order accuracy in time would be to use the trapezoidal method to discretize the system (14.7), as was done to derive the Crank-Nicolson method for the heat equation. But this is an implicit method and for hyperbolic equations there is generally no need to introduce this complication and expense.

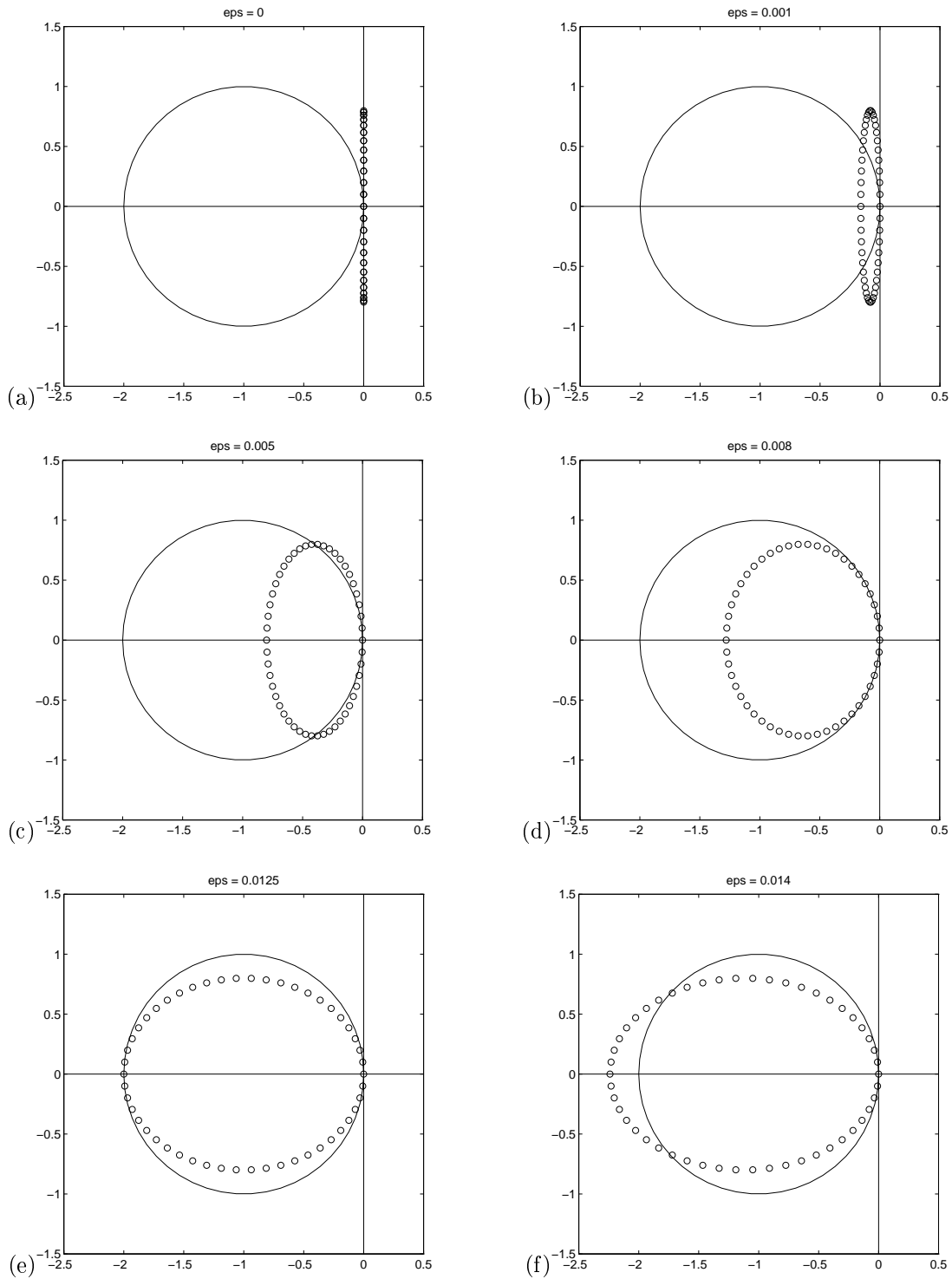


Figure 14.1: Eigenvalues of the matrix  $B$  in (14.13), for various values of  $\epsilon$ , in the case  $h = 1/50$  and  $k = 0.8h$ ,  $a = 1$ , so  $ak/h = 0.8$ .

(a) shows the case  $\epsilon = 0$  which corresponds the forward Euler method (14.3).

(d) shows the case  $\epsilon = a^2k/2$ , the Lax-Wendroff method (14.16).

(e) shows the case  $\epsilon = h^2/2k$  the Lax-Friedrichs method (14.4).

The method is stable for  $\epsilon$  between  $a^2k/2$  and  $h^2/2k$ .

Another possibility is to use a 2-stage Runge-Kutta method such as the one in Example 6.15 for the time discretization. This can be done, though some care must be exercised near boundaries and the use of a multi-stage method again typically requires additional storage.

One simple way to achieve a 2-level explicit method with higher accuracy is to use the idea of Taylor series methods, as described in Section 6.6. Applying this directly to the linear system of ODE's  $U'(t) = AU(t)$  (and using  $U'' = AU' + A^2U$ ) gives the second order method

$$U^{n+1} = U^n + kAU^n + \frac{1}{2}k^2A^2U^n.$$

Here  $A$  is the matrix (14.8) and computing  $A^2$  and writing the method at the typical grid point then gives

$$U_j^{n+1} = U_j^n - \frac{ak}{2h}(U_{j+1}^n - U_{j-1}^n) + \frac{a^2k^2}{8h^2}(U_{j-2}^n - 2U_j^n + U_{j+2}^n). \quad (14.15)$$

This method is second order accurate and explicit, but has a 5-point stencil involving the points  $U_{j-2}^n$  and  $U_{j+2}^n$ . With periodic boundary conditions this is not a problem, but with other boundary conditions this method needs more numerical boundary conditions than a 3-point method.

Note that the last term in (14.15) is an approximation to  $\frac{1}{2}a^2k^2u_{xx}$  using a centered difference based on stepsize  $2h$ . A simple way to achieve a second-order accurate 3-point method is to replace this term by the more standard 3-point formula. We then obtain the standard *Lax-Wendroff method*:

$$U_j^{n+1} = U_j^n - \frac{ak}{2h}(U_{j+1}^n - U_{j-1}^n) + \frac{a^2k^2}{2h^2}(U_{j-1}^n - 2U_j^n + U_{j+1}^n). \quad (14.16)$$

A cleaner way to derive this method is to use Taylor series expansions directly on the PDE  $u_t + au_x = 0$ , to obtain

$$u(x, t + k) = u(x, t) + ku_t(x, t) + \frac{1}{2}k^2u_{tt}(x, t) + \dots$$

Replacing  $u_t$  by  $-au_x$  and  $u_{tt}$  by  $a^2u_{xx}$  gives

$$u(x, t + k) = u(x, t) + kau_x(x, t) + \frac{1}{2}k^2a^2u_{xx}(x, t) + \dots$$

If we now use the standard centered approximations to  $u_x$  and  $u_{xx}$  and drop the higher order terms, we obtain the Lax-Wendroff method (14.16). It is also clear how we could obtain higher-order accurate explicit 2-level methods by this same approach, by retaining more terms in the series and approximating the spatial derivatives (including the higher-order spatial derivatives that will then arise) by suitably high order accurate finite difference approximations. The same approach can also be used with other PDE's. The key is to replace the time derivatives arising in the Taylor series expansion with spatial derivatives, using expressions obtained by differentiating the original PDE.

Note that with periodic boundary conditions, the method (14.16) can be written in the same form (14.13) as the Lax-Friedrichs method, but with  $\epsilon = a^2k/2$  instead of the value  $\epsilon = h^2/2k$  used in Lax-Friedrichs.

### 14.2.1 Stability analysis

We can analyze the stability of Lax-Wendroff following the same approach used for Lax-Friedrichs in Section 14.1. We can view it as Euler's method applied to the linear system of ODE's with

$$k\mu_p = -i \left( \frac{ak}{h} \right) \sin(p\pi h) + \left( \frac{ak}{h} \right)^2 (\cos(p\pi h) - 1).$$

These values all lie on an ellipse centered at  $-(ak/h)^2$  with semiaxes of length  $(ak/h)^2$  and  $ak/h$ . If  $|ak/h| \leq 1$  then all of these eigenvalues lie inside the stability region of Euler's method. Figure 14.1(d) shows an example in the case  $ak/h = 0.8$

### 14.3 Upwind methods

So far we have considered methods based on symmetric approximations to derivatives. Alternatively, one might use a nonsymmetric approximation to  $u_x$  in the advection equation, e.g.,

$$u_x(x_j, t) \approx \frac{1}{h}(U_j - U_{j-1}) \quad (14.17)$$

or

$$u_x(x_j, t) \approx \frac{1}{h}(U_{j+1} - U_j). \quad (14.18)$$

These are both *one-sided approximations*, since they use data only to one side or the other of the point  $x_j$ . Coupling one of these approximations with forward differencing in time gives the following methods for the advection equation:

$$U_j^{n+1} = U_j^n - \frac{ak}{h}(U_j^n - U_{j-1}^n) \quad (14.19)$$

or

$$U_j^{n+1} = U_j^n - \frac{ak}{h}(U_{j+1}^n - U_j^n). \quad (14.20)$$

These methods are first order accurate in both space and time. One might wonder why we would want to use such approximations, since centered approximations are more accurate. For the advection equation, however, there is an asymmetry in the equations due to the fact that the equation models translation at speed  $a$ . If  $a > 0$  then the solution moves to the right, while if  $a < 0$  it moves to the left. We will see that there are situations where it is best to acknowledge this asymmetry and use one-sided differences in the appropriate direction. In practice we often want to also use more accurate approximations however, and we will see later how to extend these methods to higher order accuracy.

The choice between the two methods (14.19) and (14.20) should be dictated by the sign of  $a$ . Note that the true solution over one time step can be written as

$$u(x_j, t + k) = u(x_j - ak, t)$$

so that the solution at the point  $x_j$  at the next time level is given by data to the *left* of  $x_j$  if  $a > 0$  whereas it is determined by data to the *right* of  $x_j$  if  $a < 0$ . This suggests that (14.19) might be a better choice for  $a > 0$  and (14.20) for  $a < 0$ .

In fact the stability analysis below shows that (14.19) is stable only if

$$0 \leq \frac{ak}{h} \leq 1. \quad (14.21)$$

Since  $k$  and  $h$  are positive, we see that this method can only be used if  $a > 0$ . This method is called the *upwind method* when used on the advection equation with  $a > 0$ . If we view the equation as modeling the concentration of some tracer in air blowing past us at speed  $a$ , then we are looking in the correct upwind direction to judge how the concentration will change with time. (This is also referred to as an *upstream differencing* method in the literature.)

Conversely, (14.20) is stable only if

$$-1 \leq \frac{ak}{h} \leq 0, \quad (14.22)$$

and can only be used if  $a < 0$ . In this case (14.20) is the proper upwind method to use.

### 14.3.1 Stability analysis

The method (14.19) can be written as

$$U_j^{n+1} = U_j^n - \frac{ak}{2h}(U_{j+1}^n - U_{j-1}^n) + \frac{ak}{2h}(U_{j+1}^n - 2U_j^n + U_{j-1}^n), \quad (14.23)$$

which puts it in the form (14.13) with  $\epsilon = ah/2$ . We have seen previously that methods of this form are stable provided  $|ak/h| \leq 1$  and also  $-2 < -2\epsilon k/h^2 < 0$ . Since  $k, h > 0$ , this requires in particular that  $\epsilon > 0$ . For Lax-Friedrichs and Lax-Wendroff, this condition was always satisfied, but for upwind the value of  $\epsilon$  depends on  $a$  and we see that  $\epsilon > 0$  only if  $a > 0$ . If  $a < 0$  then the eigenvalues of the MOL matrix lie on a circle that lies entirely in the right half plane, and the method will certainly be unstable. If  $a > 0$  then the above requirements lead to the stability restriction (14.21).

If we think of (14.23) as modeling an advection-diffusion equation, then we see that  $a < 0$  corresponds to a negative diffusion coefficient. This leads to an ill-posed equation, as in the “backward heat equation” (see Chapter 12).

The method (14.20) can also be written in a form similar to (14.23), but the last term will have a minus sign in front of it. In this case we need  $a < 0$  for any hope of stability and then easily derive the stability restriction (14.22).

The three methods Lax-Wendroff, upwind, and Lax-Friedrichs, can all be written in the same form (14.13) with different values of  $\epsilon$ . If we call these values  $\epsilon_{LW}$ ,  $\epsilon_{up}$ , and  $\epsilon_{LF}$  respectively, then we have

$$\epsilon_{LW} = \frac{a^2k}{2}, \quad \epsilon_{up} = \frac{ah}{2}, \quad \epsilon_{LF} = \frac{h^2}{2k}.$$

Note that

$$\epsilon_{LW} = \frac{ak}{h}\epsilon_{up} \quad \text{and} \quad \epsilon_{up} = \frac{ak}{h}\epsilon_{LF}.$$

If  $0 < \frac{ak}{h} < 1$  then  $\epsilon_{LW} < \epsilon_{up} < \epsilon_{LF}$  and the method is stable for any value of  $\epsilon$  between  $\epsilon_{LW}$  and  $\epsilon_{LF}$ .

### 14.3.2 The Beam-Warming method

A second-order accurate method with the same one-sided character can be derived by following the derivation of the Lax-Wendroff method, but using one-sided approximations to the spatial derivatives. This results in the *Beam-Warming* method, which for  $a > 0$  takes the form

$$U_j^{n+1} = U_j^n - \frac{ak}{2h}(3U_j^n - 4U_{j-1}^n + U_{j-2}^n) + \frac{ak}{2h}(U_j^n - 2U_{j-1}^n + U_{j-2}^n). \quad (14.24)$$

**Exercise 14.3** Compute the local truncation error for *Beam-Warming*.

**Exercise 14.4** Show that *Beam-Warming* is stable for  $0 \leq ak/h \leq 2$ .

## 14.4 Characteristic tracing

The solution to the advection equation is given by (14.1). The value of  $u$  is constant along each characteristic, which for this example are straight lines with constant slope. Over a single time step we have

$$u(x_j, t_{n+1}) = u(x_j - ak, t_n). \quad (14.25)$$

Tracing this characteristic back over time step  $k$  from the grid point  $x_j$  results in the picture shown in Figure 14.2(a). Note that if  $0 < ak/h < 1$  then the point  $x_j - ak$  lies between  $x_{j-1}$  and  $x_j$ . If we carefully choose  $k$  and  $h$  so that  $ak/h = 1$  exactly, then  $x_j - ak = x_{j-1}$  and we would find that

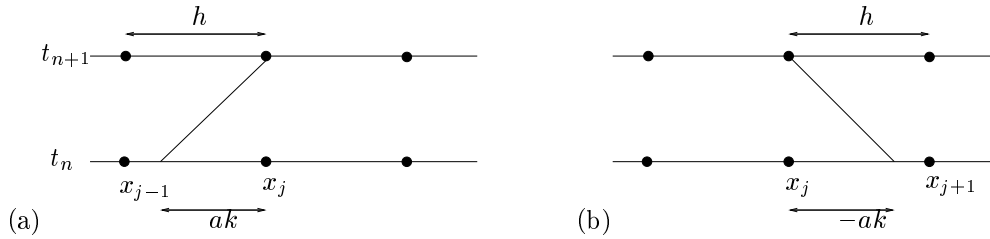


Figure 14.2: Tracing the characteristic of the advection equation back in time from the point  $(x_j, t_{n+1})$  to compute the solution according to (14.25). Interpolating the value at this point from neighboring grid values gives the upwind method (for linear interpolation) or the Lax-Wendroff or Beam-Warming methods (quadratic interpolation). (a) shows the case  $a > 0$ , (b) shows the case  $a < 0$ .

$u(x_j, t_{n+1}) = u(x_{j-1}, t_n)$ . The solution should just shift one grid cell to the right in each time step. We could compute the *exact* solution numerically with the method

$$U_j^{n+1} = U_{j-1}^n. \quad (14.26)$$

Actually, all of the 2-level methods that we have considered so far reduce to the formula (14.26) in this special case  $ak = h$ , and each of these methods happens to be exact in this case.

If  $ak/h < 1$  then the point  $x_j - ak$  is not exactly at a grid point, as illustrated in Figure 14.2. However, we might attempt to use the relation (14.25) as the basis for a numerical method by computing an approximation to  $u(x_j - ak, t_n)$  based on interpolation from the grid values  $U_i^n$  at nearby grid points. For example, we might perform simple linear interpolation between  $U_{j-1}^n$  and  $U_j^n$ . Fitting a linear function to these points gives the function

$$p(x) = U_j^n + (x - x_j) \left( \frac{U_j^n - U_{j-1}^n}{h} \right). \quad (14.27)$$

Evaluating this at  $x_j - ak$  and using this to define  $U_j^{n+1}$  gives

$$U_j^{n+1} = p(x_j - ak) = U_j^n - \frac{ak}{h}(U_j^n - U_{j-1}^n).$$

This is precisely the first-order upwind method (14.19). Note that this can also be interpreted as a linear combination of the two values  $U_{j-1}^n$  and  $U_j^n$ :

$$U_j^{n+1} = \left( 1 - \frac{ak}{h} \right) U_j^n + \frac{ak}{h} U_{j-1}^n. \quad (14.28)$$

Moreover this is a *convex combination* (i.e., the coefficients of  $U_j^n$  and  $U_{j-1}^n$  are both nonnegative and sum to 1) provided the stability condition (14.21) is satisfied, which is also the condition required to insure that  $x_j - ak$  lies between the two points  $x_{j-1}$  and  $x_j$ . In this case we are *interpolating* between these points with the function  $p(x)$ . If the stability condition is violated then we would be using  $p(x)$  to *extrapolate* outside of the interval where the data lies. It is easy to see that this sort of extrapolation can lead to instability — consider what happens if the data  $U_j^n$  is oscillatory with  $U_j^n = (-1)^j$ , for example.

To obtain better accuracy, we might try using a higher order interpolating polynomial based on more data points. If we define a quadratic polynomial  $p(x)$  by interpolating the values  $U_{j-1}^n$ ,  $U_j^n$ , and  $U_{j+1}^n$ , and then define  $U_j^{n+1}$  by evaluating  $p(x_j - ak)$ , we simply obtain the Lax-Wendroff method (14.16). Note that in this case we are properly interpolating provided that the stability restriction  $|ak/h| \leq 1$  is satisfied. If we instead base our quadratic interpolation on the three points  $U_{j-2}^n$ ,  $U_{j-1}^n$ , and  $U_j^n$ , then we obtain the Beam-Warming method (14.24).

**Exercise 14.5** Verify that Lax-Wendroff and Beam-Warming can be obtained as described above.

## Chapter 15

# The CFL Condition

The discussion of Section 14.4 suggests that for the advection equation, the point  $x_j - ak$  must be bracketed by points used in the stencil of the finite difference method if the method is to be stable and convergent. This turns out to be a *necessary* condition in general for any method developed for the advection equation: If  $U_j^{n+1}$  is computed based on values  $U_{j+p}^n, U_{j+p+1}^n, \dots, U_{j+q}^n$  with  $p \leq q$  (negative values are allowed for  $p$  and  $q$ ), then we must have  $x_{j+p} \leq x_j - ak \leq x_{j+q}$  or the method cannot be convergent. Since  $x_i = ih$ , this requires

$$-q \leq \frac{ak}{h} \leq -p.$$

This result for the advection equation is one special case of a much more general principle that is called the *CFL condition*. This condition is named after Courant, Friedrichs, and Lewy, who wrote a fundamental paper in 1928 that was essentially the first paper on the stability and convergence of finite difference methods for partial differential equations. (The original paper is in German[CFL28] but an English translation is available in [CFL67].)

To understand this general condition, we must discuss the *domain of dependence* of a time-dependent PDE. (See, e.g., [Kev90] for more details.) For the advection equation, the solution  $u(X, T)$  at some fixed point  $(X, T)$  depends on the initial data  $\eta$  at only a single point:  $u(X, T) = u(X - aT)$ . We say that the domain of dependence of the point  $(X, T)$  is the point  $X - aT$ :

$$\mathcal{D}(X, T) = \{X - aT\}.$$

If we modify the data  $\eta$  at this point then the solution  $u(X, T)$  will change, while modifying the data at any other point will have no effect on the solution at this point.

This is a rather unusual situation for a PDE. More generally we might expect the solution at  $(X, T)$  to depend on the data at several points or over a whole interval. In Chapter ?? we consider hyperbolic systems of equations of the form  $u_t + Au_x = 0$ , where  $u \in \mathbb{R}^s$  and  $A \in \mathbb{R}^{s \times s}$  is a matrix with real eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_s$ . If these values are distinct then we will see that the solution  $u(X, T)$  depends on the data at the  $s$  distinct points  $X - \lambda_1 T, \dots, X - \lambda_s T$  and hence

$$\mathcal{D}(X, T) = \{X - \lambda_p T \text{ for } p = 1, 2, \dots, s\}.$$

The heat equation  $u_t = u_{xx}$  has a much larger domain of dependence. For this equation the solution at any point  $(X, T)$  depends on the data *everywhere* and the domain of dependence is the whole real line,

$$\mathcal{D}(X, T) = (-\infty, \infty).$$

This equation is said to have infinite propagation speed, since data at any point is felt everywhere at any small time in the future (though its effect of course decays exponentially away from this point).



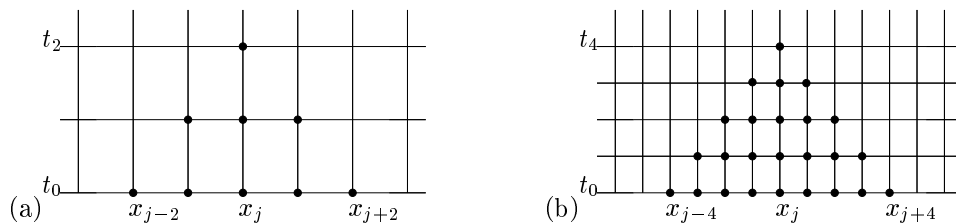


Figure 15.1: (a) Numerical domain of dependence of a grid point when using a 3-point explicit method. (b) On a finer grid.

A finite difference method also has a domain of dependence. On a particular fixed grid we define the domain of dependence of a grid point  $(x_j, t_n)$  to be the set of grid points  $x_i$  at the initial time  $t = 0$  with the property that the data  $U_i^0$  at  $x_i$  has an effect on the solution  $U_j^n$ . For example, with the Lax-Wendroff method (14.16) or any other 3-point method, the value  $U_j^n$  depends on  $U_{j-1}^{n-1}$ ,  $U_j^{n-1}$ , and  $U_{j+1}^{n-1}$ . These values depend in turn on  $U_{j-2}^{n-2}$  through  $U_{j+2}^{n-2}$ . Tracing back to the initial time we obtain a triangular array of grid points as seen in Figure 15.1(a), and we see that  $U_j^n$  depends on the initial data at the points  $x_{j-n}, \dots, x_{j+n}$ .

Now consider what happens if we refine the grid, keeping  $k/h$  fixed. Figure 15.1(b) shows the situation when  $k$  and  $h$  are reduced by a factor of 2, focusing on the same value of  $(X, T)$  which now corresponds to  $U_{2j}^{2n}$  on the finer grid. This value depends on twice as many values of the initial data, but these values all lie within the same interval and are merely twice as dense.

If the grid is refined further with  $k/h \equiv r$  fixed, then clearly the numerical domain of dependence of the point  $(X, T)$  will fill in the interval  $[X - T/r, X + T/r]$ . As we refine the grid, we hope that our computed solution at  $(X, T)$  will converge to the true solution  $u(X, T) = \eta(X - aT)$ . Clearly this can only be possible if

$$X - T/r \leq X - aT \leq X + T/r. \quad (15.1)$$

Otherwise, the true solution will depend only on a value  $\eta(X - aT)$  that is never seen by the numerical method, no matter how fine a grid we take. We could change the data at this point and hence change the true solution without having any effect on the numerical solution, so the method cannot be convergent for general initial data.

Note that the condition (15.1) translates into  $|a| \leq 1/r$  and hence  $|ak/h| \leq 1$ . This can also be written as  $|ak| \leq h$ , which just says that over a single time step the characteristic we trace back must lie within one grid point of  $x_j$  (recall the discussion of interpolation vs. extrapolation in Section 14.4).

The CFL condition generalizes this idea:

**The CFL Condition:** *A numerical method can be convergent only if its numerical domain of dependence contains the true domain of dependence of the PDE, at least in the limit as  $k$  and  $h$  go to zero.*

For the Lax-Wendroff method the condition on  $k$  and  $h$  required by the CFL condition is exactly the stability restriction we derived by von Neumann analysis. It is important to note that in general the CFL condition is only a *necessary* condition. If it is violated then the method cannot be convergent. If it is satisfied, then the method *might* be convergent, but a proper consistency and stability analysis is required to prove this and determine the proper stability restriction on  $k$  and  $h$ .

**Example 15.1.** The 3-point method (14.3) has the same numerical domain of dependence as Lax-Wendroff, but is unstable for any fixed value of  $k/h$  even though the CFL condition is satisfied for  $|ak/h| \leq 1$ .

For the first-order upwind and Beam-Warming methods, the stability restriction agrees exactly with what is required by the CFL condition.

**Example 15.2.** For the heat equation the true domain of dependence is the whole real line. It appears that any 3-point explicit method violates the CFL condition, and indeed it does if we fix  $k/h$

as the grid is refined. However, recall from Section 14.1.1 that the 3-point explicit method (13.5) is convergent as we refine the grid provided we have  $k/h^2 \leq 1/2$ . In this case when we make the grid finer by a factor of 2 in space it will become finer by a factor of 4 in time, and hence the numerical domain of dependence will cover a wider interval at time  $t = 0$ . As  $k \rightarrow 0$  the numerical domain of dependence will spread to cover the entire real line, and hence the CFL condition is satisfied in this case.

An implicit method such as the Crank-Nicolson method (13.7) satisfies the CFL condition for any time step  $k$ . In this case the numerical domain of dependence is the entire real line because the tridiagonal linear system couples together all points in such a manner that the solution at each point depends on the data at all points (*i.e.*, the inverse of a tridiagonal matrix is dense).



## Chapter 16

# Modified Equations — Numerical Diffusion and Dispersion

Our standard tool for estimating the accuracy of a finite difference method has been the “local truncation error”. Seeing how well the true solution of the PDE satisfies the difference equation gives an indication of the accuracy of the difference equation. Now we will study a slightly different approach which can be very illuminating as it reveals much more about the structure and behavior of the numerical solution rather than just its size.

The idea is to ask the following question: Is there a PDE  $v_t = \dots$  such that our numerical approximation  $U_j^n$  is actually the *exact* solution to this PDE,  $U_j^n = v(x_j, t_n)$ ? Or, less ambitiously, can we at least find a PDE that is better satisfied by  $U_j^n$  than the original PDE we were attempting to model? If so, then studying the behavior of solutions to this PDE should tell us much about how the numerical approximation is behaving. This can be advantageous because it is often easier to study the behavior of PDE’s than of finite difference formulas.

