

Essential Partial Differential Equations:  
Analytical and Computational Aspects  
Solutions to odd numbered exercises

David F. Griffiths, John W. Dold and David J. Silvester  
Springer International Publishing Switzerland, 2015

Solutions to all exercises are available to approved instructors  
by contacting the publishers.

Springer Undergraduate Mathematics Series

ISBN: 978-3-319-22568-5, e-ISBN: 978-3-319-22569-2

Essential Partial Differential Equations:  
Analytical and Computational Aspects  
Solutions to odd numbered exercises

David F. Griffiths, John W. Dold and David J. Silvester

**Exercises**

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Boundary and initial data</b>	<b>4</b>
<b>3</b>	<b>Origins of PDEs</b>	<b>6</b>
<b>4</b>	<b>Classification of PDEs</b>	<b>7</b>
<b>5</b>	<b>Boundary value problems in <math>\mathbb{R}^1</math></b>	<b>12</b>
<b>6</b>	<b>Finite difference methods in <math>\mathbb{R}^1</math></b>	<b>18</b>
<b>7</b>	<b>Maximum principles and energy methods</b>	<b>27</b>
<b>8</b>	<b>Separation of variables</b>	<b>29</b>
<b>9</b>	<b>The method of characteristics</b>	<b>36</b>
<b>10</b>	<b>Finite difference methods for elliptic PDEs</b>	<b>47</b>
<b>11</b>	<b>Finite difference methods for parabolic PDEs</b>	<b>55</b>
<b>12</b>	<b>Finite difference methods for hyperbolic PDEs</b>	<b>63</b>

## Exercises 1 Introduction

### 1.1

Function	Comment	Conclusion
$u(x, y) = A(y)$	$u_y = A'(y)$	False
$u(x, y) = A(y)$	$u_{xy} = 0$	True
$u(x, t) = A(x)B(t)$	$u_{xy} = 0$ concerns different independent variables!	False
$u(x, t) = A(x)B(t)$	$uu_{xt} = ABA'B' = u_x u_t$	True
$u(x, t, y) = A(x, y)$	$u_t = \partial_t A(x, y) = 0$	True
$u(x, t) = A(x+ct) + B(x-ct)$	$u_{tt} + c^2 u_{xx} = 2c^2(A'' + B'')$	False

### 1.3

These are not the only possible cases; you might find other PDEs:

- (a)  $u(x, t) = A(x+ct) + B(x-ct)$ :  $u_t = cA'(x+ct) - cB'(x-ct)$ ,  $u_x = A'(x+ct) + B'(x-ct)$ ,  $u_{tt} = c^2 A''(x+ct) + c^2 B''(x-ct)$ ,  $u_{xx} = A''(x+ct) + B''(x-ct)$  so  $u_{tt} - c^2 u_{xx} = 0$  (wave equation).
- (b)  $u(x, t) = A(x) + B(t)$ :  $u_t = B'(t)$  so  $u_{tx} = 0$ .
- (c)  $u(x, t) = A(x)/B(t)$ :  $\ln u = \ln A(x) - \ln B(t)$  so  $(\ln u)_{tx} = 0$  or  $uu_{tx} - u_t u_x = 0$ .
- (d)  $u(x, t) = A(xt)$ :  $u_t = xA'(xt)$ ,  $u_x = tA'(xt)$ , so  $tu_t - xu_x = 0$ .
- (e)  $u(x, t) = A(x^2 t)$ :  $u_t = x^2 A'(x^2 t)$ ,  $u_x = 2xtA'(x^2 t)$  so  $2tu_t - xu_x = 0$ .
- (f)  $u(x, t) = A(x^2/t)$ :  $u_t = -\frac{x^2}{t^2} A'(x^2/t)$ ,  $u_x = -\frac{2x}{t} A'(x^2/t)$  so  $2tu_t + xu_x = 0$ .

### 1.5

Suppose that  $u(x, t) = \frac{1}{2}c \operatorname{sech}^2(z)$ , where  $z = \frac{1}{2}\sqrt{c}(x-ct-x_0)$ . Then  $z_t = -\frac{1}{2}c^{3/2}$  and  $z_x = \frac{1}{2}c^{1/2}$  so, by the chain rule, we find

$$\begin{aligned}\partial_t u(x, t) &= \frac{1}{2}c^{5/2} \frac{\sinh z}{\cosh^3 z}, & \partial_x u(x, t) &= -\frac{1}{2}c^{3/2} \frac{\sinh z}{\cosh^3 z} \\ \partial_t u + 6u\partial_x u &= \frac{1}{2}c^{5/2}(\sinh z) \frac{\cosh^2(z) - 3}{\cosh^5 z} = -\partial_x^3 u(x, t)\end{aligned}$$

and so  $u_t + 6uu_x + u_{xxx} = 0$ .

### 1.7

Since  $u = -2\partial_x \phi = -2\phi_x/\phi$ , the partial derivatives are

$$\begin{aligned}u_t &= -2\frac{\phi_{xt}}{\phi} - 2\frac{\phi_x \phi_t}{\phi^2}, & u_x &= -2\frac{\phi_{xx}}{\phi} + 2\frac{(\phi_x)^2}{\phi^2}, \\ u_{xx} &= -2\frac{\phi_{xxx}}{\phi} + 2\frac{\phi_{xx}\phi_x}{\phi^2} + 4\frac{\phi_x}{\phi} \left( \frac{\phi_{xx}}{\phi} - \frac{(\phi_x)^2}{\phi^2} \right)\end{aligned}$$

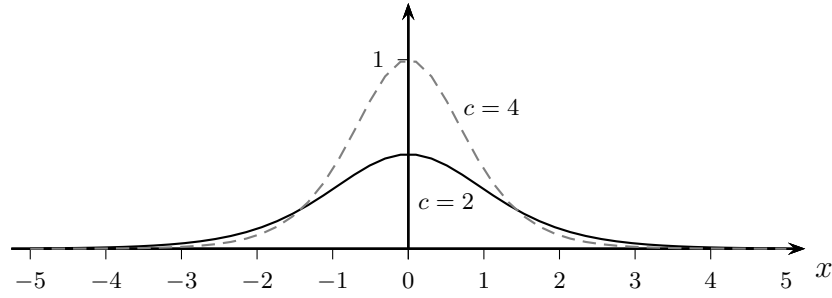


Figure 1: Soliton solutions of the KdV equation with  $c = 2$  (solid) and  $c = 4$  (dashed) for Exercise 1.5 travel with speed  $c$ .

so, with  $\phi_t = \phi_{xx}$  and  $\phi_{xt} = \phi_{xxx}$ ,

$$u_t = -2\frac{\phi_{xxx}}{\phi} - 2\frac{\phi_x\phi_{xx}}{\phi^2}$$

and  $u_t + uu_x = u_{xx}$ . With  $\phi = t^{-1/2} \exp(-x^2/(4t))$  we find  $\log \phi = -\frac{1}{2} \log(t) - x^2/(4t)$  and so  $u = x/t$ , which is a special case of equation (9.32).

## 1.9

The given function  $\phi(x, t)$  is a linear combination of solutions of the heat equation and is therefore a viable candidate for use in the Cole-Hopf transformation. We find

$$u(x, t) = 2a \frac{2e^{-2a(x-2at)} + e^{-a(x-x_0-at)}}{1 + e^{-2a(x-2at)} + e^{-a(x-x_0-at)}}$$

which is shown in Fig. 2 when  $a = 1$  and  $x_0 = 10$ . Two “step”-like solutions coalesce into one.

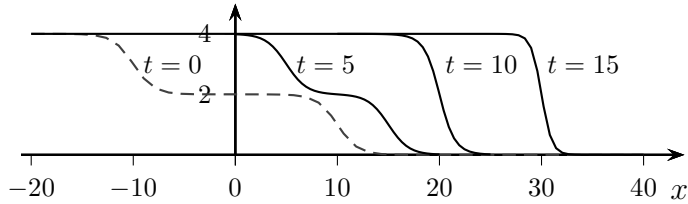


Figure 2: The solution of Burger's equation at times  $t = 0, 5, 10, 15$  for Exercise 1.9.

## Exercises 2 Boundary and initial data

### 2.1

- (a)  $u(x, 0) = A(x) + B(x) = f(x)$  so  $A(x) + B(x) = f(x)$ . There is not enough information to determine both  $A(\cdot)$  and  $B(\cdot)$ .
- (b)  $u(x, 0) = A(x) + B(0) = f(x)$  so  $A(x) = f(x) - B(0)$ . We would only need to know one value of  $B$ , namely  $B(0)$  to determine  $A$ . However, the initial conditions gives no information about  $B$ .
- (c)  $u(x, 0) = A(x)/B(0) = f(x)$  so  $A(x) = B(0)f(x)$ . We would only need to know one value of  $B$ , namely  $B(0)$  to determine  $A$ , but we have no information about  $B$ .
- (d)  $u(x, 0) = A(0) = f(x)$  so  $A(0) = f(x)$ . This tries to set a constant  $A(0)$  to something that is not constant  $f(x)$ , which is not possible!
- (e) (exactly the same)
- (f)  $u(x, 0) = A(\infty) = f(x)$  so  $A(\infty) = f(x)$ . This tries to set a constant  $A(\infty)$  to something that is not constant  $f(x)$ , which is again not possible!

### 2.3

From (1.2) with  $g(x)$  as given,

$$u(x, t) = \frac{1}{\sqrt{4\pi t}} \int_{-\infty}^{\infty} e^{-(x-s)^2/4t} g(s) ds = \frac{1}{\sqrt{4\pi t}} \int_0^{\infty} e^{-(x-s)^2/4t} ds = \frac{1}{\sqrt{\pi}} \int_{-x/\sqrt{4t}}^{\infty} e^{-z^2} dz,$$

where we have made the change of variable  $s = x + z\sqrt{4t}$  in the integrand. Thus  $u(0, t) = 1/2$  from the given result.

$$u(x, 0) = \lim_{t \rightarrow 0} \frac{1}{\sqrt{\pi}} \int_{-x/\sqrt{4t}}^{\infty} e^{-z^2} dz = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} e^{-z^2} dz = 2 \frac{1}{\sqrt{\pi}} \int_0^{\infty} e^{-z^2} dz = 1,$$

since the integrand is an even function of  $z$ .

### 2.5

First note that  $\mathcal{L}(au) = -\kappa(au)_{xx} = a(-\kappa u_{xx}) = a\mathcal{L}u$ . Similarly,  $\mathcal{B}(au) = a\mathcal{B}u$ . Then

$$(\mathcal{L}au)(x, t) = \begin{cases} (au)_t(x, t) + (\mathcal{L}au)(x, t) \\ (\mathcal{B}au)(x, t) \\ au(x, 0) \end{cases} = \begin{cases} a(u_t(x, t) + \mathcal{L}u(x, t)) & \text{for } (x, t) \in (0, 1) \times (0, T) \\ a\mathcal{B}u(x, t) & \text{for } (x, t) \in \{0, 1\} \times (0, T) \\ au(x, 0) & \text{for } t = 0, x \in [0, 1] \end{cases}.$$

and the right hand side is equal to  $a\mathcal{L}u(x, t)$ .

### 2.7

First we confirm that  $u = AT^{1/2}(T-t)^{-1/2}e^{-x^2/4(T-t)}$  is a solution of  $u_t = -u_{xx}$ :

$$\begin{aligned} u_t &= \frac{1}{2}AT^{1/2}(T-t)^{-3/2}e^{-x^2/4(T-t)} - \frac{1}{4}AT^{1/2}x^2(T-t)^{-5/2}e^{-x^2/4(T-t)}, \\ u_x &= -\frac{1}{2}AT^{1/2}x(T-t)^{-3/2}e^{-x^2/4(T-t)}, \\ u_{xx} &= -\frac{1}{2}AT^{1/2}(T-t)^{-1/2}e^{-x^2/4(T-t)} + \frac{1}{4}AT^{1/2}x^2(T-t)^{-5/2}e^{-x^2/4(T-t)} = -u_t \end{aligned}$$

so that  $u_t = -u_{xx}$ .

Note that with this solution  $|u(x, 0)| \leq A$  and  $u(0, t) \rightarrow \infty$  as  $t \rightarrow T$ .

Hence, given any  $\varepsilon > 0$  and any position  $(X, T)$ , the solution  $u = \varepsilon T^{1/2} (T-t)^{-1/2} e^{-(x-X)^2/4(T-t)}$ , satisfying the initial condition  $u(x, 0) = \varepsilon e^{-(x-X)^2/4T}$ , for which  $|u(x, 0)| \leq \varepsilon$ , becomes infinite as  $(x, t) \rightarrow (X, T)$ .

Because of this, the negative heat equation is ill-posed for  $t > 0$  when subjected to initial conditions at  $t = 0$ ; we can always find solutions that are arbitrarily small initially but that become infinite at any chosen time later on.

Suppose that the negative heat equation is subjected to *final conditions*  $u(x, t_f) = g(x)$  at some time  $t = t_f$ , say. Then it is well posed for times *before* the final time,  $t < t_f$ .

## 2.9

The highest derivatives take the form  $au_{tt} + bu_{tx} + cu_{xx}$  (or different subscripts for different independent variables):

- (a)  $u_t + u_{tx} - u_{xx} + u_x^2 = \sin u$ : semilinear.
- (b)  $u_x + u_{xx} + u_y + u_{yy} = \sin(xy)$ : linear and inhomogeneous.
- (c)  $u_x + u_{xx} - u_y - u_{yy} = \cos(xyu)$ : semilinear.
- (d)  $u_{tt} + xu_{xx} + u_t = f(x, t)$ : linear and inhomogeneous.
- (e)  $u_t + uu_{xx} + u^2 u_{tt} - u_{tx} = 0$ : quasilinear.

### Exercises 3   Origins of PDEs

#### 3.1

Since  $H$  does not depend on  $t$ , we can differentiate  $h_t + Hv_x$  with respect to  $t$ :

$$h_{tt} + Hv_{xt} = 0 \quad \Rightarrow \quad h_{tt} + H\partial_x v_t = 0.$$

But  $v_t = -gh_x$ , so  $\partial_x v_t = -gh_{xx}$  and  $h_{tt} - c^2 h_{xx} = 0$ , where  $c^2 = gH$ .

#### 3.3

This simply requires the vector form of differentiation of a product:

$$\vec{\nabla} \cdot (\vec{v}T) = T\vec{\nabla} \cdot \vec{v} + \vec{v} \cdot \vec{\nabla}T.$$

For example, suppose that  $\vec{v} = (u, v)$ , then in  $\mathbb{R}^2$ ,

$$\begin{aligned}\vec{\nabla} \cdot (\vec{v}T) &= (uT)_x + (vT)_y = u_xT + uT_x + v_yT + vT_y \\ &= T(u_x + v_y) + (uT_x + vT_y) = T\vec{\nabla} \cdot \vec{v} + \vec{v} \cdot \vec{\nabla}T.\end{aligned}$$

## Exercises 4 Classification of PDEs

### 4.1

If  $\varphi(x) = x/(1 + |x|)$ , then  $\varphi$  is a continuous function for  $x \in \mathbb{R}$  and

$$\varphi(x) = \begin{cases} \frac{x}{1+x} & x \geq 0 \\ \frac{x}{1-x} & x < 0 \end{cases} \Rightarrow \varphi'(x) = \begin{cases} \frac{1}{(1+x)^2} & x \geq 0 \\ \frac{1}{(1-x)^2} & x < 0 \end{cases}.$$

Therefore

$$\lim_{x \rightarrow 0^-} \varphi'(x) = \lim_{x \rightarrow 0} \frac{1}{(1-x)^2} = 1 = \lim_{x \rightarrow 0} \frac{1}{(1+x)^2} = \lim_{x \rightarrow 0^+} \varphi'(x).$$

That is,  $\varphi'(0^-) = \varphi'(0^+)$  showing that  $\varphi$  is continuously differentiable at the origin.

### 4.3

$u(x, t) = F(x - ct) + G(x + ct)$  and so the initial conditions give

$$\begin{aligned} u(x, 0) &= g_0(x) = F(x) + G(x) \\ u_t(x, 0) &= g_1(x) = -cF'(x) + cG'(x) \end{aligned}$$

Integrating the second of these over the interval  $(x - ct, x + ct)$  leads to

$$-F(x + ct) + F(x - ct) + G(x + ct) - G(x - ct) = \frac{1}{c} \int_{x-ct}^{x+ct} g_1(s) ds$$

while, from the first,

$$\begin{aligned} -F(x + ct) + G(x + ct) &= g_0(x + ct) \\ -F(x - ct) + G(x - ct) &= g_0(x - ct). \end{aligned}$$

Combining these gives d'Alembert's formula (4.20).

### 4.5

The highest derivatives take the form  $au_{tt} + bu_{tx} + cu_{xx}$  (or different subscripts for different independent variables):

- (a)  $u_t + u_{tx} - u_{xx} + u_x^2 = \sin u$ .  $b^2 - 4ac = 1 > 0$  so hyperbolic.
- (b)  $u_x + u_{xx} + u_y + u_{yy} = \sin(xy)$ .  $b^2 - 4ac = -4 < 0$  so elliptic.
- (c)  $u_x + u_{xx} - u_y - u_{yy} = \cos(xyu)$ .  $b^2 - 4ac = 4 > 0$  so hyperbolic.
- (d)  $u_{tt} + xu_{xx} + u_t = f(x, t)$ .  $b^2 - 4ac = -4x$  so: elliptic for  $x > 0$ , hyperbolic for  $x < 0$ , parabolic for  $x = 0$ .
- (e)  $u_t + uu_{xx} + u^2 u_{tt} - u_{tx} = 0$ .  $b^2 - 4ac = 1 - 4u^3$  so elliptic for  $u^3 > \frac{1}{4}$ , hyperbolic for  $u^3 < \frac{1}{4}$ , parabolic for  $u^3 = \frac{1}{4}$ .

### 4.7

With the change of variables  $s = x + y$ ,  $t = x - y$ , it follows from the previous answer that  $u_{xx} - u_{yy} = x + y$  becomes  $4u_{st} = s$ . Integrating with respect to  $t$  gives  $u_s = \frac{1}{4}st + f(s)$  (where  $f$  is an arbitrary function) then integrating with respect to  $s$ ,  $u = \frac{1}{8}s^2t + F(s) + G(t)$ , where  $G$



is an arbitrary function and  $F(s) = \int_s f(s) ds$  is also an arbitrary function. The general solution of  $u_{xx} - u_{yy} = x + y$  is, therefore  $u(x, y) = \frac{1}{8}(x + y)^2(x - y) + F(x + y) + G(x - y)$ . The boundary conditions give:

$$\begin{aligned} u(x, 0) &= \frac{1}{8}x^3 + F(x) + G(x) = x \\ u(0, y) &= -\frac{1}{8}y^3 + F(y) + G(-y) = -\frac{1}{2}y^3 \end{aligned}$$

and, when we replace  $y$  by  $x$  in the second of these, we find that  $F(x) + G(-x) = -\frac{3}{8}x^3$  and so

$$G(x) - G(-x) = x + \frac{1}{4}x^3.$$

With the identity<sup>1</sup>  $G(x) = \frac{1}{2}(G(x) - G(-x)) + \frac{1}{2}(G(x) + G(-x))$  we have

$$G(x) = \frac{1}{2}x + \frac{1}{8}x^3 + E(x)$$

where we have used  $E(x)$  for the even part of  $G$  (which is, as yet, undetermined). Then

$$F(x) = \frac{1}{2}x - \frac{1}{4}x^3 - E(x)$$

so that the solution is (after some algebra)

$$u(x, y) = x(1 - xy) - \frac{1}{2}(x + y)y^2 - E(x + y) + E(x - y)$$

and involves an arbitrary even function  $E(\cdot)$ .

#### 4.9

Comparing the equation  $2u_{xx} + 5u_{xt} + 3u_{tt} = 0$  with the template (4.12) we see that  $a = 2$ ,  $b = 5/2$  and  $c = 3$  so  $b^2 - ac = \frac{1}{4} > 0$ , so the equation is hyperbolic. The factorisation

$$2u_{xx} + 5u_{xt} + 3u_{tt} = (2\partial_x + 3\partial_y)(\partial_x + \partial_y)u$$

suggests the change of variables  $y = x - t$ ,  $s = 3x - 2t$  and the chain rule gives

$$\left. \begin{aligned} \partial_x &= (\partial_x y)\partial_y + (\partial_x s)\partial_s = 3\partial_y + \partial_s \\ \partial_t &= (\partial_t y)\partial_y + (\partial_t s)\partial_s = \partial_y - \partial_s \end{aligned} \right\} \Rightarrow \begin{aligned} \partial_x + \partial_t &= 4\partial_y \\ \partial_x - 3\partial_t &= 4\partial_s. \end{aligned}$$

Hence  $2u_{xx} + 5u_{xt} + 3u_{tt} = (2\partial_x + 3\partial_t)(\partial_x + \partial_t)u = 16u_{ys}$  which, on integrating twice gives the general solution  $u = F(y) + G(s) = F(x - t) + G(3x - 2t)$ .

The initial conditions give

$$F(x) + G(3x) = 0, \quad -F'(x) - 2G'(3x) = xe^{-x^2}$$

and integrating the second of these we find  $F(x) + \frac{2}{3}G(3x) = \frac{1}{2}e^{-x^2} + C$ , where  $C$  is the constant of integration. It follows that

$$F(x) = \frac{3}{2}e^{-x^2} + 3C, \quad G(3x) = -\frac{3}{2}e^{-x^2} - 3C \Rightarrow G(x) = -\frac{3}{2}e^{-x^2/9} - 3C.$$

Hence  $u(x, y) = \frac{3}{2}(e^{-(x-t)^2} - e^{-(3x-2t)^2/9})$  in which there is no arbitrary constant.

#### 4.11

Substituting  $u(r, t) = r^m f(t - r)$  into the PDE and collecting terms leads to

$$r^{m-2}f(t - r)m(m + n - 2) + r^{m-1}f'(t - r)(2m + n - 1) = 0.$$

---

<sup>1</sup>Every function can be written as the sum of an odd and an even function.

This will hold for all differentiable functions  $f$  if, and only if,

$$m(m+n-2) = 0 \text{ and } 2m+n-1 = 0.$$

Thus, either  $m = 0$  and  $n = 1$  or  $n = 3$  and  $m = -1$  giving the solutions  $u(r, t) = f(t - r)$  when  $n = 1$  and  $u(r, t) = f(t - r)/r$  when  $n = 3$  (see the previous question).

#### 4.13

The change of independent variables  $x = s \cos \alpha - t \sin \alpha$ ,  $y = s \sin \alpha + t \cos \alpha$  gives

$$\left. \begin{aligned} \partial_s &= (\partial_s x) \partial_x + (\partial_s y) \partial_y = \cos \alpha \partial_x + \sin \alpha \partial_y \\ \partial_t &= (\partial_t x) \partial_x + (\partial_t y) \partial_y = -\sin \alpha \partial_x + \cos \alpha \partial_y \end{aligned} \right\}$$

$$(\partial_s^2 + \partial_t^2)u = (\cos \alpha \partial_x + \sin \alpha \partial_y)^2 u + (-\sin \alpha \partial_x + \cos \alpha \partial_y)^2 u = u_{xx} + u_{yy}$$

since  $\alpha$  is constant and  $\cos^2 \alpha + \sin^2 \alpha = 1$ .

#### 4.15

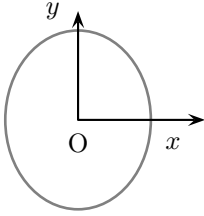


Figure 3: With  $a = 2$ ,  $b = 0$  and  $c = 3$  then  $\hat{Q} = \frac{1}{6}(3x^2 + 2y^2)$  in the solution (4.29). The figure shows a typical elliptical curve  $\mathbf{x}^\top \hat{Q} \mathbf{x} = \text{constant}$ .

#### 4.17

With  $G(x, y, t) = \frac{1}{4\pi} (\log(x^2 + (y+t)^2) - \log(x^2 + (y-t)^2))$  we find

$$\partial_t G(x, y, t) = \frac{1}{4\pi} \left( \frac{2(y+t)}{(x^2 + (y+t)^2)} - \frac{-2(y-t)}{x^2 + (y-t)^2} \right)$$

$$\partial_t G(x, y, t)|_{t=0} = \frac{1}{\pi} \left( \frac{y}{x^2 + y^2} \right) = k(x, y)$$

from (4.32).

With the change of variables  $s = x + y \tan \theta$  the interval  $-\infty < s < \infty$  becomes  $-\frac{1}{2}\pi < \theta < \frac{1}{2}\pi$  and  $ds = y \sec^2 \theta d\theta$  so that

$$\int_{-\infty}^{\infty} k(x-s, y) ds = \int_{-\pi/2}^{\pi/2} \frac{y}{(x-s)^2 + y^2} ds = \int_{-\pi/2}^{\pi/2} d\theta = \pi,$$

where we have used the identity  $1 + \tan^2 \theta = \sec^2 \theta$ .

#### 4.19

The identity

$$\tan(a-b) = \frac{\tan a - \tan b}{1 + \tan a \tan b}$$

with  $\tan a = (x+2)/y$  and  $\tan b = (x-2)/y$  gives, after a little manipulation,

$$\tan(a-b) = \frac{4y}{x^2 + y^2 - 4} = \tan(\pi u)$$

and, in the limit  $u \rightarrow 1/2$ , we have  $x^2 + y^2 \rightarrow 4$ . In particular, when  $y \rightarrow 0$  and  $x \rightarrow \pm 2$  we see that  $u(\pm 2, 0) = 1/2$ . However, the boundary condition is discontinuous at  $(\pm 2, 0)$ :  $g(2^-, 0) = 0$ ,  $g(2^+, 0) = 1$ ,  $g(-2^-, 1) = 0$ ,  $g(-2^+, 0) = 0$  from which we see that

$$u(2, 0) = \frac{1}{2}(g(2^-, 0) + g(2^+, 0)) \text{ and } u(-2, 0) = \frac{1}{2}(g(-2^-, 0) + g(-2^+, 0)).$$

#### 4.21

(a) From Exercise (4.22) with  $u = F(r)$ , we find  $\nabla^2 u = \frac{1}{r}(rF'(r))'$  so  $-\nabla^2 u = 1$  leads to

$$(rF'(r))' = -r \Rightarrow rF'(r) = -\frac{1}{2}r^2 + A \Rightarrow F'(r) = -\frac{1}{2}r + \frac{A}{r} \Rightarrow F(r) = -\frac{1}{4}r^2 + A \log r + B,$$

which is the general solution with arbitrary constants  $A$  and  $B$ . The boundary condition  $F(1) = 0$  requires  $B = \frac{1}{4}$  and, since  $\log r \rightarrow -\infty$  as  $r \rightarrow 0$ , a bounded solution on any region contain the origin, requires  $A = 0$ . Thus  $u(r, \theta) = \frac{1}{4}(1 - r^2)$ .

(b) We look for a solution of the form  $u(r, \theta) = F(r) \cos \theta$ . Substituting into Laplace's equation gives

$$\nabla^2 u = \left( F''(r) + \frac{1}{r}F'(r) - \frac{1}{r^2}F(r) \right) \cos \theta = 0.$$

We try a solution in the form  $F(r) = Ar^n$  and find  $\nabla^2 u = A(n^2 - 1) \cos \theta$ . This will be a solution if  $n = \pm 1$ . This gives two linearly independent solutions  $Ar \cos \theta$  and  $(B/r) \cos \theta$  but the second of these is unbounded when  $r \rightarrow 0$ . Therefore  $u(r, \theta) = Ar \cos \theta$  and the boundary condition  $u(1, \theta) = \cos \theta$  requires  $A = 1$ . Hence  $u = r \cos \theta$ , i.e.,  $u = x$ .

#### 4.23

Under the change of variables  $z = x + iy$ ,  $z^* = x - iy$ , the chain rule gives

$$\partial_x = \partial_z + \partial_{z^*}, \quad \partial_y = i\partial_z - i\partial_{z^*}$$

and so,

$$\begin{aligned} u_{xx} &= (\partial_z + \partial_{z^*})^2 u = u_{zz} + 2u_{zz^*} + u_{z^*z^*} \\ u_{yy} &= -(\partial_z - \partial_{z^*})^2 u = -(u_{zz} - 2u_{zz^*} + u_{z^*z^*}) \end{aligned}$$

giving  $u_{xx} + u_{yy} = 4u_{zz^*}$ . Integrating the PDE  $u_{zz^*} = 0$  with respect to  $z^*$  gives  $u_z = f(z)$  (where  $f$  is an arbitrary function) then integrating with respect to  $z$ ,  $u = F(z) + G(z^*)$ , where  $G$  is an arbitrary function and  $F(s) = \int_s f(s) ds$  is also an arbitrary function. The general solution of  $u_{xx} + u_{yy} = 0$  is, therefore,  $u(x, y) = F(x + iy) + G(x - iy)$ .

If  $F$  is a real function then, since  $\Re F(z) = \frac{1}{2}(F(z) + (F(z))^*) = \frac{1}{2}(F(z) + F(z^*))$ . This will be a solution if we choose  $G(s) = F(s)$ . (The PDE is homogeneous, so  $u = c(F(x + iy) + G(x - iy))$  is also a solution for any constant  $c$ .)

Since  $\Im F(z) = \frac{1}{2}(F(z) - (F(z))^*) = \frac{1}{2}(F(z) - F(z^*))$ . This will be a solution if we choose  $G(s) = -F(s)$ .

#### 4.25

Since  $-\nabla^2 u = v$  and  $-\nabla^2 v = f$ ,

$$-\nabla^2(-\nabla^2 u) = -\nabla^2 v = f$$

and so  $\nabla^4 u = f$ .

With  $u = Cr^2 \log(r)$  and  $\nabla^2 u = (1/r)(ru_r)_r$

$$\begin{aligned}u_r &= 2Cr \log r + Cr, \Rightarrow ru_r = 2Cr^2 \log r + Cr^2 \\(ru_r)_r &= 4Cr \log r + 4Cr \Rightarrow \nabla^2 u = 4C \log r + 4C\end{aligned}$$

so  $v = -\nabla^2 u = -4C \log r - 4C$  which satisfies  $\nabla^2 v = 0$  (see Exercise 4.20). Hence  $u$  satisfies the biharmonic equation.

## Exercises 5 Boundary value problems in $\mathbb{R}^1$

### 5.1

First we look for a particular solution  $u(x) = C$ , where  $C$  is constant. Substituting into the ODE gives  $-4C = 1$  and so  $C = -1/4$ .

Next we look for a solution of the homogeneous equation  $u'' + 3u' - 4u = 0$  in the form  $u = Ae^{\lambda x}$ , where  $A$  and  $\lambda$  are constants. Substituting into the ODE gives

$$Ae^{\lambda x}(\lambda^2 + 3\lambda - 4) = 0.$$

Thus, either  $A = 0$  (leading to the trivial solution  $u(x) = 0$ ) or  $\lambda^2 + 3\lambda - 4 = (\lambda + 4)(\lambda - 1) = 0$ , that is,  $\lambda = -4$  or  $\lambda = 1$ . The ODE is linear and homogeneous so the principle of superposition applies and the general solution of homogeneous ODE (often referred to as the complementary function) is  $u = Ae^{-4x} + Be^x$ . The general solution of the given ODE is therefore

$$u(x) = Ae^{-4x} + Be^x - \frac{1}{4}.$$

Applying the boundary conditions  $u(0) = 1$  and  $u(1) = 3$  leads to

$$\left. \begin{aligned} A + B - \frac{1}{4} &= 1 \\ Ae^{-4} + Be - \frac{1}{4} &= 3 \end{aligned} \right\} \Rightarrow 4A = \frac{5e - 13}{e - e^{-4}}, \quad 4B = \frac{13 - 5e^{-4}}{e - e^{-4}}$$

which uniquely determine  $u(x)$ .

