

---

# Finite Difference Computing with Partial Differential Equations

---

Hans Petter Langtangen<sup>1,2</sup>  
Svein Linge<sup>3,1</sup>

<sup>1</sup>Center for Biomedical Computing, Simula Research Laboratory

<sup>2</sup>Department of Informatics, University of Oslo

<sup>3</sup>Department of Process, Energy and Environmental Technology,  
University College of Southeast Norway

This easy-to-read book introduces the basics of solving partial differential equations by finite difference methods. The emphasis is on constructing finite difference schemes, formulating algorithms, implementing algorithms, verifying implementations, analyzing the physical behavior of the numerical solutions, and applying the methods and software to solve problems from physics and biology.

Mar 23, 2016

# Preface

## **hpl 1:** Recall to cite Randy's book [6].

There are so many excellent books on finite difference methods for ordinary and partial differential equations that writing yet another one requires a different view on the topic. The present book is not so concerned with the traditional academic presentation of the topic, but is focused at teaching the practitioner how to obtain reliable computations involving finite difference methods. This focus is based on a set of learning outcomes:

1. understanding of the ideas behind finite difference methods,
2. understanding how to transform an algorithm to a well-designed computer code,
3. understanding how to test (verify) the code,
4. understanding potential artifacts in simulation results.

Compared to other textbooks, the present one has a particularly strong emphasis on computer implementation and verification. It also has a strong emphasis on an intuitive understanding of constructing finite difference methods. To learn about the potential non-physical artifacts of various methods, we study exact solutions of finite difference schemes as these give deeper insight into the physical behavior of the numerical methods than the traditional (and more general) asymptotic error analysis. However, asymptotic results regarding convergence rates, typically truncation errors, are crucial for testing implementations, so an extensive appendix is devoted to the computation of truncation errors.

Various pedagogical elements are utilized to reach the learning outcomes, and these are commented upon next.

**Simplify, understand, generalize.** The book's overall pedagogical philosophy is the three-step process of first *simplifying* the problem to something we can *understand* in detail, and when that understanding is in place, we can *generalize* and hopefully address real-world applications with a sound scientific problem-solving approach. For example, in the chapter on a particular family of equations we first simplify the problem in question to a 1D, constant-coefficient equation with simple boundary conditions. We learn how to construct a finite difference method, how to implement it, and how to understand the behavior of the numerical solution. Then we can generalize to higher dimensions, variable coefficients, a source term, and more complicated boundary conditions. The solution of a compound problem is in this way an assembly of elements that are well understood in simpler settings.

**Constructive mathematics.** This text favors a constructive approach to mathematics. Instead of a set of definitions followed by popping up a method, we emphasize how to think about the construction of a method. The aim is to obtain a good intuitive understanding of the mathematical methods.

The text is written in an easy-to-read style much inspired by the following quote.

*Some people think that stiff challenges are the best device to induce learning, but I am not one of them. The natural way to learn something is by spending vast amounts of easy, enjoyable time at it. This goes whether you want to speak German, sight-read at the piano, type, or do mathematics. Give me the German storybook for fifth graders that I feel like reading in bed, not Goethe and a dictionary. The latter will bring rapid progress at first, then exhaustion and failure to resolve.*

*The main thing to be said for stiff challenges is that inevitably we will encounter them, so we had better learn to face them boldly. Putting them in the curriculum can help teach us to do so. But for teaching the skill or subject matter itself, they are overrated.* [9, p. 86] Lloyd N. Trefethen, Applied Mathematician, 1955-.

This book assumes some basic knowledge of finite difference approximations, differential equations, and scientific Python or MATLAB programming, as often met in an introductory numerical methods course. Readers without this background may start with the light companion book “Finite Difference Computing with Exponential Decay Models” [4]. That book will in particular be a useful resource for the programming parts of the present book. Since the present book deals with partial differential equations, the reader is assumed to master multi-variable calculus and linear algebra.

Fundamental ideas and their associated scientific details are first introduced in the simplest possible differential equation setting, often an ordinary differential equation, but in a way that easily allows reuse in more complex settings with partial differential equations. With this approach, new concepts are introduced with a minimum of mathematical details. The text should therefore have a potential for use early in undergraduate student programs.

**All nuts and bolts.** Many has experienced that “vast amounts of easy, enjoyable time”, as stated in the quote above, arises when mathematics is implemented on a computer. The implementation process triggers understanding, creativity, and curiosity, but many students find the transition from a mathematical algorithm to a working code difficult and spend a lot of time on “programming issues”.