In fact it is possible to find a PDE that is exactly satisfied by the  $U_j^n$ , by doing Taylor series expansions as we do to compute the local truncation error. However, this PDE will have an infinite number of terms involving higher and higher powers of  $k$  and  $h$ . By truncating this series at some point we will obtain a PDE that is simple enough to study and yet gives a good indication of the behavior of the  $U_j^n$ .

### 16.1 Upwind

This is best illustrated with an example. Consider the upwind method (14.19) for the advection equation  $u_t + au_x = 0$  in the case  $a > 0$ ,

$$U_j^{n+1} = U_j^n - \frac{ak}{h}(U_j^n - U_{j-1}^n). \quad (16.1)$$

The process of deriving the modified equation is very similar to computing the local truncation error, only now we insert the formula  $v(x, t)$  into the difference equation. This is supposed to be a function that agrees exactly with  $U_j^n$  at the grid points and so, unlike  $u(x, t)$ , the function  $v(x, t)$  satisfies (16.1) exactly:

$$v(x, t + k) = v(x, t) - \frac{ak}{h}(v(x, t) - v(x - h, t)).$$

Expanding these terms in Taylor series about  $(x, t)$  and simplifying gives

$$\left( v_t + \frac{1}{2}kv_{tt} + \frac{1}{6}k^2v_{ttt} + \dots \right) + a \left( v_x - \frac{1}{2}hv_{xx} + \frac{1}{6}h^2v_{xxx} + \dots \right) = 0.$$

We can rewrite this as

$$v_t + av_x = \frac{1}{2}(ahv_{xx} - kv_{tt}) + \frac{1}{6}(ah^2v_{xxx} - k^2v_{ttt}) + \dots$$

This is the PDE that  $v$  satisfies. If we take  $k/h$  fixed, then the terms on the right hand side are  $O(k)$ ,  $O(k^2)$ , etc. so that for small  $k$  we can truncate this series to get a PDE that is quite well satisfied by the  $U_j^n$ .

If we drop all the terms on the right hand side we just recover the original advection equation. Since we have then dropped terms of  $O(k)$ , we expect that  $U_j^n$  satisfies this equation to  $O(k)$ , as we know to be true since this upwind method is first order accurate.

If we keep the  $O(k)$  terms then we get something more interesting:

$$v_t + av_x = \frac{1}{2}(ahv_{xx} - kv_{tt}) \quad (16.2)$$

This involves second derivatives in both  $x$  and  $t$ , but we can derive a slightly different modified equation with the same accuracy by differentiating (16.2) with respect to  $t$  to obtain

$$v_{tt} = -av_{xt} + \frac{1}{2}(ahv_{xxt} - kv_{ttt})$$

and with respect to  $x$  to obtain

$$v_{tx} = -av_{xx} + \frac{1}{2}(ahv_{xxx} - kv_{ttx}).$$

Combining these gives

$$v_{tt} = a^2v_{xx} + O(k).$$

Inserting this in (16.2) gives

$$v_t + av_x = \frac{1}{2}(ahv_{xx} - a^2kv_{xx}) + O(k^2).$$

Since we have already decided to drop terms of  $O(k^2)$ , we can drop these terms here also to obtain

$$v_t + av_x = \frac{1}{2}ah \left(1 - \frac{ak}{h}\right) v_{xx}. \quad (16.3)$$

This is now a familiar advection-diffusion equation. The grid values  $U_j^n$  can be viewed as giving a *second order accurate* approximation to the true solution of this equation (whereas they only give first order accurate approximations to the true solution of the advection equation).

**Exercise 16.1** View (16.1) as a numerical method for the equation (16.2). Compute the local truncation error and verify that it is  $O(k^2)$ .

The fact that the modified equation is an advection-diffusion equation tells us a great deal about how the numerical solution behaves. Solutions to the advection-diffusion equation translate at the proper speed  $a$  but also diffuse and are smeared out. This is clearly visible in Figure 16.1.

Note that the diffusion coefficient in (16.2) is  $\frac{1}{2}(ah - a^2k)$ , which vanishes in the special case  $ak = h$ . In this case we already know that the exact solution to the advection equation is recovered by the upwind method.

Also note that the diffusion coefficient is positive only if  $0 < ak/h < 1$ . This is precisely the stability limit of upwind. If this is violated, then the diffusion coefficient in the modified equation is negative, giving an ill-posed problem with exponentially growing solutions. Hence we see that even some information about stability can be extracted from the modified equation.

**Exercise 16.2** Determine the modified equation for (14.2) and show that the diffusion coefficient is always negative.

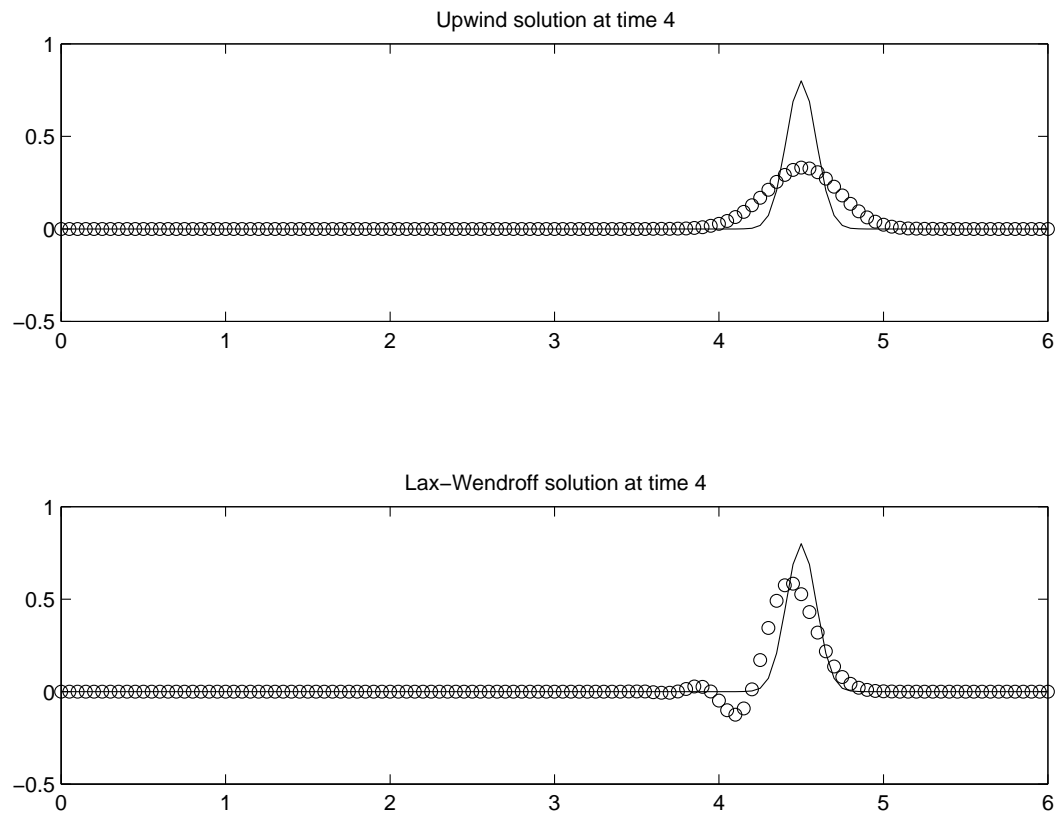


Figure 16.1: Numerical solution using upwind (diffusive) and Lax-Wendroff (dispersive) methods.

## 16.2 Lax-Wendroff

If the same procedure is followed for the Lax-Wendroff method, we find that all  $O(k)$  terms drop out of the modified equation, as is expected since this method is second order accurate on the advection equation. The modified equation obtained by retaining the  $O(k^2)$  term and then replacing time derivatives by spatial derivatives is

$$v_t + av_x = -\frac{1}{6}ah^2 \left(1 - \left(\frac{ak}{h}\right)^2\right) v_{xxx}. \quad (16.4)$$

The Lax-Wendroff method produces a *third order* accurate solution to this equation. This equation has a very different character from (16.2). The  $v_{xxx}$  term leads to *dispersive* behavior rather than diffusion. (See Chapter 12 and [Str89], [Whi74] for further details.) This is clearly seen in Figure 16.1, where the  $U_j^n$  computed with Lax-Wendroff are compared to the true solution of the advection equation. The magnitude of the error is smaller than with the upwind method for a given set of  $k$  and  $h$ , since it is a higher order method, but the dispersive term leads to an oscillating solution and also a shift in the location of the main peak, a *phase error*.

The group velocity for wave number  $\xi$  under Lax-Wendroff is

$$c_g = a - \frac{1}{2}ah^2 \left(1 - \left(\frac{ak}{h}\right)^2\right) \xi^2$$

which is less than  $a$  for all wave numbers. As a result the numerical result can be expected to develop a train of oscillations behind the peak, with the high wave numbers lagging farthest behind the correct location.

If we retain one more term in the modified equation for Lax-Wendroff, we would find that the  $U_j^n$  are fourth order accurate solutions to an equation of the form

$$v_t + av_x = \frac{1}{6}ah^2 \left(\left(\frac{ak}{h}\right)^2 - 1\right) v_{xxx} - \epsilon v_{xxxx}, \quad (16.5)$$

where the  $\epsilon$  in the fourth order dissipative term is  $O(h^3)$  and positive when the stability bound holds. This higher order dissipation causes the highest wave numbers to be damped, so that there is a limit to the oscillations seen in practice.

The fact that this method can produce oscillatory approximations is one of the reasons that the first-order upwind method is sometimes preferable in practice. In some situations nonphysical oscillations may be disastrous, for example if the value of  $u$  represents a concentration that cannot go negative or exceed some limit without difficulties arising elsewhere in the modeling process.

## 16.3 Beam-Warming

The Beam-Warming method (14.24) has a similar modified equation,

$$v_t + av_x = \frac{1}{6}ah^2 \left(2 - \frac{3ak}{h} + \left(\frac{ak}{h}\right)^2\right) v_{xxx}. \quad (16.6)$$

In this case the group velocity is greater than  $a$  for all wave numbers in the case  $0 < ak/h < 1$ , so that the oscillations move ahead of the main hump. If  $1 < ak/h < 2$  then the group velocity is less than  $a$  and the oscillations fall behind.

## Chapter 17

# Hyperbolic Systems and High-Resolution Methods

This chapter has been adapted directly from some recent lecture notes [LMDMar] on finite volume methods for hyperbolic systems. The full text of my contribution to these notes is available on-line (see the reference or the 586 webpage). More details on some of this material can be found in [LeV90].

The notation in this chapter is somewhat different from elsewhere, in particular  $q$  is used for the solution instead of  $u$ , with  $u$  now being used for velocity.

### 17.1 Scalar equations

We begin our study of conservation laws by considering the scalar case. Many of the difficulties encountered with systems of equations are already encountered here, and a good understanding of the scalar equation is required before proceeding.

#### 17.1.1 The linear advection equation

We first consider the linear advection equation,

$$q_t + aq_x = 0. \quad (17.1)$$

The Cauchy problem is defined by this equation on the domain  $-\infty < x < \infty$ ,  $t \geq 0$  together with initial conditions

$$q(x, 0) = q_0(x). \quad (17.2)$$

The solution to this problem is simply

$$q(x, t) = q_0(x - at) \quad (17.3)$$

for  $t \geq 0$ , as can be easily verified. As time evolves, the initial data simply propagates unchanged to the right (if  $a > 0$ ) or left (if  $a < 0$ ) with velocity  $a$ . The solution  $q(x, t)$  is constant along each of the rays  $x - at = x_0$ , which are known as the **characteristics** of the equation.

Note that the characteristics are curves in the  $x$ - $t$  plane satisfying the ordinary differential equations  $x'(t) = a$ ,  $x(0) = x_0$ . If we differentiate  $q(x, t)$  along one of these curves to find the rate of change of  $q$  along the characteristic, we find that

$$\begin{aligned} \frac{d}{dt}q(x(t), t) &= \frac{\partial}{\partial t}q(x(t), t) + \frac{\partial}{\partial x}q(x(t), t) x'(t) \\ &= q_t + aq_x \\ &= 0, \end{aligned} \quad (17.4)$$



confirming that  $q$  is constant along these characteristics.

### 17.1.2 Burgers' equation

Now consider the nonlinear scalar equation

$$q_t + f(q)_x = 0, \quad (17.5)$$

where  $f(q)$  is a nonlinear function of  $q$ . We will assume for the most part that  $f(q)$  is a convex function,  $f''(q) > 0$  for all  $q$  (or, equally well,  $f$  is concave with  $f''(q) < 0 \forall q$ ). The convexity assumption corresponds to a “genuine nonlinearity” assumption for systems of equations that holds in many important cases, such as the Euler equations. The nonconvex case is also important in some applications (e.g., in MHD, see Section ??) but is more complicated mathematically.

By far the most famous model problem in this field is **Burgers' equation**, in which  $f(q) = \frac{1}{2}q^2$ , so (17.5) becomes

$$q_t + qq_x = 0. \quad (17.6)$$

Actually this should be called the “inviscid Burgers' equation”, since the equation originally studied by Burgers also includes a viscous term:

$$q_t + qq_x = \epsilon q_{xx}. \quad (17.7)$$

This is about the simplest model that includes the nonlinear and viscous effects of fluid dynamics.

Consider the inviscid equation (17.6) with smooth initial data. For small time, a solution can be constructed by following characteristics. Note that (17.6) looks like an advection equation, but with the advection velocity  $q$  equal to the value of the advected quantity. The characteristics satisfy

$$x'(t) = q(x(t), t) \quad (17.8)$$

and along each characteristic  $q$  is constant, since

$$\begin{aligned} \frac{d}{dt}q(x(t), t) &= \frac{\partial}{\partial t}q(x(t), t) + \frac{\partial}{\partial x}q(x(t), t)x'(t) \\ &= q_t + qq_x \\ &= 0. \end{aligned} \quad (17.9)$$

Moreover, since  $q$  is constant on each characteristic, the slope  $x'(t)$  is constant by (17.8) and so the characteristics are straight lines, determined by the initial data.

If the initial data is smooth, then this can be used to determine the solution  $q(x, t)$  for small enough  $t$  that characteristics do not cross: for each  $(x, t)$  we can solve the equation

$$x = \xi + q(\xi, 0)t \quad (17.10)$$

for  $\xi$  and then

$$q(x, t) = q(\xi, 0). \quad (17.11)$$

### 17.1.3 Shock formation

For larger  $t$  the equation (17.10) may not have a unique solution. This happens when the characteristics cross, as will eventually happen if  $q_x(x, 0)$  is negative at any point. At the time  $T_b$  where the characteristics first cross, the function  $q(x, t)$  has an infinite slope — the wave “breaks” and a shock forms. Beyond this point there is no classical solution of the PDE, and the weak solution we wish to determine becomes discontinuous.

For times beyond the breaking time some of the characteristics have crossed and so there are points  $x$  where there are three characteristics leading back to  $t = 0$ . One can view the “solution”  $q$  at such a time as a triple-valued function. However, the density of a gas cannot possibly be triple valued at a point. We can determine the correct physical behavior by adopting the **vanishing viscosity** approach. The equation (17.6) is a model of (17.7) valid only for small  $\epsilon$  and smooth  $q$ . When it breaks down, we must return to (17.7). If the initial data is smooth and  $\epsilon$  very small, then before the wave begins to break the  $\epsilon q_{xx}$  term is negligible compared to the other terms and the solutions to the two PDEs look nearly identical. However, as the wave begins to break, the second derivative term  $q_{xx}$  grows much faster than  $q_x$ , and at some point the  $\epsilon q_{xx}$  term is comparable to the other terms and begins to play a role. This term keeps the solution smooth for all time, preventing the breakdown of solutions that occurs for the hyperbolic problem. This smooth solution has a steep transition zone where the viscous term is important. As  $\epsilon \rightarrow 0$  this zone becomes sharper and approaches the discontinuous solution known as a shock. It is this vanishing-viscosity solution that we hope to capture by solving the inviscid equation.

### 17.1.4 Weak solutions

A natural way to define a generalized solution of the inviscid equation that does not require differentiability is to go back to the integral form of the conservation law, and say that  $q(x, t)$  is a generalized solution if

$$\int_{x_1}^{x_2} q(x, t_2) dx = \int_{x_1}^{x_2} q(x, t_1) dx + \int_{t_1}^{t_2} f(q(x_1, t)) - f(q(x_2, t)) dt. \quad (17.12)$$

is satisfied for all  $x_1, x_2, t_1, t_2$ .

There is another approach that results in a different integral formulation that is often more convenient to work with. This is a mathematical technique that can be applied more generally to rewrite a differential equation in a form where less smoothness is required to define a “solution”. The basic idea is to take the PDE, multiply by a smooth “test function”, integrate one or more times over some domain, and then use integration by parts to move derivatives off the function  $q$  and onto the smooth test function. The result is an equation involving fewer derivatives on  $q$ , and hence requiring less smoothness.

In our case we will use test functions  $\phi \in C_0^1(\mathbb{R} \times \mathbb{R})$ . Here  $C_0^1$  is the space of functions that are continuously differentiable with compact support. If we multiply  $q_t + f_x = 0$  by  $\phi(x, t)$  and then integrate over space and time, we obtain

$$\int_0^\infty \int_{-\infty}^{+\infty} [\phi q_t + \phi f(q)_x] dx dt = 0. \quad (17.13)$$

Now integrate by parts, yielding

$$\int_0^\infty \int_{-\infty}^{+\infty} [\phi_t q + \phi_x f(q)] dx dt = - \int_{-\infty}^\infty \phi(x, 0) q(x, 0) dx. \quad (17.14)$$

Note that nearly all the boundary terms which normally arise through integration by parts drop out due to the requirement that  $\phi$  have compact support, and hence vanishes at infinity. The remaining boundary term brings in the initial conditions of the PDE, which must still play a role in our weak formulation.

**Definition 17.1.1** *The function  $q(x, t)$  is called a weak solution of the conservation law if (17.14) holds for all functions  $\phi \in C_0^1(\mathbb{R} \times \mathbb{R})$ .*

The vanishing-viscosity generalized solution we defined above is a weak solution in the sense of (17.14), and so this definition includes the solution we are looking for. Unfortunately, weak solutions

are often not unique, and so an additional problem is to identify *which* weak solution is the physically correct vanishing-viscosity solution. Again, one would like to avoid working with the viscous equation directly, but it turns out that there are other conditions one can impose on weak solutions that are easier to check and will also pick out the correct solution. These are usually called **entropy conditions** by analogy with the gas dynamics case, where a discontinuity is physically realistic only if the entropy of the gas *increases* as it crosses the shock.

### 17.1.5 The Riemann problem

The conservation law together with piecewise constant data having a single discontinuity is known as the Riemann problem. As an example, consider Burgers' equation  $q_t + qq_x = 0$  with piecewise constant initial data

$$q(x, 0) = \begin{cases} q_l & \text{if } x < 0 \\ q_r & \text{if } x > 0. \end{cases} \quad (17.15)$$

The form of the solution depends on the relation between  $q_l$  and  $q_r$ .

**Case I.**  $q_l > q_r$ .

In this case there is a unique weak solution,

$$q(x, t) = \begin{cases} q_l & \text{if } x < st \\ q_r & \text{if } x > st, \end{cases} \quad (17.16)$$

where

$$s = (q_l + q_r)/2 \quad (17.17)$$

is the **shock speed**, the speed at which the discontinuity travels. A general expression for the shock speed will be derived below. Note that characteristics in each of the regions where  $q$  is constant go *into* the shock (see Figure 17.1) as time advances.

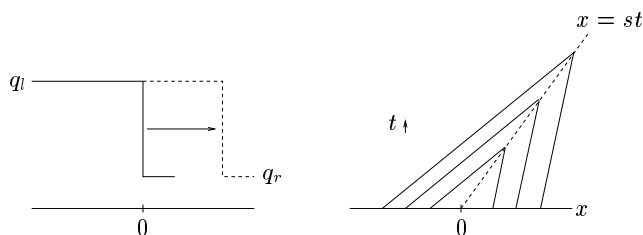


Figure 17.1: Shock wave

**Case II.**  $q_l < q_r$ .

In this case there are infinitely many weak solutions. One of these is again (17.16), (17.17) in which the discontinuity propagates with speed  $s$ . Note that characteristics now go *out* of the shock (Figure 17.2(a)) and that this solution is not stable to perturbations. If the data is smeared out slightly, or if a small amount of viscosity is added to the equation, the solution changes completely.

Another weak solution is the **rarefaction wave**

$$q(x, t) = \begin{cases} q_l & \text{if } x < q_l t \\ x/t & \text{if } q_l t \leq x \leq q_r t \\ q_r & \text{if } x > q_r t. \end{cases} \quad (17.18)$$

This solution is stable to perturbations and is in fact the vanishing-viscosity generalized solution (Figure 17.2(b)).

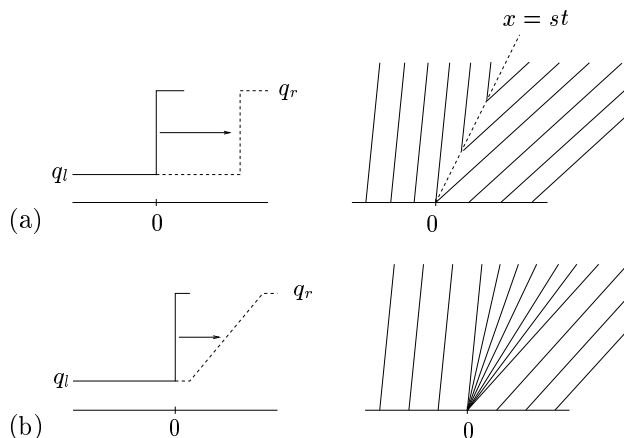


Figure 17.2: (a) Entropy-violating shock. (b) Rarefaction wave

### 17.1.6 Shock speed

The propagating shock solution (17.16) is a weak solution to Burgers' equation only if the speed of propagation is given by (17.17). The same discontinuity propagating at a different speed would not be a weak solution.

The speed of propagation can be determined by conservation. The relation between the shock speed  $s$  and the states  $q_l$  and  $q_r$  is called the **Rankine-Hugoniot jump condition**:

$$f(q_l) - f(q_r) = s(q_l - q_r) . \tag{17.19}$$

For scalar problems this gives simply

$$s = \frac{f(q_l) - f(q_r)}{q_l - q_r} = \frac{[f]}{[q]} , \tag{17.20}$$

where  $[\cdot]$  indicates the jump in some quantity across the discontinuity. Note that any jump is allowed, provided the speed is related via (17.20).

For systems of equations,  $q_l - q_r$  and  $f(q_r) - f(q_l)$  are both vectors while  $s$  is still a scalar. Now we cannot always solve for  $s$  to make (17.19) hold. Instead, only certain jumps  $q_l - q_r$  are allowed, namely those for which the vectors  $f(q_l) - f(q_r)$  and  $q_l - q_r$  are linearly dependent.

For a linear system with  $f(q) = Aq$ , (17.19) gives

$$A(q_l - q_r) = s(q_l - q_r) , \tag{17.21}$$

*i.e.*,  $q_l - q_r$  must be an eigenvector of the matrix  $A$  and  $s$  is the associated eigenvalue. For a linear system, these eigenvalues are the characteristic speeds of the system. Thus discontinuities can propagate only along characteristics, just as for the scalar advection equation.

## 17.2 Linear hyperbolic systems

We now begin to investigate systems of equations. We start with constant coefficient linear systems. Here we can solve the equations explicitly by transforming to characteristic variables. We will also obtain explicit solutions of the Riemann problem and introduce a “phase space” interpretation that will be very useful in our study of nonlinear systems.

Consider the linear system

$$\begin{aligned} q_t + Aq_x &= 0 , \\ q(x, 0) &= q_0(x) , \end{aligned} \tag{17.22}$$

where  $q : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^m$  and  $A \in \mathbb{R}^{m \times m}$  is a constant matrix. This is a system of conservation laws with the flux function  $f(q) = Aq$ . This system is called **hyperbolic** if  $A$  is diagonalizable with real eigenvalues, so that we can decompose

$$A = R\Lambda R^{-1}, \quad (17.23)$$

where  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$  is a diagonal matrix of eigenvalues and  $R = [r_1 | r_2 | \dots | r_m]$  is the matrix of right eigenvectors. Note that  $AR = R\Lambda$ , i.e.,

$$Ar_p = \lambda_p r_p \quad \text{for } p = 1, 2, \dots, m. \quad (17.24)$$

The system is called **strictly hyperbolic** if the eigenvalues are distinct. We will always make this assumption as well.

### 17.2.1 Characteristic variables

We can solve (17.22) by changing to the “characteristic variables”

$$v = R^{-1}q. \quad (17.25)$$

Multiplying (17.22) by  $R^{-1}$  and using (17.23) gives

$$v_t + \Lambda v_x = 0. \quad (17.26)$$

Since  $\Lambda$  is diagonal, this decouples into  $m$  independent scalar equations. Each of these is a constant coefficient linear advection equation, with solution

$$v_p(x, t) = v_p(x - \lambda_p t, 0). \quad (17.27)$$

Since  $v = R^{-1}q$ , the initial data for  $v_p$  is simply the  $p$ th component of the vector

$$v(x, 0) = R^{-1}q_0(x). \quad (17.28)$$

The solution to the original system is finally recovered via (17.25):

$$q(x, t) = Rv(x, t). \quad (17.29)$$

Note that the value  $v_p(x, t)$  is the coefficient of  $r_p$  in an eigenvector expansion of the vector  $q(x, t)$ , i.e., (17.29) can be written out as

$$q(x, t) = \sum_{i=1}^m v_p(x, t) r_p. \quad (17.30)$$

Combining this with the solutions (17.27) of the decoupled scalar equations gives

$$q(x, t) = \sum_{p=1}^m v_p(x - \lambda_p t, 0) r_p. \quad (17.31)$$

Note that  $q(x, t)$  depends only on the initial data at the  $m$  points  $x - \lambda_p t$ , so the domain of dependence is given by  $\mathcal{D}(\bar{x}, \bar{t}) = \{x = \bar{x} - \lambda_p \bar{t}, p = 1, 2, \dots, m\}$ .

The curves  $x = x_0 + \lambda_p t$  satisfying  $x'(t) = \lambda_p$  are the “characteristics of the  $p$ th family”, or simply “ $p$ -characteristics”. These are straight lines in the case of a constant coefficient system. Note that for a strictly hyperbolic system,  $m$  distinct characteristic curves pass through each point in the  $x$ - $t$  plane. The coefficient  $v_p(x, t)$  of the eigenvector  $r_p$  in the eigenvector expansion (17.30) of  $q(x, t)$  is constant along any  $p$ -characteristic.

### 17.2.2 The Riemann Problem

For the constant coefficient linear system, the Riemann problem can be explicitly solved. We will see shortly that the solution to a nonlinear Riemann problem has a simple structure which is quite similar to the structure of this linear solution, and so it is worthwhile studying the linear case in some detail.