### 5.3

The key identities are

$$\cosh A - \cosh B = 2 \sinh \frac{1}{2}(A+B) \sinh \frac{1}{2}(A-B), \quad \cos A - \cos B = -2 \sin \frac{1}{2}(A+B) \sin \frac{1}{2}(A-B).$$

The first of these with  $A = \frac{1}{2}\sqrt{b}$  and  $B = \sqrt{b}(x - \frac{1}{2})$  allows us to rewrite (5.6) as

$$u(x) = -\frac{2\varepsilon \sinh(\frac{1}{2}\sqrt{b}x) \sinh \frac{1}{2}\sqrt{b}(1-x)}{b \cosh \frac{1}{2}\sqrt{b}} = -2 \frac{\varepsilon}{\cosh \frac{1}{2}\sqrt{b}} \left( \frac{\sinh \frac{1}{2}\sqrt{b}x}{\sqrt{b}} \right) \left( \frac{\sinh \frac{1}{2}\sqrt{b}(1-x)}{\sqrt{b}} \right)$$

Now as  $z \rightarrow 0$ ,  $\cosh z \rightarrow 1$ ,  $\sinh(az)/z \rightarrow a$  (for any constant  $a$ ) and consequently  $u(x) \rightarrow -\frac{1}{2}\varepsilon x(1-x)$ .

Similarly, (5.8) may be rewritten as

$$u(x) = -\frac{2\varepsilon \sin \frac{1}{2}\sqrt{|b|x} \sin \frac{1}{2}\sqrt{|b|}(1-x)}{b \cos \frac{1}{2}\sqrt{|b|}} = -2 \frac{\varepsilon}{\cos \frac{1}{2}\sqrt{|b|}} \left( \frac{\sin(\frac{1}{2}\sqrt{|b|x})}{\sqrt{|b|}} \right) \left( \frac{\sin \frac{1}{2}\sqrt{|b|}(1-x)}{\sqrt{|b|}} \right)$$

and, since  $\cos z \rightarrow 1$ ,  $\sin(az)/z \rightarrow a$  as  $z \rightarrow 0$  we again obtain  $u(x) \rightarrow -\frac{1}{2}\varepsilon x(1-x)$ .

### 5.5

It follows from the inverse monotonicity of  $\mathcal{L}$  that there is a comparison function  $\varphi(x)$  such that  $\mathcal{L}\varphi \geq 1$ . Then

$$\mathcal{L}u = \mathcal{F} \geq -\|\mathcal{F}\| = -\|\mathcal{F}\| \times 1 \geq -\|\mathcal{F}\| \times (\mathcal{L}\varphi)$$

and so, using the linearity of  $\mathcal{L}$ ,  $\mathcal{L}(u + \|\mathcal{F}\|\varphi) \geq 0$  which, by inverse monotonicity, implies that  $u + \|\mathcal{F}\|\varphi \geq 0$ , as required.

### 5.7

The given equations imply  $\frac{a_0}{b_0} = \frac{u'(0)}{u(0)} = \frac{v'(0)}{v(0)}$  from which  $u'(0)v(0) - u(0)v'(0) = 0$  follows.

### 5.9

Multiplying both sides of  $-2xu'' - u' = 2f(x)$  by  $\frac{1}{2}x^{-1/2}$  we obtain

$$-x^{1/2}u'' - \frac{1}{2}x^{-1/2}u' = x^{-1/2}f(x) \Rightarrow -(2x^{1/2}u')' = x^{-1/2}f(x)$$

which has the form (5.28) with  $p(x) = x^{1/2}$ ,  $q(x) = 0$  and  $g(x) = x^{-1/2}f(x)$ . The change of independent variable  $\xi = \xi(x)$  given by (5.30), that is,

$$\xi = \int_0^x \frac{1}{p(s)} ds = 2x^{1/2}$$

then leads to (5.31), i.e.,

$$-\frac{d^2u}{d\xi^2} = \tilde{g}(\xi) = p(x)g(x) = 2f(x) = 2f(\frac{1}{4}\xi^2).$$

in which  $\tilde{q}(\xi) = 0$ . This boundary value problem with Dirichlet boundary conditions may be written in the form  $\mathcal{L}u = \mathcal{F}$ , where  $\mathcal{L}$  is identical to the operator defined in Example 5.2 except that the independent variable is  $\xi$  rather than  $x$ . The arguments used there establish that this equation has a unique solution.

### 5.11

With  $\langle u, v \rangle = \int_0^L u(x)v^*(x) dx$  :

- (a)  $\langle v, u \rangle^* = \left( \int_0^L v(x)u^*(x) dx \right)^* = \int_0^L v^*(x)u(x) dx = \langle u, v \rangle$  since  $(u^*)^* = u$ .
- (b)  $\langle u, u \rangle = \int_0^L |u(x)|^2 dx \geq 0$  since the integrand is non-negative.
- (c) Suppose that  $\langle u, u \rangle = 0$  but that  $u(x)$  is not identically zero. There must therefore be point  $a \in (0, L)$  where  $u(a) \neq 0$ . By continuity there must be an  $\varepsilon > 0$  such that  $u$  is non-zero in the interval  $(a - \varepsilon, a + \varepsilon)$  so  $\langle u, u \rangle = \int_{a-\varepsilon}^{a+\varepsilon} |u(x)|^2 dx > 0$  giving a contradiction.
- (d)  $\langle c_1u_1 + c_2u_2, v \rangle = \int_0^L (c_1u_1 + c_2u_2)v^* dx = c_1 \int_0^L u_1v^* dx + c_2 \int_0^L u_2v^* dx$  by the property of integrals and so  $\langle c_1u_1 + c_2u_2, v \rangle = c_1\langle u_1, v \rangle + c_2\langle u_2, v \rangle$ .
- (e) Suppose that  $u$  and  $v$  are orthogonal with respect to the inner product  $\langle u, v \rangle$  and are linearly dependent. Thus  $\langle u, v \rangle = 0$  and there are non-zero constants  $a, b$  such that  $au(x) + bv(x) = 0$ . However, since  $v(x) = -(a/b)u(x)$ ,

$$\langle u, v \rangle = \langle u, -(a/b)u(x) \rangle = -(a/b)u(x)\langle u, u \rangle \neq 0$$

(having used properties (d) with  $c_1 = -b/a$ ,  $c_2 = 0$  and (c) 5) giving a contradiction.

### 5.13

$\phi(x) = \sin 2\pi x$  satisfies the differential equation and boundary conditions.

Multiplying  $-u'' - 4\pi^2 u = f(x)$  by  $\phi$  and integrating over  $(0, 1)$ , then integrating by parts twice, gives

$$\begin{aligned}\int_0^1 \phi(x)f(x) \, dx &= \int_0^1 \phi(x)(-u'' - 4\pi^2 u) \, dx \\ &= -\phi(x)u'(x)\Big|_0^1 + \phi'(x)u(x)\Big|_0^1 + \int_0^1 u(x)(-\phi'' - 4\pi^2 \phi) \, dx = 0.\end{aligned}$$

The ODE cannot be satisfied with the given BCs unless  $\int_0^1 \phi(x)f(x) \, dx = 0$ .

This condition is clearly violated when  $f(x) = \phi(x) = \sin 2\pi x$ . The ODE  $-u'' - 4\pi^2 u = \sin(2\pi x)$  has the general solution

$$u(x) = x \frac{\cos 2\pi x}{4\pi} + A \sin 2\pi x + B \cos 2\pi x$$

and applying the BCs we find  $u(0) = B = 0$  and  $u(1) = \frac{1}{4\pi} + B = 0$  which are contradictory. There cannot therefore be a solution.

When  $f(x) = 1$  we find  $\int_0^1 \phi(x)f(x) \, dx = 0$ . In this case the ODE  $-u'' - 4\pi^2 u = 1$  has the general solution

$$u(x) = -\frac{1}{4\pi^2} + A \sin 2\pi x + B \cos 2\pi x$$

and, applying the BCs we find  $u(0) = -\frac{1}{4\pi^2} + B$  and  $u(1) = -\frac{1}{4\pi^2} + B$

Thus,

$$u(x) = -\frac{1}{4\pi^2} + A \sin 2\pi x - \frac{1}{4\pi^2} \cos 2\pi x$$

satisfies the differential equation and boundary conditions for any choice of constant  $A$  so we have a non-unique solution: the solution is unique up to an arbitrary multiple of  $\phi(x)$ .

### 5.15

(a) The standard argument shows that  $\phi$  cannot satisfy the boundary conditions unless  $\lambda > 0$ . So, with  $\lambda = \mu^2$ , the general solution of the ODE is  $\phi(x) = A \sin \mu x + B \cos \mu x$ . The boundary condition  $\phi(0) = 0$  requires  $B = 0$  and  $\phi'(1) = 0$  requires  $A \cos \mu = 0$ . Since  $A \neq 0$  (since it would lead to a trivial solution), it follows that  $\lambda = \lambda_n := (n - \frac{1}{2})^2 \pi^2$ ,  $n = 1, 2, \dots$ , with corresponding eigenfunctions  $\phi_n(x) = \sin(n - \frac{1}{2})\pi x$ .

(b) Suppose that  $\omega^2 = \lambda_n$  for some value of  $n$ . If  $\phi(x)$  is the corresponding eigenfunction, then

$$\int_0^1 \phi(x)(-u''(x) + \omega^2 u) \, dx = \int_0^1 \phi(x)f(x) \, dx.$$

Integrating the left hand side by parts twice and using the boundary conditions on both  $u$  and  $\phi$ , we find

$$\begin{aligned}-\phi(x)u'(x)\Big|_0^1 + \int_0^1 (\phi'(x)u'(x) + \omega^2 \phi u) \, dx &= \int_0^1 \phi(x)f(x) \, dx, \\ (-\phi(x)u'(x) + \phi'(x)u(x))\Big|_0^1 + \int_0^1 (-\phi''(x) + \omega^2 \phi)u(x) \, dx &= \int_0^1 \phi(x)f(x) \, dx, \\ 2\phi(1) - \phi'(0) &= \int_0^1 \phi(x)f(x) \, dx,\end{aligned}$$

since  $-\phi''(x) + \omega^2\phi = 0$ . The data for the problem are inconsistent unless this condition is satisfied.

When  $\omega^2 = \lambda_1 = \frac{1}{2}\pi$ , then  $\phi(x) = \sin \frac{1}{2}\pi x$  and, with  $f(x) = c$ , we find that  $c = \frac{1}{2}\pi(2 - \frac{1}{2}\pi)$ . The general solution of  $-u''(x) = \omega^2 u(x) + c$  is  $u(x) = A \sin \frac{1}{2}\pi x + B \cos \frac{1}{2}\pi x - c/\omega^2$ . Applying the BCs:

$$u(0) = 1 = B - c/\omega^2, \quad u'(1) = -2 = -\frac{1}{2}\pi B$$

both of which give the same value  $B = 4/\pi$  because the value of  $c$  was carefully chosen. The solution is, therefore,  $u(x) = 1 + A \sin \frac{1}{2}\pi x + (4/\pi)(\cos(\frac{1}{2}\pi x) - 1)$  which is unique up to an arbitrary multiple of  $\phi$ .

### 5.17

When  $u(x) = M(x)w(x)$  we find

$$u'(x) = M'(x)w(x) + M(x)w'(x), \quad u''(x) = M''(x)w(x) + 2M'(x)w'(x) + M(x)w''(x)$$

so that

$$u'' - au' - bu = Mw'' + (2M' - aM)w' + (M'' - aM' - bM)w.$$

The coefficient of  $w'$  can be made to vanish by choosing  $M$  such that  $2M' = aM$  in which case  $2M'' = a'M + aM' = (a' + \frac{1}{2}a^2)M$  and

$$M'' - aM' - bM = -\frac{1}{2}(2b + \frac{1}{2}a^2 - a')M.$$

Thus,  $u'' - au' - bu = f$  becomes

$$-w'' + Q(x)w(x) = G(x),$$

where  $Q(x) = -(M'' - aM' - bM)/M = \frac{1}{2}(2b + \frac{1}{2}a^2 - a')$ ,  $G(x) = f(x)/M(x)$  and  $M(x) = A \exp(\int^x a(s) ds)$ .

### 5.19

With  $\phi_n(x) = e^{2\pi i n x/L}$  then, for  $m \neq n$ ,

$$\begin{aligned} \langle \phi_n, \phi_m \rangle &= \int_0^L e^{2\pi i n x/L} e^{-2\pi i m x/L} dx = \int_0^L e^{2\pi i (n-m)x/L} dx \\ &= \frac{L}{2\pi i (n-m)} e^{2\pi i (n-m)x/L} \Big|_0^L = \frac{L}{2\pi i (n-m)} (e^{2\pi i (n-m)} - 1) = 0 \end{aligned}$$

since  $n - m$  is an integer and so  $e^{2\pi i (n-m)} = 1$ . Also, when  $m = n$ ,

$$\langle \phi_n, \phi_n \rangle = \int_0^L e^{2\pi i n x/L} e^{-2\pi i n x/L} dx = \int_0^L dx = L.$$

### 5.21

Suppose that the eigenvalue problem is defined by  $\mathcal{L}u = \lambda u$  with boundary conditions  $\mathcal{B}u = 0$ . Let  $v = cu$ , where  $c$  is a constant. By the linearity of  $\mathcal{L}$  and  $\mathcal{B}$ ,  $\mathcal{L}v = \mathcal{L}(cu) = c\mathcal{L}u = c\lambda u = \lambda v$  and  $\mathcal{B}(v) = \mathcal{B}(cu) = c\mathcal{B}u = 0$ . Thus,  $v$  satisfies the same equations as  $u$ :  $\mathcal{L}v = \lambda v$  and  $\mathcal{B}v = 0$ .



### 5.23

With  $\phi_n = \cos(n - \frac{1}{2})x$  for  $0 < x < \pi$  and  $m \neq n$ ,

$$\begin{aligned}\langle \phi_n, \phi_m \rangle &= \int_0^\pi \cos(n - \tfrac{1}{2})x \cos(m - \tfrac{1}{2})x \, dx \\ &= \tfrac{1}{2} \int_0^\pi (\cos(n + m - 1)x + \cos(n - m)x) \, dx \\ &= \left( \frac{\sin(n + m - 1)x}{2(n + m - 1)} + \frac{\sin 2(n - m)x}{(n - m)} \right) \Big|_0^\pi = 0.\end{aligned}$$

### 5.25

When  $\lambda \leq 4$  the argument from the previous exercise is easily adapted to show that only trivial solutions are possible for the ODE  $-u'' + 4u = \lambda u$  with BCs  $u(0) = u(\pi) = 0$ . When  $\lambda > 4$  we have the general solution

$$u(x) = A \sin \sqrt{\lambda - 4}x + B \cos \sqrt{\lambda - 4}x$$

and the BCs give

$$u(0) = 0 = B \text{ and } u(\pi) = 0 = A \sin \sqrt{\lambda - 4}\pi = 0.$$

Since  $A$  cannot be zero (this would imply that  $u(x) = 0$ , the trivial solution) so  $\lambda$  must satisfy  $\sin \sqrt{\lambda - 4}\pi = 0$ . Hence  $\sqrt{\lambda - 4}\pi = n\pi$ ,  $n = \pm 1, \pm 2, \dots$  giving  $\lambda_n = 4 + n^2$  with corresponding eigenfunction  $u_n(x) = \sin n\pi x$ ,  $n = 1, 2, 3, \dots$  (we do not include negative values of  $n$ , since  $\sin(-n\pi x) = -\sin(n\pi x)$  and we would have linearly dependent eigenfunctions).

### 5.27

With  $u = Mw$  the ODE  $-x^2 u'' + 2xu' - 2u = \lambda x^2 u$  becomes (see Exercise 5.17)

$$\begin{aligned}-x^2(M'' - aM' - bM) + 2x(M'w + Mw') - 2Mw &= \lambda x^2 Mw \\ -x^2 Mw'' + (2xM - 2x^2 M')w' + (-x^2 M'' + 2xM' - 2M)w &= \lambda x^2 Mw\end{aligned}$$

and the coefficient of  $w'$  vanishes when  $M - xM' = 0$ . Therefore we may choose  $M(x) = x$  and the eigenvalue problem becomes  $-w'' = \lambda w$ . Since  $u' = w + xw'$ , the BC  $u'(0) = 0$  becomes  $w(0) = 0$  and  $u(1) = u'(1)$  becomes  $w'(1) = 0$ .

The standard argument can be applied to show that  $\lambda = 0$  and  $\lambda < 0$  both lead to trivial solutions. For the case  $\lambda > 0$ , let  $\lambda = \mu^2$  then  $-w'' = \mu^2 w$  has general solution  $w = A \sin \mu x + B \cos \mu x$ . The BC  $w(0) = 0$  implies  $B = 0$  and  $w'(1) = 0$  then leads to  $A\mu \cos \mu = 0$ . Choosing  $A = 0$  gives immediately the trivial solution, choosing  $\mu = 0$  leads to the earlier case  $\lambda = 0$  so we are left with  $\cos \mu = 0$ . Therefore  $\mu = (n - \frac{1}{2})\pi$ ,  $n = 1, 2, \dots$  and the eigenvalues are  $\lambda_n = (n - \frac{1}{2})^2 \pi^2$  with corresponding eigenfunction  $w_n = \sin(n - \frac{1}{2})\pi x$ , i.e.,  $u_n(x) = x \sin(n - \frac{1}{2})\pi x$ ,  $n = 1, 2, \dots$

### 5.29

We shall work from first principles. Multiplying the differential equation  $-u''(x) = \lambda w(x)u(x)$  by  $u^*(x)$  and integrating by parts gives, on applying the boundary conditions,

$$\begin{aligned}\int_0^1 -u''(x)u^*(x) \, dx &= \lambda \int_0^1 w(x)u(x)u^*(x) \, dx \\ -u'(x)u^*(x) \Big|_0^1 + \int_0^1 u'(x)(u'(x))^* \, dx &= \lambda \int_0^1 w(x)|u(x)|^2 \, dx \\ |u(1)|^2 + \int_0^1 |u'(x)|^2 \, dx &= \lambda \int_0^1 w(x)|u(x)|^2 \, dx\end{aligned}$$

and so

$$\lambda = \frac{|u(1)|^2 + \int_0^1 |u'(x)|^2}{\int_0^1 w(x)|u(x)|^2 dx}$$

in which both numerator and denominator are real and positive.

### 5.31

This requires on the insertion of a factor  $w(x)$  in each of the integrands in the solution of Exercise 5.11.

### 5.33

Choosing  $u = \phi_m(x)$  and  $v = \phi_n(x)$  so that

$$\mathcal{L}u = \lambda_m w \phi_m \text{ and } \mathcal{L}v = \lambda_n w \phi_n$$

then by Lagrange's identity (5.25)  $\int_0^1 (v^* \mathcal{L}u - u \mathcal{L}v^*) dx = 0$ . However, the eigenfunctions are real by Exercise 5.30, so

$$\begin{aligned} 0 &= \int_0^1 (v^* \mathcal{L}u - u \mathcal{L}v^*) dx = \int_0^1 (\phi_n \lambda_m w \phi_m - \phi_m \lambda_n w \phi_n) dx \\ &= (\lambda_m - \lambda_n) \langle \phi_n, \phi_m \rangle_w \end{aligned}$$

and therefore  $\langle \phi_n, \phi_m \rangle_w = 0$  provided  $\lambda_m \neq \lambda_n$ .

## Exercises 6 Finite difference methods in $\mathbb{R}^1$

### 6.1

Using Taylor series expansions with remainder terms, (6.13) with  $x = x_m$  becomes

$$v(x_m \pm h) = v(x_m) \pm hv'(x_m) + \frac{1}{2}h^2v''(x_m) \pm \frac{1}{6}h^3v'''(x_m) + \frac{1}{24}h^4v''''(\xi_m^\pm),$$

where  $x_m - h < \xi_m^- < x_m < \xi_m^+ < x_m + h$ . Adding these series together we obtain

$$v(x_{m+1}) + v(x_{m-1}) = 2v(x_m) + h^2v''(x_m) + \frac{1}{24}h^4(v''''(\xi_m^-) + v''''(\xi_m^+)).$$

but, by the Intermediate Value Theorem, there must be a point  $\xi_m \in (\xi_m^-, \xi_m^+)$  such that  $\frac{1}{2}(v''''(\xi_m^-) + v''''(\xi_m^+)) = v''''(\xi_m)$ . Consequently, on rearranging,

$$v''(x_m) = h^{-2}(v(x_{m+1}) - 2v(x_m) + v(x_{m-1})) - \frac{1}{12}h^2v''''(\xi_m).$$

### 6.3

$$\begin{aligned} \Delta^+ \Delta^+ v_m &= \Delta^+(\Delta^+ v_m) = \Delta^+(v_{m+1} - v_m) = \Delta^+ v_{m+1} - \Delta^+ v_m \\ &= (v_{m+2} - v_{m+1}) - (v_{m+1} - v_m) = v_{m+2} - 2v_{m+1} + v_m = \delta^2 v_{m+1}. \end{aligned}$$

Then solution to Exercise 6.1 gives

$$h^{-2} \Delta^+ \Delta^+ v_m = h^{-2} \delta^2 v_{m+1} = v''_{m+1} + \mathcal{O}(h^2)$$

but  $v''_{m+1} = v''(x_m + h) = v''_m + hv'''_m + \mathcal{O}(h^3) = v''_m + \mathcal{O}(h)$ , Hence,  $h^{-2} \Delta^+ \Delta^+ v_m = v''_m + \mathcal{O}(h)$ . The corresponding results for  $\Delta^- \Delta^- v_m$  are:

$$h^{-2} \Delta^- \Delta^- v_m = h^{-2} \delta^2 v_{m-1} = v''_{m-1} + \mathcal{O}(h^2) = v''_m + \mathcal{O}(h).$$

### 6.5

At  $m = 0$ , the equation  $\mathcal{L}_h U_m = \mathcal{F}_{h,m}$  gives  $U_0 = \alpha$ .

For  $0 < m < M$ ,

$$-a_m U_{m-1} + b_m U_m - c_m U_{m+1} = f_m$$

and at  $m = M$ ,  $U_M = \beta$ . With  $\mathbf{u} = [U_1, U_2, \dots, U_{M-1}]^T$ , these equations may be combined to give  $A\mathbf{u} = \mathbf{f}$ , where

$$A = \frac{1}{h^2} \begin{bmatrix} b_1 & -c_1 & & & \\ -a_2 & b_2 & -c_2 & & \\ & \ddots & \ddots & \ddots & \\ & & -a_{M-2} & b_{M-2} & -c_{M-2} \\ & & & -a_{M-1} & b_{M-1} \end{bmatrix}, \quad \mathbf{f} = \begin{bmatrix} f_1 + \alpha a_1 \\ f_2 \\ \vdots \\ f_{M-2} \\ f_{M-1} + \beta c_{M-1} \end{bmatrix}.$$

### 6.7

$a_m = h^{-2}$ ,  $b_m = 2h^{-2} + x_m^2$ ,  $c_m = h^{-2}$ ,  $d_m = x_m$  for  $m = 1, 2, \dots, M-1$ .

Referring to Definition 6.9, these coefficients are all positive and  $b_m = a_m + c_m + x_m^2 \geq 0$  and the corresponding operator  $\mathcal{L}_h$  is of positive type.

### 6.9

Since  $\Phi(x)$  is a quadratic function of  $x$ ,  $h^{-2}\delta^2\Phi_m = \Phi''(x_m) = 2c$ . Also,

$$\Phi_M - \Phi_{M-1} = \Phi(1) - \Phi(1-h) = h\Phi'(1) - \frac{1}{2}h^2\Phi''(1) = -ch(1-a) + ch^2$$

Consequently,

$$-h^{-2}\delta^2\Phi_m = -2c \text{ and } h^{-1}(\Phi(1) - \Phi(1-h)) = -c(1-a) + ch.$$

Since  $h \leq \frac{1}{2}$  we choose  $a = \frac{1}{4}$  so that  $1-a-h = \frac{3}{4}-h \geq \frac{1}{4}$  and therefore  $c(1-a) - ch \geq 1$  with  $c = -4$ . Then  $-h^{-2}\delta^2\Phi_m = -2c = 8 \geq 1$ . Thus  $\Phi(x) = x(4-x)$  is a possible comparison function for Exercise 6.8.

When  $U_m = \frac{1}{2}x_m(3-x_{m-1})$  we find  $U_0 = 0$ ,  $-h^{-2}\delta^2U_m = 1$  and  $2h^{-1}(U_M - U_{M-1}) = 1$  and so  $U$  satisfies exactly the finite difference equations from Exercise 6.8. This is the only solution since  $\mathcal{L}_h$  is inverse monotone.

The general solution of the differential equation  $-u''(x) = 1$  is  $u = A + Bx - \frac{1}{2}x^2$ . The BC  $u(0) = 0$  implies  $A = 0$  while  $u'(1) = \frac{1}{2}$  requires  $B = 3/2$ . Thus  $u(x) = \frac{1}{2}x(3-x)$  and the global error is

$$E_m = u(x_m) - U_m = \frac{1}{2}x_m(x_m - x_{m-1}) = \frac{1}{2}hx_m.$$

and so  $\|E\|_{h,\infty} = \mathcal{O}(h)$ —convergence is at a first order rate.

### 6.11

The ODE is approximated by  $\mathcal{L}_h U_m = f_m$  for  $m = 1, 2, \dots, M-1$ , where

$$\begin{aligned} \mathcal{L}_h U_m &:= -h^{-2}\delta^2 U_m + 20h^{-1}\Delta U_m \\ &= -h^{-2}(1+10h)U_{m-1} + 2h^{-2}U_m - h^{-2}(1-10h)U_{m+1} \end{aligned}$$

and  $f_m = (mh)^2$ . Comparing the coefficients of  $\mathcal{L}_h$  with those of Definition 6.9,  $a_m = h^{-2}(1+10h) \geq 0$  for all  $h > 0$ ,  $c_m = h^{-2}(1-10h) \geq 0$  for  $h \leq 1/10$  and  $b_m = 2h^{-2} = a_m + c_m$ . Hence  $\mathcal{L}_h$  is a positive type operator for  $h \leq 1/10$ , i.e.,  $M \geq 10$ .

With  $\varphi(x) = Ax + B$  and  $\mathcal{L}u(x) = -u''(x) + 20u'(x)$  we find  $\mathcal{L}\varphi(x) = 20A \geq 1$  if  $A \geq 1/20$ . Also  $\varphi(0) \geq 1$  and  $\varphi(1) \geq 1$  if  $B \geq 1$  and  $A+B \geq 1$ , respectively. All these conditions can be met by choosing  $A = 1/20$  and  $B = 1$  so  $\varphi(x) = 1 + x/20$  is a comparison function for  $\mathcal{L}$  with Dirichlet BCs. Since  $\varphi$  is a linear function, then  $\mathcal{L}_h\varphi(x_m) \geq 1$  so it also acts as a comparison function for  $\mathcal{L}_h$ .  $C = \max_{0 \leq x \leq 1} \varphi(x) = 21/20$  and therefore  $\mathcal{L}_h$  with Dirichlet BCs is stable by Lemma 6.8 for  $h \leq 1/10$ .

### 6.13

The local truncation error is  $\mathcal{R}_h = \mathcal{L}_h U - \mathcal{F}_h$ , where  $\mathcal{L}_h U_m := -\varepsilon h^{-2}\delta^2 U_m + 2h^{-1}\Delta^- U_m$  and  $\mathcal{F}_{h,m} = f_m$ . Using the results in Table 6.1,

$$\begin{aligned} \mathcal{R}_{h,m} &= -\varepsilon h^{-2}\delta^2 m_m + 2h^{-1}\Delta^- u_m - f_m \\ &= -\varepsilon(u_m'' + \frac{1}{12}h^2 u_m'''' + \mathcal{O}(h^4)) + 2(u_m' + \frac{1}{2}hu_m'' + \mathcal{O}(h^2)) - f_m \\ &= (-\varepsilon u_m'' + 2u_m' - f_m) + hu_m'' + \mathcal{O}(h^2) = \mathcal{O}(h) \end{aligned}$$

since  $-\varepsilon u'' + 2u' = f$ . The order of consistency is therefore first order. Also

$$\mathcal{L}_h U_m = -(\varepsilon h^{-2} + 2h)U_{m-1} + (2\varepsilon h^{-2} + h)U_m - \varepsilon h^{-2}U_{m+1}$$

so, according to Definition 6.9, is of positive type for all  $\varepsilon > 0$  and  $h > 0$ . The comparison function  $\varphi(x) = \frac{1}{2}x + 1$ , being linear in  $x$ , also satisfies  $\mathcal{L}_h \varphi(x_m) \geq 1$ . Therefore  $\mathcal{L}_h$  is stable by Lemma 6.8 (with stability constant  $C = 3/2$ ). Finally, convergence at a first order rate is a consequence of Theorem 6.7.

### 6.15

From (6.51) the local truncation error  $\mathcal{R}_h = \widehat{\mathcal{L}}_h u - \widehat{\mathcal{F}}_h$  and, since  $\Delta^- u_m = hu'_m - \frac{1}{2}h^2 u''_m + \mathcal{O}(h^3)$  (see Table 6.1),

$$\begin{aligned} \mathcal{R}_{h,M} &= (a + \frac{1}{2}bhr_M)u_M + bh^{-1}\Delta^- u_M - (\beta + \frac{1}{2}bhf_M) \\ &= (a + \frac{1}{2}bhr(1))u(1) + b(u'(1) - \frac{1}{2}hu''(1) + \mathcal{O}(h^2)) - (\beta + \frac{1}{2}bhf(1)) \\ &= au(1) + bu'(1) - \beta + \frac{1}{2}hb[-u''(1) + r(1)u(1) - f(1) + \mathcal{O}(h^2)] \end{aligned}$$

and so  $\mathcal{R}_{h,M} = \mathcal{O}(h^2)$  since  $au(1) + bu'(1) = \beta$  and  $-u''(1) + r(1)u(1) = f(1)$ . Had we used  $\Delta^- u_m = hu'_m - \frac{1}{2}h^2 u''_m + \frac{1}{6}h^3 u'''(\xi_m)$ , where  $x_m - h < \xi_m < x_m$ , we would have found that  $\mathcal{R}_{h,M} = bh^2 u'''(\xi_m)$ . In both cases the local truncation error is of second order.

### 6.17

The leading term in the local truncation error from the previous question is  $-\frac{1}{2}bhu''(0)$  and, using the ODE  $-u'' + ru = f$  at  $x = 0$ , this becomes  $\frac{1}{2}bh(f(0) - r(0)u(0))$ . Thus, subtracting this term from the left of the BC gives the modified condition

$$aU_0 - bh^{-1}\Delta^+ U_0 - \frac{1}{2}bh(f_0 - r_0 U_0) = \alpha.$$

The corresponding local truncation error is

$$\begin{aligned} \mathcal{R}_{h,M} &= au_0 - bh^{-1}\Delta^+ u_0 - \frac{1}{2}bh(f_0 - r_0 u_0) - \alpha \\ &= au_0 - b(u'(0) + \frac{1}{2}hu''(0) + \mathcal{O}(h^2)) - \frac{1}{2}bh(f_0 - r_0 u_0) - \alpha \\ &= (au(0) - bu'(0) - \alpha) - \frac{1}{2}h(-u''(0) + r_0 u_0 - f_0) + \mathcal{O}(h^2) = \mathcal{O}(h^2) \end{aligned}$$

since  $au(0) - bu'(0) = \alpha$  and  $-u''(0) + r_0 u_0 = f_0$ .

### 6.19

The proof of Theorem 6.10 can be used to prove that  $U_m$  cannot have a negative minimum for  $1 \leq m \leq M-1$ . It remains to prove that  $U_m$  cannot have a negative minimum for either  $m = 0$  or  $m = M$ .

We begin with the left end-point. Suppose that, contrary to the statement of the theorem,  $U_m$  has a negative minimum at  $m = 0$  so that  $U_1 \geq U_0$  (which implies  $-c_0 U_1 \leq -c_0 U_0$ ) and

$$\mathcal{L}_h U_0 = b_0 U_0 - c_0 U_1 \leq (b_0 - c_0)U_0 \leq 0.$$

If this inequality were strict (because  $b_0 > c_0$ ) it would contradict the assumption  $\mathcal{L}_h U \geq 0$  and prove that a negative minimum at  $m = 0$  could not occur. Suppose therefore that equality holds which means that  $b_0 = c_0$  and  $U_1 = U_0 < 0$ .

The argument at the right end-point is essentially the same—either a negative minimum cannot occur at  $m = M$  or  $b_M = c_M$  and  $U_{M-1} = U_M < 0$ .