Most books on numerical methods concentrate on the mathematics of the subject while details on going from the mathematics to a computer implementation are less in focus. A major purpose of this text is therefore to help the practitioner by providing *all nuts and bolts* necessary for safely going from the mathematics to a well-designed and well-tested computer code. A significant portion of the text is consequently devoted to programming details.

**Python as programming language.** While MATLAB enjoys widespread popularity in books on numerical methods, we have chosen to use the Python programming language. Python is very similar to MATLAB, but contains a lot of modern software engineering tools that have become standard in the software industry and that should be adopted also for numerical computing projects. Python is at present also experiencing an exponential growth in popularity within the scientific computing community. One of the book’s goals is to present an up-to-date Python eco system for implementing finite difference methods.

**Program verification.** Program testing, called *verification*, is a key topic of the book. Good verification techniques are indispensable when debugging computer code, but also fundamental for achieving reliable simulations. Two verification techniques saturate the book: exact solution of discrete equations (where the approximation error vanishes) and empirical estimation of convergence rates in problems with exact (analytical or manufactured) solutions of the differential equation(s).

**Vectorized code.** Finite difference methods lead to code with loops over large arrays. Such code in plain Python is known to run slowly. We demonstrate, especially in Appendix C, how to port loops to fast,

compiled code in C or Fortran. However, an alternative is to vectorize the code to get rid of the explicit loops Python, and this technique is throughout the book. Vectorization becomes closely connected to the underlying array library, here `numpy`, and is often thought of as a difficult subject by students. Through numerous examples in different contexts, we hope that the present book provides a substantial contribution to explaining how algorithms can be vectorized. Not only will this speed up serial code, but with a library that can produce parallel code from `numpy` commands (such as `Numba`), vectorized code can be automatically turned into parallel code and utilize multi-core processors and GPUs. Also when creating tailored parallel code for today’s supercomputers, vectorization is useful as it emphasizes splitting up an algorithm into plain and simple array operations, where each operation is trivial to parallelize efficiently, rather than trying to develop a “smart” overall parallelization strategy.

**Analysis via exact solutions of discrete equations.** Traditional asymptotic analysis of errors is important for verification of code using convergence rates, but gives a limited understanding of how and why a correctly implemented numerical method may give non-physical results. By developing exact solutions, usually based on Fourier methods, of the discrete equations, one can obtain a physical understanding of the behavior of a numerical method. This approach is favored for analysis of methods in this book.

**Code-inspired mathematical notation.** Our primary aim is to have a clean and easy-to-read computer code, and we want a close one-to-one relationship between the computer code and mathematical description of the algorithm. This principle calls for a mathematical notation that is governed by the natural notation in the computer code. The unknown is mostly called  $u$ , but the meaning of the symbol  $u$  in the mathematical description changes as we go from the exact solution fulfilling the differential equation problem to the symbol  $\mathbf{u}$  that is naturally used in the code.

**Limited scope.** The aim of this book is not to give an overview of a lot of methods for a wide range of mathematical models. Such information can be found in numerous existing, more advanced books. The aim is rather to introduce basic concepts and a thorough understanding of how to think about computing with finite difference methods. We therefore go in depth with only the most fundamental methods and equations.

However, we have a multi-disciplinary scope and address the interplay of mathematics, numerics, computer science, and physics.

**Independent chapters.** Most book authors are careful with avoiding repetitions of material. The chapters in this book, however, contain some overlap, because I want the chapters to be meaningful on their own. Modern publishing technology makes it easy to take selected chapters from different books to make a new book tailored to a specific course. The more a chapter builds on details in other chapters, the more difficult it is to reuse chapters in new contexts. Also, most readers find it convenient that the most important information is explicitly stated, even if it was already met in another chapter.

**Supplementary materials.** All program and data files referred to in this book are available from the book's primary web site: URL: <http://hplgit.github.io/fdm-book/doc/web/>.

**Acknowledgments.** Many students have provided lots of useful feedback on the exposition and found many errors in the text. Special efforts in this regard were made by Imran Ali, Shirin Fallahi, Anders Hafreager, Daniel Alexander Mo Søreide Houshmand, Kristian Gregorius Hustad, Mathilde Nygaard Kamperud, and Fatemeh Miri.