The Riemann problem consists of the equation  $q_t + Aq_x = 0$  together with piecewise constant initial data of the form

$$q(x, 0) = \begin{cases} q_l & x < 0 \\ q_r & x > 0. \end{cases} \tag{17.32}$$

Recall that the general solution to the linear problem is given by (17.31). For the Riemann problem we can simplify the notation if we decompose  $q_l$  and  $q_r$  as

$$q_l = \sum_{p=1}^m v_p^l r_p, \quad q_r = \sum_{p=1}^m v_p^r r_p. \tag{17.33}$$

Then

$$v_p(x, 0) = \begin{cases} v_p^l & x < 0 \\ v_p^r & x > 0 \end{cases} \tag{17.34}$$

and so

$$v_p(x, t) = \begin{cases} v_p^l & \text{if } x - \lambda_p t < 0 \\ v_p^r & \text{if } x - \lambda_p t > 0. \end{cases} \tag{17.35}$$

If we let  $P(x, t)$  be the maximum value of  $p$  for which  $x - \lambda_p t > 0$ , then

$$q(x, t) = \sum_{p=1}^{P(x,t)} v_p^r r_p + \sum_{p=P(x,t)+1}^m v_p^l r_p. \tag{17.36}$$

The determination of  $q(x, t)$  at a given point is illustrated in Figure 17.3. In the case shown,  $v_1 = v_1^r$  while  $v_2 = v_2^l$  and  $v_3 = v_3^l$ . The solution at the point illustrated is thus

$$q(x, t) = v_1^r r_1 + v_2^l r_2 + v_3^l r_3. \tag{17.37}$$

Note that the solution is the same at any point in the wedge between the  $x' = \lambda_1$  and  $x' = \lambda_2$  characteristics. As we cross the  $p$ th characteristic, the value of  $x - \lambda_p t$  passes through 0 and the corresponding  $v_p$  jumps from  $v_p^l$  to  $v_p^r$ . The other coefficients  $v_i$  ( $i \neq j$ ) remain constant.

The solution is constant in each of the wedges as shown in Figure 17.4. Across the  $p$ th characteristic the solution jumps with the jump given by

$$[q] = (v_p^r - v_p^l)r_p. \tag{17.38}$$

Note that these jumps satisfy the Rankine-Hugoniot conditions (17.19), since  $f(q) = Aq$  leads to

$$[f] = A[q] = (v_p^r - v_p^l)Ar_p = \lambda_p[q]$$

and  $\lambda_p$  is precisely the speed of propagation of this jump. The solution  $q(x, t)$  in (17.36) can alternatively be written in terms of these jumps as

$$q(x, t) = q_l + \sum_{\lambda_p < x/t} (v_p^r - v_p^l)r_p \tag{17.39}$$

$$= q_r - \sum_{\lambda_p > x/t} (v_p^r - v_p^l)r_p. \tag{17.40}$$

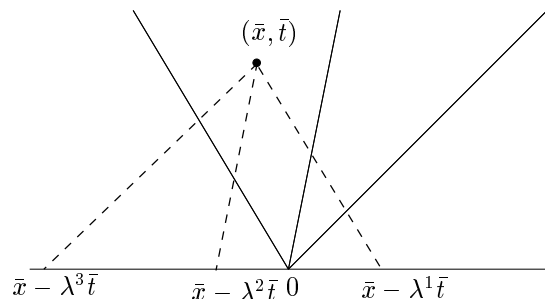


Figure 17.3: Construction of the solution to the Riemann problem at  $(\bar{x}, \bar{t})$ . We trace back along the  $p$ 'th characteristic to determine the coefficient of  $r_p$  depending on whether this lies in the left state or right state initially

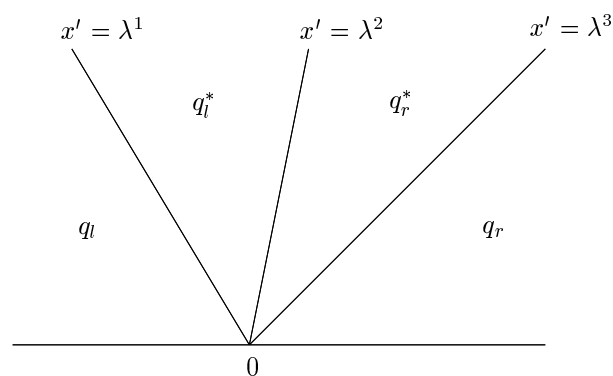


Figure 17.4: Values of solution  $q$  in each wedge of  $x-t$  plane:

$$\begin{aligned} q_l &= v_1^l r_1 + v_2^l r_2 + v_3^l r_3, & q_r^* &= v_1^r r_1 + v_2^r r_2 + v_3^r r_3, \\ q_l^* &= v_1^r r_1 + v_2^l r_2 + v_3^l r_3, & q_r &= v_1^r r_1 + v_2^r r_2 + v_3^r r_3. \end{aligned}$$

Note that the jump across each discontinuity in the solution is an eigenvector of  $A$

It might happen that the initial jump  $q_r - q_l$  is already an eigenvector of  $A$ , if  $q_r - q_l = (v_i^r - v_i^l)r_i$  for some  $i$ . In this case  $v_p^l = v_p^r$  for  $p \neq i$ . Then this discontinuity simply propagates with speed  $\lambda_i$ , and the other characteristics carry jumps of zero strength.

In general this is not the case, however, and the jump  $q_r - q_l$  cannot propagate as a single discontinuity with any speed without violating the Rankine-Hugoniot condition. We can view “solving the Riemann problem” as finding a way to split up the jump  $q_r - q_l$  into a sum of jumps,

$$q_r - q_l = (v_1^r - v_1^l)r_1 + \cdots + (v_m^r - v_m^l)r_m \equiv \alpha_1 r_1 + \cdots + \alpha_m r_m, \quad (17.41)$$

each of which *can* propagate at an appropriate speed  $\lambda_i$  with the Rankine-Hugoniot condition satisfied.

For nonlinear systems we solve the Riemann problem in much the same way: The jump  $q_r - q_l$  will usually not have the property that  $[f]$  is a scalar multiple of  $[g]$ , but we can attempt to find a way to split this jump up as a sum of jumps, across each of which this property does hold. (Although life is complicated by the fact that we may need to introduce rarefaction waves as well as shocks.) In studying the solution of the Riemann problem, the jump in the  $p$ th family, traveling at constant speed  $\lambda_p$ , is often called the  $p$ -wave.



Figure 17.5: (a) The Hugoniot locus of the state  $q_l$  consists of all states that differ from  $q_l$  by a scalar multiple of  $r_1$  or  $r_2$ . (b) Solution to the Riemann problem in the  $x-t$  plane

### 17.2.3 The phase plane

For systems of two equations, it is illuminating to view this splitting in the phase plane. This is simply the  $q_1-q_2$  plane, where  $q = (q_1, q_2)$ . Each vector  $q(x, t)$  is represented by a point in this plane. In particular,  $q_l$  and  $q_r$  are points in this plane and a discontinuity with left and right states  $q_l$  and  $q_r$  can propagate as a single discontinuity only if  $q_r - q_l$  is an eigenvector of  $A$ , which means that the line segment from  $q_l$  to  $q_r$  must be parallel to the eigenvector  $r_1$  or  $r_2$ . Figure 17.5 shows an example. For the state  $q_l$  illustrated there, the jump from  $q_l$  to  $q_r$  can propagate as a single discontinuity if and only if  $q_r$  lies on one of the two lines drawn through  $q_l$  in the direction  $r_1$  and  $r_2$ . These lines give the locus of all points that can be connected to  $q_l$  by a 1-wave or a 2-wave. This set of states is called the **Hugoniot locus**. We will see that there is a direct generalization of this to nonlinear systems in the next chapter.

Similarly, there is a Hugoniot locus through any point  $q_r$  that gives the set of all points  $q_l$  that can be connected to  $q_r$  by an elementary  $p$ -wave. These curves are again in the directions  $r_1$  and  $r_2$ .

For a general Riemann problem with arbitrary  $q_l$  and  $q_r$ , the solution consists of two discontinuities travelling with speeds  $\lambda_1$  and  $\lambda_2$ , with a new constant state in between that we will call  $q^*$ . By the discussion above,

$$q^* = v_1^r r_1 + v_2^l r_2, \tag{17.42}$$

so that  $q^* - q_l = \alpha_1 r_1$  and  $q_r - q^* = \alpha_2 r_2$ . The location of  $q^*$  in the phase plane must be where the 1-wave locus through  $q_l$  intersects the 2-wave locus through  $q_r$ . This is illustrated in Figure 17.6a.

Note that if we interchange  $q_r$  and  $q_l$  in this picture, the location of  $q^*$  changes as illustrated in Figure 17.6(b). In each case we travel from  $q_l$  to  $q_r$  by first going in the direction  $r_1$  and then in the direction  $r_2$ . This is required by the fact that  $\lambda_1 < \lambda_2$  since clearly the jump between  $q_l$  and  $q^*$  must travel slower than the jump between  $q^*$  and  $q_r$  if we are to obtain a single-valued solution.

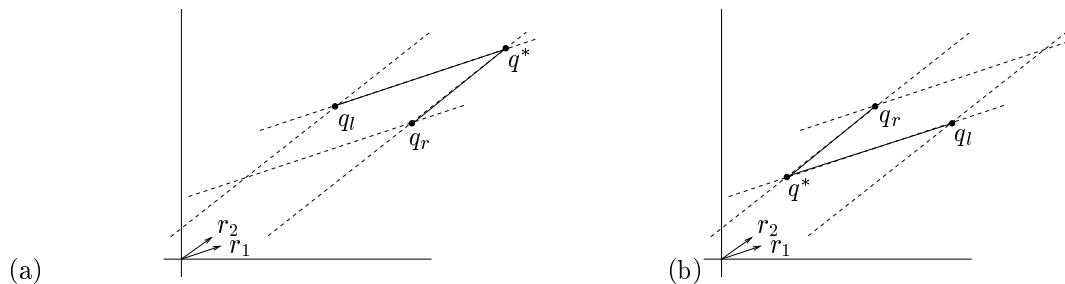


Figure 17.6: The new state  $q^*$  arising in the solution to the Riemann problem for two different choices of  $q_l$  and  $q_r$ . In each case the jump from  $q_l$  to  $q^*$  lies in the direction of the eigenvalue  $r_1$  corresponding to the slower speed, while the jump from  $q^*$  to  $q_r$  lies in the direction of the eigenvalue  $r_2$



For systems with more than two equations, the same interpretation is possible but becomes harder to draw since the phase space is now  $m$  dimensional. Since the  $m$  eigenvectors  $r_p$  are linearly independent, we can decompose any jump  $q_r - q_l$  into the sum of jumps in these directions, obtaining a piecewise linear path from  $q_l$  to  $q_r$  in  $m$ -dimensional space.

### 17.3 Finite volume methods

Many high-resolution methods for shock capturing are based on solving Riemann problems between states in neighboring grid cells. In this chapter we will develop one particular set of methods of this type. The development of such methods has a long and rich history, and numerous related methods can be found in the literature. Books such as [ATP84], [Fle88], [GR96], [Hir88], [?], [LeV90], [OB87], [PT83], [Sod85], [Tor97] contain descriptions of these methods and pointers to the literature.

Rather than viewing  $Q_i^n$  as an approximation to the single value  $q(x_i, t_n)$ , we will now view it as approximating the average value of  $q$  over an interval of length  $h = \Delta x = (b - a)/N$ . We will split the physical domain  $[a, b]$  into  $N$  intervals denoted by

$$C_i = [x_i, x_{i+1}],$$

where now  $x_i = a + (i - 1)h$ . The value  $Q_i^n$  will approximate the average value over the  $i$ 'th interval at time  $t_n$ :

$$Q_i^n \approx \frac{1}{h} \int_{x_i}^{x_{i+1}} q(x, t_n) dx \equiv \frac{1}{h} \int_{C_i} q(x, t_n) dx. \quad (17.43)$$

Notationally it might be better to denote the endpoints of the  $i$ 'th interval by  $x_{i-1/2}$  and  $x_{i+1/2}$ , which would be more symmetric and remind us that  $Q_i^n$  is an approximation to the average value between these points. However, the formulas are less cluttered if we stick to integer subscripts.

If  $q(x, t)$  is a smooth function, then the integral in (17.43) agrees with the value of  $q$  at the midpoint of the interval to  $O(h^2)$ . By working with cell averages, however, it is easier to use important properties of the conservation law in deriving numerical methods. In particular, we can insure that the numerical method is **conservative** in a way that mimics the true solution, and this is extremely important in accurately calculating shock waves. This is because  $h \sum_{i=1}^N Q_i^n$  approximates the integral of  $q$  over the entire interval  $[a, b]$ , and if we use a method that is in *conservation form* (as described below), then this discrete sum will change only due to fluxes at the boundaries  $x = a$  and  $x = b$ . The total mass within the computational domain will be preserved, or at least will vary correctly provided the boundary conditions are properly imposed.

The integral form of the conservation law (17.12), when applied to one grid cell over a single time step, gives

$$\begin{aligned} \int_{C_i} q(x, t_{n+1}) dx - \int_{C_i} q(x, t_n) dx &= \int_{t_n}^{t_{n+1}} f(q(x_i, t)) dt \\ &\quad - \int_{t_n}^{t_{n+1}} f(q(x_{i+1}, t)) dt. \end{aligned}$$

Rearranging this and dividing by  $h$  gives

$$\begin{aligned} \frac{1}{h} \int_{C_i} q(x, t_{n+1}) dx &= \frac{1}{h} \int_{C_i} q(x, t_n) dx \\ &\quad - \frac{1}{h} \left[ \int_{t_n}^{t_{n+1}} f(q(x_i, t)) dt - \int_{t_n}^{t_{n+1}} f(q(x_{i+1}, t)) dt \right]. \end{aligned} \quad (17.44)$$

This tells us exactly how the cell average of  $q$  from (17.43) should be updated in one time step. In general, however, we cannot evaluate the time integrals on the right-hand side of (17.44) exactly since

$q(x_i, t)$  varies with time along each edge of the cell, and we don't have the exact solution to work with. But this does suggest that we should develop numerical methods in the **flux-differencing form**

$$Q_i^{n+1} = Q_i^n - \frac{k}{h}(F_{i+1}^n - F_i^n), \quad (17.45)$$

where  $F_i^n$  is some approximation to the average flux along  $x = x_i$ :

$$F_i^n \approx \frac{1}{k} \int_{t_n}^{t_{n+1}} f(q(x_i, t)) dt. \quad (17.46)$$

If we can approximate this average flux based on the values  $Q^n$ , then we will have a fully-discrete method.

Since information propagates with finite speed, it is reasonable to first suppose that we can obtain  $F_i^n$  based only on the values  $Q_{i-1}^n$  and  $Q_i^n$ , the cell averages on either side of this interface. Then we might use a formula of the form

$$F_i^n = F(Q_{i-1}^n, Q_i^n),$$

where  $F$  is some **numerical flux function**. The method (17.45) then becomes

$$Q_i^{n+1} = Q_i^n - \frac{k}{h}(F(Q_i^n, Q_{i+1}^n) - F(Q_{i-1}^n, Q_i^n)). \quad (17.47)$$

The specific method obtained depends on how we choose the formula  $F$ , but in general any method of this type is an explicit method with a 3-point stencil. Moreover, it is said to be in **conservation form**, since it mimics the property (17.44) of the exact solution. Note that if we sum  $hQ_i^{n+1}$  from (17.45) over any set of cells we obtain

$$h \sum_{i=I}^J Q_i^{n+1} = h \sum_{i=I}^J Q_i^n - \frac{k}{h} (F_{J+1}^n - F_I^n). \quad (17.48)$$

The sum of the flux differences cancels out except for the fluxes at the extreme edges. Over the full domain we have exact conservation except for fluxes at the boundaries. (Boundary conditions are discussed at the end of this chapter.)

Note that (17.47) can be viewed as a direct finite difference approximation to the conservation law  $q_t + f(q)_x = 0$ , since rearranging it gives

$$\frac{Q_i^{n+1} - Q_i^n}{k} + \frac{F(Q_{i+1}^n, Q_i^n) - F(Q_i^n, Q_{i-1}^n)}{h} = 0.$$

Many methods can be equally well viewed as finite difference approximations to this equation or as finite volume methods. In obtaining a method in conservation form, the above discussion suggests that we should always discretize the conservation law in this form, rather than in the quasi-linear form  $q_t + f'(q)q_x = 0$ , for example.

## 17.4 Importance of conservation form — incorrect shock speeds

Using methods in conservation form is particularly important when solving problems with shocks or other discontinuities in the solution, as a nonconservative method may give a numerical solution that looks reasonable but is entirely wrong. For example, Burgers' equation

$$q_t + \left(\frac{1}{2}q^2\right)_x = 0 \quad (17.49)$$

can be discretized by the upwind conservative method

$$Q_i^{n+1} = Q_i^n - \frac{k}{h} \left( \frac{1}{2}(Q_i^n)^2 - \frac{1}{2}(Q_{i-1}^n)^2 \right).$$

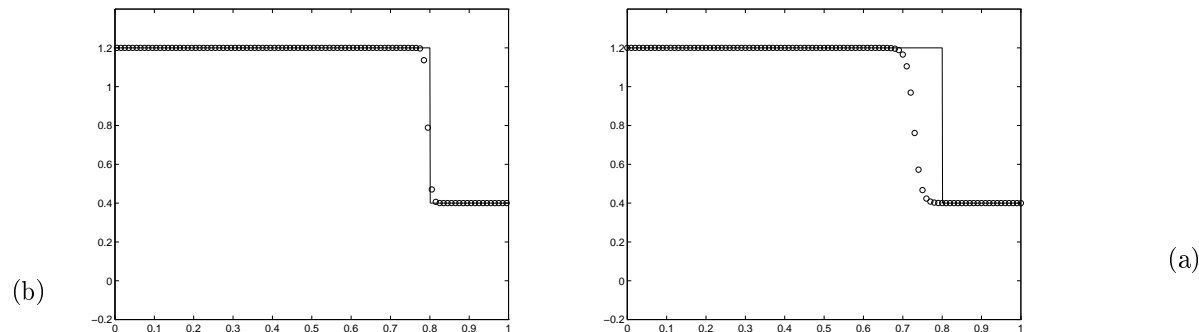


Figure 17.7: (a) True and computed solutions to Burgers' equation using a conservative upwind method. (b) True and computed solutions to Burgers' equation using a nonconservative upwind method

Figure 17.7(a) shows the resulting numerical solution for a shock wave resulting between states 1.2 and 0.4 propagating at speed 0.8. The numerical solution is slightly smeared about the correct location.

On the other hand if we discretize the quasilinear form of Burgers' equation

$$q_t + qq_x = 0 \quad (17.50)$$

using the nonconservative upwind method

$$Q_i^{n+1} = Q_i^n - \frac{k}{h} Q_i^n (Q_i^n - Q_{i-1}^n)$$

we obtain the results seen in Figure 17.7(b). The shock is moving at the wrong speed! This happens because the equations (17.49) and (17.50) are equivalent for smooth solutions but **not** for problems with shock waves. This example is discussed in more detail in [LeV90].

## 17.5 Numerical flux functions

Given that we want to use a method in conservation form, how should we define  $F(q_l, q_r)$ , the average flux at a point based on data  $q_l$  and  $q_r$  to the left and right of this point? A first attempt might be the simple average

$$F(q_l, q_r) = \frac{1}{2}(f(q_l) + f(q_r)) .$$

Using this in (17.47) would give

$$Q_i^{n+1} = Q_i^n - \frac{k}{2h}(f(Q_{i+1}^n) - f(Q_{i-1}^n)) .$$

In general, however, this method turns out to be *unconditionally unstable* for any value of  $k/h$ .

If we instead use the modified flux

$$F(q_l, q_r) = \frac{1}{2}(f(q_l) + f(q_r)) - \frac{h}{2k}(q_r - q_l) , \quad (17.51)$$

then we obtain the **Lax-Friedrichs method**,

$$Q_i^{n+1} = \frac{1}{2}(Q_{i-1}^n + Q_{i+1}^n) - \frac{k}{2h}(f(Q_{i+1}^n) - f(Q_{i-1}^n)) . \quad (17.52)$$

Note that the additional term we have added in (17.51) is a *diffusive flux* based on an approximation to  $\frac{h^2}{2k}q_x$ , and hence this modification amounts to adding some *artificial viscosity* to the centered flux formula.

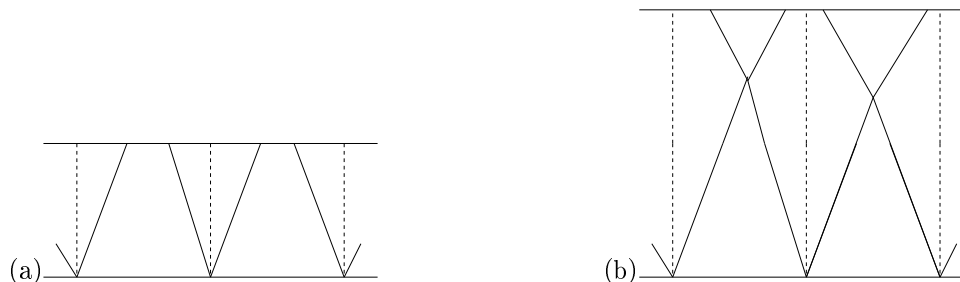


Figure 17.8: Solving the Riemann problems at each interface for Godunov’s method. (a) With Courant number less than  $1/2$  there is no interaction of waves. (b) With Courant number less than  $1$  the interacting waves do not reach the cell interfaces, so the fluxes are still constant in time

### 17.6 Godunov’s method

Many of the methods we will explore in detail are based on solving the *Riemann problem* between the states  $q_l$  and  $q_r$  in order to define the numerical flux  $F(q_l, q_r)$ . To see how this comes about, it is useful to view the data  $Q^n$  at time  $t_n$  as defining a piecewise constant function  $\tilde{q}^n(x, t_n)$  which has the value  $Q_i^n$  for all  $x$  in the interval  $C_i$ . Suppose we could solve the conservation law *exactly* over the time interval  $[t_n, t_{n+1}]$  with initial data  $\tilde{q}^n(x, t_n)$ . Call the resulting function  $\tilde{q}^n(x, t)$  for  $t_n \leq t \leq t_{n+1}$ . Then we might consider defining the numerical flux  $F_i^n$  in (17.45) by

$$F_i^n = \frac{1}{k} \int_{t_n}^{t_{n+1}} f(\tilde{q}^n(x_i, t)) dt . \tag{17.53}$$

This integral is trivial to compute compared to the integral (17.46), at least provided the time step  $k$  is small enough, because of the fact that with piecewise constant initial data we can find the exact solution easily by simply piecing together the solutions to each Riemann problem defined by the jump at each interface. Figure 17.8 illustrates this for the case of a linear hyperbolic system of 2 equations.

The crucial fact now is that the solution to the Riemann problem at  $x_i$  is a similarity solution, which is constant along each ray  $(x - x_i)/t = \text{constant}$ . In general, let  $q^*(q_l, q_r)$  denote the exact solution to the Riemann problem along the ray  $x/t = 0$ , obtained when we use data

$$q(x, 0) = \begin{cases} q_l & \text{if } x < 0 \\ q_r & \text{if } x > 0 . \end{cases}$$

Then we have

$$\tilde{q}^n(x_i, t) \equiv q^*(Q_{i-1}^n, Q_i^n) \tag{17.54}$$

for all  $t \in [t_n, t_{n+1}]$ , provided that the time step is small enough that waves from the Riemann problems do not travel farther than distance  $h$  in this time step. If this condition is violated, then the value along  $x = x_i$  may change after waves from neighboring Riemann problems pass this point. For a linear system, the maximum wave speed is  $\max_p |\lambda^p|$ , where  $\lambda^p$  are the eigenvalues of  $A$ , so this condition requires that

$$k \max_p |\lambda^p| \leq h .$$

We recognize this as being simply the CFL condition for a 3-point method, a condition which we know must be satisfied anyway for stability.

The method obtained by the procedure outlined above is known as **Godunov’s method**, and was introduced in [God59] as an approach to solving the Euler equations of gas dynamics in the presence of shock waves.

For the simplest case of scalar advection, solving the Riemann problem between states  $q_l$  and  $q_r$  gives

$$q^*(q_l, q_r) = \begin{cases} q_l & \text{if } u > 0 \\ q_r & \text{if } u < 0. \end{cases}$$

If  $u = 0$  then  $q^*$  is not well defined, as the discontinuity is stationary along the ray  $x/t = 0$ . This is no cause for concern, however, since all we really require is the flux value  $f(q^*) = uq^*$ , and if  $u = 0$  then  $f(q^*) = 0$  regardless of how we define  $q^*$ . So we obtain the numerical flux

$$F(q_l, q_r) = uq^*(q_l, q_r) = \begin{cases} uq_l & \text{if } u \leq 0 \\ uq_r & \text{if } u \geq 0. \end{cases}$$

This can also be written in the compact form

$$F(q_l, q_r) = u^+ q_l + u^- q_r \quad (17.55)$$

using the notation

$$u^+ = \max(u, 0), \quad u^- = \min(u, 0). \quad (17.56)$$

Using this in the conservative method (17.47) gives the **upwind method**. Note that solving the Riemann problem at the interface gives a flux that is defined by the value of  $Q^n$  on the upwind side of the interface, so that the method reduces to one-sided differencing in the proper direction. The method takes the form

$$Q_i^{n+1} = \begin{cases} Q_i^n - \frac{k}{h} u (Q_i^n - Q_{i-1}^n) & \text{if } u > 0 \\ Q_i^n - \frac{k}{h} u (Q_{i+1}^n - Q_i^n) & \text{if } u < 0 \end{cases} \quad (17.57)$$

as introduced in Section 14.3.

This method is easily generalized to nonlinear systems if we can solve the nonlinear Riemann problem at each cell interface, and this method gives the natural generalization of the first order upwind method to general systems of conservation laws.

Recall that  $Q_i^n$  represents an approximation to the cell average of  $q(x, t_n)$  over cell  $C_i$ ,

$$Q_i^n \approx \frac{1}{h} \int_{x_i}^{x_{i+1}} q(x, t_n) dx,$$

and the idea is to use the piecewise constant function defined by these cell values as initial data  $\tilde{q}^n(x, t_n)$  for the conservation law. Solving over time  $k$  with this data gives a function  $\tilde{q}^n(x, t_{n+1})$  that is then averaged over each cell to obtain

$$Q_i^n = \frac{1}{h} \int_{x_i}^{x_{i+1}} \tilde{q}^n(x, t_{n+1}) dx, \quad (17.58)$$

If the time step  $k$  is sufficiently small, then the exact solution  $\tilde{q}^n(x, t)$  can be determined by piecing together the solutions to the Riemann problem arising from each cell interface, as indicated in Figure 17.8(a).