We now turn to the interior grid points. Suppose that, contrary to the statement of the theorem,  $U_m$  has a negative minimum at  $m$ , where  $0 < m < M$ . The argument in Theorem 6.10 proves that either there is a contradiction of  $\mathcal{L}_h U_m \geq 0$  or both  $U_{m-1} = U_m = U_{m+1} < 0$  and  $b_m = a_m + c_m$ . Combining all three cases we see that either there is a contradiction of  $\mathcal{L}_h U_m \geq 0$  or the same negative minimum is attained at *all* grid points *and*  $b_m = a_m + c_m$  for all  $m$  (recall  $a_0 = c_M = 0$ ). However, this last possibility is ruled out by Definition 6.17 which requires that  $b_m > a_m + c_m$  for at least one value of  $m$ .

## 6.21

We find that

$$\begin{aligned} a_0 = 0, b_0 = c_0 = 1 & \quad b_0 = a_0 + c_0 \\ a_m = 1, b_m = 2, c_m = 1 & \quad b_m = a_m + c_m, \quad 0 < m < M \\ a_M = b_M = 1, c_M = 0 & \quad b_M = a_M + c_M \end{aligned}$$

so that  $b_m = a_m + c_m$  for *all*  $m$  and we do not have strict inequality for any  $m$ . Thus  $\mathcal{L}_h$  is not of positive type.

If  $V_m = C$  for  $m = 0, 1, \dots, M$  then  $\Delta^+ V_0 = 0$ ,  $\delta^2 V_m = 0$  for  $m = 1, 2, \dots, M-1$  and  $\Delta^- V_M = 0$ . Consequently  $\mathcal{L}_h V_m = 0$  for  $m = 0, 1, \dots, M$ . By linearity of  $\mathcal{L}_h$ ,

$$\mathcal{L}_h(U + V) = \mathcal{L}_h U + \mathcal{L}_h V = \mathcal{F}_h$$

so, if a solution exists then there are infinitely many solutions.

With  $\mathcal{F}_h$  given by (6.22), the equation  $\mathcal{L}_h U = \mathcal{F}_h$  corresponds to

$$\mathcal{L}_h U_m = \begin{cases} -h^{-1}(U_1 - U_0), \\ -h^{-2}\delta^2 U_m, \\ h^{-1}(U_M - U_{M-1}) \end{cases}, \quad \mathcal{F}_h = \begin{cases} \alpha, & m = 0, \\ f_m, & m = 1, 2, \dots, M-1, \\ \beta, & m = M, \end{cases}$$

which is consistent with the BVP  $-u''(x) = f(x)$  ( $0 < x < 1$ ) with BCs  $-u'(0) = \alpha$ ,  $u'(1) = \beta$ . If we integrate this ODE over the interval  $(0, 1)$  and apply the BCs, we find

$$-u'(1) + u'(0) = \int_0^1 f(x) dx \Rightarrow -\beta - \alpha = \int_0^1 f(x) dx.$$

Thus, no solution is possible unless the data satisfy this compatibility condition. When they do,  $u(x) + c$  is a solution for any constant  $c$  whenever  $u(x)$  is a solution.

Using the identity  $\delta^2 = \Delta^+ \Delta^-$  (see Exercise 6.2) we see that

$$\begin{aligned} \sum_{m=1}^{M-1} \delta^2 U_m &= \sum_{m=1}^{M-1} \Delta^+ (\Delta^- U_m) \\ &= (\Delta^- U_M - \Delta^- U_{M-1}) + (\Delta^- U_{M-1} - \Delta^- U_{M-2}) + \dots + (\Delta^- U_2 - \Delta^- U_1) \\ &= \Delta^- U_M - \Delta^+ U_0 \end{aligned}$$

by virtue of a telescoping series and  $\Delta^- U_1 = U_1 - U_0 = \Delta^+ U_0$ .

Summing the equations  $\mathcal{L}_h U_m = \mathcal{F}_{h,m}$  over  $m = 1, 2, \dots, M-1$  gives, using the BCs  $-h^{-1}\Delta^+ U_0 = \alpha$  and  $h^{-1}\Delta^- U_M = \beta$ ,

$$\begin{aligned} -h^{-2} [\Delta^- U_M - \Delta^+ U_0] &= \sum_{m=1}^{M-1} f_m \\ -h^{-1}(\beta + \alpha) &= \sum_{m=1}^{M-1} f_m \quad \Rightarrow \quad -\beta - \alpha = h \sum_{m=1}^{M-1} f_m \end{aligned}$$

which is a discrete analogue of the compatibility condition found earlier for the ODE. Since

$$h \sum_{m=1}^{M-1} f_m \rightarrow \int_0^1 f(x) dx$$

for any continuous function  $f$  as  $h \rightarrow 0$ , the continuous and discrete satisfy the same compatibility condition in the limit, but not necessarily for finite values of  $h$ .

### 6.23

Comparing the coefficients of  $\mathcal{L}_h$  with those in Definition 6.9,  $a_m = c_m = 8 > 0$  and  $b_m = 65 \geq a_m + c_m$  and so  $\mathcal{L}_h$  is of positive type.

$$\begin{aligned} \mathcal{L}_h A 8^m &= A \mathcal{L}_h 8^m = A 8^{m-1} (8 \times 8^2 - 65 \times 8 + 8) = 0 \\ \mathcal{L}_h B 8^{-m} &= B \mathcal{L}_h 8^{-m} = B 8^{-m+1} (8 - 65 \times 8 + 8 \times 8^2) = 0 \end{aligned}$$

so both sequences satisfy  $\mathcal{L}_h U = 0$ . Since  $\mathcal{L}_h$  is a linear operator,

$$\mathcal{L}_h (A 8^m + B 8^{-m}) = A \mathcal{L}_h 8^m + B \mathcal{L}_h 8^{-m} = 0$$

and  $U_m = A 8^m + B 8^{-m}$  is a solution for any  $A, B$ . The BCs lead to the equations

$$A + B = \alpha, \quad A 8^M + B 8^{-M} = \beta$$

which are readily solved to give

$$A = 8^{-M} \frac{\beta - \alpha 8^{-M}}{1 - 8^{-2M}}, \quad B = \frac{\alpha - \beta 8^{-M}}{1 - 8^{-2M}}$$

leading to

$$U_m = \frac{1}{1 - 8^{-2M}} (8^{m-M} (\beta - \alpha 8^{-M}) + 8^m (\alpha - \beta 8^{-M})).$$

When  $M = 10$  we find  $8^{-10} \approx 9 \times 10^{-10}$  and so, when  $\alpha = \pm 1$  and  $\beta = \pm 2$  we have  $U_m \approx \alpha 8^{-m} + \beta 8^{M-m}$ . The solutions with  $\alpha = \pm 1, \beta = 2$  are shown in Fig. 4. On the left  $\alpha = 1, \beta = 2$  so that  $\min(0, \alpha, \beta) = 0 \leq U_m \leq \max(0, \alpha, \beta) = 2$ . On the right  $\alpha = -1, \beta = 2$  so that  $\min(0, \alpha, \beta) = -1 \leq U_m \leq \max(0, \alpha, \beta) = 2$ .

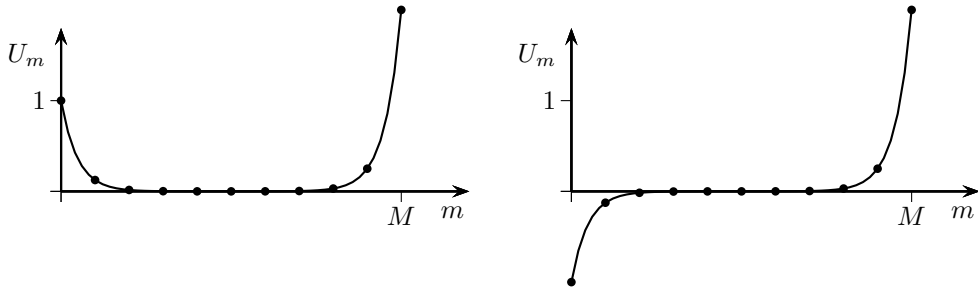


Figure 4: The points  $U_m = \alpha 8^{-m} + \beta 8^{M-m}$  for  $m = 0, 1, \dots, M$  and  $M = 10$  with  $\alpha = 1, \beta = 2$  (left) and  $\alpha = -1, \beta = 2$  (right).

### 6.25

The operator  $\mathcal{L}_h$  is of positive type by Example 6.12 (with  $r(x) = \sigma^2 > 0$ )

$$\begin{aligned}\mathcal{L}_h U_m &= \begin{cases} -h^{-1}\Delta^+ U_0 + \frac{1}{2}h\sigma^2 U_0 \\ \mathcal{L}_h U_m = -h^{-2}\delta^2 U_m + \sigma^2 U_m \\ (\sigma + \frac{1}{2}\sigma^2 h)U_M + h^{-1}\Delta^- U_M \end{cases} \\ &= \begin{cases} (h^{-1} + \frac{1}{2}h\sigma^2)U_0 - h^{-1}U_1, & m = 0, \\ h^{-2}(-U_{m-1} + (2 + \sigma^2 h^2)U_m - U_{m+1}), & m = 1, 2, \dots, M-1, \\ (\sigma + \frac{1}{2}\sigma^2 h + h^{-1})U_M - h^{-1}U_{M-1}, & m = M, \end{cases}\end{aligned}$$

which satisfy the criteria of Definition 6.17 for  $\sigma > 0$ .

### 6.27

The argument used in (6.29) and (6.30) established a second order local truncation error for  $0 < m < M$ . The local truncation error is identically zero at  $m = M$ . It remains to examine the situation for  $m = 0$ . We use the expansion  $h^{-1}\Delta^+ u_0 = (u(h) - u(0))/h = u'(0) + \frac{1}{2}hu''(0) + \mathcal{O}(h^2)$  (see Table 6.1)

$$\mathcal{R}_{h,0} = -h^{-1}\Delta^+ u_0 - 1 = (-u'(0) - 1) - \frac{1}{2}hu''(0) + \mathcal{O}(h^2) = -\frac{1}{2}hu''(0) + \mathcal{O}(h^2)$$

since  $-u'(0) = 1$ . The order of consistency is therefore first order.

For the second case,

$$\mathcal{R}_{h,0} = -h^{-1}\Delta^+ u_0 - 1 + \frac{1}{2}h = (-u'(0) - 1) - \frac{1}{2}hu''(0) + \frac{1}{2}h + \mathcal{O}(h^2) = \frac{1}{2}h(-u''(0) - 1) + \mathcal{O}(h^2).$$

However, the ODE  $-u''(x) + xu(x) = 1$  evaluated at  $x = 0$  gives  $-u''(0) = 1$  so that  $\mathcal{R}_{h,0} = \mathcal{O}(h^2)$ : the order of consistency is second order.

### 6.29

Since  $h^{-2}\delta^2 u_m = u_m'' + \mathcal{O}(h^2)$  and  $h^{-1}\Delta u_m = u_m' + \mathcal{O}(h^2)$  (see Table 6.1), the given scheme is consistent of second order.

At  $x = 1$ , we use a backward difference:

$$h^{-1}\Delta^- u_M = h^{-1}(u(1) - u(1-h)) = u'(1) - \frac{1}{2}hu''(1) + \mathcal{O}(h^2)$$

and, since the ODE evaluated at  $x = 1$  gives  $-u''(1) + 4u'(1) = 0$ , the BC  $u'(1) = 2$  gives  $u''(1) = 8$ . Consequently,

$$u'(1) = h^{-1}\Delta^- u_M + 4h + \mathcal{O}(h^2)$$

and the finite difference replacement  $h^{-1}\Delta^- u_M = 2 - 4h$  of the BC  $u'(1) = 2$  is consistent of second order.

### 6.31

$$LL^T = h^{-2} \begin{bmatrix} -1 & 1 & & & & \\ & \ddots & \ddots & & & \\ & & -1 & 1 & & \\ & & & -1 & \rho & \\ & & & & & \end{bmatrix} \begin{bmatrix} -1 & & & & & \\ 1 & \ddots & & & & \\ & \ddots & -1 & & & \\ & & 1 & -1 & & \\ & & & \rho & & \end{bmatrix} = \begin{bmatrix} 2 & -1 & & & & \\ -1 & \ddots & \ddots & & & \\ & \ddots & 2 & -1 & & \\ & & -1 & 1 + \rho^2 & & \end{bmatrix}.$$



When  $r(x) \equiv 0$  and  $a_{M,M} = 1 + ah/b$ , the matrix  $A$  in (6.47) becomes

$$A = \begin{bmatrix} 2 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & 2 & & -1 \\ & & -1 & 1 + ah/b & \end{bmatrix}.$$

Thus  $LL^\top = A$  by choosing  $\rho^2 = ah/b$ . This leads to a real value of  $\rho$  provided that  $a$  and  $b$  have the same sign.

It follows that  $\mathbf{v}^\top A \mathbf{v} = \mathbf{v}^\top LL^\top \mathbf{v} = (L^\top \mathbf{v})^\top (L^\top \mathbf{v}) = \mathbf{w}^\top \mathbf{w} \geq 0$  when  $\mathbf{w} = L^\top \mathbf{v}$  (as in Lemma 6.1). This proves that  $A$  is positive semi-definite. We need to check that it is not possible to choose a nonzero vector  $\mathbf{v}$  such that  $\mathbf{v}^\top A \mathbf{v} = 0$ . This can only occur if  $\mathbf{w} = \mathbf{0}$  but it is easily verified that the only solution of  $L^\top \mathbf{v} = \mathbf{0}$  is  $\mathbf{v} = \mathbf{0}$ .

### 6.33

Using (6.60) and the fact that  $\delta^2 f_m = 12\delta^2 x_m = 0$  leads to

$$\begin{aligned} h^{-2}(-U_{m-1} + 2U_m - U_{m+1}) + (U_{m-1} + 10U_m + U_{m+1}) &= 12x_m, \\ \text{i.e., } -(M^2 - 1)U_{m-1} + (2M^2 + 10)U_m - (M^2 - 1)U_{m+1} &= 12x_m \end{aligned}$$

since  $h^{-1} = M$ . With BCs  $u(0) = 3$  and  $u(1) = -5$ , these can be written when  $M = 4$  as

$$\begin{bmatrix} 42 & -15 & 0 \\ -15 & 42 & -15 \\ 0 & -15 & 42 \end{bmatrix} \begin{bmatrix} U_1 \\ U_2 \\ U_3 \end{bmatrix} = \begin{bmatrix} 45 + 12x_1 \\ 12x_2 \\ -75 + 12x_3 \end{bmatrix},$$

where  $x_m = m/4$ .

### 6.35

For indices  $m$  for which  $\mathcal{F}_{h,m} \neq 0$  we have the standard inequalities:

$$\mathcal{L}_h U_m = \mathcal{F}_{h,m} \leq \|\mathcal{F}_h\|_{h,\infty} \leq \|\mathcal{F}_{h,m}\|_{h,\infty} \mathcal{L}_h \Phi_m.$$

The same end result also holds for indices  $m$  such that  $\mathcal{F}_{h,m} = 0$  since  $\mathcal{L}_h U_m = 0 \leq \|\mathcal{F}_{h,m}\|_{h,\infty} \mathcal{L}_h \Phi_m$  (because the right hand side is automatically non-negative).

Thus,  $\mathcal{L}_h U_m \leq \|\mathcal{F}_{h,m}\|_{h,\infty} \mathcal{L}_h \Phi_m$  for all  $m$  from which we have

$$\mathcal{L}_h (U_m - \|\mathcal{F}_{h,m}\|_{h,\infty} \Phi_m) \leq 0$$

and then  $U_m - \|\mathcal{F}_{h,m}\|_{h,\infty} \Phi_m \leq 0$  by inverse monotonicity.

This result is particularly useful when dealing with the local truncation error since it is frequently identically zero at some grid points (see the following example).

### 6.37

Suppose that  $\mathcal{L}_h U_m := -a_m U_{m-1} + b_m U_m - c_m U_{m+1}$  and the coefficients have to be chosen so that this is to be consistent with the operator defined by  $\mathcal{L}v(x) := -v''(x) + r(x)v(x)$ .

We require  $\mathcal{L}_h v_m = (\mathcal{L}v(x))|_{x=x_m}$  for  $v(x) = 1, (x - x_m)$  and  $(x - x_m)^2$ . Since  $\mathcal{M}_h := \mathcal{L}_h - \mathcal{L}$  is a linear operator (i.e.,  $\mathcal{M}_h(u + v) = \mathcal{M}_h u + \mathcal{M}_h v$  and  $\mathcal{M}_h(cu) = c\mathcal{M}_h u$  for any sufficiently smooth functions  $u$  and  $v$  and any constant  $c$ ) it follows that

$$\mathcal{M}_h(A + B(x - x_m) + C(x - x_m)^2) = A(\mathcal{M}_h 1) + B(\mathcal{M}_h(x - x_m)) + C(\mathcal{M}_h(x - x_m)^2) = 0$$

for arbitrary constants  $A$ ,  $B$  and  $C$ . Thus  $\mathcal{M}_h p(x) = 0$  for any quadratic polynomial  $p(x)$ .

Applying the conditions  $\mathcal{L}_h v_m = (\mathcal{L}v(x))|_{x=x_m}$  for  $v(x) = 1$ ,  $(x - x_m)$  and  $(x - x_m)^2$  we find

$$\begin{aligned} \mathcal{L}_h 1 &= -a_m + b_m - c_m = r_m = (\mathcal{L}1)|_{x=x_m} \\ \mathcal{L}_h(x - x_m) &= ha_m - hc_m = 0 = (\mathcal{L}(x - x_m))|_{x=x_m} \\ \mathcal{L}_h(x - x_m)^2 &= -h^2 a_m - h^2 c_m = -2 = (\mathcal{L}(x - x_m)^2)|_{x=x_m} \end{aligned}$$

The first and third combine to give  $b_m = -2h^{-2} + r_m$  and the second and third give  $a_m = c_m = h^{-2}$ . Thus  $\mathcal{L}_h$  coincides with the approximation used in standard finite difference approximation (6.20).

### 6.39

Case (a) is entirely standard:  $u(x) = 1 - 16(x - \frac{1}{2})^4$ .

In case (b),  $u''(x) = 0$  for  $0 < x < \frac{1}{2}$  and  $u(0) = 0$  which gives  $u(x) = Ax$  for an arbitrary constant  $A$ .

For  $\frac{1}{2} < x < 1$ ,  $-u''(x) = 384(x - \frac{1}{2})^2$  which, on integrating twice gives  $u(x) = ax + b - 32(x - \frac{1}{2})^4$ . The three conditions  $u(1) = a + b - 2 = 0$ , continuity of  $u$  at  $x = \frac{1}{2}$ :  $u(\frac{1}{2}-) = u(\frac{1}{2}+)$ , i.e.,  $\frac{1}{2}A = \frac{1}{2}a + B$  and continuity of  $u'$  at  $x = \frac{1}{2}$ :  $u'(\frac{1}{2}-) = u'(\frac{1}{2}+)$ , i.e.,  $A = a$  lead to the solution

$$u(x) = 2x + \begin{cases} 0, & 0 \leq x \leq 1/2 \\ -32(x - \frac{1}{2})^4, & 1/2 < x \leq 1 \end{cases}.$$

A similar procedure in case (c) leads to

$$u(x) = 2x + \begin{cases} 0, & 0 \leq x \leq 1/2 \\ -16(x - \frac{1}{2})^3, & 1/2 < x \leq 1 \end{cases}.$$

These are shown in Fig. 5. These boundary value problems are solved numerically using the

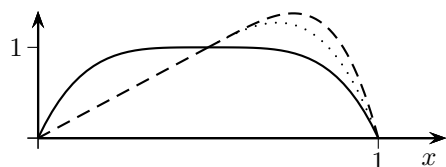


Figure 5: The solutions for Exercise 6.39: case (a) solid line, case (b) dashed line and case (c) dotted line.

standard “second order” scheme

$$-h^{-2}\delta^2 U_m = f_m, \quad m = 1, 2, \dots, M-1$$

with BCs  $U_0 = U_M = 0$ . The values of  $M$  chosen for the experiments are  $8^j$ ,  $9^j$  and  $5^j$  ( $j = 1, 2, \dots, 7$ ) and the results are summarised in Fig. 6. On the left are shown graphs of  $M^2 \times$  global error for  $M = 9$  (dots) and  $M = 16$  (crosses) in cases (a), (b) and (c). We also include the case labelled “Test” where  $f(x) = 8$  and the exact solution  $u(x) = 4x(1 - x)$  is a polynomial of degree two, for which the global error should be identically zero. The graphical results show an error of less than  $10^{-14}$  attributable to roundoff error. Including such a test in numerical experiments is recommended to test the integrity of the code. In cases (a) and (b) the graphs with different values of  $M$  are indistinguishable, in keeping with the global error  $E \propto 1/M^2$ . In case (c) the graphs of  $M^2 \times E$  appear to be two distinct continuous piecewise linear functions. Note that a grid point lies exactly on the discontinuity in  $f'(x)$  at  $x = 1/2$  when  $M = 16$  but this is not the case when  $M = 9$ .

On the right we show loglog plots of the global error as a function of  $h$ . For cases (a) and (b) the results are almost identical and lie on a line having slope two, again in keeping with  $E \propto h^2$ . The global error in case (c) is more erratic, especially for the larger values of  $h$ . However, the error itself is never larger than in cases (a) and (b) and appears to behave more smoothly as  $h \rightarrow 0$ . The theory developed in this chapter requires that the first four derivatives of the exact solution be bounded in order that the local truncation error be bounded by a multiple of  $h^2$ , and for the global error to converge to zero at a second order rate. The results of these experiments suggest that a second order convergence rate is attainable under less onerous constraints, but that goes beyond the scope of this book.

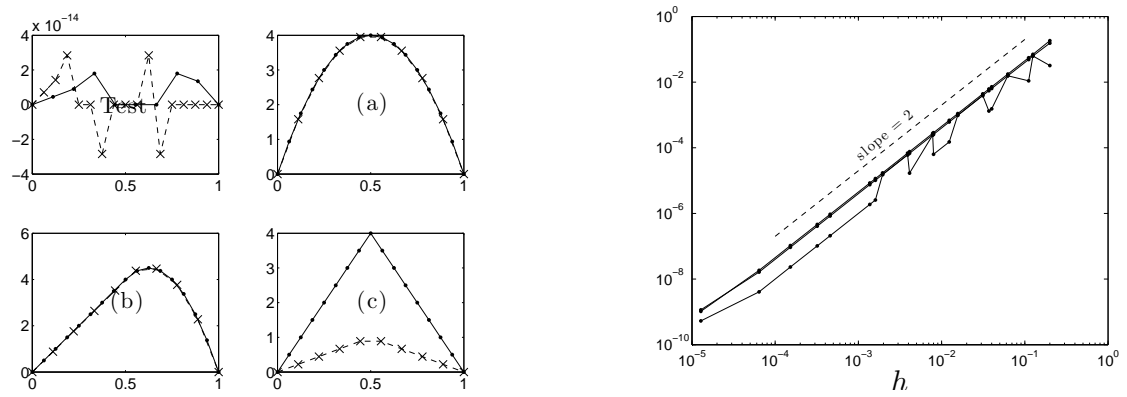


Figure 6: The graphs show  $M^2 \times E$  ( $E$  is the global error) for  $M = 9$  (dots) and  $M = 16$  (crosses) in the test case and cases (a), (b) and (c).

## Exercises 7 Maximum principles and energy methods

### 7.1

With  $v(x, t) = u(x, t) + \varepsilon(\tau - t)$  in place of (7.2) and it follows that  $-\kappa v_{xx} + v_t = -\kappa u_{xx} + u_t - \kappa\varepsilon$  so that  $-\kappa v_{xx} + v_t < 0$  for all positive values of  $\varepsilon$  and the proof proceeds as before (but restricted to  $0 \leq t \leq \tau$ ).

### 7.3

If  $\mathcal{L}u(x, t) \geq 0$  then  $-\kappa u_{xx} + u_t \geq 0$  and therefore, by Theorem 7.1 applied to  $-u$  (or part (b) of the previous solution),  $u(x, t)$  is either constant or else attains its minimum value on  $\Gamma_\tau$ . But  $\mathcal{L}u(x, t) \geq 0$  also implies that  $u(x, t) \geq 0$  for  $(x, t) \in \Gamma_\tau$ . Hence  $u(x, t) \geq 0$  for  $(x, t) \in \Omega_\tau$ .

### 7.5

The PDE  $u_t = xu_{xx} + u_x$  may be written as  $u_t = (xu_x)_x$  so, multiplying by  $u$  and integrating over the interval  $(0, 1)$  gives

$$\begin{aligned} \int_0^1 uu_t dx &= \int_0^1 u(xu_x)_x dx = (xu_x u)|_{x=0}^1 - \int_0^1 x(u_x)^2 dx \\ \frac{1}{2} \frac{d}{dt} \int_0^1 u^2 dx &= - \int_0^1 x(u_x)^2 dx \leq 0. \end{aligned}$$

The boundary terms vanish by virtue of the BCs  $u(0, t) = u(1, t) = 0$ . Thus, the energy  $E(t) := \int_0^1 u^2 dx$  is a decreasing function of  $t$  so  $E(t) \leq E(0) = \int_0^1 \sin^2 \pi x dx = 1/2$ .

### 7.7

Differentiating under the integral sign we find

$$\frac{d}{dt} \int_0^1 (u_x)^2 dx = 2 \int_0^1 u_x u_{xt} dx$$

but  $u_{xt} = u_{xxx}$ , and so

$$\frac{d}{dt} \int_0^1 (u_x)^2 dx = 2 \int_0^1 u_x u_{xxx} dx = -2 \int_1^2 (u_{xx})^2 dx + 2(u_x u_{xx}) \Big|_{x=0}^1.$$

The boundary terms clearly vanish if homogeneous Neumann conditions  $u_x = 0$  are applied at  $x = 0$  and  $x = 1$ . When homogeneous Dirichlet conditions  $u = 0$  are applied we deduce, from the PDE, that  $u_{xx} = u_t$  and boundary terms again vanish because  $u_t(0, t) = u_t(1, t) = 0$ . Thus  $\int_0^1 (u_x)^2 dx$  is a decreasing function of  $t$  and

$$\int_0^1 (u_x(x, t))^2 dx \leq \int_0^1 (u_x(x, 0))^2 dx = \int_0^1 (g'(x))^2 dx.$$

### 7.9

With homogeneous Neumann BCs:  $u_n := \vec{n} \cdot \text{grad } u = 0$  on  $\partial\Omega$ . Thus, after application of the Divergence Theorem the boundary terms still vanish so  $v$  again satisfies

$$\int_{\Omega} ((v_x)^2 + (v_y)^2) d\Omega = 0$$

from the previous solution. However, the BCs no longer allow us to conclude that  $v = 0$  from  $v_x = v_y = 0$  in  $\Omega$ . In fact the solution is only unique up to an arbitrary constant.

### 7.11

Differentiating  $m(t)$  and using the PDE gives

$$\begin{aligned} m'(t) &= \int_{-\infty}^{\infty} u_t \, dx = \int_{-\infty}^{\infty} (u_{xxx} - 6uu_x) \, dx = \int_0^1 (u_{xxx} - 3(u^2)_x) \, dx \\ &= (u_{xx} - 3u^2) \Big|_{x=-\infty}^{\infty} = 0 \end{aligned}$$

since  $u$  and its derivatives tend to zero at  $\pm\infty$ . Therefore  $m(t)$  is constant in time. Similarly for  $M(t)$ :

$$M'(t) = \int_{-\infty}^{\infty} 2uu_t \, dx = 2 \int_{-\infty}^{\infty} (uu_{xxx} - 6u^2u_x) \, dx$$

and, integrating the term involving  $uu_{xxx}$  by parts,

$$\begin{aligned} M'(t) &= 2(uu_{xx}) \Big|_{-\infty}^{\infty} - 2 \int_{-\infty}^{\infty} (u_x u_{xx} + 2(u^3)_x) \, dx = 2(uu_{xx}) \Big|_{x=-\infty}^{\infty} - \int_{-\infty}^{\infty} (\partial_x (u_x)^2 + 4(u^3)_x) \, dx \\ &= (2uu_{xx} - (u_x)^2 - 4u^3) \Big|_{-\infty}^{\infty} = 0 \end{aligned}$$

since  $u$  and its derivatives tend to zero at  $\pm\infty$ . Therefore  $M(t)$  is constant in time.

## Exercises 8 Separation of variables

### 8.1

As in Example 8.1,  $u(x, t) = X(x)T(t)$ , where  $X$  satisfies  $-X'' = \lambda X$  for  $0 < x < 1$  but, in this case, the BCs are  $X'(0) = X'(1) = 0$ .

When  $\lambda = -\mu^2 < 0$ , the general solution is  $X = Ae^{\mu x} + Be^{-\mu x}$  and the BC  $X'(0) = 0$  gives  $\mu(A - B) = 0$ . Thus, since  $\mu = 0$  is not possible (since  $\lambda < 0$ ) we must have  $A = B$ . The second BC then gives  $A\mu(e^\mu - e^{-\mu}) = 0$ . But  $e^\mu \neq e^{-\mu}$  for  $\mu \neq 0$  so we are left with  $A = 0$ , leading to the trivial solution  $X(x) = 0$ .

When  $\lambda = 0$ , the general solution is  $X = A + Bx$  and the BCs  $X'(0) = X'(1) = 0$  both require  $B = 0$  with no restriction on  $A$ . Thus  $X(x) = A$  is a nontrivial solution corresponding to an eigenvalue  $\lambda = 0$ . Since  $T'(t) = -\lambda T(t)$ , we have  $T(t) = \text{constant}$  and the corresponding solution of the heat equation is  $u(x, t) = A$ .

When  $\lambda = \mu^2 > 0$ , the general solution is  $X = A \sin \mu x + B \cos \mu x$  and the BC  $X'(0) = 0$  gives  $\mu A = 0$ . Thus, since  $\mu = 0$  is not possible (since  $\lambda > 0$ ) we must have  $A = 0$ . The second BC then gives  $B\mu \sin \mu = 0$ . To avoid the trivial solution,  $\mu$  must be chosen so that  $\sin \mu = 0$ , thus  $\mu = n\pi$  for  $n = 1, 2, \dots$ . Correspondingly,  $T'(t) = -\lambda T(t)$ , so  $T(t) = Ce^{-(n\pi)^2 t}$ , leading to fundamental solutions  $e^{-(n\pi)^2 t} \cos n\pi x$ .

The general solution is a linear combination of all fundamental solutions and so takes the form

$$u(x, t) = u(x, t) = \frac{1}{2}A_0 + \sum_{n=1}^{\infty} A_n e^{-\kappa n^2 \pi^2 t} \cos n\pi x.$$

The factor  $1/2$  in the leading term allows all the coefficients to be determined by the same formula

$$A_n = 2 \int_0^1 g(x) \cos n\pi x \, dx, \quad n = 0, 1, 2, \dots$$

(see Example 8.2). When  $g(x) = x$  we find  $A_0 = 2 \int_0^1 x \, dx = 1$  and, using integration by parts,

$$\begin{aligned} A_n &= 2 \int_0^1 x \cos n\pi x \, dx = 2x \frac{\sin n\pi x}{n\pi} \Big|_0^1 - \frac{2}{n\pi} \int_0^1 \sin n\pi x \, dx \\ &= \frac{2}{n^2 \pi^2} \cos n\pi x \Big|_0^1 = \frac{2}{n^2 \pi^2} ((-1)^n - 1). \end{aligned}$$

Thus,  $A_n = 0$  when  $n$  is even and  $A_n = 4/(n\pi)^2$  when  $n$  is odd.

### 8.3

The mean value theorem states that if  $g$  is continuous on an interval  $[a, b]$  and differentiable on the open interval  $(a, b)$  then there is a point  $c \in (a, b)$  such that

$$g'(c) = \frac{g(b) - g(a)}{b - a}.$$

If we choose  $b = x$ ,  $a = 1 - x$ , then, for  $x > \frac{1}{2}$ ,

$$g'(c) = \frac{g(x) - g(1-x)}{2x-1} = 0, \quad 1-x < c < x.$$

It follows that  $g'(\frac{1}{2}) = 0$  by taking the limit  $x \rightarrow \frac{1}{2}$ .