Recall from Section 17.6 that we do not need to perform the integration in (17.58) explicitly, which might be difficult since  $\tilde{q}^n(x, t_{n+1})$  may be very complicated as a function of  $x$ . Instead, we can use the fact that  $\tilde{q}^n(x_i, t)$  is constant in time along each cell interface so that the integral (17.53) can be computed exactly. Hence the cell average is updated by (17.47) with

$$F(q_l, q_r) = f(q^*(q_l, q_r)), \quad (17.59)$$

where  $q^*(q_l, q_r)$  is the solution to the Riemann problem between  $q_l$  and  $q_r$ , evaluated along  $x/t = 0$ .

In Figure 17.8(a) the time step was taken to be small enough that there was no interaction of waves from neighboring Riemann problems. This would be necessary if we wanted to construct the solution

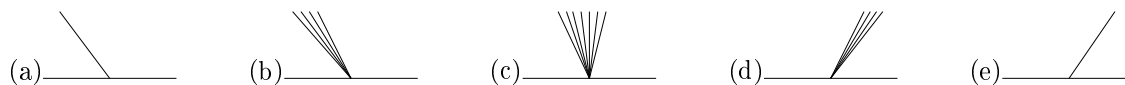


Figure 17.9: Five possible configurations for the solution to a scalar conservation law, in the  $x-t$  plane. (a) Left-going shock,  $q^* = q_r$ . (b) Left-going rarefaction,  $q^* = q_r$ . (c) Transonic rarefaction,  $q^* = q_0$ . (d) Right-going rarefaction,  $q^* = q_l$ . (e) Right-going shock,  $q^* = q_l$ .

at  $\bar{q}^n(x, t_{n+1})$  in order to explicitly calculate the cell averages (17.58). However, in order to use the flux formula (17.59) it is only necessary that  $\bar{q}^n(x_i, t)$  remain constant in time over the entire time step, which allows a time step roughly twice as large, as indicated in Figure 17.8(b). If  $s_{\max}$  represents the largest wave speed that is encountered then on a uniform grid with the cell interfaces distance  $h$  apart, we must require

$$\frac{ks_{\max}}{h} \leq 1 \quad (17.60)$$

in order to insure that the formula (17.59) is valid. Note that this is precisely the CFL condition required for stability of this 3-point method, as discussed in Chapter 15. In general  $s_{\max}k/h$  is called the **Courant number**. Figure 17.8(a) shows a case where the Courant number is less than  $1/2$  while Figure 17.8(b) shows the Courant number close to 1. Note that for a linear system of equations,  $s_{\max} = \max_p |\lambda^p|$  and this agrees with our previous definition of the Courant number.

### 17.6.1 Godunov's method on scalar equations

On a convex scalar equation with  $f'(q)$  an increasing function of  $q$ , the solution to the Riemann problem between  $q_l$  and  $q_r$  is either a shock traveling at speed  $s = [f]/[q]$  (if  $q_l > q_r$ ) or a rarefaction wave (if  $q_l < q_r$ ) bounded by  $x/t = f'(q_l)$  on the left and  $x/t = f'(q_r)$  on the right. Five possible configurations in the  $x-t$  plane are shown in Figure 17.9. In most cases the solution  $q^*$  along  $x/t = 0$  will be either  $q_r$  (if the solution is a shock or rarefaction wave moving entirely to the left, Figure 17.9(a) or (b)), or  $q_l$  (if the solution is a shock or rarefaction wave moving entirely to the right, Figure 17.9(d) or (e)).

The only case where  $q^*$  has a different value than  $q_l$  or  $q_r$  is if  $q_l < q_0 < q_r$ , where  $q_0$  is the value for which  $f'(q_0) = 0$ . This is called the **stagnation point** since the value  $q_0$  propagates with speed 0. It is also called the **sonic point** since in gas dynamics the eigenvalue  $u \pm c$  takes the value 0 only when the fluid speed is equal to the sound speed. The solution to the Riemann problem in this case, shown in Figure 17.9(c), consists of a rarefaction wave that is partly left-going and partly right-going. This is called a **transonic rarefaction** since in gas dynamics the fluid is accelerated from a subsonic velocity to a supersonic velocity through such a rarefaction. In a transonic rarefaction the value along  $x/t = 0$  is simply  $q_0$ .

We thus see that the Godunov flux function for a convex scalar conservation law is

$$F(q_l, q_r) = \begin{cases} f(q_l) & \text{if } q_l > q_0 \text{ and } s > 0 \\ f(q_r) & \text{if } q_r < q_0 \text{ and } s < 0 \\ f(q_0) & \text{if } q_l < q_0 < q_r. \end{cases} \quad (17.61)$$

## 17.7 High-resolution methods

Godunov's method is at best first order accurate on smooth solutions and generally gives very smeared approximations to shock waves or other discontinuities. In this section we will see how this method can be extended to a method that gives second-order accuracy on smooth flow, but which avoids nonphysical oscillations near discontinuities. The key is to use a better representation of the solution, say piecewise linear instead of the piecewise constant representation used in Godunov's method, but to form this reconstruction carefully by paying attention to how the data behaves nearby. In smooth regions the finite-difference approximation to the slope can be used to obtain better accuracy, but near

a discontinuity the “slope” computed by subtracting two values of  $Q$  and dividing by  $h$  may be huge and meaningless. Using it blindly in a difference approximation will introduce oscillations into the numerical solution.

One particular method will be developed here in a framework that can be interpreted as a correction phase following the solution of Riemann problems and construction of Godunov fluxes. Many other approaches can be found in the literature and a couple of these are briefly described in Section ??.

To introduce these ideas we will first consider the scalar advection equation

$$q_t + uq_x = 0$$

with  $u > 0$ , in which case Godunov’s method reduces to the simple first-order upwind method

$$Q_i^{n+1} = Q_i^n - \frac{k}{h}u(Q_i^n - Q_{i-1}^n).$$

After developing a high-resolution version of this method, the ideas can be extended to systems of equations and nonlinear problems.

### 17.7.1 Reconstruct–Solve–Average

There is another interpretation of Godunov’s method that will be useful in developing higher-order accurate methods of this type. Recall that  $\tilde{q}^n(x, t_n)$  denotes the piecewise constant function with value  $Q_i^n$  in cell  $\mathcal{C}_i$ . We defined the numerical flux of Godunov’s method by advancing the solution with this data to obtain the interface value  $\tilde{q}^n(x, t)$  over the time interval  $[t_n, t_{n+1}]$ . Another way to describe Godunov’s method is to take the advanced solution  $\tilde{q}^n(x, t_{n+1})$  at the end of the time step, and average this function over grid cell  $\mathcal{C}_i$  to obtain

$$Q_i^{n+1} \equiv \frac{1}{h} \int_{\mathcal{C}_i} \tilde{q}^n(x, t_{n+1}) dx.$$

It follows from the integral form of the conservation law that this gives exactly the same value as the flux-differencing method, though implementing it in this form would be more difficult since  $\tilde{q}^n(x, t_{n+1})$  is not constant over  $\mathcal{C}_i$ , and this integral would be difficult to evaluate directly in general. The beauty of the flux-differencing approach is that we do not need to evaluate this integral, but can find it by integrating the flux function, which is constant on the time interval of integration.

But in generalizing Godunov’s method to higher-order methods, it is useful to consider what would happen if we took this approach with a different choice of  $\tilde{q}^n(x, t_n)$ , that better approximates a smooth function. We can think of  $\tilde{q}^n(x, t_n)$  as a *reconstruction* of a function from the discrete values  $Q_i^n$ , the cell averages of the function. Instead of a piecewise constant function we might reconstruct a piecewise linear function or some other function  $\tilde{q}^n(x, t_n)$ . We can then generalize Godunov’s method to an algorithm that takes the following general form in each time step:

#### Algorithm RSA (Reconstruct-Solve-Average):

1. **Reconstruct** a function  $\tilde{q}^n(x, t_n)$  defined for all  $x$  from the cell averages  $Q_i^n$ .
2. **Solve** the hyperbolic equation exactly (or approximately) with this initial data to obtain  $\tilde{q}^n(x, t_{n+1})$  a time  $\Delta t$  later.
3. **Average** this function over each grid cell to obtain

$$Q_i^{n+1} = \frac{1}{h} \int_{\mathcal{C}_i} \tilde{q}^n(x, t_{n+1}) dx.$$

With a piecewise constant reconstruction we can solve the problem in Step 2 exactly, giving Godunov’s method. For linear systems we can solve this problem exactly even with more complicated initial data, such as the piecewise linear function considered in the next section. For nonlinear problems we may not be able to solve the problem in Step 2 exactly, but we will still be able to improve the accuracy by using an approximate solution together with piecewise linear data.

### 17.7.2 Piecewise linear reconstruction

To achieve better than first-order accuracy, we must use a better reconstruction than a piecewise constant function. From the cell averages  $Q_i^n$  we can construct a piecewise linear function of the form

$$\tilde{q}^n(x, t_n) = Q_i^n + \sigma_i^n(x - \bar{x}_i) \quad \text{for } x_i \leq x < x_{i+1}, \quad (17.62)$$

where

$$\bar{x}_i = \frac{1}{2}(x_i + x_{i+1}) = x_i + \frac{1}{2}h \quad (17.63)$$

is the center of the  $i$ 'th grid cell and  $\sigma_i^n$  is the slope on the  $i$ 'th cell. The linear function defined by (17.62) on the  $i$ 'th cell is defined in such a way that its value at the cell center  $\bar{x}_i$  is  $Q_i^n$ . More importantly, the average value of  $\tilde{q}^n(x, t_n)$  over cell  $C_i$  is  $Q_i^n$  (regardless of the slope  $\sigma_i^n$ ), so that the reconstructed function has the cell average  $Q_i^n$ . This is crucial in developing conservative methods for conservation laws. Note that Steps 2 and 3 are conservative in general, and so Algorithm RSA is conservative provided we use a *conservative reconstruction* in Step 1, as we have in (17.62). Later we will see how to write such methods in the standard conservation form (17.45).

For the scalar advection equation  $q_t + uq_x = 0$ , we can easily solve the equation with this data, and compute the new cell averages as required in Step 3 of Algorithm RSA. We have

$$\tilde{q}^n(x, t_{n+1}) = \tilde{q}^n(x - uk, t_n).$$

Until further notice we will assume that  $u > 0$  and present the formulas for this particular case. The corresponding formulas for  $u < 0$  should be easy to derive, and in Section 17.7.7 we will see a better way to formulate the methods in the general case.

Suppose also that  $|uk/h| \leq 1$ . Then it is straightforward to compute (see also Section 17.7.7) that

$$\begin{aligned} Q_i^{n+1} &= \frac{uk}{h} \left( Q_{i-1}^n + \frac{1}{2}(h - uk)\sigma_{i-1}^n \right) + \left( 1 - \frac{uk}{h} \right) \left( Q_i^n - \frac{1}{2}uk\sigma_i^n \right) \\ &= Q_i^n - \frac{uk}{h}(Q_i^n - Q_{i-1}^n) - \frac{1}{2} \frac{uk}{h} (h - uk)(\sigma_i^n - \sigma_{i-1}^n). \end{aligned} \quad (17.64)$$

### 17.7.3 Choice of slopes

Choosing  $\sigma_i^n \equiv 0$  gives Godunov's method (the upwind method for the advection equation). To obtain a second-order accurate method we want to choose nonzero slopes in such a way that  $\sigma_i^n$  approximates the derivative  $q_x$  over the  $i$ 'th grid cell. Three obvious possibilities are:

$$\text{Centered slope: } \sigma_i^n = \frac{Q_{i+1}^n - Q_{i-1}^n}{2h} \quad (\text{Fromm}), \quad (17.65)$$

$$\text{Upwind slope: } \sigma_i^n = \frac{Q_i^n - Q_{i-1}^n}{h} \quad (\text{Beam-Warming}), \quad (17.66)$$

$$\text{Downwind slope: } \sigma_i^n = \frac{Q_{i+1}^n - Q_i^n}{h} \quad (\text{Lax-Wendroff}). \quad (17.67)$$

The centered slope might seem like the most natural choice to obtain second order accuracy, but in fact all three choices give the same formal order of accuracy, and it is the other two choices that give methods we have already derived in other ways. Only the downwind slope results in a centered 3-point method, and this choice gives the **Lax-Wendroff** method (14.16). The upwind slope gives a fully-upwinded 3-point method, which is simply **Beam-Warming**.

The centered slope may seem the most symmetric choice at first glance, but due to the fact that the reconstructed function is then advected in the positive direction, the final updating formula turns out to be a non-symmetric 4-point formula, which is known as **Fromm's method**.





Figure 17.10: (a) Grid values  $Q^n$  and reconstructed  $\tilde{q}^n(\cdot, t_n)$  using Lax-Wendroff slopes. (b) After advection with  $ku = h/2$ . The dots show the new cell averages  $Q^{n+1}$ . Note the overshoot

### 17.7.4 Oscillations

Second-order methods such as Lax-Wendroff (or Beam-Warming or Fromm's method) give oscillatory approximations to discontinuous solutions. This can be easily understood using the interpretation of Algorithm RSA.

Consider the Lax-Wendroff method, for example, applied to piecewise constant data with values

$$Q_i^n = \begin{cases} 1 & \text{if } i \leq J \\ 0 & \text{if } i > J. \end{cases}$$

Choosing slopes in each grid cell based on the Lax-Wendroff prescription (17.67) gives the piecewise linear function shown in Figure 17.10(a). The slope  $\sigma_i^n$  is nonzero only for  $i = J$ .

The function  $\tilde{q}^n(x, t_n)$  has an *overshoot* with a maximum value of  $3/2$  regardless of  $h$ . When we advect this profile a distance  $uk$ , and then compute the average over the  $J$ 'th cell, we will get a value that is greater than 1 for any  $k$  with  $0 < uk < h$ . The worst case is when  $uk = h/2$ , in which case  $\tilde{q}^n(x, t_{n+1})$  is shown in Figure 17.10(b) and  $Q_J^{n+1} = 9/8$ . In the next time step this overshoot will be accentuated, while in cell  $J - 1$  we will now have a positive slope, leading to a value  $Q_{J-1}^{n+1}$  that is less than 1. This oscillation then grows with time.

The slopes proposed in the previous section were based on the assumption that the solution is smooth. Near a discontinuity there is no reason to believe that introducing this slope will improve the accuracy. On the contrary, if one of our goals is to avoid nonphysical oscillations, then in the above example we must set the slope to zero in the  $J$ 'th cell. Any  $\sigma_J^n < 0$  will lead to  $Q_J^{n+1} > 1$ , while a positive slope wouldn't make much sense. On the other hand we don't want to set all slopes to zero all the time, or we simply have the first-order upwind method. Where the solution is smooth we want second order accuracy. Moreover, we will see below that even near a discontinuity, once the solution is somewhat smeared out over more than one cell, introducing nonzero slopes can help keep the solution from smearing out too far, and hence will significantly increase the resolution and keep discontinuities fairly sharp, as long as care is taken to avoid oscillations.

This suggests that we must pay attention to *how the solution is behaving* near the  $i$ 'th cell in choosing our formula for  $\sigma_i^n$ . (And hence the resulting updating formula will be *nonlinear* even for the linear advection equation!). Where the solution is smooth we want to choose something like the Lax-Wendroff slope. Near a discontinuity we may want to "limit" this slope, using a value that is smaller in magnitude in order to avoid oscillations. Methods based on this idea are known as **slope-limiter** methods. This approach was introduced by van Leer in a series of papers [van73] through [van79], where he developed the **MUSCL** scheme for nonlinear conservation laws (Monotonic Upstream-centered Scheme for Conservation Laws). The same idea in the context of **flux limiting**, reducing the magnitude of the numerical flux to avoid oscillations, was introduced in the **flux-corrected transport (FCT)** algorithms of Boris and Book [BB73]. We can view this as creating a hybrid algorithm that is second order accurate in smooth regions but which reduces to a more robust first-order algorithm near discontinuities. This idea of hybridization was also used in early work of Harten and Zwas [HZ72]. An enormous variety of methods based on these principles have been developed in the past two decades. One of the algorithms of this type that is best known in the astrophysics community is the **piecewise parabolic method (PPM)** of Woodward and Colella [CW84], which uses a piecewise quadratic reconstruction, with appropriate limiting.

### 17.7.5 Total variation

How much should we limit the slope in a piecewise linear reconstruction? Ideally we would like to have a mathematical prescription that will allow us to use the Lax-Wendroff slope whenever possible, for second-order accuracy, while guaranteeing that no non-physical oscillations will arise. To achieve this we need a way to measure “oscillations” in the solution. This is provided by the notion of the **total variation** of a function. For a grid function  $Q$  we define

$$TV(Q) = \sum_{i=-\infty}^{\infty} |Q_i - Q_{i-1}|. \quad (17.68)$$

For an arbitrary function  $q(x)$  we can define

$$TV(q) = \sup \sum_{j=1}^N |q(\xi_j) - q(\xi_{j-1})|, \quad (17.69)$$

where the supremum is taken over all subdivisions of the real line  $-\infty = \xi_0 < \xi_1 < \dots < \xi_N = \infty$ . Note that for the total variation to be finite  $Q$  or  $q$  must approach constant values  $v^\pm$  as  $x \rightarrow \pm\infty$ .

The true solution to the advection equation simply propagates at speed  $u$  with unchanged shape, so that the total variation  $TV(q(\cdot, t))$  must be constant in time. A numerical solution to the advection equation may not have constant total variation, however. If the method introduces oscillations, then we would expect the total variation of  $Q^n$  to *increase* with time. We can thus attempt to avoid oscillations by requiring that the method does not increase the total variation:

**Definition 17.7.1** A 2-level method is called **total variation diminishing (TVD)** if, for any set of data  $Q^n$ , the values  $Q^{n+1}$  computed by the method satisfy

$$TV(Q^{n+1}) \leq TV(Q^n). \quad (17.70)$$

For a scalar conservation law, the exact solution has nonincreasing variation and so this is a reasonable condition to impose on a numerical method. Harten [Har83] introduced the use of this criterion in developing and analyzing numerical methods. For a scalar equation, steps 2 and 3 of Algorithm RSA are TVD, and so the overall method is TVD provided that the reconstruction step does not increase the variation, *i.e.*, as long as

$$TV(\bar{q}^n) \leq TV(Q^n). \quad (17.71)$$

### 17.7.6 Slope-limiter methods

Now let’s return to the derivation of numerical methods based on piecewise linear reconstruction, and consider how to “limit” the slopes so that (17.70) is satisfied. Note that setting  $\sigma_i^n \equiv 0$  works, since the piecewise constant function has the same TV as the discrete data. Hence *the first-order upwind method is TVD* for the advection equation. Hence upwind may smear solutions but cannot introduce oscillations, a familiar result.

One choice of slope that gives second-order accuracy for smooth solutions while still satisfying the TVD property is the **minmod slope**

$$\sigma_i^n = \text{minmod} \left( \frac{Q_i^n - Q_{i-1}^n}{h}, \frac{Q_{i+1}^n - Q_i^n}{h} \right), \quad (17.72)$$

where the minmod function of two arguments is defined by

$$\text{minmod}(a, b) = \begin{cases} a & \text{if } |a| < |b| \text{ and } ab > 0 \\ b & \text{if } |b| < |a| \text{ and } ab > 0 \\ 0 & \text{if } ab \leq 0. \end{cases} \quad (17.73)$$

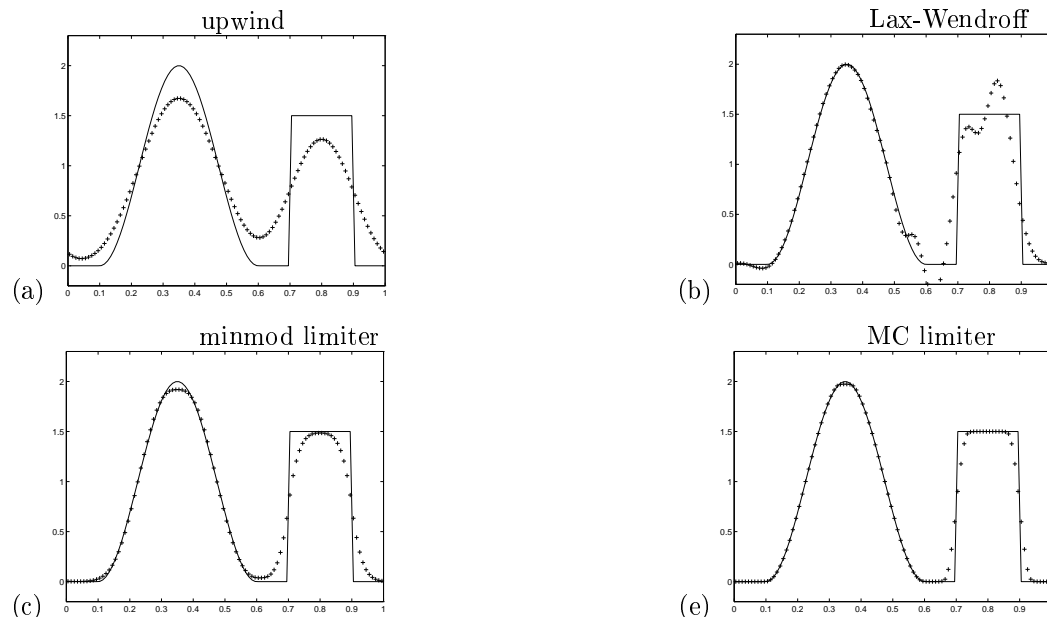


Figure 17.11: Tests on the advection equation with different limiters. All results are at time  $t = 1$ , after one revolution with periodic boundary conditions

Note that if  $a$  and  $b$  have the same sign then this selects the one which is smaller in modulus, else it returns zero.

Rather than defining the slope on the  $i$ 'th cell by always using the downwind difference (which would give Lax-Wendroff), or by always using the upwind difference (which would give Beam-Warming), the minmod method compares the two slopes and chooses the one which is smaller in magnitude. If the two slopes have different sign, then the value  $Q_i^n$  must be a local maximum or minimum, and it is easy to check in this case that we must set  $\sigma_i^n = 0$  in order to satisfy (17.71).

Figure 17.11 shows a comparison of the upwind, Lax-Wendroff, and minmod methods for an advection problem with initial data consisting of both a smooth hump and a square wave. The advection velocity is  $u = 1$  and periodic boundary conditions are used on  $[0, 1]$  so that at integer times  $t = 0, 1, 2, \dots$  the solution should agree with the initial data. The figure shows solutions at  $t = 1$  on a grid with  $h = 0.01$  and  $k = 0.005$  (so the Courant number is 0.5 — better results with all methods would be obtained with a Courant number closer to 1). We see that the minmod method does a fairly good job of maintaining good accuracy in the smooth hump and also sharp discontinuities in the square wave, with no oscillations.

Sharper resolution of discontinuities can be achieved with other limiters that do not reduce the slope as severely as minmod near a discontinuity. One popular choice is the **monotonized central-difference limiter (MC-limiter)**, which was proposed by van Leer [van77]:

$$\sigma_i^n = \text{minmod} \left( \left( \frac{Q_{i+1}^n - Q_{i-1}^n}{2h} \right), 2 \left( \frac{Q_i^n - Q_{i-1}^n}{h} \right), 2 \left( \frac{Q_{i+1}^n - Q_i^n}{h} \right) \right).$$

This compares the central-difference of Fromm's method with *twice* the one-sided slope to either side. In smooth regions this reduces to the centered slope of Fromm's method but near discontinuities it gives sharper resolution than minmod while remaining TVD.

### 17.7.7 Flux formulation with piecewise linears

The slope-limiter methods described above can be written as flux-differencing methods of the form (17.45). The updating formulas derived above can be manipulated algebraically to determine what the

numerical flux function must be. Alternatively, we can derive the numerical flux by computing the exact flux through the interface  $x_i$  using the piecewise linear solution  $\tilde{q}^n(x, t)$ , by integrating  $u\tilde{q}^n(x_i, t)$  in time from  $t_n$  to  $t_{n+1}$ . For the advection equation this is easy to do and we find that

$$\begin{aligned} F_i^n &= \frac{1}{k} \int_{t_n}^{t_{n+1}} u\tilde{q}^n(x_i, t) dt \\ &= uQ_{i-1}^n + \frac{1}{2}u(h - ku)\sigma_{i-1}^n. \end{aligned}$$

Using this in the flux-differencing formula (17.45) gives

$$Q_i^{n+1} = Q_i^n - \frac{ku}{h}(Q_i^n - Q_{i-1}^n) - \frac{1}{2}\frac{ku}{h}(h - ku)(\sigma_i^n - \sigma_{i-1}^n),$$

which agrees with (17.64).

If we also consider the case  $u < 0$ , then we will find that in general the numerical flux for a slope-limiter method is

$$F_i^n = \begin{cases} uQ_{i-1}^n + \frac{1}{2}u(h - ku)\sigma_{i-1}^n & \text{if } u \geq 0 \\ uQ_i^n - \frac{1}{2}u(h + ku)\sigma_i^n & \text{if } u \leq 0, \end{cases} \quad (17.74)$$

where  $\sigma_i^n$  is the slope in the  $i$ 'th cell  $C_i$ , chosen by one of the formulas discussed previously.

Rather than associating a slope  $\sigma_i^n$  with the  $i$ 'th cell, the idea of writing the method in terms of fluxes between cells suggests that we should instead associate our approximation to  $q_x$  with the cell interface at  $x_i$  where  $F_i^n$  is defined. Across the interface  $x_i$  we have a jump

$$\Delta Q_i^n = Q_i^n - Q_{i-1}^n \quad (17.75)$$

and this jump divided by  $h$  gives an approximation to  $q_x$ . This suggests that we write the flux (17.74) as

$$F_i^n = u^- Q_i^n + u^+ Q_{i-1}^n + \frac{1}{2}|u| \left(1 - \left|\frac{ku}{h}\right|\right) \delta_i^n, \quad (17.76)$$

where  $u^\pm$  are defined in (17.56). If  $\delta_i^n$  is the jump  $\Delta Q_i^n$  itself then this gives the Lax-Wendroff method. We see that the Lax-Wendroff flux can be interpreted as a modification to the upwind flux (17.55). This observation is crucial in the development of high-resolution methods.

### 17.7.8 Flux limiters

From the above discussion it is natural to view Lax-Wendroff as the basic second-order method based on piecewise linear reconstruction. Other second-order methods have fluxes of the form (17.76) with different choices of  $\delta_i^n$ . The slope-limiter methods can then be reinterpreted as **flux-limiter methods** by choosing  $\delta_i^n$  to be a limited version of (17.75). In general we will set

$$\delta_i^n = \phi(\theta_i^n)\Delta Q_i^n, \quad (17.77)$$

where

$$\theta_i^n = \frac{\Delta Q_I^n}{\Delta Q_i^n}. \quad (17.78)$$

The index  $I$  here is used to represent the interface on the *upwind* side of  $x_i$ :

$$I = \begin{cases} i - 1 & \text{if } u > 0 \\ i + 1 & \text{if } u < 0. \end{cases} \quad (17.79)$$

The ratio  $\theta_i^n$  can be thought of as a measure of the smoothness of the data near  $x_i$ . Where the data is smooth we expect  $\theta_i^n \approx 1$  (except at extrema). Near a discontinuity we expect that  $\theta_i^n$  may be far from 1.

The function  $\phi(\theta)$  is the **flux-limiter function**, whose value depends on the smoothness. Setting  $\phi(\theta) \equiv 1$  for all  $\theta$  gives the Lax-Wendroff method, while setting  $\phi(\theta) \equiv 0$  gives upwind. More generally we might want to devise a limiter function  $\phi$  that has values near 1 for  $\theta$  near 1, but which reduces (or perhaps increases) the slope where the data is not smooth.

There are many other ways one might choose to measure the smoothness of the data besides the variable  $\theta$  defined in (17.78). However, the framework proposed above results in very simple formulas for the  $\phi$  function corresponding to many standard methods, including all the methods discussed so far.

In particular, note the very nice feature that choosing

$$\phi(\theta) = \theta \tag{17.80}$$

results in (17.77) becoming

$$\delta_i^n = \left( \frac{\Delta Q_i^n}{\Delta Q_i^n} \right) \Delta Q_i^n = \Delta Q_i^n.$$

Hence this choice results in the jump at the interface *upwind* from  $x_i$  being used to define  $\delta_i^n$  instead of the jump at this interface. As a result, the method (17.76) with the choice of “limiter” (17.80) reduces to the Beam-Warming method.

Since the centered difference (17.65) is the average of the one-sided slopes (17.66) and (17.67), we also find that Fromm’s method can be obtained by choosing

$$\phi(\theta) = \frac{1}{2}(1 + \theta). \tag{17.81}$$

Also note that  $\phi(\theta) = 2$  corresponds to using  $\delta_i^n = 2\Delta Q_i^n$ , i.e., twice the jump at this interface, while  $\phi(\theta) = 2\theta$  results in using twice the jump at the upwind interface. Recall that these are necessary ingredients in some of the slope limiters discussed in Section 17.7.6.

Translating the various slope limiters into flux-limiter functions, we find the following expressions for some standard methods:

<b>Linear methods:</b>		
upwind :	$\phi(\theta) = 0$	
Lax-Wendroff :	$\phi(\theta) = 1$	
Beam-Warming :	$\phi(\theta) = \theta$	
Fromm :	$\phi(\theta) = \frac{1}{2}(1 + \theta)$	
<b>High-resolution limiters:</b>		(17.82)
minmod :	$\phi(\theta) = \text{minmod}(1, \theta)$	
superbee :	$\phi(\theta) = \max(0, \min(1, 2\theta), \min(2, \theta))$	
MC :	$\phi(\theta) = \max(0, \min((1 + \theta)/2, 2, 2\theta))$	
van Leer :	$\phi(\theta) = \frac{\theta +  \theta }{1 +  \theta }$	

A wide variety of other limiters have also been proposed in the literature.

### 17.7.9 TVD limiters

For simple limiters such as minmod, it is clear from the derivation as a slope-limiter (Section 17.7.6) that the resulting method is TVD, since it is easy to check that (17.71) is satisfied. For more complicated limiters we would like to have an algebraic proof that the resulting method is TVD. A fundamental tool

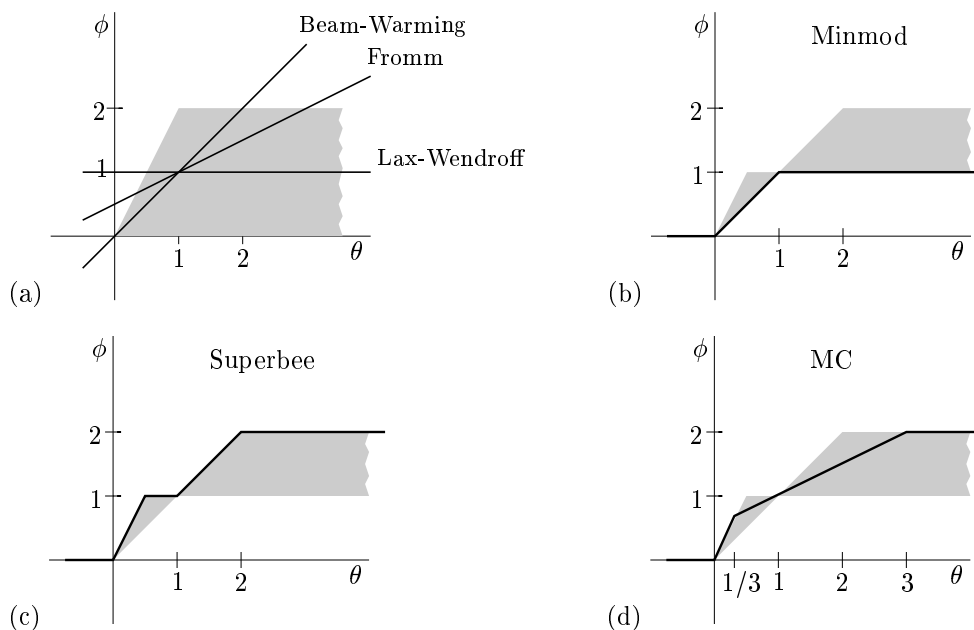


Figure 17.12: Limiter functions  $\phi(\theta)$ . (a) The shaded region shows where function values must lie for the method to be TVD. The second-order linear methods have functions  $\phi(\theta)$  that leave this region. (b) The shaded region is the Sweby region of second-order TVD methods. The minmod limiter lies along the lower boundary. (c) The superbee limiter lies along the upper boundary. (d) The MC limiter is smooth at  $\phi = 1$

in this direction is a theorem of Harten [Har83], which can be used to derive explicit algebraic conditions on the function  $\phi$  required for a TVD method. Sweby [Swe84] derived explicit constraints for limiter functions and shows that we require (see also [LeV90]):

$$0 \leq \phi(\theta) \leq \min\text{mod}(2, 2\theta). \quad (17.83)$$

This defines the **TVD region** in the  $\theta$ - $\phi$  plane: the curve  $\phi(\theta)$  must lie in this region, which is shown as the shaded region in Figure 17.12(a). This figure also shows the functions  $\phi(\theta)$  from (17.82) for the Lax-Wendroff, Beam-Warming, and Fromm methods. All of these functions lie outside the TVD region for some values of  $\theta$ , and these methods are not TVD. This graphical analysis of  $\phi$  was first presented by Sweby [Swe84], who analyzed a wide class of flux-limiter methods.

Note that for any second-order accurate method we must have  $\phi(1) = 1$ . Sweby found, moreover, that it is best to take  $\phi$  to be a convex combination of  $\phi = 1$  (Lax-Wendroff) and  $\phi = \theta$  (Beam-Warming). Other choices apparently give too much compression, and smooth data such as a sine wave tends to turn into a square wave as time evolves. Imposing this additional restriction gives the “second order TVD” region of Sweby which is shown in Figure 17.12(b).

The “high-resolution” limiter functions from (17.82) are all seen to satisfy the constraints (17.83), and these limiters all give TVD methods. The  $\phi$  functions are graphed in Figure 17.12. Note that minmod lies along the lower boundary of the Sweby region while superbee lies along the upper boundary. The fact that these functions are not smooth at  $\theta = 1$  corresponds to the fact that there is a switch in the choice of one-sided approximation used as  $\theta$  crosses this point. This can lead to a loss in accuracy near inflection points. For full second order accuracy we would like the function  $\phi$  to be smooth near  $\theta = 1$ , as for the MC limiter. The van Leer limiter is an even smoother version of this.

### 17.7.10 Linear systems

The slope-limiter or flux-limiter methods can also be extended to systems of equations. This is most easily done in the flux-limiter framework, and will be illustrated first for the linear system  $q_t + Aq_x = 0$ .

We can write  $A = A^+ + A^-$ , where  $A = R\Lambda R^{-1}$  is decomposed based on the sign of each eigenvalue:  $A^+ = R\Lambda^+R^{-1}$  and  $A^- = R\Lambda^-R^{-1}$  where  $\Lambda^+$ , for example, has the positive part of each  $\lambda$  on the diagonal, so negative eigenvalues are replaced by zero. Using this, Godunov's method for a linear system (the generalization of the upwind method) can be written in terms of the flux function

$$F(Q_{i-1}, Q_i) = (A^+Q_{i-1} + A^-Q_i). \quad (17.84)$$

The Lax-Wendroff method (14.16) can also be written in flux-differencing form (17.45) if we define the flux by

$$F(Q_{i-1}, Q_i) = \frac{1}{2}A(Q_{i-1} + Q_i) + \frac{1}{2}\frac{k}{h}A^2(Q_i - Q_{i-1}). \quad (17.85)$$

We can rewrite this as

$$F(Q_{i-1}, Q_i) = (A^+Q_{i-1} + A^-Q_i) + \frac{1}{2}|A| \left(1 - \frac{k}{h}|A|\right) (Q_i - Q_{i-1}), \quad (17.86)$$

where  $|A| = A^+ - A^-$ .

In the form (17.86), we see that the Lax-Wendroff flux can be viewed as being composed of the Godunov flux (17.84) plus a correction term, just as for the scalar advection equation. To define a flux-limiter method we must limit the magnitude of this correction term based on how the data is varying. But for a system of equations,  $\Delta Q_i = Q_i - Q_{i-1}$  is a *vector* and it is not so clear how to compare this vector with the neighboring jump  $\Delta Q_{i-1}$  or  $\Delta Q_{i+1}$  to generalize (17.77), nor which neighboring jump to consider, since the "upwind" direction is different for each eigen-component. The solution, of course, is that we must decompose the correction term in (17.86) into eigen-components and limit each scalar eigen-coefficient separately based on the algorithm for scalar advection.

We can rewrite the correction term as

$$\frac{1}{2}|A| \left(1 - \frac{k}{h}|A|\right) (Q_i - Q_{i-1}) = \frac{1}{2}|A| \left(1 - \frac{k}{h}|A|\right) \sum_{i=1}^m \alpha_i^p r^p,$$

where

$$(Q_i - Q_{i-1}) = \sum_{i=1}^m \alpha_i^p r^p$$

gives the decomposition of the jump in  $Q$  across this interface into eigenvectors of  $A$ , i.e., the solution of the Riemann problem.

The flux-limiter method is defined by replacing the scalar coefficient  $\alpha_i^p$  by a limited version, based on the scalar formulas of Section 17.7.8. We set

$$\tilde{\alpha}_i^p = \alpha_i^p \phi(\theta_i^p), \quad (17.87)$$

where

$$\theta_i^p = \frac{\alpha_I^p}{\alpha_i^p} \quad \text{with } I = \begin{cases} i-1 & \text{if } \lambda^p > 0 \\ i+1 & \text{if } \lambda^p < 0 \end{cases} \quad (17.88)$$

and  $\phi$  is one of the limiter functions of Section 17.7.8. The flux function for the flux-limiter method is then

$$F_i = A^+Q_{i-1} + A^-Q_i + \tilde{F}_i, \quad (17.89)$$

where the first term is the upwind flux and the correction flux  $\tilde{F}_i$  is defined by

$$\tilde{F}_i = \frac{1}{2}|A| \left(1 - \frac{k}{h}|A|\right) \sum_{i=1}^m \tilde{\alpha}_i^p r^p. \quad (17.90)$$

Note that in the case of a scalar equation, we can take  $r^1 = 1$  as the eigenvector of  $A = u$ , so that  $\Delta Q_i = \alpha_i^1$  which is what we called  $\delta_i$  in Section 17.7.8. The formula (17.89) then reduces to (17.76). For a system of equations, the method just presented can also be obtained by diagonalizing the linear system and applying the scalar flux limiter method to each resulting advection equation. This is what we are doing at each cell interface by solving the Riemann problem.

Also note that the flux  $F_i$  depends not only on  $Q_{i-1}$  and  $Q_i$ , but also on  $Q_{i-2}$  and  $Q_{i+1}$  in general, because neighboring jumps are used in defining the limited values  $\tilde{\alpha}_i^p$  in (17.90). The flux-limiter method thus has a 5-point stencil rather than the 3-point stencil of Lax-Wendroff. This is particularly important in specifying boundary conditions (see Section 17.8).

### 17.7.11 Implementation and CLAWPACK

For the constant coefficient linear system, we could compute the matrices  $A^+$ ,  $A^-$ , and  $|A|$  once and for all and compute the fluxes directly from the formulas given above. However, with limiters we must solve the Riemann problem at each interface to obtain a decomposition of  $\Delta Q_i$  into waves  $\alpha_i^p r^p$  and wave speeds  $\lambda^p$  and these can be used directly in the computation of  $Q_i^{n+1}$  without ever forming the matrices. This approach also generalizes directly to nonlinear systems of conservation laws, where we do not have a single matrix  $A$  but can still solve a Riemann problem at each interface for waves and wave speeds. This generalization is discussed briefly in the next section.

To accomplish this most easily, note that if we use the flux (17.89) in the flux-differencing formula (17.45) and then rearrange the upwind terms, we can write the formula for  $Q_i^{n+1}$  as

$$Q_i^{n+1} = Q_i^n - \frac{k}{h}(A^+ \Delta Q_i + A^- \Delta Q_{i+1}) - \frac{k}{h}(\tilde{F}_{i+1} - \tilde{F}_i), \quad (17.91)$$

where  $\tilde{F}_i$  is given by (17.90). Here we drop the superscript  $n$  from the current time step since we will need to use superscript  $p$  below to denote the wave family. Each of the terms in this expression can be written in terms of the waves  $\alpha_i^p r^p$  and wave speeds  $\lambda^p$ :

$$A^+ \Delta Q_i = \sum_{p=1}^m (\lambda^p)^+ \alpha_i^p r^p, \quad (17.92)$$

$$A^- \Delta Q_i = \sum_{p=1}^m (\lambda^p)^- \alpha_i^p r^p, \quad (17.93)$$

and

$$\tilde{F}_i = \frac{1}{2} \sum_{p=1}^m |\lambda^p| \left(1 - \frac{k}{h} |\lambda^p|\right) \tilde{\alpha}_i^p r^p.$$

Once we have solved the Riemann problem for the waves  $\alpha_i^p r^p$ , limited waves  $\tilde{\alpha}_i^p r^p$ , and speeds  $\lambda^p$ , we can compute everything needed for the high-resolution method (17.91).

These methods are implemented in a very general form in the CLAWPACK (Conservation LAWS PACKage) software which is available on the web [LeV].

## 17.8 Boundary conditions

So far we have only studied methods for updating the cell average  $Q_i^n$  assuming that we have neighboring cell values  $Q_{i-1}^n$  and  $Q_{i+1}^n$  and perhaps values further away as needed in order to compute the fluxes



$F_i^n$  and  $F_{i+1}^n$ . In practice we must always compute on some finite set of grid cells covering a bounded domain, and in the first and last cells we will not have the required neighboring information. Instead we have some set of physical boundary conditions that must be used in updating these cell values. One approach is to develop special formulas for use near the boundaries, which will depend both on what type of boundary condition is specified and on what sort of method we are trying to match. However, in general it is much easier to think of extending the computational domain to include a few additional cells on either end, called **ghost cells**, whose values are set at the beginning of each time step in some manner that depends on the boundary conditions and perhaps the interior solution. Then these values provide the neighboring cell values needed in updating the cells near the physical domain. The updating formula is then exactly the same in all cells, and there is no need to develop a special flux-limiter method, say, that works with boundary data instead of initial data. Instead the boundary conditions must be used in deciding how to set the values of the ghost cells, but this can generally be done in a way that depends only on the boundary conditions and is decoupled entirely from the choice of numerical method that is then applied.

Suppose the problem is on the physical domain  $[a, b]$ , which is subdivided into cells  $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_N$  with  $x_1 = a$  and  $x_{N+1} = b$ , so that  $h = (b - a)/N$ . If we use a method for which  $F_i$  depends only on  $Q_{i-1}$  and  $Q_i$ , then we need only one ghost cell on either end. The ghost cell  $\mathcal{C}_0 = [a - h, a)$  allows us to calculate the flux  $F_1$  at the left boundary while the ghost cell  $\mathcal{C}_{N+1} = [b, b + h)$  is used to calculate  $F_{N+1}$  at  $x = b$ . With a flux-limiter method of the type developed above, we will generally need two ghost cells at each boundary since, for example, the jump  $Q_0 - Q_{-1}$  will be needed in limiting the flux correction in  $F_1$ . For a method with an even wider stencil, additional ghost cells would be needed.

We will refer to the solution in the original domain  $[a, b]$  as the **interior solution**, which is computed in each time step by the numerical method. At the start of each time step we have the interior values  $Q_1^n, \dots, Q_N^n$  obtained from the previous time step (or from the initial conditions if  $n = 0$ ), and we apply a **boundary condition procedure** to fill the ghost cells with values  $Q_0^n, Q_{N+1}^n$ , etc. before applying the method on the next time step. We will look at several examples to see how the ghost cell values might be set in order to implement various physical boundary conditions.

### 17.8.1 Periodic boundary conditions

Periodic boundary conditions of the form  $q(a, t) = q(b, t)$  are very easy to impose with any numerical method. In updating  $Q_1$  we need values  $Q_0$  to the left and  $Q_2$  to the right (for a 3-point method). By periodicity the value  $Q_0$  should agree with the value  $Q_N$  in the last cell. One could code the formula for updating  $Q_1$  separately to use  $Q_N$  in place of the value  $Q_{i-1}$  that would normally be used for  $i > 1$ , but it is simpler to use the ghost-cell approach and simply set  $Q_0^n = Q_N^n$  before computing fluxes and updating the cell values, so that the same formula can then be used everywhere. With a 5-point stencil we need to fill two ghost cells at each boundary, and we set

$$Q_{-1}^n = Q_{N-1}^n, \quad Q_0^n = Q_N^n, \quad Q_{N+1}^n = Q_1^n, \quad Q_{N+2}^n = Q_2^n \quad (17.94)$$

at the start of each time step.

### 17.8.2 Outflow boundaries

Often we have artificial computational boundaries that arise simply because we can only solve the problem on a bounded domain. At such boundaries we often want to have no incoming signal, while there may be out-going waves that should leave the domain cleanly without generating spurious reflections at the artificial boundary. We thus want **nonreflecting** boundary conditions. At such boundaries we can often set ghost cell values by **extrapolation** from the interior solution. If the ghost cell value  $Q_{N+1}^n$  is set based on  $Q_N^n, Q_{N-1}^n, \dots$ , then the new value  $Q_{N+1}^{n+1}$  will effectively be computed on the basis of values to the left alone, even if the formula depends on  $Q_{N+1}^n$ , and hence this reduces to some sort of upwind method. The simplest approach is to use a **zero-order extrapolation**, meaning extrapolation

by a constant function. We simply set

$$Q_{N+1}^n = Q_N^n, \quad Q_{N+2}^n = Q_N^n \quad (17.95)$$

at the start of each time step. The idea of extrapolation at outflow boundaries turns out to be extremely powerful in conjunction with methods based on solving Riemann problems. If there is no jump in the values at the boundary, there are no waves in the Riemann solution and in particular no incoming waves.



## Chapter 18

# Mixed Equations and Fractional Step Methods

We have now studied the solution of various types of time-dependent equations — ODE's such as those arising in chemical kinetics, diffusion equations, and advection equations. In practice several processes may be happening simultaneously, and the PDE model will not be a pure equation of any of the types already discussed, but rather a mixture. In this chapter we discuss mixed equations such as reaction-diffusion equations or advection-diffusion equations. There are various ways to handle mixed equations, and we will consider two basic approaches:

- Unsplit methods, in which a single finite-difference formula is developed to advance the full mixed equation over one time step.
- Fractional step (splitting) methods, in which the problem is broken down into pieces corresponding to the different processes, and a numerical method appropriate for each separate piece is applied independently. This approach is also often used to split multi-dimensional problems into a sequence of one-dimensional problems (*dimensional splitting*). We have seen an example of this for the heat equation in the LOD method of Section 13.8.

### 18.1 Advection-reaction equations

**Example 18.1.** We begin with a simple advection-reaction equation of the form

$$u_t + au_x = -\lambda u, \quad (18.1)$$

with data  $u(x, 0) = \eta(x)$ . This would model, for example, the transport of a radioactive material in a fluid flowing at constant speed  $a$  down a pipe. The material decays as it flows along, at rate  $\lambda$ . We can easily compute the exact solution of (18.1), since along the characteristic  $dx/dt = a$  we have  $du/dt = -\lambda u$ , and hence

$$u(x, t) = e^{-\lambda t} \eta(x - at). \quad (18.2)$$

#### 18.1.1 Unsplit methods

It is easy to develop unsplit methods for (18.1). For example, an obvious extension of the upwind method for advection would be (assuming  $a > 0$ ),

$$U_j^{n+1} = U_j^n - \frac{ak}{h}(U_j^n - U_{j-1}^n) - k\lambda U_j^n. \quad (18.3)$$

This method is first order accurate and stable for  $0 < ak/h < 1$ .

A second order Lax-Wendroff style method can be developed by using the Taylor series

$$u(x, t + k) \approx u(x, t) + ku_t(x, t) + \frac{1}{2}k^2 u_{tt}(x, t). \quad (18.4)$$

As in the derivation of Lax-Wendroff, we must compute  $u_{tt}$  from the PDE, obtaining

$$u_{tt} = -au_{xt} - \lambda u_t.$$

Since

$$u_{tx} = -au_{xx} - \lambda u_x,$$

we obtain

$$u_{tt} = a^2 u_{xx} + 2a\lambda u_x + \lambda^2 u. \quad (18.5)$$

Note that this is more easily obtained by using

$$\partial_t u = (-a\partial_x - \lambda)u,$$

and hence

$$\partial_t^2 u = (-a\partial_x - \lambda)^2 u = (a^2 \partial_x^2 - 2a\lambda \partial_x + \lambda^2)u. \quad (18.6)$$

Using this expression for  $u_{tt}$  in (18.4) gives

$$\begin{aligned} u(x, t + k) &\approx u - k(au_x + \lambda u) + \frac{1}{2}k^2(a^2 u_{xx} + 2a\lambda u_x + \lambda^2 u) \\ &= \left(1 - k\lambda + \frac{1}{2}k^2 \lambda^2\right) u - ka \left(1 - \frac{1}{2}k\lambda\right) u_x + \frac{1}{2}k^2 a^2 u_{xx}. \end{aligned} \quad (18.7)$$

We can now approximate  $x$ -derivatives by finite differences to obtain the second-order method

$$U_j^{n+1} = \left(1 - k\lambda + \frac{1}{2}k^2 \lambda^2\right) U_j^n - \frac{ka}{2h} \left(1 - \frac{1}{2}k\lambda\right) (U_{j+1}^n - U_{j-1}^n) + \frac{k^2 a^2}{2h^2} (U_{j-1}^n - 2U_j^n + U_{j+1}^n). \quad (18.8)$$

Note that in order to correctly model the equation (18.1) to second order accuracy, we must properly model the interaction between the  $au_x$  and the  $\lambda u$  terms, which brings in the mixed term  $-\frac{1}{2}k^2 a\lambda u_x$  in the Taylor series expansion.

For future use we also note that the full Taylor series expansion can be written as

$$u(x, t + k) = \sum_{j=0}^{\infty} \frac{k^j}{j!} (-a\partial_x - \lambda)^j u(x, t), \quad (18.9)$$

which can be written formally as

$$u(x, t + k) = e^{-k(a\partial_x + \lambda)} u(x, t) \quad (18.10)$$

The operator  $e^{-k(a\partial_x + \lambda)}$ , which is defined via the Taylor series in (18.9), is called the *solution operator* for the PDE (18.1). Note that for any value of  $k$  we have

$$u(x, k) = e^{-k(a\partial_x + \lambda)} u(x, 0).$$

### 18.1.2 Fractional step methods

A fractional step method for (18.1) is applied by first splitting the equation into two *subproblems*

$$\text{Problem A: } u_t + au_x = 0, \quad (18.11)$$

$$\text{Problem B: } u_t = -\lambda u. \quad (18.12)$$

We have developed methods for each of these two equations separately. The idea with the fractional step method is to combine these by applying the two methods in an alternating manner. As a simple example, suppose we use the upwind method for the A-step and forward Euler for the ODE in the B-step. Then the simplest fractional step method over one time step would consist of the following 2 stages:

$$\text{A-step: } U_j^* = U_j^n - \frac{ak}{h}(U_j^n - U_{j-1}^n), \quad (18.13)$$

$$\text{B-step: } U_j^{n+1} = U_j^* - k\lambda U_j^*. \quad (18.14)$$

Note that we first take a time step of length  $k$  with upwind, starting with initial data  $U_j^n$  to obtain the intermediate value  $U_j^*$ . Then we take a time step of length  $k$  using forward Euler, starting with the data  $U^*$  obtained from the first stage.

It may seem that we have advanced the solution by time  $2k$  after taking these two steps of length  $k$ . However, in each stage we used only some of the terms in the original PDE, and the two stages combined give a consistent approximation to solving the original PDE (18.1) over a single time step of length  $k$ .

To check this consistency, we can combine the two stages by eliminating  $U^*$  to obtain a method in a more familiar form:

$$\begin{aligned} U_j^{n+1} &= (1 - k\lambda)U_j^* \\ &= (1 - k\lambda) \left[ U_j^n - \frac{ak}{h}(U_j^n - U_{j-1}^n) \right] \\ &= U_j^n - \frac{ak}{h}(U_j^n - U_{j-1}^n) - k\lambda U_j^n + \frac{ak^2\lambda}{h}(U_j^n - U_{j-1}^n). \end{aligned} \quad (18.15)$$

The first three terms on the right-hand side agree with the unsplit method (18.3). The final term is  $O(k^2)$  (since  $(U_j^n - U_{j-1}^n)/h \approx u_x = O(1)$ ) and so a local truncation error analysis will show that this method, though slightly different from (18.3), is also consistent and first order accurate on the original equation (18.1).

A natural question is whether we could improve the accuracy by using a more accurate method in each step. For example, suppose we use Lax-Wendroff in the A-step and the trapezoidal method, or the 2-stage Runge-Kutta method from Example 6.15, in the B-step. Would we then obtain a second order accurate method for the original equation? For this particular equation, the answer is yes. In fact if we use  $p$ 'th order accurate methods for each step, the result will be a  $p$ 'th order accurate method for the full original equation. But this equation is very special in this regard, and this claim should seem surprising. One would think that splitting the equation into pieces in this manner would introduce some error that depends on the size of the time step  $k$  and is independent of how well we then approximate the subproblem in each step. In general this is true — there is a “splitting error” that in general would be  $O(k)$  for the type of splitting used above, and so the resulting fractional step method will be only first order accurate, no matter how well we then approximate each step. This will be analyzed in more detail below.

For the case of equation (18.1) there is no splitting error. This follows from the observation that we can solve (18.1) over any time period  $k$  by first solving equation (18.11) over time  $k$ , and then using the result as data to solve the equation (18.12) over time  $k$ . To verify this, let  $u^*(x, k)$  be the exact solution to the A-problem,

$$\begin{aligned} u_t^* + au_x^* &= 0 \\ u^*(x, 0) &= \eta(x). \end{aligned} \quad (18.16)$$

We use a different symbol  $u^*(x, t)$  for the solution to this problem rather than  $u(x, t)$ , which we reserve for the exact solution to the original problem.