With the change of variable  $x = 1 - s$  and  $v(s, t) = u(1 - s, t)$ ,

$$v_t = u_t(1 - s, t), \quad v_s = -u_x(1 - s, t), \quad v_{ss} = u_{xx}(1 - s, t)$$

and therefore

$$v_t - \kappa v_{ss} = u_t(1 - s, t) - \kappa u_{xx}(1 - s, t) = u_t(x, t) - \kappa u_{xx}(x, t) = 0.$$

Also  $v(0, t) = u(1, t) = 0$ ,  $v(1, t) = u(0, t) = 0$  and  $v(s, 0) = u(1 - s, 0) = g(1 - s) = g(s)$  (since  $g$  is symmetric about  $x = 1/2$ ).

### 8.5

From Example 8.1 the general solution of the heat equation with BCs  $u = 0$  at both  $x = 0$  and  $x = 1$  is given by (8.6) and the coefficients by (8.8). When  $g(1 - x) = g(x)$  we have, making a change of variable  $x = 1 - s$ ,

$$\begin{aligned} \langle g, X_n \rangle &= \int_0^1 g(x) \sin n\pi x \, dx = \int_0^1 g(1 - x) \sin n\pi x \, dx = \int_0^1 g(s) \sin n\pi(1 - s) \, ds \\ &= - \int_0^1 g(s) \cos n\pi \sin n\pi s \, ds = (-1)^{n+1} \int_0^1 g(s) \sin n\pi s \, ds = (-1)^{n+1} \langle g, X_n \rangle. \end{aligned}$$

It follows that  $\langle g, X_n \rangle = -\langle g, X_n \rangle$  when  $n$  is even. Consequently  $\langle g, X_n \rangle = 0$ , and so  $A_n = 0$  when  $n$  is even.

### 8.7

From Exercise 4.20

$$u_x = \cos \theta u_r - \frac{1}{r} \sin \theta u_\theta = \cos \theta f'(r),$$

since  $u = f(r)$ . Allowing  $x, y \rightarrow 0$  would lead to  $u_x(x, y)$  being multivalued at the origin (because it would depend on the angle at which the origin was approached) unless  $f'(0) = 0$ .

### 8.9

When  $g(r) = 1 - (r/b)^2$  for  $0 \leq r \leq b$  and  $g(r) = 0$  for  $b < r \leq a$

$$\int_0^a r g(r) J_0\left(\frac{r\xi_n}{a}\right) \, dr = \int_0^b r(1 - (r/b)^2) J_0\left(\frac{r\xi_n}{a}\right) \, dr$$

From Exercise D.4 and writing  $\omega = b\xi_n/a$ ,

$$\int_0^b r g(r) J_0\left(\frac{r\xi_n}{a}\right) \, dr = \left(\frac{a}{\xi_n}\right)^2 \int_0^\omega x J_0(x) \, dx = \frac{b^2}{\omega} J_1(\omega)$$

and, using  $\int x^3 J_0(x) \, dx = 2x^2 J_0(x) + x(x^2 - 4)J_1(x)$ , we find

$$\int_0^b r \left(\frac{r}{b}\right)^2 J_0\left(\frac{r\xi_n}{a}\right) \, dr = \frac{1}{b^2} \left(\frac{a}{\xi_n}\right)^4 \int_0^\omega x^3 J_0(x) \, dx = \frac{b^2}{\omega^4} (2\omega^2 J_0(\omega) + \omega(\omega^2 - 4)J_1(\omega)).$$

These combine to give the required result.

**8.11**

When  $a = 1/2$ ,  $b = 1/4$ , and  $g(r) = 1$  for  $0 < r < b$  and is otherwise zero, we find on integration by parts,

$$A_n = 4 \int_0^{\frac{1}{4}} r \sin(2n\pi r) \, dr = \frac{1}{n^2 \pi^2} (\sin \frac{1}{2} n\pi - \frac{1}{2} n\pi \cos \frac{1}{2} n\pi)$$

Hence

$$A_n = \begin{cases} -\frac{1}{4m\pi} \cos m\pi = \frac{1}{4m\pi} (-1)^{m-1}, & \text{when } n = 2m \\ \frac{1}{(2m-1)^2 \pi^2} \sin \frac{1}{2} (2m-1)\pi = \frac{1}{(2m-1)^2 \pi^2} (-1)^{m-1}, & \text{when } n = 2m-1. \end{cases}$$

**8.13**

The solution is given by (8.26) with the coefficients in Exercise 8.11 with  $a = 1$  and  $g(r) = 1$ . Thus,

$$A_n = 2 \int_0^1 r \sin(n\pi r) \, dr = -\frac{2}{n\pi} \cos(n\pi) = \frac{2}{n\pi} (-1)^{n-1}$$

and the solution is

$$u(r, t) = 2 \sum_{n=1}^{\infty} (-1)^{n-1} e^{-n^2 \pi^2 t} \frac{\sin n\pi r}{n\pi r}.$$

The result follows by taking the limit  $r \rightarrow 0$  and using the fact that  $\frac{\sin x}{x} \rightarrow 1$ . At large times the leading term dominates and  $u(0, t) \approx 2e^{-\pi^2 t}$ .

**8.15**

When  $g_0$  is given by (8.43)

$$g_0(x+ct) + g_0(x-ct) = \sum_{n=1}^{\infty} A_n (\sin n\pi(x+ct) + \sin n\pi(x-ct)) = 2 \sum_{n=1}^{\infty} A_n \sin n\pi x \cos n\pi ct.$$

When  $g_1$  is given by (8.44)

$$\begin{aligned} \int_{x-ct}^{x+ct} g_1(z) \, dz &= \sum_{n=1}^{\infty} B_n c n \pi \int_{x-ct}^{x+ct} \sin n\pi z \, dz \\ &= -c \sum_{n=1}^{\infty} B_n (\cos n\pi(x+ct) - \cos n\pi(x-ct)) = 2c \sum_{n=1}^{\infty} B_n \sin n\pi x \sin n\pi ct. \end{aligned}$$

Combining these in d'Alembert's solution (4.20) leads to (8.42).

**8.17**

Substituting  $u(x, t) = X(x)T(t)$  into the PDE gives, on dividing by  $X(x)T(t)$ ,

$$\frac{X''(x)}{X(x)} = \frac{T''(t)}{T(t)} = -\lambda$$

leading to the eigenvalue problem  $-X''(x) = \lambda X(x)$ ,  $-a < x < a$  with  $X(-a) = X(a) = 0$ . The standard arguments can be used to show that only trivial solutions are possible when  $\lambda \leq 0$ .



When the origin is not one of the endpoints the calculations are simplified by writing the general solution for  $\lambda = \mu^2 > 0$  as

$$X(x) = A \sin \mu(x + \beta)$$

for arbitrary constants  $A$  and  $\beta$ . The BC  $X(-a) = 0$  immediately sets  $\beta = a$  and  $X(a) = 0$  leads to  $A \sin 2\mu a = 0$ . Thus, to avoid trivial solutions we must choose  $\mu = \frac{1}{2}n\pi/a$  and the eigenvalues are  $\lambda_n = (\frac{1}{2}n\pi)^2$  with corresponding eigenfunctions  $X_n(x) = \sin(\frac{1}{2}n\pi(x + a)/a)$ . These eigenfunctions are closely related to those in Example 8.1.

The ODE  $T''(t) + \mu^2 T(t) = 0$  has the general solution  $T(t) = C \sin \mu t + D \cos \mu t$ , for arbitrary constants  $C$  and  $D$ . This leads to the general solution

$$u(x, t) = \sum_{n=1}^{\infty} \left( C_n \sin \frac{n\pi t}{2a} + D_n \cos \frac{n\pi t}{2a} \right) \sin \frac{n\pi(x + a)}{2a}$$

that satisfies the PDE and all BCs.

The initial conditions  $u(x, 0) = g_0(x)$  and  $u_t(x, 0) = 0$  give

$$g_0(x) = \sum_{n=1}^{\infty} D_n \sin \frac{n\pi(x + a)}{2a}, \quad 0 = \sum_{n=1}^{\infty} C_n \frac{n\pi}{2a} \sin \frac{n\pi(x + a)}{2a}.$$

Since the eigenfunctions are orthogonal on the interval  $(-a, a)$ ,

$$D_n = \frac{\langle g_0, X_n \rangle}{\langle X_n, X_n \rangle}, \quad C_n = 0.$$

and  $\langle X_n, X_n \rangle = \frac{1}{2}a$ .

When  $g_0(x) = \max(0, 1 - (x/b)^2)$  ( $b < a$ ) a change of variable  $x = bs$  and integrating by parts twice gives

$$\begin{aligned} D_n &= \frac{2}{a} \int_{-b}^b \left(1 - \left(\frac{x}{b}\right)^2\right) \sin \frac{1}{2}n\pi(x + a)/a \, dx = \frac{2b}{a} \int_{-1}^1 (1 - s^2) \sin \omega(s + \alpha) \, ds, \quad \omega = \frac{n\pi b}{a}, \quad \alpha = \frac{a}{b} \\ &= 8 \frac{1}{\omega} \frac{\sin \alpha \omega}{\alpha \omega} \left( \frac{\sin \omega}{\omega} - \cos \omega \right). \end{aligned}$$

The solution is shown in Fig. 8.7 (dashed line) at time  $t = .75$  with  $a = 1/2$ ,  $b = 1/4$  (see Example 8.11).

### 8.19

Suppose  $u = \Phi(x, y, a, b)$  represents the solution (8.56) of subproblem  $P_1$  in Example 8.12, so that

$$\Phi(x, 0, a, b) = \sum_{n=1}^{\infty} A_n \sin \left( \frac{n\pi x}{a} \right)$$

on  $E_1$  and  $\Phi = 0$  on the remaining edges. The PDE is unchanged under the linear change of variable  $y \mapsto b - y$  ( $x$  unchanged) and maps the rectangle into itself with the edges  $E_1$  being interchanged  $E_3$ . This gives the solution

$$u(x, y) = \Phi(x, b - y, a, b) = \sum_{n=1}^{\infty} C_n \frac{\sinh(n\pi y/(ab))}{\sinh(n\pi/a)} \sin \left( \frac{n\pi x}{a} \right), \quad u(x, b) = \sum_{n=1}^{\infty} C_n \sin \left( \frac{n\pi x}{a} \right)$$

and  $u = 0$  on edges  $E_j$ ,  $j = 1, 2, 4$ . This is the solution to  $P_3$ .

The PDE is unchanged under the interchange  $x \leftrightarrow y$  and the domain is also unaltered if we make the interchange  $a \leftrightarrow b$ . Thus

$$u(x, y) = \Phi(y, x, b, a) = \sum_{n=1}^{\infty} D_n \frac{\sinh(n\pi(1-x/a)/b)}{\sinh(n\pi/b)} \sin\left(\frac{n\pi y}{b}\right), \quad u(0, y) = \sum_{n=1}^{\infty} D_n \sin\left(\frac{n\pi y}{b}\right)$$

is the solution to  $P_4$ .

Finally, applying the linear change of variable  $x \mapsto a - x$  ( $y$  unchanged) gives the solution to  $P_2$ :

$$u(x, y) = \Phi(y, a - x, b, a) = \sum_{n=1}^{\infty} B_n \frac{\sinh(n\pi x/(ab))}{\sinh(n\pi/b)} \sin\left(\frac{n\pi y}{b}\right), \quad u(a, y) = \sum_{n=1}^{\infty} B_n \sin\left(\frac{n\pi y}{b}\right).$$

### 8.21

By virtue of Exercise 8.19 the solution to problem  $P_2$  is given by

$$u(x, y) = \sum_{n=1}^{\infty} B_n \frac{\sinh(n\pi x)}{\sinh n\pi} \sin(n\pi y), \quad u(1, y) = \sum_{n=1}^{\infty} B_n \sin(n\pi y)$$

and, with  $u(1, y) = 1 - y$ ,

$$B_n = 2 \int_0^1 (1 - y) \sin(n\pi y) dy = \frac{2}{n\pi} - 2 \int_0^1 \frac{1}{n\pi} \cos n\pi y dy = \frac{2}{n\pi} - 2 \frac{\sin n\pi}{n^2 \pi^2} = \frac{2}{n\pi}.$$

### 8.23

Substituting  $u(x, y) = X(x)Y(y)$  into the PDE and dividing by  $X(x)Y(y)$  gives

$$\frac{-X''(x) + 2X'(x)}{X(x)} = \frac{Y''(y)}{Y(y)}$$

and, the left hand side being a function of  $x$  only, while the right hand side is a function of  $y$  only, we deduce that both must be constant. The homogeneous BCs  $u(x, 0) = u(x, 1) = 0$  imply that  $Y(0) = Y(1) = 0$ . We therefore look for eigenfunctions in the  $y$ -variable and set the separation constant to  $-\lambda$  so that  $-Y'' = \lambda Y$ .

This is the eigenvalue problem solved in Example 8.1 (for  $X(x)$ ) and so the eigenvalues are  $\lambda_n = (n\pi)^2$  with corresponding eigenfunctions  $Y_n = \sin(n\pi y)$ ,  $n = 1, 2, \dots$ . The corresponding ODE for  $X$  is  $-X'' + 2X' + \lambda X = 0$ . The general solution is a linear combination of  $e^{x+\sigma x}$  and  $e^{x-\sigma x}$ , where  $\sigma = \sqrt{1+\lambda^2}$ . This may be written in several equivalent ways but, in view of the BC  $X'(1) = 0$  being applied at  $x = 1$ , the most convenient form is

$$X(x) = e^x (C \cosh \sigma(1-x) + D \sinh \sigma(1-x))$$

which satisfies  $X'(1) = 0$  if  $C = \sigma D$ . Thus the fundamental solutions are

$$u_n = e^x (\sigma \cosh \sigma(1-x) + \sinh \sigma(1-x)) \sin n\pi y, \quad \sigma = \sqrt{1+\lambda_n^2},$$

for  $n = 1, 2, \dots$  and the general solution is  $u(x, y) = \sum_{n=1}^{\infty} D_n u_n(x, y)$ .

To match the additional BC  $u(0, y) = y(1 - y)$ , we have

$$y(1 - y) = \sum_{n=1}^{\infty} A_n \sin n\pi y, \quad A_n = D_n(\sigma \cosh \sigma + \sinh \sigma),$$

so that (see the discussion from equation (8.7) to (8.8))

$$A_n = 2 \int_0^1 y(1 - y) \sin n\pi y \, dy = \frac{4}{(n\pi)^3} (1 - (-1)^n).$$

Hence,

$$\begin{aligned} u(x, y) &= e^x \sum_{n=1}^{\infty} \frac{4}{(n\pi)^3} (1 - (-1)^n) \frac{\sigma \cosh \sigma(1 - x) + \sinh \sigma(1 - x)}{\sigma \cosh \sigma + \sinh \sigma} \sin n\pi y \\ &= e^x \sum_{\substack{n=1 \\ n \text{ odd}}}^{\infty} \frac{8}{(n\pi)^3} \frac{\sigma \cosh \sigma(1 - x) + \sinh \sigma(1 - x)}{\sigma \cosh \sigma + \sinh \sigma} \sin n\pi y. \end{aligned}$$

### 8.25

Using Exercise 4.20 we find that Laplace's equation in polar coordinates becomes

$$u_{rr} + \frac{1}{r}u_r + \frac{1}{r^2}u_{\theta\theta} = 0$$

and substituting  $u(r, \theta) = R(r)\Theta(\theta)$  into the PDE leads to

$$\frac{r^2 R''(r) + rR'(r)}{R(r)} = -\frac{\Theta''(\theta)}{\Theta(\theta)}.$$

The left hand side is a function of  $r$  only, while the right hand side is a function of  $\theta$  only so we deduce that both must be constant. The BCs  $u(r, 0) = u(r, \pi/4)$  suggest that we look for eigenfunctions in the  $\theta$ -variable and set the separation constant to  $\lambda$  so that  $-\Theta'' = \lambda\Theta$ .

(a)  $\lambda = -\mu^2 < 0$ , then  $-\Theta'' = -\mu^2\Theta$  has general solution  $\Theta = Ae^{\mu\theta} + Be^{-\mu\theta}$ . This cannot satisfy  $\Theta(0) = \Theta(\pi/4) = 0$  for any  $\mu \neq 0$  unless  $A = B = 0$ .

(b)  $\lambda = 0$ , then  $-\Theta'' = 0$  has general solution  $\Theta = A + B\theta$  which cannot satisfy  $\Theta(0) = \Theta(\pi/4) = 0$  unless  $A = B = 0$ .

(c)  $\lambda = \mu^2 > 0$ , then  $-\Theta'' = \mu^2\Theta$  has general solution  $\Theta = A \sin \mu\theta + B \cos \mu\theta$ . The BC  $\Theta(0) = 0$  requires  $B = 0$ . Then  $\Theta(\pi/4) = 0$  requires  $A \sin \mu\pi/4 = 0$ . This will lead to the trivial solution  $A = 0$  unless  $\mu = 4n$ ,  $n = 1, 2, \dots$ . The eigenvalues are therefore given by  $\lambda_n = 16n^2$  and corresponding eigenfunctions  $\Theta_n = \sin(4n\theta)$ . Then  $R$  satisfies

$$r^2 R''(r) + rR'(r) = 16n^2 R.$$

Looking for solutions of the form  $R = Ar^\alpha$  we find that  $r^2 R''(r) + rR'(r) = \alpha^2 Ar^\alpha$  and so  $\alpha = \pm 4n$ . There are two solutions and, by linearity of the ODE, a linear combination is also a solution and therefore

$$R(r) = Dr^{4n} + Er^{-4n}$$

for arbitrary constants  $D$  and  $E$ . The BC  $u(1, \theta)$  implies that  $R(1) = 0$  so that  $E = -D$  and the general solution satisfying the PDE and the homogeneous BCs is

$$u(r, \theta) = \sum_{n=1}^{\infty} D_n (r^{4n} - r^{-4n}) \sin 4n\theta.$$

The BC  $u(2, \theta) = g(\theta)$  then means that the coefficients are given by

$$D_n (2^{4n} - 2^{-4n}) = \frac{8}{\pi} \int_0^{\pi/4} g(\theta) \sin 4n\theta \, d\theta.$$

due to the mutual orthogonality of the eigenfunctions over  $(0, \pi/4)$  and

$$\int_0^{\pi/4} \sin^2 4n\theta \, d\theta = \frac{1}{8}\pi.$$

### 8.27

Written in polar coordinates, we require the eigenvalues  $\lambda$  and eigenfunctions  $u$  (not identically zero) such that

$$-\left(u_{rr} + \frac{1}{r}u_r + \frac{1}{r^2}u_{\theta\theta}\right) = \lambda u.$$

Substituting  $u(r, \theta) = R(r)\Theta(\theta)$  into the PDE leads to

$$-\frac{r^2 R''(r) + rR'(r) + \lambda r^2 R(r)}{R(r)} = \frac{\Theta''(\theta)}{\Theta(\theta)}.$$

The left hand side is a function of  $r$  only, while the right hand side is a function of  $\theta$  only so we deduce that both must be constant. The periodic BCs  $u(r, \theta) = u(r, \theta + 2k\pi)$  for any integer  $k$  suggests that we look for eigenfunctions in the  $\theta$ -variable and set the separation constant to  $\alpha$  so that  $-\Theta'' = \alpha\Theta$ .

It is readily shown that there are no periodic solutions for  $\alpha < 0$  so we set  $\alpha = \nu^2 \geq 0$ , so that  $\Theta = A \sin \nu\theta + B \cos \nu\theta$ . This will be periodic of period  $2\pi$  if  $\nu = n$ ,  $n = 0, 1, \dots$  (Note that  $n = 0$  is permissible). This being the case, then

$$r^2 R'' + rR'(\lambda r^2 - n^2)R = 0.$$

The simple change of variable  $x = r\sqrt{\lambda}$  converts this into Bessel's equation (D.1) with  $\nu = n$ . The general solution that remains bounded at the centre of the circle is

$$R_n(r) = C J_n(r\sqrt{\lambda}),$$

where  $J_n$  is the Bessel function of the first kind of order  $n$  (see Appendix D). It is necessary to have  $R(a) = 0$  in order for  $u(a, \theta) = 0$ . Therefore, since  $C = 0$  leads to the trivial solution, we have to choose  $\lambda$  in such a way that  $a\sqrt{\lambda} = \xi_{n,m}$ , the  $m$ th nonnegative zero of  $J_n$ ,  $m = 1, 2, \dots$ . Hence the eigenvalues are  $\lambda_n = (\xi_{n,m}/a)^2$  and the eigenfunctions are as given in the question.

## Exercises 9 The method of characteristics

### 9.1

(a)  $u_t + tu_x = u$ :

$$\frac{dt}{1} = \frac{dx}{t} = \frac{du}{u} \Rightarrow \frac{dx}{dt} = t, \quad \frac{du}{dt} = u$$

giving, in terms of  $(t, k)$ :  $x = k + \frac{1}{2}t^2$ ,  $u = A(k)e^t$ .  
In terms of  $(x, t)$ :  $u = A(x - \frac{1}{2}t^2)e^t$ .

(b)  $tu_t - u_x = 1$ :

$$\frac{dt}{t} = \frac{dx}{-1} = \frac{du}{1} \Rightarrow \frac{dx}{dt} = -t, \quad \frac{du}{dx} = -1$$

giving, in terms of  $(x, k)$ :  $t = ke^{-x}$  and  $u = A(k) - x$ .  
In terms of  $(x, t)$ :  $u = A(te^x) - x$ .

(c)  $u_t + xu_x = -u$ :

$$\frac{dt}{1} = \frac{dx}{x} = \frac{du}{-u} \Rightarrow \frac{dx}{dt} = x, \quad \frac{du}{dt} = -u$$

giving, in terms of  $(t, k)$ :  $x = ke^t$  and  $u = A(k)e^{-t}$ .  
In terms of  $(x, t)$ :  $u = A(xe^{-t})e^{-t}$ .

(d)  $xu_t - u_x = t$ :

$$\frac{dt}{x} = \frac{dx}{-1} = \frac{du}{t} \Rightarrow \frac{dx}{dt} = -x, \quad \frac{du}{dx} = -t$$

giving,  $t = k - \frac{1}{2}x^2$  and so  $\frac{du}{dx} = -k + \frac{1}{2}x^2$ .

Hence in terms of  $(x, k)$ :  $u = A(k) - kx + \frac{1}{6}x^3$ .  
In terms of  $(x, t)$ :  $u = A(t + \frac{1}{2}x^2) - (t + \frac{1}{2}x^2)x + \frac{1}{6}x^3$ .

(e)  $tu_t + xu_x = x$ :

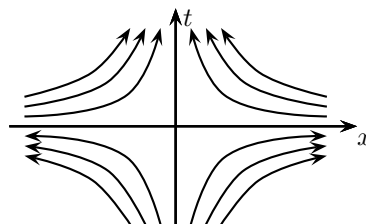
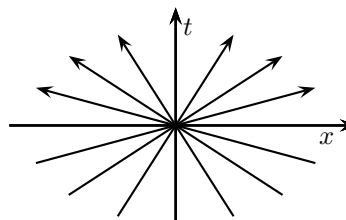
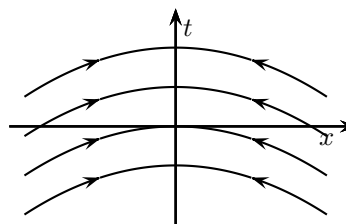
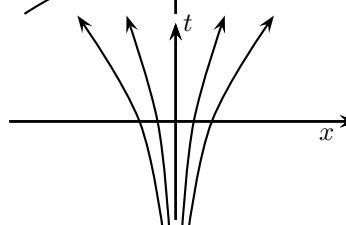
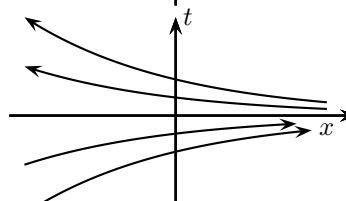
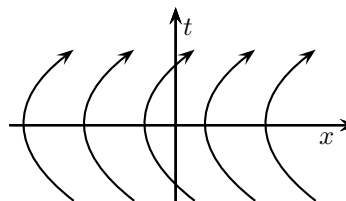
$$\frac{dt}{t} = \frac{dx}{x} = \frac{du}{x} \Rightarrow \int \frac{dt}{t} = \int \frac{dx}{x}, \quad \frac{du}{dx} = 1$$

giving, in terms of  $(x, k)$ :  $t = kx$  and  $u = A(k) + x$ .  
In terms of  $(x, t)$ :  $u = A(t/x) + x$ .

(f)  $tu_t - xu_x = t$ :

$$\frac{dt}{t} = \frac{dx}{-x} = \frac{du}{t} \Rightarrow \int \frac{dt}{t} = - \int \frac{dx}{x}, \quad \frac{du}{dt} = 1$$

giving, in terms of  $(t, k)$ :  $x = k/t$  and  $u = A(k) + t$ .  
In terms of  $(x, t)$ :  $u = A(tx) + t$ .



(g)  $xu_t - tu_x = xt$ :

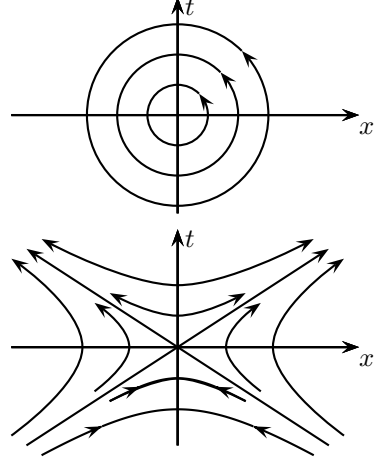
$$\frac{dt}{x} = \frac{dx}{-t} = \frac{du}{xt} \Rightarrow \int t \, dt = - \int x \, dx, \quad \frac{du}{dt} = t$$

giving, in terms of  $(t, k)$ :  $x^2 = k - t^2$  and  $u = A(k) + \frac{1}{2}t^2$ .  
In terms of  $(x, t)$ :  $u = A(x^2 + t^2) + \frac{1}{2}t^2$ .

(h)  $xu_t + tu_x = -xu$ :

$$\frac{dt}{x} = \frac{dx}{t} = \frac{du}{-xu} \Rightarrow \int t \, dt = \int x \, dx, \quad \frac{du}{dt} = -u$$

giving, in terms of  $(t, k)$ :  $x^2 = k + t^2$  and  $u = A(k)e^{-t}$ .  
In terms of  $(x, t)$ :  $u = A(x^2 - t^2)e^{-t}$ .



### 9.3

The nonzero component of  $\mathbf{f}_1$  is equal to the first component of  $V^T \mathbf{g}(X - \lambda_1 T)$ . Thus  $\mathbf{f}_1 = D_1 V^T \mathbf{g}(X - \lambda_1 T)$ , where  $D_1$  is the  $d \times d$  matrix zero matrix except that its  $(1, 1)$  entry is equal to one. Hence,

$$\mathbf{u}_1 = V^{-T} D_1 V^T \mathbf{g}(X - \lambda_1 T) \Rightarrow \|\mathbf{u}_1\|_2 \leq \|V^{-T}\|_2 \|D_1\|_2 \|V^T\|_2 \|\mathbf{g}(X - \lambda_1 T)\|_2.$$

However,  $\|V^{-T}\|_2 = \|V^{-1}\|_2$ ,  $\|V^T\|_2 = \|V\|_2$ ,  $\|D_1\|_2 = 1$  and  $\|\mathbf{g}(X - \lambda_1 T)\|_2 \leq M_2$  so the result follows.

A similar argument holds for  $V^T \mathbf{u}_j = \mathbf{f}_j$ , where the only nonzero component of  $\mathbf{f}_j$  is its  $j$ th component, which is  $\mathbf{v}_j^T \mathbf{g}(X - \lambda_j T)$ . This time we can write  $\mathbf{f}_j = D_j V^T \mathbf{g}(X - \lambda_j T)$ , where  $D_j$  is the  $d \times d$  matrix zero matrix except that its  $(j, j)$  entry is equal to one.

The bound on the solution of (9.6) follows by applying the triangle inequality to  $\mathbf{u} = \sum_{j=1}^d \mathbf{u}_j$ :  
 $\|\mathbf{u}\|_2 \leq \sum_{j=1}^d \|\mathbf{u}_j\|_2$ .

When  $A$  is symmetric its eigenvalues are mutually orthogonal and  $\kappa_2(V) = 1$ .

### 9.5

At  $P_1$ , where  $X_1 > 2T_1$ , (see Fig. 7, Left)

- (a)  $AP_1$  is a  $\Gamma_1$ -characteristic:  $x+t = X_1+T_1$  along which  $u+w = \text{constant}$ . Thus  $x_A = X_1+T_1$  and

$$u(P) + w(P) = u(A) + w(A). \quad (9.5a)$$

- (b)  $BP_1$  is a  $\Gamma_2$ -characteristic:  $x-t = X_1-T_1$  along which  $v+w = \text{constant}$ . Thus  $x_B = X_1-T_1$  and

$$v(P) + w(P) = v(B) + w(B). \quad (9.5b)$$

- (c)  $CP_3$  is a  $\Gamma_3$ -characteristic:  $x-2t = X_1-2T_1$  along which  $u+v = \text{constant}$ . Thus  $x(C) = X_1-2T_1$  and

$$u(P) + v(P) = u(C) + v(A). \quad (9.5c)$$

Equations (9.5a-c) are clearly of the form (9.13), where the columns of  $V$  are the eigenvectors of  $A^T$  and

$$\mathbf{f} = [\mathbf{v}_3^T \mathbf{g}(A), \mathbf{v}_2^T \mathbf{g}(B), \mathbf{v}_1^T \mathbf{g}(C)]^T.$$

The initial-conditions provide values for  $u$ ,  $v$  and  $w$  at A, B, and C and these equations may be solved to give  $u$ ,  $v$  and  $w$  at  $P_1$ .

At  $P_3$ , where  $X_3 < T_3$ , (see Fig. 7, Centre)

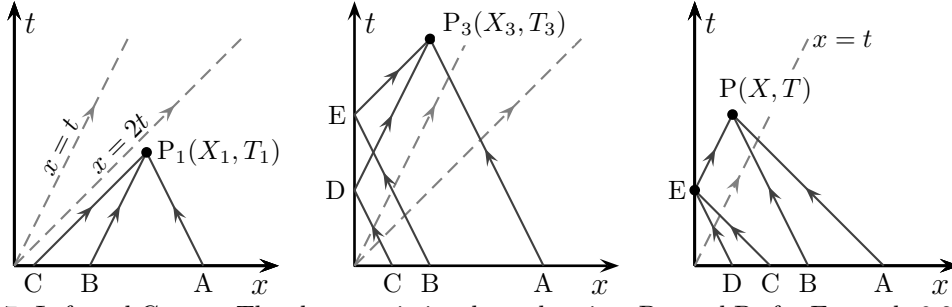


Figure 7: Left and Centre: The characteristics through points  $P_1$ , and  $P_3$  for Example 9.2 drawn backwards in time, with reflections when they intersect the  $t$ -axis. The dashed lines show the characteristics  $x - t = 0$  and  $x - 2t = 0$  that pass through the origin. Right: The characteristics for Exercise 9.7.

- (a)  $AP_3$  is a  $\Gamma_1$ -characteristic:  $x + t = X_3 + T_3$  along which  $u + w = \text{constant}$ . Thus  $x_A = X_3 + T_3$  and

$$u(P) + w(P) = u(A) + w(A). \quad (9.5A)$$

- (b) The  $\Gamma_2$ -characteristic through  $P_3$  intersects the  $t$ -axis at D:  $x - t = X_3 - T_3$  along which  $v + w = \text{constant}$ . Thus  $t_D = T_3 - X_3$  and  $v(P) + w(P) = v(D) + w(D)$ . One of the BCs on  $x = 0$  is  $u(0, t) = 0$  and so  $v(P) + w(P) = v(D)$ . Now CD is a  $\Gamma_1$ -characteristic so  $x_C = t_D$ ,  $u + w$  is constant and, applying the BCs, gives  $u(C) + w(C) = v(D)$ . Combining these results leads to

$$v(P) + w(P) = u(C) + w(C). \quad (9.5B)$$

- (c) The  $\Gamma_3$ -characteristic through  $P_3$  intersects the  $t$ -axis at E:  $x - 2t = X_3 - 2T_3$  along which  $u + v = \text{constant}$ . Thus  $t_E = T_3 - \frac{1}{2}X_3$  and  $u(P) + v(P) = u(E) + v(E)$ . On of the BCs on  $x = 0$  is  $u(0, t) = v(0, t)$  and so  $u(P) + v(P) = 2u(E)$ . Now BE is a  $\Gamma_1$ -characteristic so  $x_B = t_E$ ,  $u + w$  is constant and, applying the BCs, gives  $u(B) + w(B) = u(E)$ . Combining these results leads to

$$u(P) + v(P) = 2u(B) + 2w(B). \quad (9.5C)$$

Equations (9.5A-C) are clearly of the form (9.13), where the columns of  $V$  are the eigenvectors of  $A^T$  and the components of  $\mathbf{f}$  are linear combinations of the values of  $\mathbf{u}$  evaluated at the points where the characteristics intersect the  $x$ -axis.