Then we have

$$u^*(x, k) = \eta(x - ak).$$

If we now use this as data in solving the B-problem (18.12), we will be solving

$$u_t^{**} = -\lambda u^{**} \quad (18.17)$$

with initial data

$$u^{**}(x, 0) = \eta(x - ak).$$

This is just an ODE at each point  $x$ , and the solution is

$$u^{**}(x, k) = e^{-\lambda k} \eta(x - ak).$$

Comparing this with (18.2), we see that we have indeed recovered the solution to the original problem by this 2-stage procedure.

Physically we can interpret this as follows. Think of the original equation as modeling a radioactive tracer that is advecting with constant speed  $a$  (carried along in a fluid, say) and also decaying with rate  $\lambda$ . Since the decay properties are independent of the position  $x$ , we can think of first advecting the tracer over time  $k$  without allowing any decay, and then holding the fluid and tracer stationary while we allow it to decay for time  $k$ . We will get the same result, and this is what we have done in the fractional step method.

Another way to examine the splitting error, which must be used more generally when we do not know the exact solution to the equations involved, is to use Taylor series expansions. If we look at a time step of length  $k$ , then solving the A-equation gives

$$\begin{aligned} u^*(x, k) &= u^*(x, 0) + k u_t^*(x, 0) + \frac{1}{2} k^2 u_{tt}^*(x, 0) + \dots \\ &= u^*(x, 0) - a k u_x^*(x, 0) + \frac{1}{2} a^2 k^2 u_{xx}^*(x, 0) - \dots \end{aligned} \quad (18.18)$$

Similarly, if we solve the problem (18.17) with general initial data we obtain

$$\begin{aligned} u^{**}(x, k) &= u^{**}(x, 0) + k u_t^{**}(x, 0) + \frac{1}{2} k^2 u_{tt}^{**}(x, 0) + \dots \\ &= \left( 1 - k\lambda + \frac{1}{2} k^2 \lambda^2 + \dots \right) u^{**}(x, 0). \end{aligned} \quad (18.19)$$

If we now use the result from (18.18) as the initial data in (18.19), we obtain

$$\begin{aligned} u^{**}(x, k) &= \left( 1 - k\lambda + \frac{1}{2} k^2 \lambda^2 - \dots \right) \left( u^*(x, 0) - a k u_x^*(x, 0) + \frac{1}{2} a^2 k^2 u_{xx}^*(x, 0) + \dots \right) \\ &= u^* - (a k u_x^* + \lambda u^*) + \frac{1}{2} k^2 (a^2 u_{xx}^* + 2a\lambda u_x^* + \lambda^2 u^*) + \dots \end{aligned} \quad (18.20)$$

Comparing this with the Taylor series expansion (18.7) that we used in deriving the unsplit Lax-Wendroff method shows that this agrees with  $u(x, k)$ , at least for the 3 terms shown, and in fact to all orders.

Note that the mixed term  $k^2 a \lambda u_x$  needed in the  $u_{tt}$  term from (18.5) now arises naturally from taking the product of the two Taylor series (18.18) and (18.19). In fact, we see that for this simple equation we can write (18.19) as

$$u^{**}(x, k) = e^{-k\lambda} u^{**}(x, 0)$$

while (18.18) can be written formally as

$$u^*(x, k) = e^{-a k \partial_x} \eta(x).$$

If we now use  $u^*(x, k)$  as the data  $u^{**}(x, 0)$  as we do in the fractional step method, we obtain

$$u^{**}(x, k) = e^{-k\lambda} e^{-ak\partial_x} \eta(x).$$

Multiplying out the Taylor series as we did in (18.20) verifies that these exponentials satisfy the usual rule, that to take the product we need only add the exponents, *i.e.*,

$$u^{**}(x, k) = e^{-k(a\partial_x + \lambda)} \eta(x).$$

The exponential appearing here is exactly the solution operator for the original equation, and so again we see that  $u^{**}(x, k) = u(x, k)$ .

The fact that there is no splitting error for the problem (18.1) is a reflection of the fact that, for this problem, the solution operator for the full problem is exactly equal to the product of the solution operators of the two subproblems (18.11) and (18.12). This is not generally the case in other problems.

**Example 18.2.** Suppose we modify the equation slightly so that the decay rate  $\lambda$  depends on  $x$ ,

$$u_t + au_x = -\lambda(x)u. \quad (18.21)$$

Then our previous argument for the lack of a splitting error breaks down — advecting the tracer a distance  $ak$  and then allowing it to decay, with rates given by the values of  $\lambda$  at the final positions, will not in general give the same result as if the decays occurs continuously as it advects, using the instantaneous rate given by  $\lambda(x)$  at each point passed.

This can be analyzed formally using Taylor series expansions again. Rather than going through this for this particular example, we will first examine the more general case and then apply it to this problem.

## 18.2 General formulation of fractional step methods

Consider a more general PDE of the form

$$u_t = (\mathcal{A} + \mathcal{B})u \quad (18.22)$$

where  $\mathcal{A}$  and  $\mathcal{B}$  may be differential operators, *e.g.*,  $\mathcal{A} = -a\partial_x$  and  $\mathcal{B} = \lambda(x)$  in the previous example. For simplicity suppose that  $\mathcal{A}$  and  $\mathcal{B}$  do not depend explicitly on  $t$ , *e.g.*,  $\lambda(x)$  is a function of  $x$  but not of  $t$ . Then we can compute that

$$u_{tt} = (\mathcal{A} + \mathcal{B})u_t = (\mathcal{A} + \mathcal{B})^2 u,$$

and in general

$$\partial_t^j u = (\mathcal{A} + \mathcal{B})^j u.$$

We have used this idea before in calculating Taylor series, *e.g.*, in (18.6).

Note that if  $\mathcal{A}$  or  $\mathcal{B}$  do depend on  $t$ , then we would have to use the product rule,

$$u_{tt} = (\mathcal{A} + \mathcal{B})u_t + (\mathcal{A}_t + \mathcal{B}_t)u$$

and everything would become more complicated.

In our simple case we can write the solution at time  $t$  using Taylor series as

$$\begin{aligned} u(x, k) &= u(x, 0) + k(\mathcal{A} + \mathcal{B})u(x, 0) + \frac{1}{2}k^2(\mathcal{A} + \mathcal{B})^2 u(x, 0) + \cdots \\ &= \left( I + k(\mathcal{A} + \mathcal{B}) + \frac{1}{2}k^2(\mathcal{A} + \mathcal{B})^2 + \cdots \right) u(x, 0) \\ &= \sum_{j=0}^{\infty} \frac{k^j}{j!} (\mathcal{A} + \mathcal{B})^j u(x, 0), \end{aligned} \quad (18.23)$$



which formally could be written as

$$u(x, k) = e^{k(\mathcal{A}+\mathcal{B})}u(x, 0).$$

With the fractional step method, we instead compute

$$u^*(x, k) = e^{k\mathcal{A}}u(x, 0),$$

and then

$$u^{**}(x, k) = e^{k\mathcal{B}}e^{k\mathcal{A}}u(x, 0),$$

and so the *splitting error* is

$$u(x, k) - u^{**}(x, k) = \left( e^{k(\mathcal{A}+\mathcal{B})} - e^{k\mathcal{B}}e^{k\mathcal{A}} \right) u(x, 0). \quad (18.24)$$

This should be calculated using the Taylor series expansions. We have (18.23) already, while

$$\begin{aligned} u^{**}(x, k) &= \left( I + k\mathcal{B} + \frac{1}{2}k^2\mathcal{B}^2 + \cdots \right) \left( I + k\mathcal{A} + \frac{1}{2}k^2\mathcal{A}^2 + \cdots \right) u(x, 0) \\ &= \left( I + k(\mathcal{A} + \mathcal{B}) + \frac{1}{2}k^2(\mathcal{A}^2 + 2\mathcal{B}\mathcal{A} + \mathcal{B}^2) + \cdots \right) u(x, 0). \end{aligned} \quad (18.25)$$

The  $I + k(\mathcal{A} + \mathcal{B})$  terms agree with (18.23). In the  $k^2$  term, however, the term from (18.23) is

$$\begin{aligned} (\mathcal{A} + \mathcal{B})^2 &= (\mathcal{A} + \mathcal{B})(\mathcal{A} + \mathcal{B}) \\ &= \mathcal{A}^2 + \mathcal{A}\mathcal{B} + \mathcal{B}\mathcal{A} + \mathcal{B}^2. \end{aligned} \quad (18.26)$$

In general this is *not* the same as

$$\mathcal{A}^2 + 2\mathcal{B}\mathcal{A} + \mathcal{B}^2,$$

and so the splitting error is

$$u(x, k) - u^{**}(x, k) = \frac{1}{2}k^2(\mathcal{A}\mathcal{B} - \mathcal{B}\mathcal{A})u(x, 0) + O(k^3). \quad (18.27)$$

The splitting error is zero only in the special case when the differential operators  $\mathcal{A}$  and  $\mathcal{B}$  commute (in which case it turns out that all the higher order terms in the splitting error also vanish).

**Example 18.3.** For the problem considered in Example 18.1

$$\mathcal{A} = -a\partial_x \quad \text{and} \quad \mathcal{B} = -\lambda.$$

We then have  $\mathcal{A}\mathcal{B}u = \mathcal{B}\mathcal{A}u = a\lambda u_x$ . These operators commute for  $\lambda$  constant and there is no splitting error.

**Example 18.4.** Now suppose  $\lambda = \lambda(x)$  depends on  $x$  as in Example 18.2. Then we have

$$\mathcal{A}\mathcal{B}u = a\partial_x(\lambda(x)u) = a\lambda(x)u_x + a\lambda'(x)u$$

while

$$\mathcal{B}\mathcal{A}u = \lambda(x)au_x.$$

These are not the same unless  $\lambda'(x) = 0$ . In general the splitting error will be

$$u(x, k) - u^{**}(x, k) = \frac{1}{2}k^2a\lambda'(x)u(x, 0) + O(k^3).$$

If we now design a fractional step method based on this splitting, we will see that the splitting error alone will introduce an  $O(k^2)$  error in each time step, which can be expected to accumulate to an  $O(k)$  error after the  $T/k$  time steps needed to reach some fixed time  $T$  (in the best case, assuming the method is stable). Hence even if we solve each subproblem *exactly* within the fractional step method, the resulting method will be only first order accurate. If the subproblems are actually solved with numerical methods that are  $p$ 'th order accurate, the solution will still only be first order accurate no matter how large  $p$  is.

### 18.3 Strang splitting

It turns out that a slight modification of the splitting idea will yield second order accuracy quite generally (assuming each subproblem is solved with a method of at least this accuracy). The idea is to solve the first subproblem  $u_t = \mathcal{A}u$  over only a half time step of length  $k/2$ . Then we use the result as data for a full time step on the second subproblem  $u_t = \mathcal{B}u$ , and finally take another half time step on  $u_t = \mathcal{A}u$ . We can equally well reverse the roles of  $\mathcal{A}$  and  $\mathcal{B}$  here. This approach is often called *Strang splitting* as it was popularized in a paper by Strang[Str68] on solving multi-dimensional problems.

To analyze the Strang splitting, note that we are now approximating the solution operator  $e^{k(\mathcal{A}+\mathcal{B})}$  by  $e^{\frac{1}{2}k\mathcal{A}}e^{k\mathcal{B}}e^{\frac{1}{2}k\mathcal{A}}$ . Taylor series expansion of this product shows that

$$\begin{aligned} e^{\frac{1}{2}k\mathcal{A}}e^{k\mathcal{B}}e^{\frac{1}{2}k\mathcal{A}} &= \left( I + \frac{1}{2}k\mathcal{A} + \frac{1}{8}k^2\mathcal{A}^2 + \cdots \right) \left( I + k\mathcal{B} + \frac{1}{2}k^2\mathcal{B}^2 + \cdots \right) \left( I + \frac{1}{2}k\mathcal{A} + \frac{1}{8}k^2\mathcal{A}^2 + \cdots \right) \\ &= I + k(\mathcal{A} + \mathcal{B}) + \frac{1}{2}k^2(\mathcal{A}^2 + \mathcal{A}\mathcal{B} + \mathcal{B}\mathcal{A} + \mathcal{B}^2) + O(k^3). \end{aligned} \quad (18.28)$$

Comparing with (18.23), we see that the  $O(k^2)$  term is now captured correctly. The  $O(k^3)$  term is not correct in general, however, unless  $\mathcal{A}\mathcal{B} = \mathcal{B}\mathcal{A}$ .

**Exercise 18.1** Compute the  $O(k^3)$  term in the splitting error for the Strang splitting.

Note that over several time steps we can simplify the expression obtained with the Strang splitting. After  $n$  steps we have

$$U^n = \left( e^{\frac{1}{2}k\mathcal{A}}e^{k\mathcal{B}}e^{\frac{1}{2}k\mathcal{A}} \right) \left( e^{\frac{1}{2}k\mathcal{A}}e^{k\mathcal{B}}e^{\frac{1}{2}k\mathcal{A}} \right) \cdots \left( e^{\frac{1}{2}k\mathcal{A}}e^{k\mathcal{B}}e^{\frac{1}{2}k\mathcal{A}} \right) U^0 \quad (18.29)$$

repeated  $n$  times. Dropping the parentheses and noting that  $e^{\frac{1}{2}k\mathcal{A}}e^{\frac{1}{2}k\mathcal{A}} = e^{k\mathcal{A}}$ , we obtain

$$U^n = e^{\frac{1}{2}k\mathcal{A}}e^{k\mathcal{B}}e^{k\mathcal{A}}e^{k\mathcal{B}}e^{k\mathcal{A}} \cdots e^{k\mathcal{B}}e^{\frac{1}{2}k\mathcal{A}}U^0. \quad (18.30)$$

This differs from the first order splitting only in the fact that we start and end with a half time step on  $\mathcal{A}$ , rather than starting with a full step and ending with  $\mathcal{B}$ .

Another way to achieve this same effect is to simply take steps of length  $k$  on each problem, as in the first-order splitting, but to alternate the order of these steps in alternate time steps, e.g.,

$$U^1 = e^{k\mathcal{B}}e^{k\mathcal{A}}U^0$$

$$U^2 = e^{k\mathcal{A}}e^{k\mathcal{B}}U^1$$

$$U^3 = e^{k\mathcal{B}}e^{k\mathcal{A}}U^2$$

$$U^4 = e^{k\mathcal{A}}e^{k\mathcal{B}}U^3$$

etc.

If we take an even number of time steps, then we obtain

$$\begin{aligned} U^n &= (e^{k\mathcal{A}}e^{k\mathcal{B}})(e^{k\mathcal{B}}e^{k\mathcal{A}})(e^{k\mathcal{A}}e^{k\mathcal{B}})(e^{k\mathcal{B}}e^{k\mathcal{A}}) \cdots (e^{k\mathcal{A}}e^{k\mathcal{B}})(e^{k\mathcal{B}}e^{k\mathcal{A}})U^0 \\ U^n &= e^{k\mathcal{A}}(e^{k\mathcal{B}}e^{k\mathcal{B}})(e^{k\mathcal{A}}e^{k\mathcal{A}})(e^{k\mathcal{B}}e^{k\mathcal{B}}) \cdots (e^{k\mathcal{B}}e^{k\mathcal{B}})e^{k\mathcal{A}}U^0. \end{aligned}$$

Since  $e^{k\mathcal{B}}e^{k\mathcal{B}} = e^{2k\mathcal{B}}$ , this is essentially the same as (18.29) but with  $\frac{1}{2}k$  replaced by  $k$ .

The fact that the Strang splitting is so similar to the first order splitting suggests that the first order splitting is not really so bad, and in fact it is not. While formally only first order accurate, the coefficient of the  $O(k)$  term may be much smaller than coefficients in the second order terms arising from discretization of  $e^{k\mathcal{A}}$  and  $e^{k\mathcal{B}}$ .



# Part II

# Appendices



# Appendix A1

## Measuring Errors

In order to discuss the accuracy of a numerical solution, or the relative virtues of one numerical method, vs. another, it is necessary to choose a manner of measuring that error. It may seem obvious what is meant by the error, but as we will see there are often many different ways to measure the error which can sometimes give quite different impressions as to the accuracy of an approximate solution.

### A1.1 Errors in a scalar value

First consider a problem in which the answer is a single value  $z \in \mathbb{R}$ . Consider, for example, the scalar ODE

$$u'(t) = f(u(t)), \quad u(0) = \eta$$

and suppose we are trying to compute the solution at some particular time  $T$ , so  $z = u(T)$ . Denote the computed solution by  $\hat{z}$ . Then the error in this computed solution is

$$E = \hat{z} - z.$$

#### A1.1.1 Absolute error

A natural measure of this error would be the absolute value of  $E$ ,

$$|E| = |\hat{z} - z|.$$

This is called the *absolute error* in the approximation.

As an example, suppose that  $z = 2.2$  while some numerical method produced a solution  $\hat{z} = 2.20345$ . Then the absolute error is

$$|\hat{z} - z| = 0.00345 = 3.45 \times 10^{-3}.$$

This seems quite reasonable — we have a fairly accurate solution with three correct digits and the absolute error is fairly small, on the order of  $10^{-3}$ . We might be very pleased with an alternative method that produced an error of  $10^{-6}$  and horrified with a method that produced an error of  $10^6$ .

But note that our notion of what is a large error or a small error might be thrown off completely if we were to choose a different set of units for measuring  $z$ . For example, suppose the  $z$  discussed above were measured in meters, so  $z = 2.2$  meters is the correct solution. But suppose that instead we expressed the solution (and the approximate solution) in nanometers rather than meters. Then the true solution is  $z = 2.2 \times 10^9$  and the approximate solution is  $\hat{z} = 2.20345 \times 10^9$ , giving an absolute error of

$$|\hat{z} - z| = 3.45 \times 10^6.$$

We have an error that seems huge and yet the solution is just as accurate as before, with three correct digits.

Conversely, if we measured  $z$  in kilometers then  $z = 2.2 \times 10^{-3}$  and  $\hat{z} = 2.20345 \times 10^{-3}$  so

$$|\hat{z} - z| = 3.45 \times 10^{-6}.$$

The error seems much smaller and yet there are still only three correct digits.

### A1.1.2 Relative error

The above difficulties arise from a poor choice of scaling of the problem. One way to avoid this is to consider the *relative error*, defined by

$$\left| \frac{\hat{z} - z}{z} \right|.$$

The size of the error is scaled by the size of the value being computed. For the above examples, the relative error in  $\hat{z}$  is equal to

$$\left| \frac{2.20345 - 2.2}{2.2} \right| = \left| \frac{2.20345 \times 10^9 - 2.2 \times 10^9}{2.2 \times 10^9} \right| = 1.57 \times 10^{-3}$$

The value of the relative error is the same no matter what units we use to measure  $z$ , a very desirable feature. Also note that in general a relative error that is on the order of  $10^{-k}$  indicates that there are roughly  $k$  correct digits in the solution, matching our intuition.

For these reasons the relative error is often a better measure of accuracy than the absolute error. Of course if we know that our problem is “properly” scaled, so that the solution  $z$  has magnitude order 1, then it is fine to use the absolute error, which is roughly the same as the relative error in this case.

In fact it is generally better to insure that the problem is properly scaled than to rely on the relative error. Poorly scaled problems can lead to other numerical difficulties, particularly if several different scales arise in the same problem so that some numbers are orders of magnitude larger than others for nonphysical reasons. Unless otherwise noted below, we will assume that the problem is scaled in such a way that the absolute error is meaningful.

## A1.2 “Big-oh” and “little-oh” notation

In discussing the rate of convergence of a numerical method we use the notation  $O(h^p)$ , the so-called “big-oh” notation. In case this is unfamiliar, here is a brief review of the proper use of this notation.

If  $f(h)$  and  $g(h)$  are two functions of  $h$  then we say that

$$f(h) = O(g(h)) \text{ as } h \rightarrow 0$$

if there is some constant  $C$  such that

$$\left| \frac{f(h)}{g(h)} \right| < C \text{ for all } h \text{ sufficiently small,}$$

or equivalently, if we can bound

$$|f(h)| < C|g(h)| \text{ for all } h \text{ sufficiently small.}$$

Intuitively, this means that  $f(h)$  decays to zero *at least as fast* as the function  $g(h)$  does. Usually  $g(h)$  is some monomial  $h^q$ , but this isn’t necessary.

It is also sometimes convenient to use the “little-oh” notation

$$f(h) = o(g(h)) \text{ as } h \rightarrow 0.$$

This means that

$$\left| \frac{f(h)}{g(h)} \right| \rightarrow 0 \text{ as } h \rightarrow 0.$$

This is slightly stronger than the previous statement, and means that  $f(h)$  decays to zero *faster* than  $g(h)$ . If  $f(h) = o(g(h))$  then  $f(h) = O(g(h))$  though the converse may not be true. Saying that  $f(h) = o(1)$  simply means that the  $f(h) \rightarrow 0$  as  $h \rightarrow 0$ .

**Examples:**

$$2h^3 = O(h^2) \text{ as } h \rightarrow 0, \text{ since } \frac{2h^3}{h^2} = 2h < 1 \text{ for all } h < 1/2.$$

$$2h^3 = o(h^2) \text{ as } h \rightarrow 0, \text{ since } 2h \rightarrow 0 \text{ as } h \rightarrow 0.$$

$$\sin(h) = O(h) \text{ as } h \rightarrow 0, \text{ since } \sin h = h - \frac{h^3}{3} + \frac{h^5}{5} + \dots < h \text{ for all } h > 0.$$

$$\sin(h) = h + o(h) \text{ as } h \rightarrow 0, \text{ since } (\sin h - h)/h = O(h^2).$$

$$\sqrt{h} = O(1) \text{ as } h \rightarrow 0, \text{ and also } \sqrt{h} = o(1), \text{ but } \sqrt{h} \text{ is not } O(h).$$

$$1 - \cos h = o(h) \text{ and } 1 - \cos h = O(h^2) \text{ as } h \rightarrow 0.$$

$$e^{-1/h} = o(h^q) \text{ as } h \rightarrow 0 \text{ for every value of } q.$$

$$\text{To see this, let } x = 1/h \text{ then } \frac{e^{-1/h}}{h^q} = e^{-x} x^q \rightarrow 0 \text{ as } x \rightarrow \infty.$$

Note that saying  $f(h) = O(g(h))$  is a statement about how  $f$  behaves in the limit as  $h \rightarrow 0$ . This notation is sometimes abused by saying, for example, that if  $h = 10^{-3}$  then the number  $3 \times 10^{-6}$  is  $O(h^2)$ . Though it is clear what is meant, this is really meaningless mathematically and may be misleading when analyzing the accuracy of a numerical method. If the error  $E(h)$  on a grid with  $h = 10^{-3}$  turns out to be  $3 \times 10^{-6}$ , we cannot conclude that the method is second order accurate. It could be, for example, that the error  $E(h)$  has the behavior

$$E(h) = 0.003 h \tag{A1.1}$$

in which case  $E(10^{-3}) = 3 \times 10^{-6}$  but it is not true that  $E(h) = O(h^2)$ . In fact the method is only first order accurate, which would become apparent as we refined the grid.

Conversely, if

$$E(h) = 10^6 h^2 \tag{A1.2}$$

then  $E(10^{-3}) = 1$  which is much larger than  $h^2$ , and yet it is still true that

$$E(h) = O(h^2) \text{ as } h \rightarrow 0.$$

Also note that there is more to the choice of a method than its asymptotic rate of convergence. While in general a second order method outperforms a first order method, if we are planning to compute on a grid with  $h = 10^{-3}$  then we would prefer a first order method with error (A1.1) over a second order method with error (A1.2).

### A1.3 Errors in vectors

Now suppose  $z \in \mathbb{R}^m$  is a vector with  $m$  components, for example the solution to a system of  $m$  ODE's at some particular fixed time  $T$ . Then  $\hat{z}$  is a vector of approximate values and the error  $e = \hat{z} - z$  is also a vector in  $\mathbb{R}^m$ . In this case we can use some vector norm to measure the error.

There are many ways to define a vector norm. In general a vector norm is simply a mapping from vectors  $x$  in  $\mathbb{R}^m$  to nonnegative real numbers, satisfying the following conditions (which generalize important properties of the absolute value for scalars):



1.  $\|x\| \geq 0$  for any  $x \in \mathbb{R}^m$ , and  $\|x\| = 0$  if and only if  $x = \vec{0}$ .
2. If  $a$  is any scalar then  $\|ax\| = |a| \|x\|$ .
3. If  $x, y \in \mathbb{R}^m$ , then  $\|x + y\| \leq \|x\| + \|y\|$ . (Triangle inequality)

One common choice is the max-norm (or infinity-norm) denoted by  $\|\cdot\|_\infty$ :

$$\|e\|_\infty = \max_{1 \leq i \leq m} |e_i|.$$

It is easy to verify that  $\|\cdot\|_\infty$  satisfies the required properties. A bound on the max-norm of the error is nice because we know that every component of the error can be no greater than the max-norm. For some problems, however, there are other norms which are either more appropriate or easier to bound using our analytical tools.

Two other norms that are frequently used are the 1-norm and 2-norm,

$$\|e\|_1 = \sum_{i=1}^m |e_i| \quad \text{and} \quad \|e\|_2 = \sqrt{\sum_{i=1}^m |e_i|^2}. \quad (\text{A1.3})$$

These are special cases of the general family of  $p$ -norms, defined by

$$\|e\|_p = \left[ \sum_{i=1}^m |e_i|^p \right]^{1/p}. \quad (\text{A1.4})$$

Note that the max-norm can be obtained as the limit as  $p \rightarrow \infty$  of the  $p$ -norm.

### A1.3.1 Norm equivalence

With so many different norms to choose from, it is natural to ask whether results on convergence of numerical methods will depend on our choice of norm. Suppose  $e^h$  is the error obtained with some step size  $h$ , and that  $\|e^h\| = O(h^q)$  in some norm, so that the method is  $q$ th order accurate. Is it possible that the rate will be different in some other norm? The answer is “no”, due to the following result on the “equivalence” of all norms on  $\mathbb{R}^m$ . (Note that this result is only valid as long as the dimension  $m$  of the vector is fixed as  $h \rightarrow 0$ . See Section A1.5 for an important case where the length of the vector depends on  $h$ .)

Let  $\|\cdot\|_\alpha$  and  $\|\cdot\|_\beta$  represent two different vector norms on  $\mathbb{R}^m$ . Then there exist two constants  $C_1$  and  $C_2$  such that

$$C_1 \|x\|_\alpha \leq \|x\|_\beta \leq C_2 \|x\|_\alpha \quad (\text{A1.5})$$

for all vectors  $x \in \mathbb{R}^m$ . For example, it is fairly easy to verify that the following relations hold among the norms mentioned above:

$$\|x\|_\infty \leq \|x\|_1 \leq m \|x\|_\infty \quad (\text{A1.6a})$$

$$\|x\|_\infty \leq \|x\|_2 \leq \sqrt{m} \|x\|_\infty \quad (\text{A1.6b})$$

$$\|x\|_2 \leq \|x\|_1 \leq \sqrt{m} \|x\|_2. \quad (\text{A1.6c})$$

Now suppose that  $\|e^h\|_\alpha \leq Ch^q$  as  $h \rightarrow 0$  in some norm  $\|\cdot\|_\alpha$ . Then we have

$$\|e^h\|_\beta \leq C_2 \|e^h\|_\alpha \leq C_2 Ch^q$$

and so  $\|e^h\|_\beta = O(h^q)$  as well. In particular, if  $\|e^h\| \rightarrow 0$  in some norm then the same is true in any other norm and so the notion of “convergence” is independent of our choice of norm. This will *not* be true in Section A1.4, where we consider approximating functions rather than vectors.

### A1.3.2 Matrix norms

For any vector norm  $\|\cdot\|$  we can define a corresponding matrix norm. The norm of a matrix  $A \in \mathbb{R}^{m \times m}$  is denoted by  $\|A\|$  and has the property that  $C = \|A\|$  is the *smallest* value of the constant  $C$  for which the bound

$$\|Ax\| \leq C\|x\| \quad (\text{A1.7})$$

holds for *every* vector  $x \in \mathbb{R}^m$ . Hence  $\|A\|$  is defined by

$$\|A\| = \max_{\substack{x \in \mathbb{R}^m \\ x \neq \vec{0}}} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \|Ax\|.$$

It would be rather difficult to calculate  $\|A\|$  from the above definitions, but for the most commonly used norms there are simple formulas for computing  $\|A\|$  directly from the matrix:

$$\|A\|_1 = \max_{1 \leq j \leq m} \sum_{i=1}^m |a_{ij}| \quad (\text{maximum column sum}) \quad (\text{A1.8a})$$

$$\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^m |a_{ij}| \quad (\text{maximum row sum}) \quad (\text{A1.8b})$$

$$\|A\|_2 = \sqrt{\rho(A^T A)}. \quad (\text{A1.8c})$$

In the definition of the 2-norm,  $\rho(B)$  denotes the spectral radius of the matrix  $B$  (the maximum modulus of an eigenvalue). In particular, if  $A = A^T$  is symmetric, then  $\|A\|_2 = \rho(A)$ .

## A1.4 Errors in functions

Now consider a problem in which the solution is a function  $u(x)$  over some interval  $a \leq x \leq b$  rather than a single value or vector. Some numerical methods, such as finite element or collocation methods, produce an approximate solution  $\hat{u}(x)$  which is also a function. Then the error is given by a function

$$e(x) = \hat{u}(x) - u(x).$$

We can measure the magnitude of this error using standard function space norms, which are quite analogous to the vector norms described above. For example, the max-norm is given by

$$\|e\|_\infty = \max_{a \leq x \leq b} |e(x)|. \quad (\text{A1.9})$$

The 1-norm and 2-norms are given by integrals over  $[a, b]$  rather than by sums over the vector elements:

$$\|e\|_1 = \int_a^b |e(x)| dx, \quad (\text{A1.10})$$

$$\|e\|_2 = \left( \int_a^b |e(x)|^2 dx \right)^{1/2}. \quad (\text{A1.11})$$

These are again special cases of the general  $p$ -norm, defined by

$$\|e\|_p = \left( \int_a^b |e(x)|^p dx \right)^{1/p}. \quad (\text{A1.12})$$

## A1.5 Errors in grid functions

Finite difference methods do not produce a function  $\hat{u}(x)$  as an approximation to  $u(x)$ . Instead they produce a set of values  $U_i$  at grid points  $x_i$ . For example, on a uniform grid with  $N + 1$  equally spaced points with spacing  $h = (b - a)/N$ ,

$$x_i = a + ih, \quad i = 0, 1, \dots, N,$$

our approximation to  $u(x)$  would consist of the  $N + 1$  values  $(U_0, U_1, \dots, U_N)$ . How can we measure the error in this approximation? We want to compare a set of discrete values with a function.

We must first decide what the values  $U_i$  are supposed to be approximating. Often the value  $U_i$  is meant to be interpreted as an approximation to the pointwise value of the function at  $x_i$ , so  $U_i \approx u(x_i)$ . In this case it is natural to define a vector of errors  $e = (e_0, e_1, \dots, e_N)$  by

$$e_i = U_i - u(x_i).$$

This is not always the proper interpretation of  $U_i$ , however. For example, some numerical methods are derived using the assumption that  $U_i$  approximates the average value of  $u(x)$  over an interval of length  $h$ , e.g.,

$$U_i \approx \frac{1}{h} \int_{x_{i-1}}^{x_i} u(x) dx, \quad i = 1, 2, \dots, N.$$

In this case it would be more appropriate to compare  $U_i$  to this cell average in defining the error. Clearly the errors will be different depending on what definition we adopt, and may even exhibit different convergence rates (see Example 3.1 below), so it is important to make the proper choice for the method being studied.

Once we have defined the vector of errors  $(e_0, \dots, e_N)$ , we can measure its magnitude using some norm. Since this is simply a vector with  $N + 1$  components, it would be tempting to simply use one of the vector norms discussed above, e.g.,

$$\|e\|_1 = \sum_{i=0}^N |e_i|. \tag{A1.13}$$

However, this choice would give a very misleading idea of the magnitude of the error. The quantity in (A1.13) can be expected to be roughly  $N$  times as large as the error at any single grid point and here  $N$  is not the dimension of some physically relevant space, but rather the number of points on our grid. If we refine the grid and increase  $N$ , then the quantity (A1.13) might well *increase* even if the error at each grid point decreases, which is clearly not the correct behavior.

Instead we should define the norm of the error by discretizing the integral in (A1.10), which is motivated by considering the vector  $(e_0, \dots, e_N)$  as a discretization of some error function  $e(x)$ . This suggests defining

$$\|e\|_1 = h \sum_{i=0}^N |e_i| \tag{A1.14}$$

with the factor of  $h$  corresponding to the  $dx$  in the integral. Note that since  $h = (b - a)/N$ , this scales the sum by  $1/N$  as the number of grid points increases, so that  $\|e\|_1$  is the average value of  $e$  over the interval (times the length of the interval), just as in (A1.10). The norm (A1.14) will be called a *grid-function norm* and is distinct from the related vector norm. The set of values  $(e_0, \dots, e_N)$  will sometimes be called a *grid function* to remind us that it is a special kind of vector that represents the discretization of a function.

Similarly, the  $p$ -norm should be scaled by  $h^{1/p}$ , so that the  $p$ -norm for grid functions is

$$\|e\|_p = \left( h \sum_{i=0}^N |e_i|^p \right)^{1/p}. \tag{A1.15}$$

Since  $h^{1/p} \rightarrow 1$  as  $p \rightarrow \infty$ , the max-norm remains unchanged:

$$\|e\|_\infty = \max_{0 \leq i \leq N} |e_i|,$$

which makes sense from (A1.9).

In two space dimensions we have analogous norms of functions and grid functions, e.g.,

$$\begin{aligned} \|e\|_p &= \left( \iint |e(x, y)|^p dx dy \right)^{1/p} && \text{for functions} \\ \|e\|_p &= \left( \Delta x \Delta y \sum_i \sum_j |e_{ij}|^p \right)^{1/p} && \text{for grid functions} \end{aligned}$$

with the obvious extension to more dimensions.

### A1.5.1 Norm equivalence

Note that we still have an equivalence of norms in the sense that, for any *fixed*  $N$  (and hence fixed  $h$ ), there are constants  $C_1$  and  $C_2$  such that

$$C_1 \|x\|_\alpha \leq \|x\|_\beta \leq C_2 \|x\|_\alpha$$

for any vector  $e \in \mathbb{R}^{N+1}$ . For example, translating (A1.6a) to the context of grid-function norms gives the bounds

$$h \|e\|_\infty \leq \|e\|_1 \leq Nh \|e\|_\infty = (b - a) \|e\|_\infty, \tag{A1.16a}$$

$$\sqrt{h} \|e\|_\infty \leq \|e\|_2 \leq \sqrt{Nh} \|e\|_\infty = \sqrt{b - a} \|e\|_\infty, \tag{A1.16b}$$

$$\sqrt{h} \|e\|_2 \leq \|e\|_1 \leq \sqrt{Nh} \|e\|_2 = \sqrt{b - a} \|e\|_2. \tag{A1.16c}$$

However, since these constants may depend on  $N$  and  $h$ , this equivalence does not carry over when we consider the behavior of the error as we refine the grid so that  $h \rightarrow 0$  and  $N \rightarrow \infty$ .

We are particularly interested in the convergence rate of a method, and would like to show that

$$\|e^h\| \leq O(h^q)$$

for some  $q$ . In the last section we saw that the rate is independent of the choice of norm if  $e^h$  is a vector in the space  $\mathbb{R}^m$  with fixed dimension  $m$ . But now  $m = N + 1$  and grows as  $h \rightarrow 0$ , and as a result the rate may be quite different in different norms. This is particularly noticeable if we approximate a discontinuous function, as the following example shows.

**Example 3.1.** Set

$$u(x) = \begin{cases} 0 & x \leq \frac{1}{2} \\ 1 & x > \frac{1}{2} \end{cases}$$

Let  $N$  be even and let

$$U_i^h = \begin{cases} 0 & i < N/2 \\ \frac{1}{2} & i = N/2 \\ 1 & i > N/2 \end{cases}$$

be the discrete approximation on the grid with spacing  $h = 1/N$  on the interval  $0 \leq x \leq 1$ . This is illustrated in Figure A1.1 for  $N = 8$ . Define the error  $e_i^h$  by

$$e_i^h = U_i^h - u(x_i) = \begin{cases} \frac{1}{2} & i = N/2 \\ 0 & \text{otherwise.} \end{cases}$$

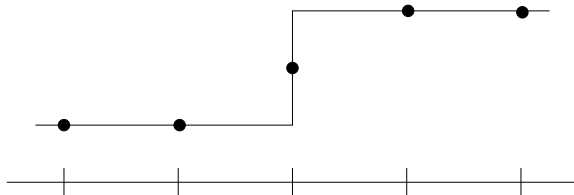


Figure A1.1: The function  $u(x)$  and the discrete approximation.

No matter how fine the grid is, there is always an error of magnitude  $1/2$  at  $i = N/2$  and hence

$$\|e^h\|_\infty = \frac{1}{2} \text{ for all } h.$$

On the other hand, in the 1-norm (A1.14) we have

$$\|e^h\|_1 = h/2 = O(h) \text{ as } h \rightarrow 0.$$

We see that the 1-norm converges to zero as  $h$  goes to zero while the max-norm does not.

How should we interpret this? Should we say that  $U^h$  is a first order accurate approximation to  $u(x)$  or should we say that it does not converge? It depends on what we are looking for. If it is really important that the maximum error over all grid points be uniformly small, then the max-norm is the appropriate norm to use and the fact that  $\|e^h\|_\infty$  does not approach zero tells us that we are not achieving our goal. On the other hand this may not really be required, and in fact this example illustrates that it is unrealistic to expect pointwise convergence in problems where the function is discontinuous. For many purposes the approximation shown in Figure A1.1 would be perfectly acceptable.

This example also illustrates the effect of choosing a different definition of the “error”. If we were to define the error by

$$e_i^h = U_i^h - \frac{1}{h} \int_{x_i-h/2}^{x_i+h/2} u(x) dx,$$

then we would have  $e_i^h \equiv 0$  for all  $i$  and  $h$  and  $\|e^h\| = 0$  in every norm, including the max-norm. With this definition of the error our approximation is not only acceptable, it is the best possible approximation.

## Appendix A2

# Estimating errors in numerical solutions

When developing a computer program to solve a differential equation, it is generally a good idea to test the code and ensure that it is producing correct results with the expected accuracy. How can we do this?

A first step is often to try the code on a problem for which the exact solution is known, in which case we can compute the error in the numerical solution exactly. Not only can we then check that the error is small on some grid, we can also refine the grid and check how the error is behaving asymptotically, to verify that the expected order of accuracy and perhaps even error constant are seen. Of course one must be aware of some of the issues raised in Appendix A1, *e.g.*, that the expected order may only appear for  $h$  sufficiently small.

It is important to test a computer program by doing grid refinement studies even if the results look quite good on one particular grid. A subtle error in programming (or in deriving the difference equations or numerical boundary conditions) can lead to a program that gives reasonable results and may even converge to the correct solution, but at less than the optimal rate. Consider, for example, the First Attempt of Section 2.11.

Of course in practice we are usually trying to solve a problem for which we do not know the exact solution, or we wouldn't bother with a numerical method in the first place. However, there are often simplified versions of the problem for which exact solutions are known, and a good place to start is with these special cases. They may reveal errors in the code that will affect the solution of the real problem as well.

This is generally not sufficient however, even when it is possible, since in going from the easy special case to the real problem there may be new errors introduced. How do we estimate the error in a numerical solution if we do not have the exact solution to compare it with?

The standard approach, when we can afford to, is to compute a numerical solution on a very fine grid and use this as a “reference solution” (or “fine-grid” solution). This can be used as a good approximation to the exact solution in estimating the error on other, much coarser, grids. When the fine grid is fine enough, we can obtain good estimates not only for the errors, but also for the order of accuracy. See Section A2.2.

Often we cannot afford to take very fine grids, especially in more than one space dimension. We may then be tempted to use a grid that is only slightly finer than the grid we are testing in order to generate a reference solution. When done properly this approach can also yield accurate estimates of the order of accuracy, but more care is required. See Section A2.3 below.

## A2.1 Estimates from the true solution

First suppose we know the true solution. Let  $E(h)$  denote the error in the calculation with grid spacing  $h$ , as computed using the true solution. In this chapter we suppose that  $E(h)$  is a scalar, typically some norm of the error over the grid, i.e.,

$$E(h) = \|q^h - \hat{q}^h\|$$

where  $q^h$  is the numerical solution vector (grid function) and  $\hat{q}^h$  is the true solution evaluated on the same grid.

If the method is  $p$ 'th order accurate then we expect

$$E(h) = Ch^p + o(h^p) \quad \text{as } h \rightarrow 0,$$

and if  $h$  is sufficiently small then

$$E(h) \approx Ch^p. \tag{A2.1}$$

If we refine the grid by a factor of 2, say, then we expect

$$E(h/2) \approx C(h/2)^p.$$

Defining the *error ratio*

$$R(h) = E(h) / E(h/2), \tag{A2.2}$$

we expect

$$R(h) \approx 2^p, \tag{A2.3}$$

and hence

$$p \approx \log_2(R(h)). \tag{A2.4}$$

Here refinement by a factor of 2 is used only as an example, since this choice is often made in practice. But more generally if  $h_1$  and  $h_2$  are any two grid spacings, then we can estimate  $p$  based on calculations on these two grids using

$$p \approx \frac{\log(E(h_1)/E(h_2))}{\log(h_1/h_2)}. \tag{A2.5}$$

Hence we can estimate the order  $p$  based on any two calculations. (This will only be valid if  $h$  is small enough that (A2.1) holds, of course.)

Note that we can also estimate the error constant  $C$  by

$$C \approx E(h)/h^p$$

once  $p$  is known.

## A2.2 Estimates from a fine-grid solution

Now suppose we don't know the exact solution but that we can afford to run the problem on a very fine grid, say with grid spacing  $\bar{h}$ , and use this as a reference solution in computing the errors on some sequence of much coarser grids. In order to compare  $q^h$  on the coarser grid with  $q^{\bar{h}}$  on the fine grid, we need to make sure that these two grids contain coincident grid points where we can directly compare the solutions. Typically we choose the grids in such a way that all grid points on the coarser grid are

also fine grid points. (This is often the hardest part of doing such grid refinement studies — getting the grids and indexing correct.)

Let  $\bar{q}^h$  be the restriction of the fine-grid solution to the  $h$ -grid, so that we can define the approximate error  $\bar{E}(h) \equiv \|q^h - \bar{q}^h\|$ , analogous to the true error  $E(h) = \|q^h - \hat{q}^h\|$ . What is the error in this approximate error? We have

$$q^h - \bar{q}^h = (q^h - \hat{q}^h) + (\hat{q}^h - \bar{q}^h)$$

If the method is supposed to be  $p$ 'th order accurate and  $\bar{h}^p \ll h^p$ , then the second term on the right hand side (the true error on the  $\bar{h}$ -grid) should be negligible compared to the first term (the true error on the  $h$ -grid) and  $\bar{E}(h)$  should give a very accurate estimate of the error.

WARNING: Estimating the error and testing the order of accuracy by this approach only confirms that the code is converging to *some* function with the desired rate. It is perfectly possible that the code is converging very nicely to the *wrong* function. Consider a second-order accurate method applied to a two-point boundary value problem, for example, and suppose that we code everything properly except that we mistype the value of one of the boundary values. Then a grid-refinement study of this type would show that the method is converging with second order accuracy, as indeed it is. The fact that it is converging to the solution of the wrong problem would not be revealed by this test. One must use other tests as well, not least of which is checking that the computed solutions make sense physically, e.g., that the correct boundary conditions are in fact satisfied.

More generally, a good understanding of the problem being solved, a knowledge of how the solution should behave, good physical intuition and common sense are all necessary components in successful scientific computing. Don't believe the numbers coming out simply because they are generated by a computer, even if the computer also tells you that they are second order accurate!

## A2.3 Estimates from coarser solutions

Now suppose that our computation is very expensive even on relatively coarse grids, and we cannot afford to run a calculation on a much finer grid in order to test the order of accuracy. Suppose, for example, that we are only willing to run the calculation on grids with spacing  $h$ ,  $h/2$  and  $h/4$ , and wish to estimate the order of accuracy from these three calculations, without using any finer grids. Since we can estimate the order from any two values of the error, we could define the errors in the two coarser grid calculations by using the  $h/4$  calculation as our reference solution. Will we get a good estimate for the order?

In the notation used above, we now have  $\bar{h} = h/4$  while  $h = 4\bar{h}$  and  $h/2 = 2\bar{h}$ . Assuming the method is  $p$ 'th order accurate and that  $h$  is small enough that (A2.1) is valid (a poor assumption, perhaps, if we are using very coarse grids!), we expect

$$\begin{aligned} \bar{E}(h) &= E(h) - E(\bar{h}) \\ &\approx Ch^p - C\bar{h}^p \\ &= (4^p - 1)C\bar{h}^p. \end{aligned}$$

Similarly,

$$\bar{E}(h/2) \approx (2^p - 1)C\bar{h}^p.$$

The ratio of approximate errors is thus

$$\bar{R}(h) = \bar{E}(h)/\bar{E}(h/2) \approx \frac{4^p - 1}{2^p - 1} = 2^p + 1.$$

This differs significantly from (A2.3). For a first-order accurate method with  $p = 1$ , we now have  $\bar{R}(h) \approx 3$  and we should expect the apparent error to decrease by a factor of 3 when we go from  $h$  to  $h/2$ , not by the factor of 2 that we normally expect. For a second-order method we expect a factor



of 5 improvement rather than a factor of 4. This increase in  $\bar{R}(h)$  results from the fact that we are comparing our numerical solutions to another approximate solution that has a similar error.

We can obtain a good estimate of  $p$  from such calculations (assuming (A2.1) is valid), but to do so we must calculate  $p$  by

$$p \approx \log_2(\bar{R}(h) - 1)$$

rather than by (A2.4). The approximation (A2.4) would overestimate the order of accuracy.

Again we have used refinement by factors of 2 only as an example. If the calculation is very expensive we might want to refine the grid more slowly, using for example  $h$ ,  $3h/4$  and  $h/2$ . One can develop appropriate approximations to  $p$  based on any three grids. The tricky part may be to estimate the error at grid points on the coarser grids if these are not also grid points on the  $\bar{h}$  grid. Interpolation can be used, but then one must be careful to insure that sufficiently accurate interpolation formulas are used that the error in interpolation does not contaminate the estimate of the error in the numerical method being studied.

## A2.4 Extrapolation methods

We have seen that we can estimate the dominant term of the error  $Ch^p$  by using the results of numerical calculations alone, without knowing the exact solution. This immediately suggests that we should be able to improve our computed solution on any of the grids used by subtracting a good estimate of the error. This is the basis of *extrapolation methods*, often called *Richardson extrapolation*, which is a very important technique for improving the accuracy with modest additional work. See [Kel76] or any introductory numerical analysis text for more information.

## Appendix A3

# Convergence of iterative methods

Consider the nonlinear equation

$$g(x) = 0 \tag{A3.1}$$

and suppose we want to solve for a root  $x^*$  numerically. Many iterative methods take the form

$$x^{[k+1]} = G(x^{[k]}) \tag{A3.2}$$

where  $G$  is some iteration function. The iteration (A3.2) is often called a *fixed point iteration* (FPI) since the root  $x^*$  should be a fixed point of  $G$ :

$$x^* = G(x^*). \tag{A3.3}$$

Sometimes the original function  $g(x)$  has a form that makes one choice of  $G$  clear (although the obvious choice may not be the best). For example, if we are using the Backward Euler method to solve the IVP  $y' = f(y)$ , then

$$y_{n+1} = y_n + hf(y_{n+1}) \tag{A3.4}$$

and we must solve for  $y_{n+1}$ . The function  $g$  is given by

$$g(x) = x - y_n - hf(x)$$

and an obvious candidate for  $G$  is

$$G(x) = y_n + hf(x). \tag{A3.5}$$

In *Newton's method*,  $G$  is chosen to enhance convergence,

$$G(x) = x - g(x)/g'(x). \tag{A3.6}$$

Note that  $G(x^*) = x^*$  provided that  $g'(x^*) \neq 0$ . For example, Newton applied to Backward Euler gives

$$G(x) = x - (x - y_n - hf(x))/(1 - hf'(x)).$$

The secant method is more complicated since it defines  $x^{[k+1]}$  based on two previous values,

$$x^{[k+1]} = H(x^{[k]}, x^{[k-1]}).$$

**Convergence of FPI.** It is easy to determine when the FPI (A3.2) converges and how fast. Subtracting (A3.3) from (A3.2) gives

$$x^{[k+1]} - x^* = G(x^{[k]}) - G(x^*)$$

and so

$$\begin{aligned} |x^{[k+1]} - x^*| &= |G(x^{[k]}) - G(x^*)| \\ &\leq K|x^{[k]} - x^*| \end{aligned} \quad (\text{A3.7})$$

provided that  $G$  is Lipschitz continuous near the root  $x^*$ . If  $G$  is continuously differentiable then

$$K \approx |G'(x^*)|$$

provided  $x^{[k]}$  is close to  $x^*$ .

Let  $e^{[k]} = x^{[k]} - x^*$  be the error in iteration  $k$ . Then (A3.7) gives

$$|e^{[k+1]}| \leq K|e^{[k]}| \quad (\text{A3.8})$$

and so

$$|e^{[k]}| \leq K^k |e^{[0]}|.$$

The iteration converges provided  $K < 1$ . This requires in particular that  $|G'(x^*)| < 1$  and that we start close enough to  $x^*$ .

*Example.* For the backward Euler method (A3.4),  $G$  is given by (A3.5) and  $y_n$  is fixed, so

$$G'(x) = hf'(x).$$

If  $f$  is smooth then  $|G'| < 1$  provided  $h$  is small enough. (For stiff problems, however, this puts a restriction on  $h$  that is similar to what would be required by explicit Euler, so this is not really attractive and something like Newton's method must be used.)

**Convergence rate.** The relation (A3.8) shows that the FPI is in general *linearly convergent* if  $0 < K < 1$ , i.e.,  $|e^{[k+1]}|$  is approximately equal to a constant times  $|e^{[k]}|$  to the *first* power. The constant is roughly  $|G'(x^*)|$  as we approach  $x^*$ .

**Convergence of Newton's method.** Now consider Newton's method where  $G$  is given by (A3.6). In this case

$$G'(x) = \frac{g(x)g''(x)}{(g'(x))^2}$$

and so  $G'(x^*) = 0$  (provided  $g'(x^*) \neq 0$ ). So in (A3.7) we can obtain a better bound on the error by expanding

$$G(x^{[k]}) = G(x^*) + G'(x^*)(x^{[k]} - x^*) + \frac{1}{2}G''(x^*)(x^{[k]} - x^*)^2 + \dots$$

We find that

$$|G(x^{[k]}) - G(x^*)| \approx \frac{1}{2}|G''(x^*)||e^{[k]}|^2$$

which shows that Newton is *quadratically convergent*:

$$|e^{[k+1]}| \approx K|e^{[k]}|^2.$$

(This is still assuming  $g'(x^*) \neq 0$ . *Exercise:* what if  $g'(x^*) = 0$  but  $g''(x^*) \neq 0$ ?)

**Determining the rate of convergence from numerical results.** Suppose we apply some iterative method to  $g(x) = 0$  and get a table of errors  $e^{[k]} = x^{[k]} - x^*$ . How can we determine the rate of convergence? (For the moment assume we are testing our program on a problem where we know  $x^*$  so we can compute  $e^{[k]}$ . This is generally a good place to start.)

If we expect

$$|e^{[k+1]}| \approx K|e^{[k]}|^p$$

for some  $p$ , then taking logarithms gives

$$\log |e^{[k+1]}| \approx \log K + p \log |e^{[k]}| \quad (\text{A3.9})$$

so we expect the relation between  $\log|e^{[k]}|$  and  $\log|e^{[k+1]}|$  to be linear. Similarly we expect

$$\log|e^{[k+2]}| \approx \log K + p \log|e^{[k+1]}| \quad (\text{A3.10})$$

Note that (A3.9) and (A3.10) form a system of two linear equations for the unknowns  $\log K$  and  $p$ . So based on three iterates  $e^{[k]}$ ,  $e^{[k+1]}$  and  $e^{[k+2]}$  we can estimate these values. Of course the values we get will depend on which iterates we look at, but we would expect that as we converge to  $x^*$  these values will converge to the correct values.