## 9.7

The three families of characteristics are:

$$\begin{aligned} \Gamma_1 : \lambda_1 &= -2, & x + 2t &= \text{constant}, & \mathbf{v}_1^T \mathbf{u} &= u + w = \text{constant}, \\ \Gamma_2 : \lambda_2 &= -1, & x + t &= \text{constant}, & \mathbf{v}_2^T \mathbf{u} &= v + w = \text{constant}, \\ \Gamma_3 : \lambda_3 &= 1, & x - t &= \text{constant}, & \mathbf{v}_3^T \mathbf{u} &= u + v = \text{constant}. \end{aligned}$$

In order to find the solution at  $P(X, T)$  with  $X < T$  (see Fig. 7, Right) we follow the characteristics through  $P$  backwards in time until they intersect the initial line  $t = 0$  ( $x \geq 0$ ). When a characteristic intersects the  $t$ -axis (at E, for example), we also have to include the characteristics through E as shown.

AP: This is a  $\Gamma_1$  characteristic  $x + 2t = X + 2T$  along which  $u + w$  is constant. Thus,

$$u(A) + w(A) = u(P) + w(P), \quad x(A) = X + 2T. \quad (9.7a)$$

BP: This is a  $\Gamma_2$  characteristic  $x + t = X + T$  along which  $v + w$  is constant. Thus,

$$v(B) + w(B) = v(P) + w(P), \quad x(B) = X + T. \quad (9.7b)$$

EP: This is a  $\Gamma_3$  characteristic  $x - t = X - T$  along which  $u + v$  is constant. Thus,

$$u(E) + v(E) = u(P) + v(P), \quad t(E) = T - X. \quad (9.7c)$$

CE: This is a  $\Gamma_1$  characteristic  $x + 2t = 2(T - X)$  along which  $u + w$  is constant. The BC specifies  $w(0, t) = 0$ , and so,

$$u(E) = u(C) + w(C), \quad x(C) = 2(T - X). \quad (9.7d)$$

DE: This is a  $\Gamma_2$  characteristic  $x + t = (T - X)$  along which  $v + w$  is constant. The BC specifies  $w(0, t) = 0$ , and so,

$$v(E) = v(D) + w(D), \quad x(D) = (T - X). \quad (9.7e)$$

Hence, from (9.7c-e),

$$u(P) + v(P) = u(C) + w(C) + v(D) + w(D).$$

This, together with (9.7a-b) constitute three linearly independent equations to determine  $u(P)$ ,  $v(P)$  and  $w(P)$ . The coefficient matrix of this system is  $V^T$  whose rows are the eigenvectors of  $A^T$ .

## 9.9

With  $\mathbf{u} = [u, v]^T$  and  $\mathbf{f} = [f, g]^T$ , the equations become

$$\mathbf{u}_t + A\mathbf{u}_x = \mathbf{f}, \quad A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

The matrix  $A$  has eigenvalue  $\lambda = 1$  with eigenvector  $\mathbf{v}_1 = [1, 1]^T$  and eigenvalue  $\lambda = -1$  with eigenvector  $\mathbf{v}_2 = [1, -1]^T$ .

Multiplying  $\mathbf{u}_t + A\mathbf{u}_x = \mathbf{f}$  by  $\mathbf{v}_1^T = [1, 1]$  and supposing that  $f = 0$ ,  $g = G_t$ , we find that  $U^+ = \mathbf{v}_1^T \mathbf{u} = u + v$  satisfies  $U_t^+ + U_x^+ = G_t$  for which the characteristic equation is (see (9.2))

$$\frac{dt}{1} = \frac{dx}{1} = \frac{dU^+}{G_t}$$

so that

$$\frac{dx}{dt} = 1, \frac{dU^+}{dt} = G_t \Rightarrow x = t + k, \quad U^+ = G(x, t) + A(k), \Rightarrow u + v = G(x, t) + A(x - t),$$

where  $k$  is an arbitrary constant and  $A(k)$  an arbitrary function. Similarly, multiplying  $\mathbf{u}_t + A\mathbf{u}_x = \mathbf{f}$  by  $\mathbf{v}_2^T = [1, -1]$ , we find that  $U^- = \mathbf{v}_2^T \mathbf{u} = u - v$  satisfies  $U_t^- + U_x^- = -G_t$  for which the characteristic equation is

$$\frac{dt}{1} = \frac{dx}{-1} = \frac{dU^-}{-G_t}$$

so that

$$\frac{dx}{dt} = -1, \frac{dU^-}{dt} = G_t \Rightarrow x = -t + k, \quad U^- = G(x, t) + B(k), \Rightarrow u - v = -G(x, t) + B(x + t)$$



where  $B$  is an arbitrary function.

Combining these we have

$$u = \frac{1}{2}(A(x-t) + B(x+t)), \quad v = G(x, t) + \frac{1}{2}(A(x-t) - B(x+t)).$$

The initial condition  $u(x, 0) = 0$  gives  $B(x) = -A(x)$  and then  $v(x, 0) = 0$  gives  $A(x) = -G(x, 0)$ . Hence the solution is

$$u(x, t) = \frac{1}{2}(G(x+t, 0) - G(x-t, 0)), \quad v(x, t) = G(x, t) - \frac{1}{2}(G(x-t, 0) + G(x+t, 0)).$$

### 9.11

Suppose that  $P$  has coordinates  $(X, Y)$ .

- (a) For  $X > 2Y > 0$  (see Fig. 8) both characteristics through  $P_1$  intersect the boundary on the  $x$ -axis. Hence  $u(X, Y)$  is determined by the initial conditions—the solution to Exercise 9.10 fulfils these conditions.

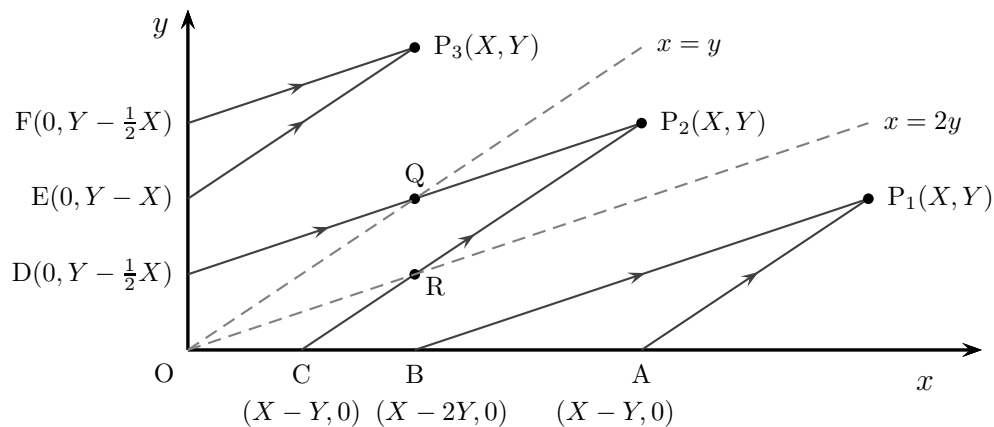


Figure 8: The characteristics through three general points  $P_j$  ( $j = 1, 2, 3$ ) for Exercise 9.11.

- (b) For  $0 < X < Y$  both characteristics through  $P_3$  intersect the boundary on the  $y$ -axis. Hence  $u(X, Y)$  is determined by the boundary conditions. We begin with the general solution  $u(x, y) = F(x-y) + G(x-2y)$  of the PDE derived in the solution to Exercise 9.10. Applying the BCs leads to

$$F(-y) + G(-2y) = f_0(y), \quad F'(-y) + G'(-2y) = f_1(y) \Rightarrow -F(-y) - \frac{1}{2}G(-2y) = \int^y f_1(s) ds$$

from which we obtain

$$\begin{aligned} G(-2y) = 2f_0(y) + 2 \int^y f_1(s) ds &\Rightarrow G(t) = 2f_0(-\frac{1}{2}t) + 2 \int^{-t/2} f_1(s) ds, \\ F(-y) = -f_0(y) - 2 \int^y f_1(s) ds &\Rightarrow F(t) = -2f_0(-t) - 2 \int^{-t} f_1(s) ds. \end{aligned}$$

The solution satisfying the BCs is, therefore,

$$u(x, y) = 2f_0(y - \frac{1}{2}x) - f_0(y - x) + 2 \int_{y-x}^{y-\frac{1}{2}x} f_1(s) ds.$$

(c) For  $2Y > X > Y > 0$  (see Fig. 8) one characteristic through  $P_2$  intersects the boundary on each axis. Hence  $u(X, Y)$  is determined by both the initial and boundary conditions.

Using  $u(x, y) = F(x - y) + G(x - 2y)$  we find that, along  $x = y$ ,  $u(y, y) = F(0) + G(-y)$  while, along  $x = 2y$ ,  $u(2y, y) = F(y) + G(0)$ . Hence

$$F(y) = u(2y, y) - G(0), \quad G(y) = u(-y, -y) - F(0).$$

and

$$u(X, Y) = u(2(X - Y), X - Y) + u(2Y - X, 2Y - X) - F(0) - G(0).$$

Setting  $X = Y = 0$  shows that  $F(0) + G(0) = u(0, 0)$ . The characteristics  $y = x + Y - X$  and  $2y = x + 2Y - X$  through  $P_3$  intersect the characteristics  $y = x$  and  $2y = x$  through the origin at  $Q$  and  $R$  whose coordinates are  $(2Y - X, 2Y - X)$  and  $(2(X - Y), X - Y)$ , respectively. Hence,  $u(P) = u(Q) + u(R) - u(0, 0)$ .

### 9.13

The operator associated with the PDE  $u_{tt} + u_{tx} - 2u_{xx} = t$  factorizes:

$$u_{tt} + u_{tx} - 2u_{xx} = (\partial_t - \partial_x)(\partial_t + 2\partial_x)u$$

Hence, with  $v = (\partial_t + 2\partial_x)u$  we have  $(\partial_t - \partial_x)v = t$  whose characteristic equations are

$$\frac{dt}{1} = \frac{dx}{-1} = \frac{dv}{t} \Rightarrow \frac{dx}{dt} = -1, \quad \frac{dv}{dt} = t$$

so the characteristics are  $x + t = k_1$  along which  $v = \frac{1}{2}t^2 + a(k_1)$ , where  $a(\cdot)$  is an arbitrary function. Then  $u$  is the solution of  $(\partial_t + 2\partial_x)u = v$  whose characteristic equations are

$$\frac{dt}{1} = \frac{dx}{2} = \frac{du}{v} \Rightarrow \frac{dx}{dt} = 2, \quad \frac{du}{dt} = v$$

so the characteristics are  $x - 2t = k_2$  along which  $u = \frac{1}{6}t^3 + \int^t a(x + s) ds + B(k_2)$ . This gives the general solution

$$u(x, t) = \frac{1}{6}t^3 + A(x + t) + B(x - 2t),$$

where  $A(x + t) = \int^t a(x + s) ds$ .

The initial conditions  $u(x, 0) = u_t(x, 0) = 0$  give

$$A(x) + B(x) = 0, \quad A'(x) - 2B'(x) = 0 \Rightarrow A(x) - 2B(x) = \text{constant} = C$$

and therefore  $A(x) = \frac{1}{3}C = -B$ . Hence  $u(x, t) = \frac{1}{6}t^3$ .

### 9.15

$$\begin{aligned} \mathcal{L}_1 \mathcal{L}_2 &= (a\partial_t + b\partial_x)(c\partial_t + d\partial_x) = a\partial_t(c\partial_t + d\partial_x) + b\partial_x(c\partial_t + d\partial_x) \\ &= a(c_t\partial_t + c\partial_t^2 + d_t\partial_x + d\partial_x\partial_t) + b(c_x\partial_t + c\partial_t\partial_x + d_x\partial_x + d\partial_x^2\partial_t) \\ &= \mathcal{L} + (ac_t + bc_x)\partial_t + (ad_t + bd_x)\partial_x \end{aligned}$$

and so  $\mathcal{L}_1 \mathcal{L}_2 = \mathcal{L}$  if  $ac_t + bc_x = 0$  and  $ad_t + bd_x = 0$ .

Thus,  $\mathcal{L}u = 0$  is equivalent to  $\mathcal{L}_1 v = 0$  where  $\mathcal{L}_2 u = v$ .

For the given PDE, comparing coefficients we find  $ac = 1$ ,  $ad + bc = t - 1$  and  $bd = -t$ . A solution of these is  $a = c = 1$ ,  $b = t$ ,  $d = -1$  so that  $\mathcal{L}_1 = \partial_t + t\partial_x$ ,  $\mathcal{L}_2 = \partial_t - \partial_x$ . The characteristic equations of  $\mathcal{L}_1 v = 0$  are

$$\frac{dt}{1} = \frac{dx}{t} = \frac{dv}{0} \Rightarrow \frac{dx}{dt} = t, \quad \frac{dv}{dt} = 0, \Rightarrow x = k_1 + \frac{1}{2}t^2, \quad v = A(k_1)$$

and so  $v = u_t - u_x = A(k_1)$  is constant along  $x - \frac{1}{2}t^2 = k_1$ . The characteristic through the point  $(x, t)$  cuts the  $x$ -axis at  $(x - \frac{1}{2}t^2, 0)$ , at which point  $u_t = g_1(x - \frac{1}{2}t^2)$  and  $u_x = g'_0(x - \frac{1}{2}t^2)$ . Hence,

$$u_t - u_x = A(x - \frac{1}{2}t^2), \quad \text{where } A(z) = g_1(z) - g'_0(z).$$

This PDE  $\mathcal{L}_2 u = v$  has characteristic equations

$$\begin{aligned} \frac{dt}{1} = \frac{dx}{-1} = \frac{du}{v} &\Rightarrow \frac{dx}{dt} = -1, \Rightarrow x = k_2 - t, \\ \frac{du}{dt} &= A(x - \frac{1}{2}t^2) = A(k_2 - t - \frac{1}{2}t^2) \Rightarrow u = \int_0^t A(x + t - s - \frac{1}{2}s^2) ds + g_0(x + t) \end{aligned}$$

using the initial condition  $u(x, 0) = g_0(x)$ .

### 9.17

The PDE  $u_t + uu_x = -2u$  has the characteristic equations

$$\frac{dt}{1} = \frac{dx}{u} = \frac{du}{-2u}$$

leading to

$$\frac{dx}{dt} = u, \quad \frac{du}{dt} = -2u.$$

The second of these has general solution  $u = Ce^{-2t}$  from which the first ODE gives  $x = k - \frac{1}{2}Ce^{-2t}$ , where  $k$  and  $C$  are related arbitrary constants—which we express by writing  $C = A(k)$ . Consider a characteristic emanating from a point on the boundary given by  $x = 0$  and  $t = t^*$ , say. We then have  $k = \frac{1}{2}A(k)e^{-2t^*}$  and the BC  $u(0, t^*) = e^{-t^*} = A(k)e^{-2t^*}$  so that  $A(k) = e^{t^*}$ ,  $k = \frac{1}{2}e^{-t^*}$  and therefore  $kA(k) = \frac{1}{2}$ . These results tie the arbitrary constants for a characteristic to its point of intersection on the  $t$ -axis. Furthermore,  $u = e^{t^* - 2t}$  and  $x = \frac{1}{2}(e^{-t^*} - e^{t^* - 2t})$ . Eliminating  $t^*$  shows that  $u$  is related to  $x$  and  $t$  via the quadratic equation

$$u^2 + 2xu - e^{-2t} = 0$$

whose roots are

$$u = -x \pm \sqrt{x^2 + e^{-2t}}.$$

Since  $u > 0$  for  $x = 0$  the positive root is appropriate and, by rationalising the result, we obtain

$$u = \frac{e^{-2t}}{x + \sqrt{x^2 + e^{-2t}}}.$$

### 9.19

Under the change of variables  $s = -x$ ,  $v(s, t) = -u(-s, t)$  the PDE becomes

$$-v_t + (-v)(-v)_{-s} = v \Rightarrow v_t + vv_s = -u$$

so remains invariant. The initial condition becomes  $v(s, 0) = -u(-s, 0) = -g(-s)$  for  $s \in \mathbb{R}$ . Thus  $v(s, 0) = g(s)$  when  $g$  is an odd function in which case  $v$  satisfies the same PDE and initial condition as  $u$  so  $v(s, t) = u(s, t) = -u(-s, t)$  so that  $u$  must be an odd function of  $x$ . With  $g(x) = x/(1 + |x|)$ , which is an odd function, in (9.27) we have, for  $k \geq 0$ ,

$$u = \frac{k}{1+k}e^{-t}, \quad x = k + \frac{k}{1+k}(e^{-t} - 1).$$

Defining  $U = ue^t$  and  $a(t) = e^{-t} - 1$  to simplify the notation, we find

$$U = \frac{k}{1+k}, \quad x = k + Ua(t) \quad \Rightarrow \quad k = \frac{U}{1-U}, \quad x = \frac{U}{1-U} + Ua(t)$$

which gives the relationship between  $U$ ,  $x$  and  $t$ . On rearranging we obtain the quadratic equation

$$a(t)U^2 - (1 + x + a(t))U + x = 0$$

whose roots are

$$U = \frac{1}{2a(t)}((1 + x + a(t)) \pm \sqrt{(1 + x + a(t))^2 - 4a(t)x}).$$

Since  $a(t) \rightarrow 0$  as  $t \rightarrow 0$ , we have to choose the root for which the numerator also vanishes at  $t = 0$ —this is the negative root. Thus, since  $u = Ue^{-t}$ ,

$$\begin{aligned} u(x, t) &= \frac{e^{-t}}{2a(t)}((1 + x + a(t)) - \sqrt{(1 + x + a(t))^2 - 4a(t)x}) \\ &= \frac{2xe^{-t}}{(1 + x + a(t)) + \sqrt{(1 + x + a(t))^2 - 4a(t)x}} \end{aligned}$$

which is valid for  $x \geq 0$ . The solution for  $x < 0$  follows since  $u$  must be an odd function of  $x$ .

## 9.21

The general solution of Burgers' equation takes the form (9.30):

$$u = g(k), \quad x = g(k)t + k.$$

The characteristics pass through the point  $x = k$  at  $t = 0$  and  $u(x, 0) = g(x)$  and  $u = \text{constant}$  along characteristics. When  $g(x) = \max\{0, 1 - |x|\}$  the initial condition is  $u = 0$  for  $|x| \geq 1$  and the characteristics are therefore  $x = \text{constant}$  if there are no shocks present.

- (a) Consider characteristics emanating from the points  $(x, t) = (k, 0)$  for  $0 \leq k \leq 1$ . Here  $g(x) = 1 - x$  and so

$$u = 1 - k, \quad x = (1 - k)t + k \quad \Rightarrow \quad u(x, t) = \frac{1 - x}{1 - t}$$

so this solution is valid only until  $t = 1$  when the entire family of these characteristics intersect at  $(1, 1)$  and a shock forms. Note that  $u(1, t) = 0$  for  $0 \leq t < 1$ .

- (b) Consider characteristics emanating from the points  $(x, t) = (k, 0)$  for  $-1 \leq k \leq 0$ . Here  $g(x) = 1 + x$  and so

$$u = 1 + k, \quad x = (1 + k)t + k \quad \Rightarrow \quad u(x, t) = \frac{1 + x}{1 + t}$$

a solution that is valid for  $t > 0$ . Note that  $u(-1, t) = 0$  for  $t \geq 0$ .

- (c) We have seen that the characteristics in case (a) intersect to form a shock at  $x = t = 1$ . However, the rightmost characteristic from case (b) (i.e.,  $k = 0^-$ ) also intersects the characteristic  $x = 1$  and the speed and location of the shock is determined by the interaction of characteristics of case (b) with those for  $x > 1$ , which are vertical lines since  $u = 0$  on each (see Fig. 9). The appropriate Rankine-Hugoniot condition (see Example 9.11) is

$$s'(t) = \frac{1}{2}(u^+ + u^-)$$

with  $u^+ = 0$  and  $u^- = 1 + k$ . Along the (b)-characteristics we have  $k = (x - t)/(1 + t)$  and therefore, when  $x = s(t)$ ,

$$s'(t) = \frac{1}{2} \frac{1 + s}{1 + t}, \quad t > 1, \quad s(1) = 1.$$

This has solution  $s(t) = \sqrt{2 + 2t} - 1$ . The solution for  $t > 1$  is therefore given by

$$u(x, t) = \begin{cases} 0, & x \leq -1 \\ \frac{1+x}{1+t}, & -1 < x < s(t) \\ 0, & x > s(t). \end{cases}$$

The solution at selected times is shown in Fig. 9.

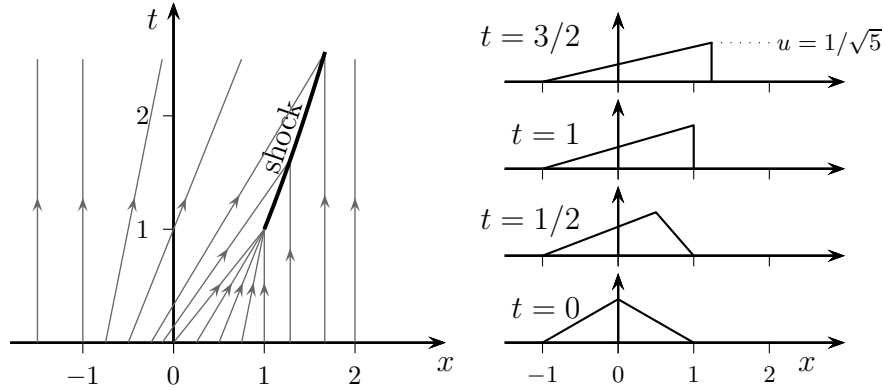


Figure 9: The characteristics (left) and solutions at selected times (right) for Exercise 9.21.

The triangular profile persists for all time. Prior to shock formation it has a fixed base of length 2 and a constant height  $u = 1$  so the area of the triangle is 1 for  $0 \leq t \leq 1$ . After shock formation the base of the triangle has length  $(1 + s(t))$  and height  $u(s(t), t) = (1 + s(t))/(1 + t)$  so the area is

$$\frac{1}{2}(1 + s(t)) \frac{1 + s(t)}{1 + t} = 1$$

since  $1 + s(t) = \sqrt{2 + 2t}$ .

### 9.23

$$g(x) = 1 - \frac{1}{1 + |x|} = \begin{cases} \frac{1}{1+x}, & x \geq 0 \\ \frac{2x-1}{x-1}, & x < 0. \end{cases}$$

Hence, for  $k = k^+ > 0$ , (9.41) leads to

$$s = k + \frac{t}{1+k} \quad \Rightarrow \quad k^2 - (s-1)k + t - s = 0$$

and, therefore,

$$k = \frac{1}{2}(s-1 \pm \sqrt{(s-1)^2 - 4(t-s)}).$$

Since  $k = k^+ > 0$  we have to choose the positive square root (this is readily checked at  $t = 0$ ). Similarly, for  $k = k^- < 0$ , (9.41) leads to

$$s = k + t \frac{2k-1}{k-1} \quad \Rightarrow \quad k^2 - (1+s-2t)k + s-t = 0$$

and, therefore,

$$k = \frac{1}{2}(1+s-2t \pm \sqrt{(1+s-2t)^2 - 4(s-t)}).$$

This time we have to choose the negative square root since  $k = k^- < 0$ . Clearly, when  $s(t) = t$  these give  $k^+ = k^- = 0 = k^*$  for  $t \geq 1$ .

### 9.25

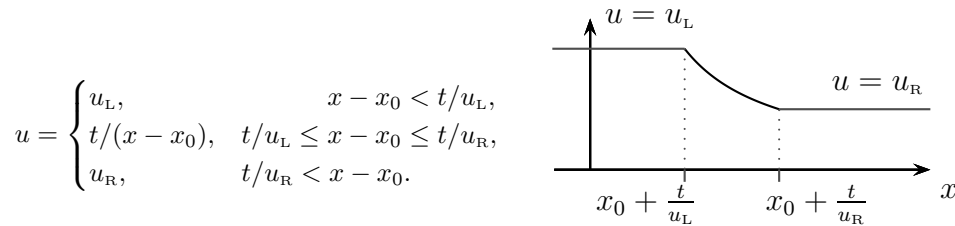
Differentiating the expression  $q'(u) = \frac{x-x_0}{t}$  with respect to  $x$  and  $t$  gives

$$\partial_x q'(u) = q''(u)u_x = \frac{1}{t}, \quad \partial_t q'(u) = q''(u)u_t = -\frac{x-x_0}{t^2}$$

so that  $q''(u)u_t = -q'(u)q''(u)u_x$ , that is,  $u_t + q'(u)u_x = 0$  provided that  $q''(u) \neq 0$ .

### 9.27

$q(u) = \log u$  so  $q'(u) = 1/u = (x-x_0)/t$  gives the expansion fan



### 9.29

The behaviour of the solution is broadly similar to that in the example—an expansion fan forms at  $x = 0$  and a shock at  $x = 1$  and they will collide at some point in time.

The expansion fan with  $u_L = 1$  and  $u_R = 2$  is given by

$$u = \begin{cases} 1, & x < t, \\ x/t, & t \leq x \leq 2t, \\ 2, & 2t < x. \end{cases}$$

A shock is already present at  $x = 1$ ,  $t = 0$  with  $u_L = 2$ ,  $u_R = 1$ . The shock speed is  $s'(t) = \frac{1}{2}(u_L + u_R) = \frac{3}{2}$  and its location is therefore  $x = s(t) = 1 + \frac{3}{2}t$  (since  $s(0) = 1$ ).

The rightmost characteristic from the expansion fan meets the shock wave when  $x = 2t = 1 + \frac{3}{2}t$ , i.e., when  $t = 2$  and  $x = 4$ .

The solution is illustrated in Fig. 10 for  $t = \frac{1}{3}$  (left) and  $t = 2$  (right).

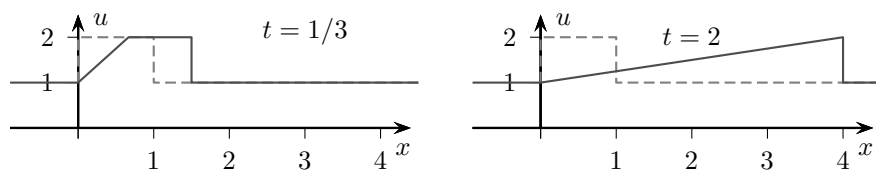


Figure 10: The solution of Exercise 9.29 when  $t = \frac{1}{3}$  (left) and  $t = 2$  (right). The dashed lines show the initial condition.

## Exercises 10 Finite difference methods for elliptic PDEs

### 10.1

We use the replacements  $u_{xx} = h_x^{-2}\delta_x^2 u - \frac{1}{12}h_x^2\partial_x^4 u + \mathcal{O}(h_x^4)$  and  $u_{yy} = h_y^{-2}\delta_y^2 u - \frac{1}{12}h_y^2\partial_y^4 u + \mathcal{O}(h_y^4)$  and so

$$\nabla^2 u = u_{xx} + u_{yy} = h_x^{-2}\delta_x^2 u + h_y^{-2}\delta_y^2 u - \frac{1}{12}(h_x^2\partial_x^4 u + h_y^2\partial_y^4 u) + \mathcal{O}(h_x^4) + \mathcal{O}(h_y^4).$$

Hence the local truncation error for Poisson's equation:  $-\nabla^2 u = f$  is

$$\begin{aligned}\mathcal{R}_{\ell,m} &= -(h_x^{-2}\delta_x^2 u + h_y^{-2}\delta_y^2 u) - f \\ &= -\nabla^2 u - \frac{1}{12}(h_x^2\partial_x^4 u + h_y^2\partial_y^4 u) + \mathcal{O}(h_x^4) + \mathcal{O}(h_y^4) - f \\ &= -\frac{1}{12}(h_x^2\partial_x^4 u + h_y^2\partial_y^4 u) + \mathcal{O}(h_x^4) + \mathcal{O}(h_y^4).\end{aligned}$$

Our finite difference replacement is

$$-(h_x^{-2}\delta_x^2 + h_y^{-2}\delta_y^2)U_{\ell,m} = f_{\ell,m}$$

and, because  $\delta_x^2 U_{\ell,m} = U_{\ell-1,m} - 2U_{\ell,m} + U_{\ell+1,m}$ ,  $\delta_y^2 U_{\ell,m} = U_{\ell,m-1} - 2U_{\ell,m} + U_{\ell,m+1}$  this becomes, on multiplying by  $h_x h_y$

$$\begin{aligned}\frac{h_y}{h_x}(U_{\ell-1,m} - 2U_{\ell,m} + U_{\ell+1,m}) + \frac{h_x}{h_y}(U_{\ell,m-1} - 2U_{\ell,m} + U_{\ell,m+1}) &= h_x h_y f_{\ell,m} \\ 2(\theta + \theta^{-1})U_{\ell,m} - \theta(U_{\ell-1,m} + U_{\ell+1,m}) - \theta^{-1}(U_{\ell,m-1} + U_{\ell,m+1}) &= h_x h_y f_{\ell,m}.\end{aligned}$$

If we use this scheme to solve Poisson's equation on the rectangle  $\{0 < x < L_x, 0 < y < L_y\}$  with  $h_x = L_x/M_x$  and  $h_y = L_y/M_y$  (where  $M_x$  and  $M_y$  are specified integers indicating the number of grid points in each direction) the totality of equations can be written in matrix vector form by defining the  $M_y \times M_y$  tridiagonal matrix

$$D = \begin{bmatrix} 2(\theta + \theta^{-1}) & -\theta^{-1} & & \\ -\theta^{-1} & 2(\theta + \theta^{-1}) & -\theta^{-1} & \\ & \ddots & \ddots & \ddots \\ & -\theta^{-1} & 2(\theta + \theta^{-1}) & -\theta^{-1} \\ & & -\theta^{-1} & 2(\theta + \theta^{-1}) \end{bmatrix}$$

and the  $M_x \times M_x$  block tridiagonal matrix

$$A = \begin{bmatrix} D & -\theta I & & \\ -\theta I & D & -\theta I & \\ & \ddots & \ddots & \ddots \\ & & -\theta I & D & -\theta I \\ & & & -\theta I & D \end{bmatrix}$$

then, in the natural column ordering of points,  $A\mathbf{u} = \mathbf{f}$ . This reduces to the standard 5-point formula (10.6) when  $\theta = 1$ .