**Example.** Applying fixed point iteration to the equation  $x = \cos(x)$  starting from  $x^{[0]} = 1$  gives the following results:

nu	x	err	K	p
0	0.1000000000000000D+01	0.2609D+00		
1	0.5403023058681397D+00	0.1988D+00		
2	0.8575532158463933D+00	0.1185D+00	0.2563D+01	1.9030
3	0.6542897904977793D+00	0.8480D-01	0.3365D+00	0.6461
4	0.7934803587425656D+00	0.5440D-01	0.1440D+01	1.3276
5	0.7013687736227565D+00	0.3772D-01	0.4163D+00	0.8248
6	0.7639596829006543D+00	0.2487D-01	0.1032D+01	1.1367
7	0.7221024250267078D+00	0.1698D-01	0.5022D+00	0.9169
8	0.7504177617637606D+00	0.1133D-01	0.8518D+00	1.0599
9	0.7314040424225099D+00	0.7681D-02	0.5703D+00	0.9615
10	0.7442373549005569D+00	0.5152D-02	0.7641D+00	1.0268
11	0.7356047404363474D+00	0.3480D-02	0.6155D+00	0.9823
12	0.7414250866101092D+00	0.2340D-02	0.7198D+00	1.0121
13	0.7375068905132428D+00	0.1578D-02	0.6424D+00	0.9920
14	0.7401473355678759D+00	0.1062D-02	0.6971D+00	1.0055
15	0.7383692041223232D+00	0.7159D-03	0.6573D+00	0.9963
16	0.7395672022122562D+00	0.4821D-03	0.6855D+00	1.0025
17	0.7387603198742113D+00	0.3248D-03	0.6653D+00	0.9983
18	0.7393038923969059D+00	0.2188D-03	0.6796D+00	1.0011
19	0.7389377567153445D+00	0.1474D-03	0.6694D+00	0.9992

Newton's method, applied to  $g(x) = x - \cos(x)$ , gives:

nu	x	err	K	p
0	0.10000000000000000D+01	0.2609D+00		
1	0.7503638678402439D+00	0.1128D-01		
2	0.7391128909113617D+00	0.2776D-04	0.1473D+00	1.9123
3	0.7390851333852841D+00	0.1701D-09	0.2163D+00	1.9980
4	0.7390851332151607D+00			

Actually for this problem we do not know what the true  $x^*$  is, and the tables above were computed by first running Newton for 4 iterations (at which point it had converged to machine precision) and then using this value  $x^{[4]}$  as  $x^*$ . For this reason we cannot compute the error in  $x^{[4]}$  in the above table for the Newton method. This shows that we can check the convergence rate even for problems where we don't know the exact solution.

Also note that even if we had not computed the errors, we could easily guess how accurate the solutions are just from the tables of  $x$  values, by seeing how many digits remain unchanged in later iterations. With Newton's method the number of correct digits roughly doubles in each iteration, as expected from quadratic convergence, while with the fixed point iteration it appears to take about 6 iterations to compute each additional digit. This is typical of linearly convergent methods — the number of iterations to compute each additional digit is independent of the accuracy already achieved. Computing one additional digit requires reducing the error by a factor of  $1/10$ , and doing so takes  $n$  iterations, where  $K^n = 1/10$ . In the above example,  $K \approx 0.68$  and so  $n \approx \log(0.1)/\log(0.68) = 5.97$ . To converge to 16 digits we would need roughly 16 times as many iterations, *i.e.*,  $n \approx \log(10^{-16})/\log(0.68) = 95.5$ .

The number of iterations needed depends on how small  $K$  is. When we start looking at iterative methods for linear systems we will see that most methods are only linearly convergent and that the constant  $K$  approaches 1 as the matrix becomes more ill-conditioned. If  $K = 0.99$  then it takes roughly 230 iterations for each digit, 3665 iterations for 16 digits.

## Appendix A4

# Some Matrix Properties

If  $A$  is a matrix that approximates a differential operator then its transpose  $A^T$  (with elements  $(A^T)_{ij} = A_{ji}$ ) approximates the adjoint operator and symmetric matrices (for which  $A = A^T$ ) correspond to self-adjoint operators. These relations will be discussed in this chapter along with the notion of positive definiteness for operators and matrices.

### A4.1 Adjoins and symmetry

If  $u, v \in \mathbb{R}^m$  are two vectors, then the inner product of  $u$  and  $v$  is

$$(u, v) = u^T v = v^T u = \sum_{i=1}^m u_i v_i \in \mathbb{R}. \quad (\text{A4.1})$$

If  $A \in \mathbb{R}^{m \times m}$  then  $v^T A u \in \mathbb{R}$  is a scalar that can be interpreted as  $(v, Au)$ . Alternatively, since  $v^T A = (A^T v)^T$ , we can write  $v^T A u = (A^T v, u)$  and so we see that

$$(v, Au) = (A^T v, u). \quad (\text{A4.2})$$

Now suppose  $u(x)$  and  $v(x)$  are two functions defined on  $0 \leq x \leq 1$  and define the inner product of these functions in the usual way by

$$(u, v) = \int_0^1 u(x)v(x) dx. \quad (\text{A4.3})$$

Note the relationship between this and (A4.1). Let  $L$  be a linear differential operator mapping functions to functions. Then for any function  $u(x)$ ,  $Lu$  is another function of  $x$  and

$$(v, Lu) = \int_0^1 v(x)(Lu)(x) dx.$$

The *adjoint* of  $L$  is the linear operator  $L^*$  that satisfies

$$(v, Lu) = (L^* v, u)$$

for all functions  $u$  in the domain of  $L$  and  $v$  in the domain of  $L^*$ . Compare this to (A4.2). The domains of functions where these operators are defined depend on the boundary conditions. For simplicity in the discussion here we assume they are defined on the space of functions satisfying homogeneous Dirichlet boundary conditions,  $u(0) = u(1) = 0$ . (For a more complete discussion of adjoints, see a text on differential equations.)

As an example, consider the linear differential operator  $L$  defined by

$$(Lu)(x) = a_2(x)u''(x) + a_1(x)u'(x) + a_0(x)u(x). \quad (\text{A4.4})$$

Then we have

$$(v, Lu) = \int_0^1 [v(x)a_2(x)u''(x) + v(x)a_1(x)u'(x) + v(x)a_0(x)u(x)] dx.$$

To compute the adjoint operator, we must rewrite this as the integral of  $u(x)$  times some differential operator applied to  $v$ . The basic tool is integration by parts, which yields

$$\begin{aligned} (v, Lu) &= v(x)a_2(x)u'(x)|_0^1 - \int_0^1 (v(x)a_2(x))'u'(x) dx \\ &\quad + v(x)a_1(x)u(x)|_0^1 - \int_0^1 (v(x)a_1(x))'u(x) dx \\ &\quad + \int_0^1 v(x)a_0(x)u(x) dx. \end{aligned}$$

The boundary terms all drop out because  $u$  and  $v$  are assumed to both satisfy the homogeneous Dirichlet boundary conditions. The first integral can be integrated by parts once again to move the remaining derivative off of  $u$  and onto  $v$ . Again the boundary terms drop out and we obtain

$$(v, Lu) = \int_0^1 [(v(x)a_2(x))''u(x) - (v(x)a_1(x))'u(x) + v(x)a_0(x)u(x)] dx$$

which can be interpreted as  $(L^*v, u)$  where  $L^*$  is defined by

$$(L^*v)(x) = (v(x)a_2(x))'' - (v(x)a_1(x))' + v(x)a_0(x).$$

This can be rewritten as

$$(L^*v)(x) = a_2(x)v''(x) + [2a_2'(x) - a_1(x)]v'(x) + [a_2''(x) - a_1'(x) + a_0(x)]v(x).$$

This defines the adjoint operator.

In general the adjoint operator is different from the original operator since there is no reason to expect that  $2a_2'(x) - a_1(x) = a_1(x)$ , for example. However, for many important problems that arise from physical principles, the adjoint operator is the same as the original operator, i.e., the equation is *self adjoint* and  $L^* = L$ . This is the analog of a symmetric matrix.

For example, consider the operator appearing in the steady-state heat equation with a varying heat conductivity  $\kappa(x)$ :

$$Lu = (\kappa u')' = \kappa u'' + \kappa' u'. \quad (\text{A4.5})$$

In this case  $a_2(x) = \kappa(x)$ ,  $a_1(x) = \kappa'(x)$ , and  $a_0(x) = 0$ , so we do have  $2a_2' - a_1 = a_1$ , for example. It is much clearer to see *why* this equation is self adjoint if we go through the integration by parts again explicitly:

$$\begin{aligned} (v, Lu) &= \int (\kappa u')' v dx \\ &= - \int \kappa u' v' dx \\ &= \int u (\kappa v')' dx \\ &= (Lv, u). \end{aligned} \quad (\text{A4.6})$$

The symmetry of the first and second integrations by parts is apparent. Recall from Section 2.14 that for this self-adjoint differential equation we can develop a discretization in terms of a symmetric matrix.

## A4.2 Positive definiteness

A symmetric matrix  $A \in \mathbb{R}^{m \times m}$  is called *positive definite* if  $u^T A u > 0$  for all vectors  $u \in \mathbb{R}^m$  other than the zero vector. ( $A$  is *positive semi-definite* if  $u^T A u \geq 0$  for all  $u \neq \vec{0}$ , *negative definite* if  $u^T A u < 0$  for all  $u \neq \vec{0}$ , and *indefinite* if the sign of  $u^T A u$  can change for different  $u$ 's.) Note that  $A$  is negative definite if and only if  $-A$  is positive definite.

A symmetric matrix always has real eigenvalues and orthogonal eigenspaces and so we can write

$$A = R \Lambda R^T$$

where  $R$  is the matrix of right eigenvectors and  $\Lambda$  is a diagonal matrix of eigenvalues. Based on this it is easy to show that  $A$  is:

positive definite	if and only if all eigenvalues are positive,
positive semi-definite	if and only if all eigenvalues are nonnegative,
negative definite	if and only if all eigenvalues are negative,
negative semi-definite	if and only if all eigenvalues are nonpositive,
indefinite	if there are eigenvalues of both signs.

The proofs follow directly from the observation that

$$u^T A u = u^T R \Lambda R^T u = w^T \Lambda w = \sum_{i=1}^m \lambda^i w_i^2$$

where  $w = R^T u$ .

In studying iterative methods for linear systems, we will see that some methods such as conjugate-gradient methods depend on  $A$  being positive definite (or negative definite since we can always negate the system). Why might we expect this property in matrices arising from physical problems?

As one example, again consider the self-adjoint operator (A4.5), and examine the second line of (A4.6) in the case where  $v(x) = u(x)$ :

$$(u, Lu) = - \int \kappa(x) (u'(x))^2 dx. \quad (\text{A4.7})$$

Since the conductivity  $\kappa(x)$  is positive everywhere we see that  $(u, Lu) < 0$  for any function  $u(x)$  that is not identically zero. In other words, the differential operator  $L$  is negative definite.

The finite difference approximation (2.50) results in a matrix that is negative definite (Exercise 2.8). This is easy to verify:

$$\begin{aligned} h^2 U^T A U &= \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_m \end{bmatrix}^T \begin{bmatrix} -(\kappa_{1/2} + \kappa_{3/2})U_1 + \kappa_{3/2}U_2 \\ \kappa_{3/2}U_1 - (\kappa_{3/2} + \kappa_{5/2})U_2 + \kappa_{5/2}U_3 \\ \vdots \\ \kappa_{m-1/2}U_{m-1} - (\kappa_{m-1/2} + \kappa_{m+1/2})U_m \end{bmatrix} \\ &= -(\kappa_{1/2} + \kappa_{3/2})U_1^2 + 2\kappa_{3/2}U_1U_2 - (\kappa_{3/2} + \kappa_{5/2})U_2^2 \\ &\quad + 2\kappa_{5/2}U_2U_3 - \cdots + 2\kappa_{m-1/2}U_{m-1}U_m - (\kappa_{m-1/2} + \kappa_{m+1/2})U_m^2. \end{aligned}$$

Rearranging this gives

$$\begin{aligned} h^2 U^T A U &= -\kappa_{1/2}U_1^2 - \kappa_{3/2}(U_2 - U_1)^2 \\ &\quad - \kappa_{5/2}(U_3 - U_2)^2 - \cdots - \kappa_{m-1/2}(U_m - U_{m-1})^2 - \kappa_{m+1/2}U_m^2. \end{aligned} \quad (\text{A4.8})$$

Since each  $\kappa$  is positive and these are multiplied by squares which are nonnegative, we see that

$$U^T A U \leq 0.$$

Moreover this can be equal to zero only if  $U_1 = U_2 = \cdots = U_m = 0$ , i.e., only if  $U = \vec{0}$ , and so  $A$  is negative definite. Note that the sum in (A4.8), when divided by  $h$ , is a discrete approximation of the integral appearing in (A4.7).



### A4.3 The maximum principle

The equation  $(\kappa u)' = 0$  with  $\kappa(x) > 0$  and Dirichlet boundary conditions  $u(0) = \alpha$ ,  $u(1) = \beta$  can be shown to satisfy a maximum principle:  $u$  can attain a maximum or minimum only on the boundary of the interval (unless it is constant throughout). In particular,  $u$  lies between  $\alpha$  and  $\beta$  everywhere on the interval.

The discrete solution to the system with matrix (2.50) can be easily shown to satisfy the same property. If we let  $U_0 = \alpha$  and  $U_{m+1} = \beta$ , then the  $i$ 'th equation of the system  $AU = 0$  with the matrix (2.50) can be rewritten as

$$U_i = \left( \frac{\kappa_{i-1/2}}{\kappa_{i-1/2} + \kappa_{i+1/2}} \right) U_{i-1} + \left( \frac{\kappa_{i+1/2}}{\kappa_{i-1/2} + \kappa_{i+1/2}} \right) U_{i+1}$$

for all  $i = 1, 2, \dots, m$ . Since  $\kappa > 0$  everywhere and

$$\left( \frac{\kappa_{i-1/2}}{\kappa_{i-1/2} + \kappa_{i+1/2}} \right) + \left( \frac{\kappa_{i+1/2}}{\kappa_{i-1/2} + \kappa_{i+1/2}} \right) = 1,$$

we see that  $U_i$  is just a convex combination of the two neighboring values, and hence must lie between the two values. This observation can be easily turned into a proof.

### A4.4 Diagonal dominance

A matrix is said to be (strictly) *diagonally dominant* if the magnitude of the diagonal element in each row is larger than the sum of the magnitudes of all the off-diagonal elements of the same row,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|.$$

$A$  is *weakly diagonally dominant* if the above holds with  $\geq$  instead of  $>$  and there is equality for some row.

The beautiful *Gerschgorin Circle Theorem* states that the eigenvalues of a matrix all lie in the union  $\mathcal{B} = \cup_i \mathcal{B}_i$  of balls  $\mathcal{B}_i$  in the complex plane defined by

$$\mathcal{B}_i = \left\{ z \in \mathbb{C} : |z - a_{ii}| \leq \sum_{j \neq i} |a_{ij}| \right\}.$$

It follows that if a real matrix is strictly diagonally dominant and its diagonal elements are all positive, then the matrix must be positive definite. If it is only weakly diagonally dominant, and has positive diagonal elements, then it is at least positive semidefinite. It is in fact positive definite if we can show by other means that it is also nonsingular (ruling out 0 as a possible eigenvalue).

This technique can be used to give another proof that the matrix (2.50) is positive definite. It is clearly weakly diagonally dominant, and to see that it is nonsingular we use the maximum principle: If  $Au = 0$  then we are solving the Dirichlet problem with  $\alpha = \beta = 0$  and the solution, by the maximum principle, must be  $u = 0$ . Since this is the only solution to  $Au = 0$ , the matrix must be nonsingular.

Diagonal dominance is also important computationally. For example, it can be shown that we can solve a linear system with a diagonally dominant matrix without pivoting and the algorithm will be stable. This can be very important for sparse matrices where we want to preserve the structure and avoid introducing any more nonzeros than necessary.

## A4.5 Quadratic forms

The fact that certain iterative methods depend on the matrix being symmetric positive definite comes from the fact that in this case we can reformulate the problem of solving the linear system  $Au = F$  as a minimization problem. These iterative methods are based on optimization techniques.

Consider the function

$$g(u) = \frac{1}{2}u^T Au - 2u^T F.$$

This function maps vectors  $u \in \mathbb{R}^m$  to real numbers and is quadratic in the elements of the vector  $u$ . It is called a *quadratic form*. Let's look at this function for two different  $2 \times 2$  matrices to illustrate its behavior:

$$A_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} : \quad g_1(u) = u_1^2 + u_2^2 - 2(u_1 F_1 + u_2 F_2),$$

$$A_2 = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} : \quad g_2(u) = u_1^2 - u_2^2 - 2(u_1 F_1 + u_2 F_2).$$

The matrix  $A_1$  is positive definite and the graph of  $g_1(u)$  is a bowl opening upwards. The function  $g_1(u)$  has a unique global minimum. The matrix  $A_2$  is indefinite and the graph of  $g_2(u)$  has a saddle point but no global minimum or maximum. As we move off towards  $\infty$  in different directions in the  $u_1$ - $u_2$  plane,  $g_2(u)$  may go towards either  $+\infty$  or  $-\infty$  depending on the direction. This is typical of indefinite matrices.

More generally, for any positive-definite matrix  $A \in \mathbb{R}^{m \times m}$  the term  $u^T Au$  in  $g(u)$  is always positive and must go to  $+\infty$  as  $u$  increases in magnitude in any direction in  $\mathbb{R}^m$ . Moreover this quadratic growth overwhelms whatever the linear term  $u^T F$  is doing, and so the function  $g(u)$  must go to  $+\infty$  in every direction and therefore must have a global minimum at some point. Hence for positive definite matrices it makes sense to solve a minimization problem. Why would we want to? Because simple calculus tells us that the minimum of  $g(u)$  must occur at the point where the gradient vector  $\nabla g(u)$  is zero. Differentiating  $g(u)$  with respect to each component  $u_j$  shows that (provided  $A$  is symmetric)

$$\nabla g(u) = Au - F$$

This is the zero vector precisely at the point  $u$  that solves the original linear system  $Au = F$ . So if we can minimize  $g(u)$ , we have solved the system.



# Bibliography

- [ALY88] L. M. Adams, R. J. LeVeque, and D. M. Young. Analysis of the SOR iteration for the 9-point Laplacian. *SIAM J. Num. Anal.*, 25:1156–1180, 1988.
- [AMR88] U. Ascher, R. Mattheij, and R. Russell. *Numerical Solution of Boundary Value Problems for Ordinary Differential Equations*. Prentice-Hall, 1988.
- [ATP84] D. A. Anderson, J. C. Tannehill, and R. H. Pletcher. *Computational Fluid Mechanics and Heat Transfer*. McGraw-Hill, 1984.
- [BB73] J. P. Boris and D. L. Book. Flux corrected transport I, SHASTA, a fluid transport algorithm that works. *J. Comput. Phys.*, 11:38–69, 1973.
- [Bri87] W. L. Briggs. *A Multigrid Tutorial*. SIAM, 1987.
- [But87] J. C. Butcher. *The Numerical Analysis of Ordinary Differential Equations : Runge-Kutta and General Linear Methods*. J. Wiley, Chichester, 1987.
- [CFL28] R. Courant, K. O. Friedrichs, and H. Lewy. Über die partiellen Differenzgleichungen der mathematischen Physik. *Math. Ann.*, 100:32–74, 1928.
- [CFL67] R. Courant, K. O. Friedrichs, and H. Lewy. On the partial difference equations of mathematical physics. *IBM Journal*, 11:215–234, 1967.
- [CH52] C. F. Curtiss and J. O. Hirschfelder. Integration of stiff equations. *Proc. Nat. Acad. Sci. U.S.A.*, 38:235–23, 1952.
- [CHQZ88] C. Canuto, M.Y. Hussaini, A. Quarteroni, and T. A. Zang. *Spectral Methods in Fluid Dynamics*. Springer, New York, 1988.
- [CL55] E. A. Coddington and N. Levinson. *Theory of Ordinary Differential Equations*. McGraw-Hill, New York, 1955.
- [CL88] T. F. Coleman and C. F. Van Loan. *Handbook for matrix computations*. SIAM, 1988.
- [CW84] P. Colella and P. Woodward. The piecewise-parabolic method (PPM) for gas-dynamical simulations. *J. Comput. Phys.*, 54:174–201, 1984.
- [DER86] I. S. Duff, A. M. Erisman, and J. K. Reid. *Direct Methods for Sparse Matrices*. Oxford University Press, 1986.
- [Fle88] C. A. J. Fletcher. *Computational Techniques for Fluid Dynamics*. Springer-Verlag, 1988.
- [For96] B. Fornberg. *A Practical Guide to Pseudospectral Methods*. Cambridge University Press, 1996.
- [FW60] G. E. Forsyth and W. R. Wasow. *Finite Difference Methods for Partial Differential Equations*. Wiley, 1960.

- [Gea71] C. W. Gear. *Numerical Initial Value Problems in Ordinary Differential Equations*. Prentice-Hall, 1971.
- [GL81] A. George and J. W. H. Liu. *Computer Solution of Large Sparse Positive-definite Systems*. Prentice-Hall, 1981.
- [GL89] G. H. Golub and C. F. Van Loan. *Matrix computations*. Johns Hopkins Press, 1989.
- [GO77] D. Gottlieb and S. A. Orszag. *Numerical Analysis of Spectral Methods*. CBMS-NSF Regional Conference Series in Applied Mathematics, #26, SIAM, Philadelphia, 1977.
- [GO92] G. H. Golub and J. M. Ortega. *Scientific computing and differential equations : an introduction to numerical methods*. Academic Press, 1992.
- [God59] S. K. Godunov. *Mat. Sb.*, 47:271, 1959.
- [GR96] E. Godlewski and P.-A. Raviart. *Numerical Approximation of Hyperbolic Systems of Conservation Laws*. Springer-Verlag, New York, 1996.
- [Gre97] A. Greenbaum. *Iterative Methods for Solving Linear Systems*. SIAM, Philadelphia, 1997.
- [Har83] A. Harten. High resolution schemes for hyperbolic conservation laws. *J. Comput. Phys.*, 49:357–393, 1983.
- [Hen62] P. Henrici. *Discrete Variable Methods in Ordinary Differential Equations*. John Wiley and Sons, New York, 1962.
- [Hir88] C. Hirsch. *Numerical Computation of Internal and External Flows*. Wiley, 1988.
- [HNW87] E. Hairer, S. P. Norsett, and G. Wanner. *Solving Ordinary Differential Equations I. Nonstiff Problems*. Springer-Verlag, Berlin, Heidelberg, 1987.
- [HNW93] E. Hairer, S. P. Norsett, and G. Wanner. *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic Problems*. Springer-Verlag, New York, 1993.
- [HY81] L. A. Hageman and D. M. Young. *Applied Iterative Methods*. Academic Press, 1981.
- [HZ72] A. Harten and G. Zwas. Self-adjusting hybrid schemes for shock computations. *J. Comput. Phys.*, 9:568, 1972.
- [Ise96] A. Iserles. *Numerical Analysis of Differential Equations*. Cambridge University Press, Cambridge, 1996.
- [Jes84] D. C. Jespersion. Multigrid methods for partial differential equations. In G. H. Golub, editor, *Studies in Numerical Analysis*, pages 270–317. MAA Studies in Mathematics, Vol. 24, 1984.
- [Joh87] C. Johnson. *Numerical solution of Partial Differential Equations by the Finite Element Method*. Cambridge University Press, 1987.
- [KC81] J. Kevorkian and J. D. Cole. *Perturbation Methods in Applied Mathematics*. Springer, New York, 1981.
- [Kel76] H. B. Keller. *Numerical Solution of Two Point Boundary Value Problems*. SIAM, 1976.
- [Kev90] J. Kevorkian. *Partial Differential Equations*. Wadsworth & Brooks/Cole, 1990.
- [Lam73] J. D. Lambert. *Computational Methods in Ordinary Differential Equations*. Wiley, 1973.
- [Lax72] P. D. Lax. *Hyperbolic Systems of Conservation Laws and the Mathematical Theory of Shock Waves*. SIAM Regional Conference Series in Applied Mathematics, #11, 1972.

- [LeV] R. J. LeVeque. CLAWPACK software. available from [netlib.bell-labs.com](http://netlib.bell-labs.com) in [netlib/pdes/claw](http://netlib/pdes/claw) or on the Web at the URL <http://www.amath.washington.edu/~rjl/clawpack.html>.
- [LeV90] R. J. LeVeque. *Numerical Methods for Conservation Laws*. Birkhäuser-Verlag, 1990.
- [LMDMar] R. J. LeVeque, D. Mihalas, E. Dorfi, and E. Müller. *Numerical Methods in Astrophysical Fluid Flow*. Twenty-seventh Saas-Fee Course, (A. Gautschy and O. Steiner, editors) Springer-Verlag, to appear. (<ftp://amath.washington.edu/pub/rjl/papers/saasfee.ps.gz>).
- [Loa97] C. F. Van Loan. *Introduction to Scientific Computing*. Prentice Hall, Upper Saddle River, NJ, 1997.
- [LT88] R. J. LeVeque and L. N. Trefethen. Fourier analysis of the SOR iteration. *IMA J. Numer. Anal.*, 8:273–279, 1988.
- [LT98] D. Levy and E. Tadmor. From semidiscrete to fully discrete: stability of Runge-Kutta schemes by the energy method. *SIAM Review*, 40:40–73, 1998.
- [McC89] S. F. McCormick. *Multilevel Adaptive Methods for Partial Differential Equations*. SIAM, 1989.
- [MG80] A. R. Mitchell and D. F. Griffiths. *The Finite Difference Method in Partial Differential Equations*. Wiley, 1980.
- [MM94] K. W. Morton and D. F. Mayers. *Numerical Solution of Partial Differential Equations*. Cambridge University Press, Cambridge, 1994.
- [OB87] E. S. Oran and J. P. Boris. *Numerical Simulation of Reactive Flow*. Elsevier, 1987.
- [PT83] R. Peyret and T. D. Taylor. *Computational Methods for Fluid Flow*. Springer, 1983.
- [RM67] R. D. Richtmyer and K. W. Morton. *Difference Methods for Initial-value Problems*. Wiley-Interscience, 1967.
- [SF73] G. Strang and G. J. Fix. *An Analysis of the Finite Element Method*. Prentice-Hall, 1973.
- [Sod85] G. Sod. *Numerical Methods for Fluid Dynamics*. Cambridge University Press, 1985.
- [Str68] G. Strang. On the construction and comparison of difference schemes. *SIAM J. Num. Anal.*, 5:506–517, 1968.
- [Str89] J. C. Strikwerda. *Finite Difference Schemes and Partial Differential Equations*. Wadsworth & Brooks/Cole, 1989.
- [Swa84] Paul N. Swarztrauber. Fast Poisson solvers. In Gene H. Golub, editor, *Studies in Numerical Analysis*, volume 24, pages 319–370, Washington, D.C., 1984. The Mathematical Association of America.
- [Swe84] P. K. Sweby. High resolution schemes using flux limiters for hyperbolic conservation laws. *SIAM J. Num. Anal.*, 21:995–1011, 1984.
- [Tor97] E. F. Toro. *Riemann solvers and numerical methods for fluid dynamics*. Springer-Verlag, Berlin, Heidelberg, 1997.
- [van73] B. van Leer. Towards the ultimate conservative difference scheme I. The quest of monotonicity. *Springer Lecture Notes in Physics*, 18:163–168, 1973.

- [van77] B. van Leer. Towards the ultimate conservative difference scheme IV. A new approach to numerical convection. *J. Comput. Phys.*, 23:276–299, 1977.
- [van79] B. van Leer. Towards the ultimate conservative difference scheme V. A second order sequel to Godunov’s method. *J. Comput. Phys.*, 32:101–136, 1979.
- [Var62] R. S. Varga. *Matrix Iterative Analysis*. Prentice-Hall, 1962.
- [Whi74] G. Whitham. *Linear and Nonlinear Waves*. Wiley-Interscience, 1974.
- [You50] D. M. Young. *Iterative Methods for Solving Partial Differential Equations of Elliptic Type*. PhD thesis, Harvard University, 1950.
- [You71] D. M. Young. *Iterative Solution of Large Linear Systems*. Academic Press, 1971.