### 10.3



We use the standard finite difference operators

$$\delta_x^2 U_{\ell,m} = U_{\ell+1,m} - 2U_{\ell,m} + U_{\ell-1,m}, \quad \delta_y^2 U_{\ell,m} = U_{\ell,m+1} - 2U_{\ell,m} + U_{\ell,m-1}$$

so

$$\begin{aligned} \delta_x^2 \delta_y^2 U_{\ell,m} &= \delta_x^2 (\delta_y^2 U)_{\ell,m} = (\delta_y^2 U)_{\ell+1,m} - 2(\delta_y^2 U)_{\ell,m} + (\delta_y^2 U)_{\ell-1,m} \\ &= (U_{\ell+1,m+1} - 2U_{\ell+1,m} + U_{\ell+1,m-1}) - 2(U_{\ell,m+1} - 2U_{\ell,m} + U_{\ell,m-1}) \\ &\quad + (U_{\ell-1,m+1} - 2U_{\ell-1,m} + U_{\ell-1,m-1}). \end{aligned}$$

The stencil of this result may be visualized as the “outer” product of a column and a row vector:

$$\begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} [1, -2, 1] = \begin{bmatrix} 1 & -2 & 1 \\ -2 & 4 & -2 \\ 1 & -2 & 1 \end{bmatrix}.$$

Then

$$2h^2 \mathcal{L}_h^+ - \delta_x^2 \delta_y^2 = \begin{bmatrix} -2 & -2 & -2 \\ -2 & 8 & -2 \\ -2 & -2 & -2 \end{bmatrix} - \begin{bmatrix} 1 & -2 & 1 \\ -2 & 4 & -2 \\ 1 & -2 & 1 \end{bmatrix} = \begin{bmatrix} -1 & -1 & -1 \\ -1 & 4 & -1 \\ -1 & -1 & -1 \end{bmatrix} = 2h^2 \mathcal{L}_h^\times$$

We know from equation (10.13) that

$$\mathcal{L}_h^+ u_{\ell,m} = -\nabla^2 u_{\ell,m} + \mathcal{R}_{\ell,m}^+,$$

where

$$\mathcal{R}_{\ell,m}^+ = -\frac{1}{12}h^2 (\partial_x^4 u + \partial_y^4 u) + \mathcal{O}(h^4)$$

so it then follows from  $\mathcal{L}_h^\times U_{\ell,m} = \mathcal{L}_h^+ U_{\ell,m} - \frac{1}{2}h^{-2}\delta_x^2 \delta_y^2 U_{\ell,m}$  that

$$\begin{aligned} \mathcal{L}_h^\times u_{\ell,m} &= \mathcal{L}_h^+ u_{\ell,m} - \frac{1}{2}h^{-2}\delta_x^2 \delta_y^2 u_{\ell,m} = \mathcal{L}_h^+ u_{\ell,m} - \frac{1}{2}h^2 (h^{-4}\delta_x^2 \delta_y^2 u_{\ell,m}) \\ &= -\nabla^2 u - \frac{1}{12}h^2 (\partial_x^4 u + \partial_y^4 u) + \mathcal{O}(h^4) - \frac{1}{2}h^2 (\partial_x^2 \partial_y^2 u_{\ell,m} + \mathcal{O}(h^2)) \\ &= -\nabla^2 u - \frac{1}{12}h^2 (\partial_x^4 u + 6\partial_x^2 \partial_y^2 u + \partial_y^4 u) + \mathcal{O}(h^4). \end{aligned}$$

## 10.5

In standard column-wise order, the equations  $\mathcal{L}_h^\times U = F$  are

$$\left[ \begin{array}{ccc|ccc|ccc} 4 & & & -1 & & & & & \\ & 4 & & & -1 & & & & \\ & & 4 & & & -1 & & & \\ \hline & -1 & & 4 & & & -1 & & \\ -1 & & -1 & & 4 & & & -1 & -1 \\ & & -1 & & & 4 & & -1 & \\ \hline & & & -1 & & & 4 & & \\ & & & -1 & & -1 & & 4 & \\ & & & & -1 & & & & 4 \end{array} \right] \mathbf{u} = \left[ \begin{array}{c} b_0 + b_2 + b_{14} \\ b_1 + b_3 \\ b_2 + b_4 + b_6 \\ \hline b_{13} + b_{15} \\ 0 \\ b_5 + b_7 \\ \hline b_{10} + b_{12} + b_{14} \\ b_9 + b_{11} \\ b_6 + b_8 + b_{10} \end{array} \right]$$

where  $b_j$  denotes the boundary value at the  $j$ th boundary node and these have been numbered clockwise from 0–15 starting at the origin. For general  $M$ , the coefficient matrix has the block tridiagonal structure:

$$A = \begin{bmatrix} 4I & T & & \\ T & 4I & T & \\ & \ddots & \ddots & \ddots \\ & & T & 4I & T \\ & & & T & 4I \end{bmatrix}, \text{ where } T = \begin{bmatrix} & -1 & & \\ -1 & & -1 & \\ & \ddots & \ddots & \ddots \\ & & -1 & -1 \\ & & & -1 \end{bmatrix}.$$

The diagonal blocks of  $A$  are a multiple of the identity  $I$  and the nonzero off-diagonal blocks  $T$  are themselves tridiagonal.

Switching to an odd-even numbering system where a node  $(x_\ell, y_m)$  is even if  $\ell + m$  is even. The numbering on a grid with  $h = 1/M$  and  $M = 4$  is shown below. The numbers outside the box refer to the boundary nodes where, unlike Exercise 10.2, the corner nodes are now involved. There are  $N_{\text{even}} = 5$  even nodes (numbered 1–5) and  $N_{\text{odd}} = 4$  nodes (numbered 6–9).

$$\mathbf{u} = [U_{1,1}, U_{1,3}, U_{2,2}, U_{3,1}, U_{3,3}, U_{1,2}, U_{2,1}, U_{2,3}, U_{3,2}]^T$$

4	5	6	7	8
3	<b>2</b>	8	<b>5</b>	9
2	6	<b>3</b>	9	10
1	<b>1</b>	7	<b>4</b>	11
0	15	14	13	12

In terms of this ordering, we define

$$\mathbf{u}_{\text{even}} = [U_1, U_2, U_3, U_4, U_5]^T, \quad \mathbf{u}_{\text{odd}} = [U_6, U_7, U_8, U_9]^T, \quad \mathbf{u} = \begin{bmatrix} \mathbf{u}_{\text{even}} \\ \mathbf{u}_{\text{odd}} \end{bmatrix}.$$

With this numbering system the equations  $\mathcal{L}_h^\times U = F$  become, in matrix-vector form,

$$\begin{bmatrix} 4 & & -1 & & \\ & 4 & -1 & & \\ -1 & -1 & 4 & -1 & -1 \\ & & -1 & 4 & \\ & & -1 & & 4 \end{bmatrix} \mathbf{u}_{\text{odd}} = \begin{bmatrix} b_0 + b_2 + b_{14} \\ b_2 + b_4 + b_6 \\ b_{10} + b_{12} + b_{14} \\ b_6 + b_8 + b_{10} \end{bmatrix}$$

$$\begin{bmatrix} 4 & -1 & -1 & \\ -1 & 4 & & -1 \\ -1 & & 4 & -1 \\ & -1 & -1 & 4 \end{bmatrix} \mathbf{u}_{\text{even}} = \begin{bmatrix} b_1 + b_3 \\ b_{13} + b_{15} \\ b_5 + b_7 \\ b_9 + b_{11} \end{bmatrix}.$$

The general pattern is difficult to detect with such a small value of  $M$  but, for even nodes (say), the structure is

$$\begin{bmatrix} 4I_1 & B & & \\ B^T & 4I_2 & B & \\ & B^T & 4I_1 & B \\ & & \ddots & \ddots & \ddots \end{bmatrix}, \quad B = \begin{bmatrix} -1 & & & \\ -1 & -1 & & \\ & -1 & -1 & \\ & & \ddots & \ddots \end{bmatrix},$$

where  $I_1$  and  $I_2$  are identity matrices whose dimensions are the number of even nodes in odd and even numbered columns, respectively. The matrix  $B$  is square when there are the same number of nodes in consecutive columns, otherwise it is rectangular. The structure is similar for odd nodes with the roles of the pairs  $I_1, B$  and  $I_2, B^T$  being interchanged.

One reason for studying different orderings of unknowns is the profound effect these can have on the time taken to solve the finite difference equations.

### 10.7

If  $Q(0) = Q(h) + Ch^2 + \mathcal{O}(h^3)$  for some constant  $C$ , then  $Q(0) = Q(h/2) + C(h/2)^2 + \mathcal{O}(h^3)$ . Subtracting these gives  $Ch^2 = 4(Q(h/2) - Q(h))/3 + \mathcal{O}(h^3)$  and so

$$Q(h) = Q(0) + \frac{4}{3}(Q(h) - Q(h/2)) + \mathcal{O}(h^3).$$

Thus  $4(Q(h) - Q(h/2))/3 \equiv 4(U_P^h - U_P^{h/2})/3$  gives an estimate of the leading term in the error in  $Q(h) \equiv U_P^h$ .

### 10.9

Suppose that the Neumann BC  $-u_x = g(0, y)$  is approximated at  $y = y_m$  by the second-order central difference formula

$$-\Delta_x U_{0,m} = g(0, mh) \quad \Rightarrow \quad -U_{1,m} + U_{-1,m} = 2hg(0, mh)$$

while (10.6) with  $\ell = 0$  gives

$$4U_{0,m} - U_{0,m+1} - U_{1,m} - U_{0,m-1} - U_{-1,m} = h^2 f_{0,m}$$

and adding these two equations gives

$$4U_{0,m} - U_{0,m+1} - 2U_{1,m} - U_{0,m-1} - U_{-1,m} = 2hg_{0,m} + h^2 f_{0,m}$$

which is (10.22b), as required.

### 10.11

The local truncation error associated with (10.24) is

$$\mathcal{R}_h|_{0,M} = -h^{-1}\Delta_x^+ u_{0,M} + h^{-1}\Delta_y^- u_{0,M} - \frac{1}{2}hf_{0,M} - g(1^+, 0) - g(0^-, 1).$$

We have, from the previous exercise,

$$h^{-1}\Delta_x^+ u_{0,M} = u_x(0, 1) + \frac{1}{2}hu_{xx}(0, 1) + \frac{1}{6}h^2u_{xxx}(\xi, 1)$$

and by making the change of variables  $x \mapsto 1-y$ ,  $y \mapsto 1-x$ ,  $h \mapsto -h$  (so that  $h^{-1}\Delta^+ \mapsto -h^{-1}\Delta^-$ ),

$$h^{-1}\Delta_y^- u_{0,M} = u_y(0, 1) - \frac{1}{2}hu_{yy}(0, 1) + \frac{1}{6}h^2u_{yyy}(0, 1-\eta), \quad 0 < \eta < h.$$

Hence,

$$\begin{aligned} \mathcal{R}_h|_{0,M} &= -\left(u_x(0, 1) + \frac{1}{2}hu_{xx}(0, 1) + \frac{1}{6}h^2u_{xxx}(\xi, 1)\right) \\ &\quad + \left(u_y(0, 1) - \frac{1}{2}hu_{yy}(0, 1) + \frac{1}{6}h^2u_{yyy}(0, 1-\eta)\right) \\ &\quad - \frac{1}{2}hf_{0,M} - g(1^+, 0) - g(0^-, 1) \\ &= (-u_x(0, 1) - g(1^+, 0)) + (u_y(0, 1) - g(0^-, 1)) - \frac{1}{2}h(-\nabla^2 u(0, 1) - f_{0,M}) \\ &\quad - \frac{1}{6}h^2(u_{xxx}(\xi, 1) + u_{yyy}(0, 1-\eta)) \end{aligned}$$

and therefore  $\mathcal{R}_h|_{0,M} = \mathcal{O}(h^2)$  if the third partial derivatives of  $u$  are bounded in the neighbourhood of the corner  $(0, 1)$ .

### 10.13

The factors are:

**Geometry** The domain  $0 < x, y < 1$  is unaffected by the interchange  $x \leftrightarrow y$ .

**Boundary values**  $u(x, y) = x^2 + y^2$  is unaffected by the interchange  $x \leftrightarrow y$ .

**PDE** Neither the source term  $f(x, y) = xy$  nor the differential operator  $-\nabla^2 \equiv \partial_x^2 + \partial_y^2$  is affected by the interchange  $x \leftrightarrow y$ .

Thus  $u(x, y) = u(y, x)$ .

Since the finite difference grid is also unaffected by the interchange  $x \leftrightarrow y$ , the numerical solution  $U$  must also share the same symmetry, that is  $U_{\ell, m} = U_{m, \ell}$ . It follows that  $U_{\ell, \ell+1} = U_{\ell+1, \ell}$  and  $U_{\ell-1, \ell} = U_{\ell, \ell-1}$  so the 5-point formula (10.6) applied at the grid point  $(\ell h, \ell h)$  gives

$$\begin{aligned}\mathcal{L}_h U_{\ell, \ell} &= h^{-2}(4U_{\ell, \ell} - U_{\ell+1, \ell} - U_{\ell, \ell+1} - U_{\ell-1, \ell} - U_{\ell, \ell-1}) \\ &= h^{-2}(4U_{\ell, \ell} - U_{\ell+1, \ell} - U_{\ell+1, \ell} - U_{\ell, \ell-1} - U_{\ell, \ell-1}) \\ &= h^{-2}(4U_{\ell, \ell} - 2U_{\ell+1, \ell} - 2U_{\ell, \ell-1})\end{aligned}$$

and  $f_{\ell, \ell} = x_{\ell} y_{\ell} = \ell^2 h^2$ .

When  $M = 4$  then symmetry reduces the number of unknowns to  $N = \frac{1}{2}M(M-1) = 6$  and these, together with the boundary nodes, are numbered as shown.

$$\begin{aligned}\mathbf{u} &= [U_{1,1}, U_{2,1}, U_{2,2}, U_{3,1}, U_{3,2}, U_{3,3}]^T \\ \mathbf{f} &= \frac{1}{16}[1, 2, 4, 2, 6, 9]^T \\ \mathbf{g} &= [2g_1, g_2, 0, g_3 + g_5, g_6, 2g_7]^T = \frac{1}{16}[2, 4, 0, 27, 20, 50]^T,\end{aligned}$$

4	5	6	7	8
3	4	5	<b>6</b>	7
2	2	<b>3</b>	<b>5</b>	6
1	<b>1</b>	<b>2</b>	<b>4</b>	5
0	1	2	3	4

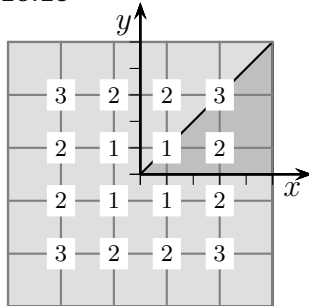
where  $g_j$  is the value of  $U$  at the  $j$ th boundary node.

The numerical solution is then found by solving  $A\mathbf{u} = \mathbf{f} + \mathbf{g}$ , where

$$A = 16 \begin{bmatrix} 4 & -2 & & & & \\ -1 & 4 & -1 & -1 & & \\ & -2 & 4 & & -2 & \\ & -1 & & 4 & -1 & \\ & & -1 & -1 & 4 & -1 \\ & & & -2 & 4 & \end{bmatrix} = 16D \begin{bmatrix} 2 & -1 & & & & \\ -1 & 4 & -1 & -1 & & \\ & -1 & 2 & & -1 & \\ & -1 & & 4 & -1 & \\ & & -1 & -1 & 4 & -1 \\ & & & & -1 & 2 \end{bmatrix},$$

where the diagonal matrix  $D$  is equal to the identity matrix except for entries corresponding to nodes on the diagonal of the grid, where  $D_{\ell, \ell} = 2$ . The purpose of introducing  $D$  being that  $A$  is the product of a diagonal matrix and a symmetric matrix. When  $M$  is large the cost of solving  $(D^{-1}A)\mathbf{u} = D^{-1}(\mathbf{f} + \mathbf{g})$  is approximately half the cost of solving the original system.

### 10.15

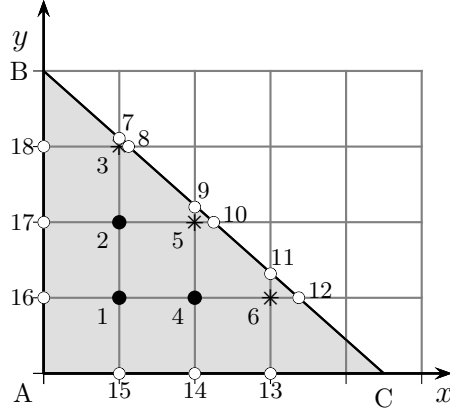


With  $M = 5$  ( $h = 2/5$ ) the implications of symmetry are shown in the figure and reveal that there are only 3 independent unknowns. Applying the 5-point approximation of the Poisson equation at nodes 1, 2 & 3 leads to the system:

$$\frac{25}{4} \begin{bmatrix} 2 & -2 & \\ -1 & 3 & -1 \\ & -2 & 4 \end{bmatrix} \begin{bmatrix} U_1 \\ U_2 \\ U_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

For general odd values of  $M$ , there are  $(M-1)^2$  unknowns (1 per grid point) of which only about one eighth are independent. We therefore have to solve for  $\frac{1}{8}(M-1)^2$  unknowns.

### 10.17



Node	$h_+$	$k_+$
3	1/32	1/36
5	2/32	2/36
7	3/32	3/36

The 6 internal grid points and the 12 active boundary nodes are numbered as shown. The standard 5-point finite difference equations

$$-h^{-2}(\delta_x^2 + \delta_y^2)U_{\ell,m} + U_{\ell,m} = 0.$$

can only be deployed at points 1, 2 and 4 where there is a uniform grid in both directions. These give:

$$\begin{aligned} P_1 : \quad & 65U_1 - 16U_2 - 16U_4 = 16(U_{15} + U_{16}) = 16[4x(9-8x)]_{x=1/4} = 112, \\ P_2 : \quad & -16U_1 + 65U_2 - 16U_3 - 16U_5 = 16U_{17} = 0, \\ P_4 : \quad & -16U_1 + 65U_4 - 16U_5 - 16U_6 = 16U_{14} = 16[4x(9-8x)]_{x=1/2} = 160. \end{aligned}$$

We focus first on the approximation of  $u_{xx}$  at the points 3, 5 and 6. At points 3, 5 and 6 we use equation (10.26) to give

$$\begin{aligned} P_3 : \quad & h_+ = \frac{1}{32}, \quad h_- = h = \frac{1}{4} : \quad (u_{xx})_3 \approx \frac{64}{9} \left( \frac{32}{1}(u_8 - u_3) - 4(u_3 - u_{18}) \right), \\ P_5 : \quad & h_+ = \frac{2}{32}, \quad h_- = h = \frac{1}{4} : \quad (u_{xx})_5 \approx \frac{64}{10} \left( \frac{32}{2}(u_{10} - u_5) - 4(u_5 - u_2) \right), \\ P_6 : \quad & h_+ = \frac{3}{32}, \quad h_- = h = \frac{1}{4} : \quad (u_{xx})_6 \approx \frac{64}{11} \left( \frac{32}{3}(u_{12} - u_6) - 4(u_6 - u_4) \right). \end{aligned}$$

The  $u_{yy}$  derivatives are treated similarly.

$$\begin{aligned} P_3 : \quad & k_+ = \frac{1}{36}, \quad k_- = k = \frac{1}{4} : \quad (u_{yy})_3 \approx \frac{72}{10} \left( \frac{36}{1}(u_7 - u_3) - 4(u_3 - u_2) \right), \\ P_5 : \quad & k_+ = \frac{2}{36}, \quad k_- = k = \frac{1}{4} : \quad (u_{yy})_5 \approx \frac{72}{11} \left( \frac{36}{2}(u_9 - u_5) - 4(u_5 - u_4) \right), \\ P_6 : \quad & k_+ = \frac{3}{36}, \quad k_- = k = \frac{1}{4} : \quad (u_{yy})_6 \approx \frac{72}{12} \left( \frac{36}{1}(u_{11} - u_6) - 4(u_6 - u_{13}) \right). \end{aligned}$$

Putting these together we find the following finite difference expressions for the approximations to  $-\nabla^2 u + u = 0$  at the 3 grid points closest to the hypotenuse.

$$\begin{aligned} P_3 : \quad & -\frac{64}{9} \left( \frac{32}{1}(U_8 - U_3) - 4(U_3 - U_{18}) \right) - \frac{72}{10} \left( \frac{36}{1}(U_7 - U_3) - 4(U_3 - U_2) \right) + U_3 = 0 \\ & \quad \quad \quad 545U_3 - \frac{144}{5}U_2 = 0 \\ P_5 : \quad & -\frac{64}{10} \left( \frac{32}{2}(U_{10} - U_5) - 4(U_5 - U_2) \right) - \frac{72}{11} \left( \frac{36}{2}(U_9 - U_5) - 4(U_5 - U_4) \right) + U_5 = 0 \\ & \quad \quad \quad 273U_5 - \frac{128}{5}U_2 - \frac{288}{11}U_4 = 0 \\ P_6 : \quad & -\frac{64}{11} \left( \frac{32}{3}(U_{12} - U_6) - 4(U_6 - U_4) \right) - \frac{72}{12} \left( \frac{36}{1}(U_{11} - U_6) - 4(U_6 - U_{13}) \right) + U_6 = 0 \\ & \quad \quad \quad \frac{979}{3}U_6 - \frac{256}{11}U_4 = 24U_{13}, \end{aligned}$$

where  $U_{13} = [4x(9 - 8x)]_{x=3/4} = 9$ . These provide 6 equations in the 6 unknowns.

### 10.19

In view of the identity

$$2(U_{\ell,1} - U_{\ell,0}) = (U_{\ell,1} - 2U_{\ell,0} + U_{\ell,-1}) + (U_{\ell,1} - U_{\ell,-1}) = \delta_\theta^2 U_{\ell,0} + 2\Delta_\theta U_{\ell,0}$$

equation (10.33b) may be written

$$\mathcal{B}_h U_{\ell,0} := -\Delta\theta^{-1}\Delta_\theta U_{\ell,0} - \frac{1}{2}r_\ell^2 \Delta\theta \mathcal{L}_h U_{\ell,0},$$

where  $\mathcal{L}_h$  (defined by (10.33)) is an approximation of  $\mathcal{L}u = (1/r)(ru_r)_r + u_{\theta\theta}/r^2$ . The local truncation error of the BC  $\mathcal{B}_h U_{\ell,0} = 0$  is

$$\begin{aligned} \mathcal{R}_h|_{\ell,0} &:= \mathcal{B}_h u_{\ell,0} = \Delta\theta^{-1}\Delta_\theta u_{\ell,0} - \frac{1}{2}r_\ell^2 \Delta\theta \mathcal{L}_h u_{\ell,0} \\ &= [-u_\theta + \mathcal{O}(\Delta\theta^2)]_{\ell,0} - \frac{1}{2}r_\ell^2 \Delta\theta [\mathcal{L}u + \mathcal{O}(h^2)]_{\ell,0} \\ &= [-u_\theta - \frac{1}{2}r_\ell^2 \Delta\theta \mathcal{L}u]_{\ell,0} + \mathcal{O}(\Delta\theta^2) + \mathcal{O}(h^2) = \mathcal{O}(\Delta\theta^2) + \mathcal{O}(h^2) \end{aligned}$$

as required (since  $u_\theta = 0$  and  $\mathcal{L}u = 0$  at  $(\ell h, 0)$ ).

### 10.21

With

$$u(r, \theta) = r^2 \left( \frac{1}{\pi} \left[ \log\left(\frac{1}{r}\right) \sin 2\theta + \left(\frac{\pi}{4} - \theta\right) r^2 \cos 2\theta \right] - \frac{\pi}{4} \right) r^2$$

we find that  $u(r, 0) = 0$ ,  $u(r, \frac{1}{2}\pi) = 0$  and  $u(1, \theta) = \frac{1}{4}(1 - 4\theta/\pi) \cos 2\theta - \frac{1}{4}$ .

Consider the functions  $f(r, \theta) = r^2 \log(\frac{1}{r}) \sin 2\theta$ ,  $g(r, \theta) = (\frac{\pi}{4} - \theta) r^2 \cos 2\theta$  and  $h(r, \theta) = \frac{1}{4} r^2$ .

$$\begin{aligned} f_r &= -r(2r \log(r) + 1) \sin(2\theta), & f_\theta &= 2r^2 \log\left(\frac{1}{r}\right) \cos 2\theta \\ f_{rr} &= -(2r \log(r) + 3) \sin(2\theta), & f_{\theta\theta} &= 4r^2 \log(r) \sin 2\theta \end{aligned}$$

so that  $\nabla^2 f := f_{rr} + \frac{1}{r}u_r + \frac{1}{r^2}f_{\theta\theta} = -4 \sin 2\theta$ . Similarly

$$\begin{aligned} g_r &= 2\left(\frac{\pi}{4} - \theta\right) r \cos 2\theta, & g_\theta &= -r^2 \cos(2\theta) - 2\left(\frac{1}{4}\pi - \theta\right) r^2 \sin 2\theta \\ g_{rr} &= 2\left(\frac{\pi}{4} - \theta\right) \cos 2\theta, & g_{\theta\theta} &= 4r^2 \sin 2\theta - 4\left(\frac{1}{4}\pi - \theta\right) r^2 \cos 2\theta \end{aligned}$$

and so  $\nabla^2 g = 4 \sin 2\theta$ . Hence  $\mathcal{L}(f + g) = 0$ :  $f + g$  is a harmonic function. Also  $\nabla^2 h = -1$ . Thus  $-\nabla^2 u = 1$  since  $u = f + g + h$ .

Also,

$$u_{rr} = -2 \log(r) \sin 2\theta + \text{terms involving } \theta \text{ only}$$

so that  $|u_{rr}| \rightarrow 0$  as  $r \rightarrow 0$  so long as  $\sin 2\theta \neq 0$ .

### 10.23

We rewrite (10.41) as

$$\begin{aligned} \mathcal{L}_h U_{\ell,m} &= -h_x^{-1} h_y^{-1} \left( \frac{1}{\rho} a(U_{\ell-1,m} - 2U_{\ell,m} + U_{\ell+1,m}) + \frac{1}{2} b(U_{\ell+1,m+1} - U_{\ell-1,m+1} + U_{\ell-1,m-1} - U_{\ell+1,m-1}) \right. \\ &\quad \left. + c\rho(U_{\ell,m-1} - 2U_{\ell,m} + U_{\ell,m+1}) \right. \\ &\quad \left. - \frac{1}{2} \gamma(4U_{\ell,m} - 2U_{\ell+1,m} + U_{\ell+1,m+1} - 2U_{\ell,m-1} + U_{\ell-1,m-1} - 2U_{\ell-1,m} + U_{\ell-1,m-1} - 2U_{\ell,m-1} + U_{\ell+1,m-1}) \right). \end{aligned}$$

This is of the form

$$\mathcal{L}_h u|_P = h_x^{-1} h_y^{-1} \left( \alpha_0 U|_P - \sum_{j=1}^8 \alpha_j U|_{Q_j} \right)$$

with (these are easiest to calculate by adjusting the coefficients in Fig. 10.15)

$$\begin{aligned} \alpha_0 &= 2a/\rho + 2c\rho - 2\gamma, & \alpha_1 &= a/\rho - \gamma, & \alpha_2 &= -\frac{1}{2}b + \frac{1}{2}\gamma \\ \alpha_3 &= c\rho - \gamma, & \alpha_4 &= \frac{1}{2}b + \frac{1}{2}\gamma, & \alpha_5 &= \alpha_1 \\ \alpha_6 &= -\frac{1}{2}b + \frac{1}{2}\gamma, & \alpha_7 &= \alpha_3, & \alpha_8 &= \frac{1}{2}b + \frac{1}{2}\gamma. \end{aligned}$$

These coefficients will be non-negative if  $\gamma$  is chosen such that  $|b| \leq \gamma$ ,  $\gamma \leq a/\rho$  and  $\gamma \leq c\rho$ . When  $\rho = \sqrt{a/c}$  these become

$$|b| \leq \gamma \leq \sqrt{ac}$$

which can always be achieved when  $b^2 \leq ac$ .

### 10.25

Using  $\Delta_x^- = \Delta - \frac{1}{2}\delta_x^2$ , then the analogue of the operator  $\mathcal{L}_h^-$  defined by (10.51) that is appropriate for the PDE  $-\varepsilon \nabla^2 u + pu_x = 0$  with  $p > 0$  may be written as

$$\begin{aligned} \mathcal{L}_h^- U_{\ell,m} &= -\varepsilon h^{-2} (\delta_x^2 + \delta_y^2) U_{\ell,m} + ph^{-1} \Delta_x^- U_{\ell,m} \\ &= -(\varepsilon h^{-2} + \frac{1}{2}ph^{-1}) \delta_x^2 U_{\ell,m} - \varepsilon h^{-2} \delta_y^2 U_{\ell,m} + ph^{-1} \Delta_x U_{\ell,m} \\ &= -\varepsilon h^{-2} (1 + \text{Pe}_h) \delta_x^2 U_{\ell,m} - \varepsilon h^{-2} \delta_y^2 U_{\ell,m} + ph^{-1} \Delta_x U_{\ell,m}, \end{aligned}$$

where  $\text{Pe}_h = ph/(2\varepsilon)$ .

When  $p < 0$ , the backward difference  $\Delta_x^-$  should be replaced by a forward difference  $\Delta_x^+$  and, with the identity  $\Delta_x^+ = \Delta + \frac{1}{2}\delta_x^2$ , the corresponding manipulations are

$$\begin{aligned} \mathcal{L}_h^+ U_{\ell,m} &= -\varepsilon h^{-2} (\delta_x^2 + \delta_y^2) U_{\ell,m} + ph^{-1} \Delta_x^+ U_{\ell,m} \\ &= -(\varepsilon h^{-2} - \frac{1}{2}ph^{-1}) \delta_x^2 U_{\ell,m} - \varepsilon h^{-2} \delta_y^2 U_{\ell,m} + ph^{-1} \Delta_x U_{\ell,m} \\ &= -\varepsilon h^{-2} (1 - \text{Pe}_h) \delta_x^2 U_{\ell,m} - \varepsilon h^{-2} \delta_y^2 U_{\ell,m} + ph^{-1} \Delta_x U_{\ell,m}. \end{aligned}$$

Here  $\text{Pe}_h < 0$  so that  $\text{Pe}_h = -|\text{Pe}_h|$  and therefore both cases are accommodated in the single formula

$$\mathcal{L}_h^\pm U_{\ell,m} := -\varepsilon h^{-2} (1 + |\text{Pe}_h|) \delta_x^2 U_{\ell,m} - \varepsilon h^{-2} \delta_y^2 U_{\ell,m} + h^{-1} p \Delta_x U_{\ell,m}.$$

The generalization to the PDE  $-\varepsilon(u_{xx} + u_{yy}) + pu_x + qu_y = 0$ , is

$$-\varepsilon h^{-2} (1 + |\text{Pe}_{h,x}|) \delta_x^2 U_{\ell,m} - \varepsilon h^{-2} (1 + |\text{Pe}_{h,y}|) \delta_y^2 U_{\ell,m} + h^{-1} p \Delta_x U_{\ell,m} + h^{-1} q \Delta_y U_{\ell,m} = 0,$$

where  $\text{Pe}_{h,x} := ph/(2\varepsilon)$  and  $\text{Pe}_{h,y} := qh/(2\varepsilon)$ .

## Exercises 11 Finite difference methods for parabolic PDEs

### 11.1

Differentiating  $u_t = u_{xx}$  with respect to  $t$  we find  $u_{tt} = u_{xxt}$  while differentiating twice with respect to  $x$  leads to  $u_{xxxx} = u_{xxt}$ . Hence (11.12) becomes

$$\begin{aligned}\mathcal{R}_m^n &= \frac{1}{2}k u_{tt}|_m^n - \frac{1}{12}h^2 u_{xxxx}|_m^n + \mathcal{O}(k^2) + \mathcal{O}(h^4) \\ &= \frac{1}{2}h^2(r - \frac{1}{6})u_{xxt}|_m^n + \mathcal{O}(k^2) + \mathcal{O}(h^4).\end{aligned}$$

When  $r = 1/6$  the order of consistency is  $\mathcal{O}(k^2) + \mathcal{O}(h^4)$ .

### 11.3

Equations (11.8b) give, for  $m = 1, 2, \dots, M-1$ ,

$$\begin{aligned}U_1^{n+1} &= r g_0(nk) + (1-2r)U_1^n + rU_2^n \\ U_2^{n+1} &= rU_1^n + (1-2r)U_2^n + rU_3^n \\ &\vdots \\ U_{M-2}^{n+1} &= rU_{M-3}^n + (1-2r)U_{M-2}^n + rU_{M-1}^n \\ U_{M-1}^{n+1} &= rU_{M-2}^n + (1-2r)U_{M-1}^n + r g_1(nk)\end{aligned}$$

which, in matrix-vector form, become

$$\mathbf{u}^{n+1} = \begin{bmatrix} 1-2r & r & 0 & \cdots & 0 \\ r & 1-2r & r & & \\ 0 & r & \ddots & \ddots & \\ & & \ddots & r & \\ 0 & & & r & 1-2r \end{bmatrix} \mathbf{u}^n + r \begin{bmatrix} g_0(nk) \\ 0 \\ \vdots \\ 0 \\ g_1(nk) \end{bmatrix}$$

where  $\mathbf{u}^n = [U_1^n, U_2^n, \dots, U_{M-1}^n]^\top$ . The coefficient matrix is clearly the same as that in the BTCS scheme (11.25) with  $r$  replaced by  $-r$ .

### 11.5

The finite difference scheme  $U_m^{n+1} = \frac{1}{3}rU_{m-2}^n + (1-r)U_m^n + \frac{2}{3}rU_{m+1}^n$  can be written in the form (11.10) with

$$\mathcal{L}_h U_m^n = -\frac{1}{3}h^{-2}(U_{m-2}^n - 3U_m^n + 2U_{m+1}^n).$$

Using the Taylor expansions of  $u_{m-2}^n$  and  $u_{m+1}^n$  about the point  $(x_m, t^n)$ , we find that (all the terms on the right hand side are evaluated at  $(x_m, t^n)$ )

$$\begin{aligned}\mathcal{L}_h u_m^n &= -\frac{1}{3}h^{-2} \left( u - 2hu_x + \frac{1}{2!}(2h)^2 u_{xx} - \frac{1}{3!}(2h)^3 u_{xxx} + \mathcal{O}(h^4) \right. \\ &\quad \left. - 3u \right. \\ &\quad \left. + 2u + 2hu_x + 2\frac{1}{2!}h^2 u_{xx} + \frac{1}{3!}h^3 u_{xxx} + \mathcal{O}(h^4) \right) \\ &= -u_{xx} + \mathcal{O}(h).\end{aligned}$$

Since  $\mathcal{L}_h u = \mathcal{L}u + \mathcal{O}(h)$ ,  $\mathcal{L}_h$  is first order consistent with  $\mathcal{L}$  and, from (11.11), the given method is therefore consistent with the heat equation. The stencil and domain of dependence are shown in Fig. 11.



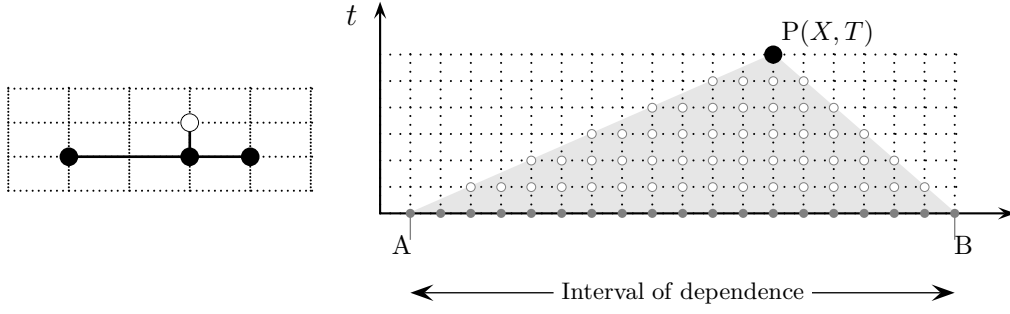


Figure 11: Stencil (left) and domain of dependence (right) for Exercise 11.5.

### 11.7

Replacing  $U$  by  $-U$  in Theorem 11.5 we have

$$\begin{aligned} -h^{-2}\delta_x^2(-U_m^n) + k^{-1}((-U_m^{n+1}) - (-U_m^n)) &\leq 0 \\ \text{i.e., } -h^{-2}\delta_x^2 U_m^n + k^{-1}(U_m^{n+1} - U_m^n) &\geq 0 \end{aligned}$$

for  $(x_m, t_{n+1}) \in \Omega_\tau$ . Thus, if  $r = k/h^2 \leq 1/2$  then  $U$  is either constant or else attains its *minimum* ( $-U$  attains its maximum) value on  $\Gamma_\tau$ .

If  $-h^{-2}\delta_x^2 U_m^n + k^{-1}(U_m^{n+1} - U_m^n) = 0$  then  $U$  is either constant or else attains its maximum and minimum value on  $\Gamma_\tau$ . The case when  $U$  is constant can only occur if  $U$  is the same constant on  $\Gamma_\tau$ . So, when  $U_0^n = U_M^n = 0$ , the maximum absolute value must occur at  $t = 0$ .

### 11.9

The local truncation error of the BTCS approximation (11.22) is, by definition,

$$\mathcal{R}_h|_m^n := \frac{1}{k}(u_m^{n+1} - u_m^n - r\delta_x^2 u_m^{n+1}) = \frac{1}{k}(u_m^{n+1} - u_m^n) - h^{-2}\delta_x^2 u_m^{n+1}$$

and, using  $u_m^n = u_m^{n+1} - ku_t|_m^{n+1} + \frac{1}{2}k^2 u_{tt}|_m^{n+1} + \mathcal{O}(k^3)$  together with  $\delta_x^2 u_m^{n+1} = u_{xx}|_m^{n+1} + \frac{1}{12}h^2 u_{xxxx}|_m^{n+1} + \mathcal{O}(h^4)$ , we find

$$\begin{aligned} \mathcal{R}_h|_m^n &= (u_t|_m^{n+1} - \frac{1}{2}ku_{tt}|_m^{n+1} + \mathcal{O}(k^2)) - (u_{xx}|_m^{n+1} + \frac{1}{12}h^2 u_{xxxx}|_m^{n+1} + \mathcal{O}(h^4)) \\ &= (u_t - u_{xx})|_m^{n+1} - \frac{1}{2}h^2(ru_{tt} + \frac{1}{6}u_{xxxx})|_m^{n+1} + \mathcal{O}(k^2) + \mathcal{O}(h^4) \end{aligned}$$

and so  $\mathcal{R}_h = \mathcal{O}(k) + \mathcal{O}(h^2)$ . Since  $u_{tt} = u_{txx} = \partial_x^2 u_t = u_{xxxx}$ , this becomes

$$\mathcal{R}_h|_m^n = -\frac{1}{2}h^2(r + \frac{1}{6})u_{txx}|_m^{n+1} + \mathcal{O}(k^2) + \mathcal{O}(h^4).$$

### 11.11

Replacing  $U$  by  $-U$  in Theorem 11.11 we deduce that

$$\begin{aligned} -h^{-2}\delta_x^2(-U_m^{n+1}) + k^{-1}((-U_m^{n+1}) - (-U_m^n)) &\leq 0 \\ \text{i.e., } -h^{-2}\delta_x^2 U_m^{n+1} + k^{-1}(U_m^{n+1} - U_m^n) &\geq 0 \end{aligned}$$

and so  $U$  is either constant or else attains its *minimum* value on  $\Gamma_\tau$  when  $r = k/h^2 \leq 1/2$ . Thus, if  $-h^{-2}\delta_x^2 U_m^{n+1} + k^{-1}(U_m^{n+1} - U_m^n) = 0$  then  $U$  is either constant or else attains its maximum and

minimum value on  $\Gamma_\tau$ . The case when  $U$  is constant can only occur if  $U$  is the same constant on  $\Gamma_\tau$ . So, when  $U_0^n = U_M^n = 0$ , the maximum absolute value must occur at  $t = 0$ .

### 11.13

With  $\mathcal{L}_h U = -h^{-2}\delta_x^2 U$ , the defining equation (11.30) of the  $\theta$ -method becomes

$$U_m^{n+1} - \theta r \delta_x^2 U_m^{n+1} = U_m^n + (1 - \theta) r \delta_x^2 U_m^n$$

i.e.,

$$-r\theta U_{m-1}^{n+1} + (1 + 2r\theta)U_m^{n+1} - r\theta U_{m+1}^{n+1} = r(1 - \theta)U_{m-1}^n + (1 - 2(1 - \theta)r)U_m^n + r(1 - \theta)U_{m+1}^n$$

so, with the given BCs, these give, for  $m = 1$  and  $m = M - 1$ ,

$$\begin{aligned} (1 + 2r\theta)U_1^{n+1} - r\theta U_2^{n+1} &= (1 - (1 - 2\theta)r)U_1^n + r(1 - \theta)U_2^n + r((1 - \theta)g_0(t_n) + \theta g_0(t_{n+1})) \\ -r\theta U_{M-2}^{n+1} + (1 + 2r\theta)U_{M-1}^{n+1} &= r(1 - \theta)U_{M-2}^n + (1 - (1 - 2\theta)r)U_{M-1}^n + r((1 - \theta)g_1(t_n) + \theta g_1(t_{n+1})). \end{aligned}$$

Thus, for  $(M - 1) \times (M - 1)$  matrices  $B$  and  $C$ :

$$B\mathbf{u}^{n+1} = C\mathbf{u}^n + \mathbf{f}^n$$

where

$$\begin{aligned} B &= \begin{bmatrix} 1 + 2r\theta & -r\theta & 0 & \cdots & 0 \\ -r\theta & 1 + 2r\theta & -r\theta & & \\ 0 & -r\theta & \ddots & \ddots & \\ & & \ddots & & -r\theta \\ 0 & & & -r\theta & 1 + 2r\theta \end{bmatrix} = A(\theta r) \\ C &= \begin{bmatrix} 1 - 2r(1 - \theta) & r(1 - \theta) & 0 & \cdots & 0 \\ r(1 - \theta) & 1 - 2r(1 - \theta) & r(1 - \theta) & & \\ 0 & r(1 - \theta) & \ddots & \ddots & \\ & & \ddots & & r\theta \\ 0 & & & r(1 - \theta) & 1 - 2r(1 - \theta) \end{bmatrix} = A((\theta - 1)r), \end{aligned}$$

where  $A \equiv A(r)$  is the coefficient matrix appearing in (11.25). Also,

$$\mathbf{f}^n = r[(1 - \theta)g_0(t_n) + \theta g_0(t_{n+1}), 0, \dots, 0, (1 - \theta)g_1(t_n) + \theta g_1(t_{n+1})]^\top.$$

### 11.15

Replacing  $U$  by  $-U$  in Theorem 11.16 we deduce that

$$\begin{aligned} -\frac{1}{2}h^{-2}\delta_x^2((-U_m^{n+1}) + (-U_m^n)) + k^{-1}((-U_m^{n+1}) - (-U_m^n)) &\leq 0 \\ \text{i.e., } -\frac{1}{2}h^{-2}\delta_x^2(U_m^{n+1} + U_m^n) + k^{-1}(U_m^{n+1} - U_m^n) &\geq 0 \end{aligned}$$

for  $(x_m, t_{n+1}) \in \Omega_\tau$ . If  $r = k/h^2 \leq 1$  then, from Theorem 11.16,  $-U$  is either constant or else attains its maximum value—therefore  $U$  attains its *minimum* value, on  $\Gamma_\tau$ .

Thus, if  $-\frac{1}{2}h^{-2}\delta_x^2(U_m^{n+1} + U_m^n) + k^{-1}(U_m^{n+1} - U_m^n) = 0$  then  $U$  is either constant or else attains both its maximum and minimum value on  $\Gamma_\tau$ . The case when  $U$  is constant can only occur if  $U$

is the same constant on  $I_\tau$ . So, when  $U_0^n = U_M^n = 0$ , the maximum absolute value must occur at  $t = 0$ .

### 11.17

We use two properties of a sum of non-negative terms: it is always smaller than the (largest term)  $\times$  (the number of terms) and it is always larger than the largest single term. The norm

$$\|U\|_{h,2} = \left( h \sum_{m=0}^M |U_m|^2 \right)^{1/2} \leq \left( h \sum_{m=0}^M 1 \right)^{1/2} \max_{0 \leq m \leq M} |U_m| = \|U\|_{h,\infty}.$$

If the largest term among  $\{U_m^n\}_{m=0}^M$  occurs at  $m = 0$  or  $m = M$ , then  $\|U\|_{h,\infty} = |U_0|$  or  $|U_M|$  and  $\|U\|_{h,2} \geq \frac{1}{2}\sqrt{h}\|U\|_{h,\infty}$ .

If  $\|U\|_{h,\infty} = |U_m|$  with  $0 < m < M$ , then  $\|U\|_{h,2} \geq \sqrt{h}\|U\|_{h,\infty}$ . Thus, in general,

$$\frac{1}{2}\sqrt{h}\|U\|_{h,\infty} \leq \|U\|_{h,2} \leq \|U\|_{h,2}.$$

For the two possible grid functions a simple calculation reveals

Case	$\ U\ _{h,\infty}$	$\ U\ _{h,2}$	
(a)	1	1	Right hand bound attained
(b)	1	$\frac{1}{2}\sqrt{h}$	Left hand bound attained

These examples show that the bounds are attained so they cannot be improved upon.

### 11.19

From the solution to Exercise 11.13, the  $\theta$ -method applied to the heat equation gives

$$U_m^{n+1} - \theta r \delta_x^2 U_m^{n+1} = U_m^n + (1 - \theta) r \delta_x^2 U_m^n.$$

We now substitute  $U_m^n = \xi^n e^{i\kappa m h}$  and use the relationships  $U_m^{n+1} = \xi U_m^n$  and (11.47b) to give (we define  $s = \sin^2(\frac{1}{2}\kappa h)$  as a convenient abbreviation)

$$[1 + 4r\theta s]\xi U_m^n = [1 - 4r(1 - \theta)s]U_m^n \Rightarrow \xi = \frac{1 - 4r(1 - \theta)s}{1 + 4r\theta s}.$$

We can rewrite this as

$$\xi = 1 - \frac{4rs}{1 + 4r\theta s}$$

so  $\xi \leq 1$  for all  $r > 0$  and  $\theta \geq 0$ . In order that  $\xi \geq -1$  we require

$$1 - 4r(1 - \theta)s \geq -1 - 4r\theta s \Rightarrow 2r(1 - 2\theta)s \leq 1.$$

This holds for all  $r > 0$  when  $\theta \geq \frac{1}{2}$  otherwise, it will hold for all  $s \in [0, 1]$  only when  $r$  is restricted by

$$r \leq \frac{1}{2(1 - 2\theta)}.$$

This reduces, as it should, to  $r \leq \frac{1}{2}$  when  $\theta = 0$  and the method becomes the FTCS scheme.

### 11.21

The FTCS scheme (11.53)

$$\begin{aligned} U_m^{n+1} &= U_m^n + r \delta_x^2 U_m^n + c \Delta_x U_m^n \\ &= (r - \frac{1}{2}\rho)U_{m-1}^n + (1 - 2r)U_m^n + (r + \frac{1}{2}\rho)U_{m+1}^n \end{aligned}$$

is of the form

$$U_m^{n+1} = \alpha_{-1}U_{m-1}^n + \alpha_0U_m^n + \alpha_1U_{m+1}^n$$

with  $\alpha_{-1} = r - \frac{1}{2}\rho$ ,  $\alpha_0 = 1 - 2r$ ,  $\alpha_1 = r + \frac{1}{2}\rho$ . Non-negativity of the coefficients requires  $\frac{1}{2}\rho \leq r \leq \frac{1}{2}$ . The left inequality simplifies to  $h \leq 2\varepsilon$  (so that stability is only possible if the spatial grid size is sufficiently small compared to  $\varepsilon$ ), that is,  $\text{Pe}_h \leq 1$ , where  $\text{Pe}_h := h/(2\varepsilon)$  is the mesh Peclet number. The right inequality leads to  $k \leq h^2/2\varepsilon$ —the stability region is shown in Fig. 12 (Left).

The conditions  $\frac{1}{2}\rho^2 \leq r \leq \frac{1}{2}$  for  $\ell_2$ -stability require  $k \leq 2\varepsilon$  and  $k \leq h^2/2\varepsilon$ , both coinciding when  $\text{Pe}_h = 1$ . The corresponding stability region is shown in Fig. 12 (Right). Thus the scheme is  $\ell_2$ -stable for any spatial grid size.

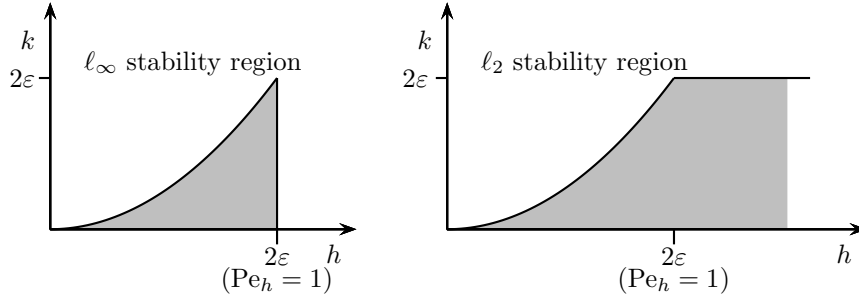


Figure 12: The stability regions for Exercise 11.21.

The proof of Theorem 11.5 may be used to establish a maximum principle under the conditions  $\frac{1}{2}\rho \leq r \leq \frac{1}{2}$ . The only change necessary is the derivation of an upper bound on  $U_m^{j+1}$ : This now reads

$$U_m^{j+1} \leq \alpha_{-1}U_{m-1}^j + \alpha_0U_m^j + \alpha_1U_{m+1}^j$$

and, because  $U_{m-1}^j, U_m^j, U_{m+1}^j \leq K_\tau$  and the non-negativity of the coefficients

$$U_m^{j+1} \leq (\alpha_{-1} + \alpha_0 + \alpha_1)K_\tau = K_\tau$$

since  $\alpha_{-1} + \alpha_0 + \alpha_1 = 1$ .

### 11.23

We substitute  $U_m^n = \xi^n e^{i\kappa m h}$  into the finite difference scheme

$$U_m^{n+1} = \frac{1}{3}rU_{m-2}^n + (1-r)U_m^n + \frac{2}{3}rU_{m+1}^n$$

and use the relationships  $U_m^{n+1} = \xi U_m^n$ ,  $U_{m-2}^n = e^{-2i\kappa m h}U_m^n$  and  $U_{m+1}^n = e^{i\kappa m h}U_m^n$  to give

$$\xi U_m^n = \frac{1}{3}r e^{-2i\kappa h} U_m^n + (1-r)U_m^n + \frac{2}{3}r e^{i\kappa h} U_m^n \Rightarrow \xi = 1 - r + \frac{1}{3}r(2e^{i\kappa h} + e^{-2i\kappa h}).$$

Extracting real and imaginary parts of the right hand side we have

$$\xi = (1 - r - r \cos \kappa h) + \frac{1}{3}ir \sin \kappa h = (1 - 2r \sin^2 \frac{1}{2}\kappa h) + \frac{1}{3}ir \sin \kappa h$$

so that, after expressing  $\sin \kappa h$  in terms of half-angles,

$$|\xi|^2 - 1 = -4rs[1 - r + \frac{1}{9}r(1 - s)], \quad s = \sin^2 \frac{1}{2}\kappa h.$$

Hence  $|\xi|^2 - 1 \leq 0$  if the bracketed term on the right hand side is positive for  $s \in [0, 1]$ . Since it is a linear expression in  $s$ , only its end-points need to be examined. It is non-negative at  $s = 0$  if  $r \leq 9/8$  and at  $s = 1$  if  $r \leq 1$ . The method is therefore  $\ell_2$  stable for  $r \leq 1$ .

### 11.25

With the Robin end condition  $-u_x(0, t) + \sigma u(0, t) = g_0(t)$ , the numerical boundary condition (11.57) is replaced by

$$-\frac{1}{2}h^{-1}(-U_{-1}^n + U_1^n) + \sigma U_0^n = g_0(nk).$$

Using (11.59) to eliminate  $U_{-1}^n$  leads to

$$\frac{1}{2hr}(U_0^{n+1} - (1 - 2r - 2hr\sigma)U_0^n - 2rU_1^n) = g_0(nk).$$

or, for computational purposes,

$$U_0^{n+1} = (1 - 2r - 2hr\sigma)U_0^n + 2rU_1^n + 2rhg_0(nk).$$

However, it is the former version that is correctly scaled for determining the local truncation error (all the terms involving  $u$  on the right hand side in the following expansions are evaluated at  $(0, nk)$ ):

$$\begin{aligned} \mathcal{R}_h|_0^n &= \frac{1}{2hr}(u_0^{n+1} - (1 - 2r - 2hr\sigma)u_0^n - 2ru_1^n) - g_0(nk) \\ &= \frac{1}{2hr}(u + ku_t + \frac{1}{2}k^2u_{tt} + \mathcal{O}(k^3) \\ &\quad - (1 - 2r - 2hr\sigma)u \\ &\quad - 2r(u + hu_x + \frac{1}{2}h^2u_{xx} + \frac{1}{6}h^3u_{xxx} + \mathcal{O}(h^4))) - g_0(nk) \\ &= \frac{1}{2hr}((1 - (1 - 2r - 2hr\sigma) - 2r)u + ku_t - 2hru_x + \frac{1}{2}k^2u_{tt} - ku_{xx} + \dots) - g_0(nk) \\ &= (\sigma u - u_x - g_0(nk)) + \frac{1}{2}h(u_t - u_{xx}) + \frac{1}{6}h(3ku_{tt} - 2hu_{xxx}) + \dots \end{aligned}$$

Hence  $\mathcal{R}_h = \mathcal{O}(hk) + \mathcal{O}(h^2)$  is consistent of second order in  $h$  since  $k = \mathcal{O}(h)$ ,  $u(0, t) - u_x(0, t) = g_0(t)$  and  $u_t = u_{xx}$ .

### 11.27

The Crank–Nicolson scheme (11.34b) for the heat equation at  $m = 0$  is

$$-\frac{1}{2}rU_{-1}^{n+1} + (1 + r)U_0^{n+1} - \frac{1}{2}rU_1^{n+1} = \frac{1}{2}rU_{-1}^n + (1 - r)U_0^n + \frac{1}{2}rU_1^n.$$

and using (11.57) at  $t = nk$  and also with  $n$  replaced by  $n + 1$  gives

$$\begin{aligned} -\frac{1}{2}h^{-1}(-U_{-1}^n + U_1^n) &= g_0(nk) &\Rightarrow U_{-1}^n &= U_1^n + 2hg_0(nk) \\ -\frac{1}{2}h^{-1}(-U_{-1}^{n+1} + U_1^{n+1}) &= g_0((n + 1)k) &\Rightarrow U_{-1}^{n+1} &= U_1^{n+1} + 2hg_0((n + 1)k). \end{aligned}$$

On substituting for  $U_{-1}^n$  and  $U_{-1}^{n+1}$  the Crank–Nicolson scheme then becomes

$$(1 + r)U_0^{n+1} - rU_1^{n+1} = (1 - r)U_0^n + rU_1^n + hr(g_0(nk) + g_0((n + 1)k)).$$

### 11.29

We prove the result by induction on  $n$ . Let  $K_0 = \|U^0\|_{h,\infty}$ , then the induction hypothesis  $\|U^n\|_{h,\infty} \leq \|U^0\|_{h,\infty}$  is satisfied at  $n = 0$ . Let us suppose that it holds up until  $n = j$  so that  $\|U^j\|_{h,\infty} \leq K_0$ . We need to prove that  $\|U^{j+1}\|_{h,\infty} \leq K_0$  when  $r \leq \frac{1}{2}$ . This is done in a two step process, exploiting the one-dimensional nature of the factors constituting the scheme. Let  $V_{\ell,m}^{j+1} = (1 + r\delta_y^2)U_{\ell,m}^j$  then

$$\begin{aligned} |V_{\ell,m}^{j+1}| &= |rU_{\ell,m-1}^j + (1 - 2r)U_{\ell,m}^j + rU_{\ell,m+1}^j| \\ &\leq (r + |1 - 2r| + r)K_0 = K_0 \end{aligned}$$

since  $1 - 2r \geq 0$ . A similar argument with  $U_{\ell,m}^{j+1} = (1 + r\delta_y^2)V_{\ell,m}^{j+1}$  gives

$$\begin{aligned} |U_{\ell,m}^{j+1}| &= |rV_{\ell-1,m}^j + (1 - 2r)V_{\ell,m}^j + rV_{\ell+1,m}^j| \\ &\leq (r + |1 - 2r| + r)K_0 = K_0 \end{aligned}$$

and so the induction hypothesis holds with  $n = j + 1$ . Thus  $\|U^n\|_{h,\infty} \leq K_0$  holds for  $n = 0, 1, \dots$  provided that  $r \leq \frac{1}{2}$ .

### 11.31

(a) The given finite difference scheme may be written

$$U_m^{n+1} = U_m^n + \frac{k}{h^2 r_m} (r_{m+1/2}(U_{m+1}^n - U_m^n) - r_{m-1/2}(U_m^n - U_{m-1}^n)),$$

Since  $r_{m\pm 1/2} = r_m \pm h/2$ , this may be written as

$$\begin{aligned} U_m^{n+1} &= U_m^n + \frac{k}{h^2 r_m} (r_m(U_{m+1}^n - 2U_m^n + U_{m-1}^n) + \frac{1}{2}h(U_{m+1}^n - U_{m-1}^n)) \\ k^{-1}(U_m^{n+1} - U_m^n) &= h^{-2}\delta_r^2 U_m^n + \frac{1}{r_m}h^{-1}\Delta_r U_m^n \end{aligned}$$

which is the same as the scheme obtained by replacing the spatial derivatives in  $u_{rr} + \frac{1}{r}u_r$  with second order accurate differences and the time derivative by a forward difference  $k^{-1}\Delta_t^+$ .

(b) We may also write the scheme as

$$U_m^{n+1} = \frac{k}{h^2}(1 - \frac{h}{2r_m})U_{m-1}^n + (1 - 2\frac{k}{h^2})U_m^n + \frac{k}{h^2}(1 + \frac{h}{2r_m})U_{m+1}^n$$

which is of positive type provided that  $k \leq h^2/2$  (the formula is valid only for  $m \geq 1$  so  $h/(2r_m) = 1/(2m) \leq \frac{1}{2}$ ).

It was shown in Exercise 8.7 that  $u_r \rightarrow 0$  as  $r \rightarrow 0$  when  $u$  possesses circular symmetry. Hence  $u_r/r$  is of the form  $0/0$  at the origin and, by l'Hôpital's rule:

$$\lim_{r \rightarrow 0} \frac{u_r}{r} = \lim_{r \rightarrow 0} \frac{\partial_r u_r}{\partial_r r} = u_{rr}(0, t)$$

so that the PDE becomes  $u_t = 2u_{rr}$  at  $r = 0$ . The standard FTCS approximation of this equation is

$$U_0^{n+1} = U_0^n + 2\frac{k}{h^2}\delta_r^2 U_0^n = U_0^n + 2\frac{k}{h^2}(U_{-1}^n - 2U_0^n + U_1^n)$$

but, because of symmetry,  $U_{-1}^n = U_1^n$  and so

$$U_0^{n+1} = U_0^n + 4\frac{k}{h^2}(U_1^n - U_0^n) = (1 - 4\frac{k}{h^2})U_0^n + 4\frac{k}{h^2}U_1^n$$

which is of positive type for  $k \leq h^2/4$ .

### 11.33

We use the expressions derived for  $\mathcal{L}_h U_{\ell,m}$  derived in the solution to Exercise 10.17—the notation, and numbering of nodes, is taken from that solution. Then

$$\begin{aligned} P_1 : U_1^{n+1} &= U_1^n - k(65U_1^n - 16U_2^n - 16U_4^n - 112), \\ P_2 : U_2^{n+1} &= U_2^n - k(-16U_1^n + 65U_2^n - 16U_3^n - 16U_5^n), \\ P_3 : U_3^{n+1} &= U_3^n - k(545U_3^n - \frac{144}{5}U_2^n), \\ P_4 : U_4^{n+1} &= U_4^n - k(-16U_1^n + 65U_4^n - 16U_5^n - 16U_6^n - 160), \\ P_5 : U_5^{n+1} &= U_5^n - k(273U_5^n - \frac{128}{5}U_2^n - \frac{288}{11}U_4^n), \\ P_6 : U_6^{n+1} &= U_6^n - k(\frac{979}{3}U_6^n - \frac{256}{11}U_4^n - 9). \end{aligned}$$

These equations will be of positive type if the coefficient of  $U_j^n$  is non-negative at the point  $P_j$ . The most restrictive condition occurs at  $P_3$  where it is required that  $k \leq 1/545$ . The corresponding limit for a regular grid is  $k \leq 1/64$  (as, for instance, at  $P_j$ ,  $j = 1, 2, 4$ ).

## Exercises 12 Finite difference methods for hyperbolic PDEs

### 12.1

Differentiating the PDE  $u_t + au_x = 0$  with respect to  $t$  gives:

$$u_{tt} + au_{tx} = 0 \quad \Rightarrow \quad u_{tt} + a\partial_x u_t = 0 \quad \Rightarrow \quad u_{tt} + a\partial_x(-au_x) = 0$$

and so  $u_{tt} = a^2 u_{xx}$ . Using this in (12.21) leads to

$$\mathcal{R}_h|_m^n = \frac{1}{2}ah(c-1)u_{xx}|_m^n + \mathcal{O}(h^2) + \mathcal{O}(k^2),$$

as required.

### 12.3

The product of the coefficients of  $U_{m\pm 1}^n$  have opposite sign so it is not possible for both to be non-negative.

The amplification factor of the scheme is

$$\xi = 1 - ic \sin \kappa h \quad \Rightarrow \quad |\xi|^2 = 1 + c^2 \sin^2 \kappa h$$

thus  $\xi(\frac{1}{2}\pi) = 1 + c^2$  and it is not possible to find a constant  $C$ , independent of  $h$  and  $k$  such that  $|\xi| \leq 1 + Ck$  (see Definition 11.20).

### 12.5

- (a) The Lax–Friedrichs scheme can be written as<sup>2</sup>

$$U_m^n = \frac{1}{2}(1+c)U_{m-1}^n + U_m^n + \frac{1}{2}(1-c)U_{m+1}^n$$

so its stencil is the same as that of the Lax–Wendroff method (see Fig. 12.6 (left)).

- (b) The local truncation error is (note the scaling), using the results from Table 6.1,

$$\begin{aligned} \mathcal{R}_h|_m^n &:= k^{-1}(u_m^{n+1} - u_m^n + c\Delta_x u_m^n - \frac{1}{2}\delta_x^2 u_m^n) \\ &= (u_t + \frac{1}{2}ku_{tt} + \mathcal{O}(k^2)) + a(u_x + \frac{1}{6}h^2 u_{xxx} + \mathcal{O}(h^2)) - \frac{1}{2}(h^2/k)u_{xx} + \mathcal{O}(h^4) \\ &= (u_t + au_x) + \frac{1}{2}ku_{tt} - \frac{1}{2}(h^2/k)u_{xx} + \mathcal{O}(k^2) + \mathcal{O}(h^2) + \mathcal{O}(h^4/k) \end{aligned}$$

so that  $\mathcal{R}_h = \mathcal{O}(h) + \mathcal{O}(k)$  if  $c$  is fixed—the method is consistent of first order.

- (c) We see from part (a) that the coefficients on the right hand side are non-negative for  $-1 \leq c \leq 1$  in which case the scheme is of positive type.
- (d) The amplification factor of the scheme is, writing  $\theta = \kappa h$ ,

$$\xi = 1 - ic \sin \theta - 2 \sin^2 \frac{1}{2}\theta = \cos \theta - ic \sin \theta.$$

Therefore,

$$|\xi|^2 - 1 = \cos^2 \theta - 1 + c^2 \sin^2 \theta = -(1 - c^2) \sin^2 \theta$$

and  $|\xi|^2 \leq 1$  for all  $\theta \in [-\pi, \pi]$  if, and only if,  $-1 \leq c \leq 1$ .

---

<sup>2</sup>Note that this can also be deduced from (12.26) by replacing  $U_m^n$  by the average  $\frac{1}{2}(U_{m-1}^n + U_{m+1}^n)$ .



- (e) The Lax-Friedrichs method has the same order of consistency as the FTBS and FTFS methods so holds no advantage in that department. One major advantage that it does hold is that it is stable for  $|c| \leq 1$  regardless of the sign of the advection speed so that it can be applied to systems of hyperbolic equations that have both right and left moving characteristics.
- (f) The end product of the calculation in part (b) may be reorganized to give

$$\mathcal{R}_h|_m^n = (u_t + au_x - \frac{1}{2} \frac{h^2}{k} (1 + c^2) u_{xx}) + \mathcal{O}(k^2) + \mathcal{O}(h^2) + \mathcal{O}(h^4/k)$$

where we have used  $u_{tt} = a^2 u_{xx}$ . Thus the local truncation error is consistent of order  $\mathcal{O}(k^2) + \mathcal{O}(h^2)$  (provided that  $c = ak/h$  is fixed) with the PDE  $u_t + au_x = \varepsilon u_{xx}$ , where  $\varepsilon = \frac{1}{2} h^2 (1 + c^2)/k$ .

The scheme is consistent with an advection-diffusion equation which is of parabolic type. The numerical solutions will therefore become smoother (damped) as time proceeds.

## 12.7

We substitute  $U_m^n = \xi^n e^{i\kappa m h}$  into the Lax-Wendroff method scheme

$$U_m^{n+1} = [1 - c\Delta_x + \frac{1}{2} c^2 \delta_x^2] U_m^n$$

and use the relationships  $U_m^{n+1} = \xi U_m^n$ , (11.47b) and (11.47c) to give

$$\xi U_m^n = U_m^n - i c \sin(\kappa h) U_m^n - 2c^2 \sin^2(\frac{1}{2} \kappa h) U_m^n \Rightarrow \xi = 1 - i c \sin \kappa h - 2c^2 \sin^2 \frac{1}{2} \kappa h.$$

Then, writing  $s = \sin^2 \frac{1}{2} \kappa h$

$$\begin{aligned} |\xi|^2 - 1 &= (1 - 2c^2 s)^2 - 1 + c^2 \sin^2(\kappa h) = -4c^2 s + 4c^4 s^2 + 4c^2 s(1 - s) \\ &= -4c^2 s^2 (1 - c^2). \end{aligned}$$

Clearly  $|\xi|^2 - 1 \leq 0$  if  $c^2 \leq 1$ .

## 12.9

With  $\mu = 2, \nu = 0$  the formula (12.17) leads to the coefficients

$$\alpha_{-2} = \prod_{\substack{\ell=-2 \\ \ell \neq -2}}^0 \frac{\ell + c}{\ell + 2} = -\frac{1}{2} c(1 - c), \quad \alpha_{-1} = \prod_{\substack{\ell=-2 \\ \ell \neq -1}}^0 \frac{\ell + c}{\ell + 1} = c(2 - c), \quad \alpha_0 = \prod_{\substack{\ell=-2 \\ \ell \neq 0}}^0 \frac{\ell + c}{\ell} = \frac{1}{2} (1 - c)(2 - c)$$

$U_m^{n+1} = \alpha_{-2} U_{m-2}^n + \alpha_{-1} U_{m-1}^n + \alpha_0 U_m^n$  gives the Warming-Beam scheme (12.64).

## 12.11

- (a) The local truncation error of Leith's method is

$$\begin{aligned} \mathcal{R}_h|_m^n &:= k^{-1} (u_m^{n+1} - u_m^n + c\Delta_x U_m^n - (r + \frac{1}{2} c^2) \delta_x^2 u_m^n), \\ &= u_t + \frac{1}{2} k u_{tt} + \mathcal{O}(k^2) + a(u_x + \frac{1}{6} h^2 u_{xxx} + \mathcal{O}(h^2)) - (\varepsilon + \frac{1}{2} a^2 k)(u_{xx} + \mathcal{O}(h^2)) \\ &= (u_t + au_x - \varepsilon u_{xx}) + \frac{1}{2} k (u_{tt} - a^2 u_{xx}) + \mathcal{O}(k^2) + \mathcal{O}(h^2) \end{aligned}$$

and so is consistent of order  $\mathcal{O}(k) + \mathcal{O}(h^2)$  with the advection-diffusion equation  $u_t + au_x = \varepsilon u_{xx}$ .

(b) Differentiating the PDE with respect to  $t$  and  $x$  gives

$$\begin{aligned} u_{tt} + au_{xt} &= \varepsilon u_{xxt} \\ u_{xt} + au_{xx} &= \varepsilon u_{xxx} \end{aligned}$$

so, subtracting  $a \times$  second from the first equation leads to  $u_{tt} = a^2 u_{xx} + \varepsilon(u_{xxt} - au_{xxx})$  and, therefore,

$$\mathcal{R}_h|_m^n = \frac{1}{2}k\varepsilon(u_{xxt} - au_{xxx}) + \mathcal{O}(k^2) + \mathcal{O}(h^2).$$

From the point of view of convergence as  $h, k \rightarrow 0$  the scheme is clearly consistent of order  $\mathcal{O}(k) + \mathcal{O}(h^2)$  with the advection-diffusion equation. However, the factor  $\varepsilon$  multiplying the leading term in the local truncation error means that in computations where  $\varepsilon \ll 1$ , the method will, on coarser grids, effectively perform as a second order scheme.

(c) We observe that the amplification factor for Leith's scheme

$$\xi = 1 - 4(r + \frac{1}{2}c^2) \sin^2(\frac{1}{2}\kappa h) + i\rho \sin(\kappa h)$$

is the same as the FTCS scheme (11.53) with  $r$  replaced by  $(r + \frac{1}{2}c^2)$  and  $\rho$  replaced by  $-c = -ak/h$ . The stability conditions (11.56) become

$$\frac{1}{2}c^2 \leq r + \frac{1}{2}c^2 \leq \frac{1}{2},$$

so the left inequality is always satisfied leaving  $2r + c^2 \leq 1$ . Written in terms of  $h$  and  $k$ , this requires

$$2\varepsilon k + a^2 k^2 \leq h^2.$$

If both sides are multiplied by  $a^2/\varepsilon^2$  we find

$$2\frac{ka^2}{\varepsilon} + \left(\frac{ka^2}{\varepsilon}\right)^2 \leq \left(\frac{ah}{\varepsilon}\right)^2 \Rightarrow 2\hat{k} + \hat{k}^2 \leq \hat{h}^2,$$

where the stability region is independent of any parameters when expressed in terms of  $\hat{k} := ka^2/\varepsilon$  and  $\hat{h} := ha/\varepsilon$  (known as non-dimensional variables because their values are independent of the units used to measure  $k$ ,  $h$ ,  $a$  and  $\varepsilon$ ). Thus, in the  $\hat{h}$ - $\hat{k}$  plane the boundary is a branch of the hyperbola

$$(1 + \hat{k})^2 - \hat{h}^2 = 1$$

whose centre is at  $\hat{k} = -1$ ,  $\hat{h} = 0$  and has asymptotes  $\hat{k} = -1 \pm \hat{h}$ . The region bounded by this curve and the  $\hat{h}$ -axis is shown shaded in Fig. 13 but the scales shown on the axes are for the grid sizes  $h$  and  $k$ . Also shown for reference is the curve  $k = \varepsilon h^2/2$  (dashed), which is the corresponding stability limit for the FTCS approximation of the heat equation  $u_t = \varepsilon u_{xx}$ . We can see that this is the relevant limit when  $\varepsilon/a$  is large (when we should focus on the stability region near the origin where  $h$  and  $k$  are small). On the other hand, when  $\varepsilon/a$  is small—the advection dominated case—we enjoy a limit which is approximately that for the FTCS approximation of the advection equation  $u_t + au_x = 0$ , i.e.,  $k \leq ah$ .

## 12.13

(a) We evaluate the coefficients supplied in (12.33a) at the Courant numbers  $c = -1, 0, 1, 2$  and display the results in the following table.

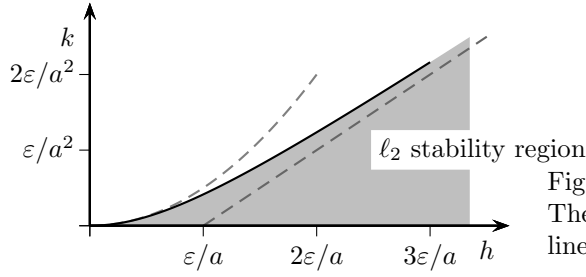


Figure 13: The  $\ell_2$ -stability region for Leith's method. The dashed curve shows  $k = \varepsilon h^2/2$  and the dashed line the asymptote  $k = (ah - \varepsilon)/a^2$ .

Coeffs.	$\alpha_{-2}$	$\alpha_{-1}$	$\alpha_0$	$\alpha_1$
$c$	$-\frac{1}{6}c(1-c^2)$	$\frac{1}{2}c(1+c)(2-c)$	$\frac{1}{2}(1-c^2)(2-c)$	$\frac{1}{6}c(c-1)(2-c)$
2	1	0	0	0
1	0	1	0	0
0	0	0	1	0
-1	0	0	0	1

and so  $U_m^{n+1} = U_{m-c}^n$  for  $c = -1, 0, 1, 2$ .

(b) Using (12.33a), the LTE is given by

$$\mathcal{R}_h|_m^n := k^{-1}(u_m^{n+1} - u_m^n + c\Delta_x u_m^n - \frac{1}{2}c^2\delta_x^2 u_m^n - \frac{1}{6}c(1-c^2)\Delta_x^- \delta_x^2 u_m^n)$$

The expansion of all terms except the last on the right hand side are available in Table 6.1. For the final term we proceed as follows (many other routes to the same destination are possible): Since  $\Delta_x^- v_m^n = hv_x - \frac{1}{2}h^2v_{xx} + \mathcal{O}(h^3)$  (Table 6.1) then, with  $v = \delta_x^2 u_m^n = h^2u_{xx} + \mathcal{O}(h^4)$ ,

$$\Delta_x^- \delta_x^2 u_m^n = h\partial_x(h^2u_{xx} + \mathcal{O}(h^4)) + \mathcal{O}(h^4) = h^3u_{xxx} - h^4u_{xxxx} + \mathcal{O}(h^5).$$

Now

$$\begin{aligned} \mathcal{R}_h|_m^n = & u_t + \frac{1}{2}ku_{tt} + \frac{1}{6}k^2u_{ttt} + \frac{1}{24}k^3u_{tttt} + \mathcal{O}(k^4) \\ & + a(u_x + \frac{1}{6}h^2u_{xxx} + \mathcal{O}(h^4)) \\ & - \frac{1}{2}ka^2(u_{xx} + \frac{1}{12}h^2u_{xxxx} + \mathcal{O}(h^4)) \\ & - \frac{1}{6}a(1-c^2)(h^2u_{xxx} - h^3u_{xxxx} + \mathcal{O}(h^4)) \end{aligned}$$

so that

$$\begin{aligned} \mathcal{R}_h|_m^n = & u_t + au_x + \frac{1}{2}k(u_{tt} - a^2u_{xx}) + \frac{1}{6}k^2(u_{ttt} + a^2u_{xxx}) \\ & + \frac{1}{24}k^3(u_{tttt} + a^2u_{xxxx}) + \frac{1}{24}a^4h^3(1+c)(1-c)(2-c)u_{xxxx} + \mathcal{O}(h^4). \end{aligned}$$

But  $\partial_t^j u = (-a\partial_x)^j u$  so that this reduces to

$$\mathcal{R}_h|_m^n = \frac{1}{24}a^4h^3(1+c)(1-c)(2-c)u_{xxxx} + \mathcal{O}(h^4).$$

Notice how the leading coefficient vanishes at  $c = -1, 1, 2$  (as do all subsequent coefficients in deference to part (a)).

(c) The CFL condition follows immediately from Example 12.3.

- (d) We substitute  $U_m^n = \xi^n e^{i\kappa m h}$  into the finite difference scheme and use the relationships  $U_m^{n+1} = \xi U_m^n$ ,  $\Delta_x U_m^n = (1 - e^{-i\kappa h})$ , (11.47b) and (11.47c) to give

$$\xi = 1 - i c \sin \theta - 2c^2 \sin^2 \frac{1}{2}\theta - \frac{2}{3}c(1 - c^2)(1 - e^{-i\theta}) \sin^2 \frac{1}{2}\theta,$$

where  $\theta = \kappa h$ . We now write  $s = \sin^2 \frac{1}{2}\theta$  and then  $1 - e^{-i\theta} = (1 - \cos \theta) + i \sin \theta = 2s + i \sin \theta$  so that, collecting real and imaginary parts

$$\xi = 1 - 2c^2 s - \frac{4}{3}c(1 - c^2)s^2 - i c \left(1 + \frac{2}{3}c(1 - c^2)s\right) \sin \theta.$$

Therefore, since  $\sin^2 \theta = 4s(1 - s)$  and defining  $A = \frac{2}{3}c(1 - c^2)s$  in order to lighten the complexity,

$$\begin{aligned} \xi &= 1 - 2c^2 s - 2As - i(A + c) \sin \theta. \\ |\xi|^2 - 1 &= (1 - 2s(A + c^2))^2 - 1 + 4(A + c)^2 s(1 - s) \\ &= -4s(A + c^2) + 4s^2(A + c^2)^2 + 4(A + c)^2 s - 4(A + c)^2 s^2 \\ &= -4s((A + c^2) - (A + c)^2) + 4s^2((A + c^2)^2 - (A + c)^2) \\ &= -4sA(1 - 2c - A) + 4s^2(c^2 - c)(c^2 + c + 2A) \\ &= -4sA(1 - 2c - A) - \underbrace{4s^2 c^2(1 - c^2) - 8s^2 A c(1 - c)}_{-6scA} \\ &= -4sA(1 - \frac{1}{2}c - \underbrace{A + 2sc(1 - c)}_{\frac{2}{3}sc(1 - c)(c - 2)}) = -\frac{8}{3}c(1 - c^2)(2 - c)s^2 \left[1 + \frac{2}{3}sc(1 - c)\right]. \end{aligned}$$

- (e) For  $\ell_2$ -stability it is necessary to have  $|\xi|^2 - 1 \leq 0$ . Since  $F(s) := (|\xi|^2 - 1)/s^2$  is linear in  $s$  and, in order for it to be non-positive for  $0 \leq s \leq 1$ , it must be non-positive at  $s = 0$  and  $s = 1$ . Now

$$\begin{aligned} F(0) &= -\frac{8}{3}c(1 - c^2)(2 - c) \leq 0 & \Rightarrow c \in [-1, -\frac{1}{2}] \cup [0, 1] \cup [\frac{3}{2}, 2] \\ F(1) &= -\frac{8}{9}c(1 - c^2)(2 - c)(1 + 2c)(3 - 2c) & \Rightarrow c \in (-\infty, -1] \cup [0, 1] \cup [2, \infty) \end{aligned}$$

and the only interval held in common is  $c \in [0, 1]$ .

## 12.15

The  $m$ th component of  $C\mathbf{v}_j$  is

$$\begin{aligned} (C\mathbf{v}_j)_m &= -c(\mathbf{v}_j)_{m-1} + (1 + c)(\mathbf{v}_j)_m = -ce^{2\pi i(m-1)j/M} + (1 + c)e^{2\pi i m j/M} \\ &= (1 + c - ce^{-2\pi i j/M})(\mathbf{v}_j)_m. \end{aligned}$$

this holds for all  $m = 1, 2, \dots, M$  (at  $M = 1$  we need to recognise that  $(\mathbf{v}_j)_M = (\mathbf{v}_j)_{M-1}$  because of the periodic nature of its components). Thus  $C\mathbf{v}_j = \lambda_j \mathbf{v}_j$ , where  $\lambda_j = 1 + c - ce^{-2\pi i j/M}$ . Comparing this with the amplification factor  $\xi(\kappa h)$  given by (12.36), it is seen that

$$\lambda_j = 1/\xi(2\pi j h)$$

corresponding to the wavenumber  $\kappa = 2\pi j$ . This result generalises readily to all *constant coefficient* finite difference approximations of parabolic and hyperbolic PDEs with periodic BCs because all circulant matrices of a given dimension share the same set of eigenvectors.

### 12.17

From (12.22) and Exercise 12.16 we find, respectively,

$$\xi_F(\theta) = 1 - c + ce^{-i\theta} \text{ and } \xi_B(\theta) = \frac{1}{1 - c + ce^{i\theta}} = \frac{1}{\xi_F^*(\theta)}$$

(where  $\xi_F^*(\theta)$  is the complex conjugate of  $\xi_F(\theta)$ ) so that  $\xi_F^*(\theta)\xi_B(\theta) = 1$ . Since  $|\xi_F(\theta)| |\xi_B(\theta)| = 1$  it follows that if  $|\xi_F(\theta)| < 1$  (the FTBS scheme is stable) for some scaled wavenumber  $\theta$ , then  $|\xi_B(\theta)| > 1$  (the BTFS scheme is unstable) at that wavenumber, and *vice versa*. Their stability regions are therefore complements of each other. Since the FTBS scheme is stable for  $0 \leq c \leq 1$  we deduce that the BTFS scheme is stable for  $c \geq 1$  and  $c \leq 0$ .

### 12.19

With  $\eta = (1 - c) + (1 + c)e^{-i\kappa h}$ , then

$$\eta^* = (1 - c) + (1 + c)e^{i\kappa h} = ((1 - c)e^{-i\kappa h} + (1 + c))e^{i\kappa h}$$

so that (12.39) becomes  $\xi = e^{-i\kappa h}\eta/\eta^*$ . It follows that  $|\xi| = 1$  (for all  $c$ ) since  $|\eta| = |\eta^*|$  and  $|e^{-i\kappa h}| = 1$ .

### 12.21

The amplification factor of the method (12.5) is

$$\xi = \sum_{j=-\mu}^{\nu} \alpha_j e^{ij\kappa h}$$

and its LTE is

$$\mathcal{R}_h|_m^n := k^{-1} (u_m^{n+1} - \sum_{j=-\mu}^{\nu} \alpha_j u_{m+j}^n).$$

The solution of the advection equation  $u_t + au_x = 0$  with initial condition  $u(x, 0) = \exp(i\kappa x)$  is  $u(x, t) = \exp(i\kappa(x - at))$  so the LTE becomes

$$\begin{aligned} \mathcal{R}_h|_m^n &= k^{-1} (e^{i\kappa(x_m - a(n+1)k)} - \sum_{j=-\mu}^{\nu} \alpha_j e^{i\kappa(x_m + j - ank)}) \\ &= k^{-1} (e^{-ai\kappa k} - \sum_{j=-\mu}^{\nu} \alpha_j e^{ij\kappa h}) e^{i\kappa(x_m - ank)} = k^{-1} (e^{-ci\theta} - \xi(\theta)) u_m^n, \end{aligned}$$

where  $\theta = \kappa h$ . This shows that equation (12.50) holds more generally than just for the Lax-Wendroff method.

Now  $h^{p+1}\partial_x^{p+1}u = (i\kappa h)^{p+1}u = (i\theta)^{p+1}u$  with a similar expression for  $h^{p+2}\partial_x^{p+2}u$  and the given error expansion gives

$$k\mathcal{R}_h|_m^n = C_{p+1}(i\theta)^{p+1}u + C_{p+2}(i\theta)^{p+2}u + \mathcal{O}(h^{p+3})$$

and, using (12.50), we obtain

$$\xi(\theta) = e^{-i\theta} - C_{p+1}(i\theta)^{p+1} - C_{p+2}(i\theta)^{p+2}u + \mathcal{O}(h^{p+3})$$

showing that the amplification factor is an order  $(p + 1)$  approximation of  $e^{-i\theta}$ .

- (a)  $p$  is odd: then  $(i\theta)^{p+1}$  is real and  $(i\theta)^{p+2}$  is imaginary. In fact, supposing that  $p = 2q - 1$  (where  $q$  is a positive integer),

$$(i\theta)^{p+1} = i^{2q}\theta^{p+1} = (-1)^q\theta^{p+1} = (-1)^{(p+1)/2}\theta^{p+1}$$

and  $(i\theta)^{p+2} = (-1)^{(p+1)/2}i\theta^{p+2}$ . Hence

$$\xi(\theta) = (\cos c\theta - (-1)^{(p+1)/2}C_{p+1}\theta^{p+1}) - i(\sin c\theta + (-1)^{(p+1)/2}i\theta^{p+2}) + \mathcal{O}(\theta^{p+3})$$

leading to

$$\begin{aligned} |\xi(\theta)|^2 &= (\cos c\theta - (-1)^{(p+1)/2}C_{p+1}\theta^{p+1})^2 + (\sin c\theta + (-1)^{(p+1)/2}C_{p+2}\theta^{p+2})^2 + \mathcal{O}(\theta^{p+3}) \\ &= \cos^2 c\theta + \sin^2 c\theta - 2(-1)^{(p+1)/2}C_{p+1}\theta^{p+1} \cos c\theta + \mathcal{O}(\theta^{p+3}) \\ &= 1 - 2(-1)^{(p+1)/2}C_{p+1}\theta^{p+1} + \mathcal{O}(\theta^{p+3}), \end{aligned}$$

where we have used  $\theta^{p+1} \cos c\theta = \theta^{p+1} + \mathcal{O}(\theta^{p+3})$  and  $\theta^{p+2} \sin c\theta = \mathcal{O}(\theta^{p+3})$ .

- (b)  $p$  is even: then  $(i\theta)^{p+1}$  is imaginary and  $(i\theta)^{p+2}$  is real. In fact, supposing that  $p = 2q$  (where  $q$  is a positive integer),

$$(i\theta)^{p+1} = i(i^{2q})\theta^{p+1} = (-1)^q i\theta^{p+1} = (-1)^{p/2} i\theta^{p+1}$$

and  $(i\theta)^{p+2} = (-1)^{p/2+1}\theta^{p+2}$ . Hence

$$\xi(\theta) = (\cos c\theta - (-1)^{p/2+1}C_{p+2}\theta^{p+2}) - i(\sin c\theta + (-1)^{p/2}C_{p+1}\theta^{p+1}) + \mathcal{O}(\theta^{p+3})$$

leading to

$$\begin{aligned} |\xi(\theta)|^2 &= (\cos c\theta - (-1)^{p/2+1}C_{p+2}\theta^{p+2})^2 + (\sin c\theta + (-1)^{p/2}C_{p+1}\theta^{p+1})^2 + \mathcal{O}(\theta^{p+3}) \\ &= \cos^2 c\theta + \sin^2 c\theta - 2(-1)^{p/2+1}C_{p+2}\theta^{p+2} \cos c\theta + 2(-1)^{p/2}C_{p+1}\theta^{p+1} \sin c\theta + \mathcal{O}(\theta^{p+3}) \\ &= 1 - 2(-1)^{p/2+1}C_{p+2}\theta^{p+2} + 2(-1)^{p/2}C_{p+1}\theta^{p+1} \sin c\theta + \mathcal{O}(\theta^{p+3}) \\ &= 1 - 2(-1)^{p/2+1}(cC_{p+1} + C_{p+2})\theta^{p+2} + \mathcal{O}(\theta^{p+3}), \end{aligned}$$

where we have used  $\theta^{p+2} \cos c\theta = \theta^{p+2} + \mathcal{O}(\theta^{p+4})$  and  $\theta^{p+1} \sin c\theta = c\theta^{p+1} + \mathcal{O}(\theta^{p+3})$ .

The bottom line is that for a method whose order of consistency  $p$  (odd), the order of dissipation is  $(p + 1)$  but, if the order of consistency  $p$  (even), the order of dissipation is  $(p + 2)$ .

### 12.23

When  $c > 0$ , characteristics travel from left-to-right (see Fig. 14, left) and the domain of dependence of the anchor point (shown as  $\circ$ ) is shaded. The method must be used as

$$U_m^{n+1} = (U_m^n + cU_{m-1}^{n+1})/(1 + c), \quad m = 1, 2, \dots, M$$

so the BC must be placed at  $x = 0$ .

The corresponding situation when  $c \leq -1$  is shown in Fig. 14 (right) and the BC needs to be at  $x = 1$ .

### 12.25

The Lax-Wendroff method is applied at the points  $(x_m, nk)$ ,  $0 < m < M$ ,  $n \geq 0$ . At these points



Figure 14: The two modes of operation of the BTBS scheme in Exercise 12.23 with target points ( $\circ$ ). Both modes are stable provided that the characteristics lie in the shaded regions.

the LTE  $\mathcal{R}_m^n = 0$  because the exact solution is a quadratic polynomial and the expression for the LTE involves the 3rd derivatives of  $u$  (see (12.30)). So

$$\begin{aligned}\mathcal{R}_m^n &= \frac{1}{k}(u_m^{n+1} - u_m^n + c\Delta_x u_m^n - \frac{1}{2}c^2\delta_x^2 u_m^n) = 0, \\ U_m^{n+1} &= (1 - c\Delta_x + \frac{1}{2}c^2\delta_x^2) U_m^n,\end{aligned}$$

and so the global error  $E = u - U$  itself is a solution of the Lax-Wendroff scheme:

$$\begin{aligned}E_m^{n+1} &= (1 - c\Delta_x + \frac{1}{2}c^2\delta_x^2) E_m^n \\ &= \frac{1}{2}c(1 + c)E_{m-1}^n + (1 - c^2)E_m^n + \frac{1}{2}c(c - 1)E_{m+1}^n\end{aligned}$$

with starting condition  $E_m^0 = 0$  and a homogeneous BC  $E_0^n = 0$  at  $x = 0$ . Thus, by a domain of dependence argument,  $E_m^n = 0$  for  $x_m + t_n \leq 1$ .

The LTE of the FTBS scheme (see (12.21)) at  $x = 1$  is, for  $u(x, t) = (x - t)^2$  and  $a = 1$ ,

$$\mathcal{R}_M^n = k^{-1}(u_M^{n+1} - cu_{M-1}^n - (1 - c)u_M^n) = -h(1 - c)$$

and so

$$E_M^{n+1} = cE_{M-1}^n - (1 - c)E_M^n - hk(1 - c)$$

for  $n = 0, 1, 2, \dots$

Assuming now that  $E_m^n \rightarrow A_m$  as  $n \rightarrow \infty$ , the Lax-Wendroff equations reduce to

$$(1 + c)A_{m-1} - 2cA_m + (c - 1)A_{m+1} = 0,$$

Setting

$$C = \frac{1}{2}(1 - c)^2 h^2 \left( -\frac{1 + c}{1 - c} \right)^{1-M}, \quad \rho = \frac{c + 1}{c - 1}$$

then the proposed solution is  $A_m = C\rho^m$ . Substituting, we find,

$$(1 + c)A_{m-1} - 2cA_m + (c - 1)A_{m+1} = C\rho^m \left( \frac{1 + c}{\rho} - 2c + (c - 1)\rho \right) = C\rho^m ((c - 1) - 2c + (c + 1)) = 0$$

so these equations are satisfied.

When  $E_m^n \rightarrow A_m$  as  $n \rightarrow \infty$ , the FTBS equations reduce to  $A_M - A_{M-1} = -h^2(1 - c)$ . Substituting the putative solution into the left hand side of this, we find

$$A_M - A_{M-1} = C\rho^{M-1}(\rho - 1) = \frac{2}{c - 1}C\rho^{M-1} = -h^2(1 - c)$$

as required.

## 12.27

The FTCS method is stable only for positive eigenvalues and the FTBS method is stable only for negative eigenvalues but the matrix  $A$  has eigenvalues  $\pm a$  of both signs.

(a) Applying the FTBS method to (12.66a) and the FTBS method to (12.66b) leads to

$$\begin{aligned}(U_m^{n+1} - V_m^{n+1}) &= (U_m^n - V_m^n) - c\Delta_x^-(U_m^n - V_m^n) = c(U_{m-1}^n - V_{m-1}^n) + (1-c)(U_m^n - V_m^n) \\ (U_m^{n+1} + V_m^{n+1}) &= (U_m^n + V_m^n) + c\Delta_x^+(U_m^n + V_m^n) = (1-c)(U_m^n + V_m^n) + c(U_{m+1}^n + V_{m+1}^n)\end{aligned}$$

so, by adding and subtracting, we find

$$\begin{aligned}U_m^{n+1} &= (1-c)U_m^n + \frac{1}{2}c(U_{m-1}^n - V_{m-1}^n) + \frac{1}{2}c(U_{m+1}^n + V_{m+1}^n) \\ &= (1 + \frac{1}{2}c\delta_x^2)U_m^n + c\Delta_x V_m^n \\ V_m^{n+1} &= (1-c)V_m^n - c(U_{m-1}^n - V_{m-1}^n) + c(U_{m+1}^n + V_{m+1}^n) \\ &= (1 + \frac{1}{2}c\delta_x^2)V_m^n + c\Delta_x U_m^n \\ \mathbf{U}_m^{n+1} &= (1 + \frac{1}{2}c\delta_x^2)\mathbf{U}_m^n - (k/h)A\Delta_x \mathbf{U}_m^n = \begin{bmatrix} (1 + \frac{1}{2}c\delta_x^2) & c\Delta_x \\ c\Delta_x & (1 + \frac{1}{2}c\delta_x^2) \end{bmatrix} \mathbf{U}_m^n.\end{aligned}$$

(b) With  $A = \begin{bmatrix} 0 & -a \\ -a & 0 \end{bmatrix}$  the matrix of eigenvectors is  $V = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$  and  $\Lambda = \text{diag}(-a, a)$ . Therefore  $|A| = |a|I$  (where  $I$  is the  $2 \times 2$  identity matrix) and so  $|A| = V|A|V^{-1} = |a|VV^{-1} = |a|I$ . The method now reads

$$\mathbf{U}_m^{n+1} = \mathbf{U}_m^n - (k/h)A\Delta_x \mathbf{U}_m^n + \frac{1}{2}|c|\delta_x^2 \mathbf{U}_m^n$$

which is identical to the method in part (a) provided  $a > 0$ .

(c) Following the steps in Example 12.18 the Lax–Friedrichs scheme is found to be

$$\mathbf{U}_m^{n+1} = [I - (k/h)A\Delta_x + \frac{1}{2}\delta_x^2]\mathbf{U}_m^n,$$

which is identical to the method in part (b) when  $|c| = 1$ .

### 12.29

Consider the critical factor in the final term on the right hand side of (12.77)

$$\Delta_x^-(\varphi_m^n \Delta_x^+ U_m^n) = \varphi_m^n \Delta_x^+ U_m^n - \varphi_{m-1}^n \Delta_x^+ U_{m-1}^n$$

but  $\Delta_x^+ U_{m-1}^n = \Delta_x^- U_m^n$  so that this becomes

$$\begin{aligned}\Delta_x^-(\varphi_m^n \Delta_x^+ U_m^n) &= \varphi_m^n \Delta_x^+ U_m^n - \varphi_{m-1}^n \Delta_x^- U_m^n \\ &= \left(\frac{\varphi_m^n}{\rho_m^n} - \varphi_{m-1}^n\right) \Delta_x^- U_m^n,\end{aligned}$$

where  $\rho_m^n = \Delta_x^- U_m^n / \Delta_x^+ U_m^n$ , and the result follows.

### 12.31

The minmod limiter is a special case of the Chakravarthy–Osher limiter. Chakravarthy–Osher:

$$\begin{aligned}\rho \leq 0 : & \quad \max\{0, \min(\rho, \psi)\} = \max\{0, \rho\} = 0 \\ 0 < \rho \leq \psi : & \quad \max\{0, \min(\rho, \psi)\} = \max\{0, \rho\} = \rho \\ \rho > \psi : & \quad \max\{0, \min(\rho, \psi)\} = \max\{0, \psi\} = \psi.\end{aligned}$$



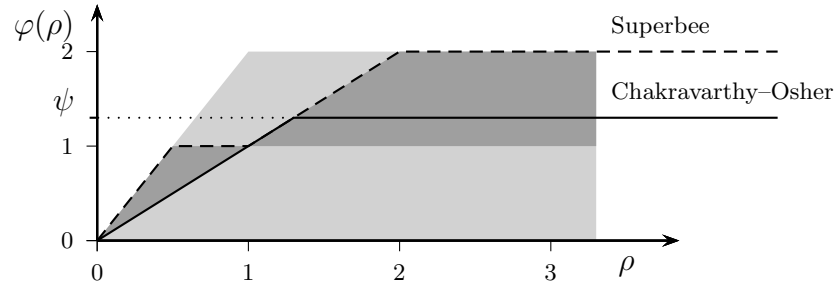


Figure 15: The flux limiters for Exercise 12.31.

Superbee:

$$\begin{aligned}
 \rho \leq 0 : & \quad \max\{0, \min(2\rho, 1), \min(\rho, 2)\} = \max\{0, \rho\} = 0 \\
 0 < \rho \leq \frac{1}{2} : & \quad \max\{0, \min(2\rho, 1), \min(\rho, 2)\} = \max\{0, 2\rho, \rho\} = 2\rho \\
 \frac{1}{2} < \rho \leq 1 : & \quad \max\{0, \min(2\rho, 1), \min(\rho, 2)\} = \max\{0, 1, \rho\} = 1 \\
 1 < \rho \leq 2 : & \quad \max\{0, \min(2\rho, 1), \min(\rho, 2)\} = \max\{0, 1, \rho\} = \rho \\
 \rho > 2 : & \quad \max\{0, \min(2\rho, 1), \min(\rho, 2)\} = \max\{0, 1, 2\} = 2.
 \end{aligned}$$

Essential Partial Differential Equations

Analytical and Computational Aspects

Griffiths, D.F.; Dold, J.W.; Silvester, D.J.

2015, XI, 368 p. 106 illus., 1 illus. in color., Softcover

ISBN: 978-3-319-22568-5