*Oslo, February 2016*

*Hans Petter Langtangen, Svein Linge*



# Contents

Preface .....	v
<b>1 Vibration ODEs</b> .....	1
1.1 Finite difference discretization .....	1
1.1.1 A basic model for vibrations .....	2
1.1.2 A centered finite difference scheme .....	2
1.2 Implementation .....	5
1.2.1 Making a solver function .....	5
1.2.2 Verification .....	7
1.2.3 Scaled model .....	10
1.3 Long time simulations .....	11
1.3.1 Using a moving plot window .....	11
1.3.2 Making animations .....	13
1.3.3 Using Bokeh to compare graphs .....	16
1.3.4 Using a line-by-line ascii plotter .....	19
1.3.5 Empirical analysis of the solution .....	20
1.4 Analysis of the numerical scheme .....	22
1.4.1 Deriving a solution of the numerical scheme .....	22
1.4.2 Exact discrete solution .....	24
1.4.3 Convergence .....	25
1.4.4 The global error .....	26
1.4.5 Stability .....	27

1.4.6	About the accuracy at the stability limit .....	28
1.5	Alternative schemes based on 1st-order equations.....	30
1.5.1	The Forward Euler scheme .....	31
1.5.2	The Backward Euler scheme .....	32
1.5.3	The Crank-Nicolson scheme.....	32
1.5.4	Comparison of schemes.....	34
1.5.5	Runge-Kutta methods .....	36
1.5.6	Analysis of the Forward Euler scheme .....	37
1.6	Energy considerations .....	39
1.6.1	Derivation of the energy expression .....	40
1.6.2	An error measure based on energy .....	42
1.7	The Euler-Cromer method .....	44
1.7.1	Forward-backward discretization.....	44
1.7.2	Equivalence with the scheme for the second-order ODE	46
1.7.3	Implementation .....	47
1.7.4	The velocity Verlet algorithm .....	49
1.8	Generalization: damping, nonlinear spring, and external excitation .....	50
1.8.1	A centered scheme for linear damping .....	51
1.8.2	A centered scheme for quadratic damping.....	52
1.8.3	A forward-backward discretization of the quadratic damping term.....	53
1.8.4	Implementation .....	54
1.8.5	Verification .....	55
1.8.6	Visualization .....	56
1.8.7	User interface .....	56
1.8.8	The Euler-Cromer scheme for the generalized model..	57
1.9	Exercises and Problems .....	59
1.10	Applications of vibration models .....	66
1.10.1	Oscillating mass attached to a spring.....	66
1.10.2	General mechanical vibrating system .....	68
1.10.3	A sliding mass attached to a spring .....	70
1.10.4	A jumping washing machine .....	71
1.10.5	Motion of a pendulum .....	71
1.10.6	Dynamic free body diagram during pendulum motion	74
1.10.7	Motion of an elastic pendulum .....	79
1.10.8	Vehicle on a bumpy road .....	85

---

1.10.9 Bouncing ball . . . . .	86
1.10.10 Electric circuits . . . . .	87
1.11 Exercises . . . . .	88
<b>2 Wave equations . . . . .</b>	<b>93</b>
2.1 Simulation of waves on a string . . . . .	93
2.1.1 Discretizing the domain . . . . .	94
2.1.2 The discrete solution . . . . .	95
2.1.3 Fulfilling the equation at the mesh points . . . . .	95
2.1.4 Replacing derivatives by finite differences . . . . .	95
2.1.5 Formulating a recursive algorithm . . . . .	97
2.1.6 Sketch of an implementation . . . . .	98
2.2 Verification . . . . .	100
2.2.1 A slightly generalized model problem . . . . .	100
2.2.2 Using an analytical solution of physical significance . .	101
2.2.3 Manufactured solution . . . . .	102
2.2.4 Constructing an exact solution of the discrete equations	104
2.3 Implementation . . . . .	106
2.3.1 Callback function for user-specific actions . . . . .	107
2.3.2 The solver function . . . . .	107
2.3.3 Verification: exact quadratic solution . . . . .	109
2.3.4 Visualization: animating the solution . . . . .	110
2.3.5 Running a case . . . . .	114
2.3.6 Working with a scaled PDE model . . . . .	115
2.4 Vectorization . . . . .	117
2.4.1 Operations on slices of arrays . . . . .	117
2.4.2 Finite difference schemes expressed as slices . . . . .	120
2.4.3 Verification . . . . .	121
2.4.4 Efficiency measurements . . . . .	122
2.4.5 Remark on the updating of arrays . . . . .	124
2.5 Exercises . . . . .	125
2.6 Generalization: reflecting boundaries . . . . .	129
2.6.1 Neumann boundary condition . . . . .	129
2.6.2 Discretization of derivatives at the boundary . . . . .	129
2.6.3 Implementation of Neumann conditions . . . . .	131
2.6.4 Index set notation . . . . .	132
2.6.5 Verifying the implementation of Neumann conditions .	134

2.6.6 Alternative implementation via ghost cells .....	136
2.7 Generalization: variable wave velocity .....	139
2.7.1 The model PDE with a variable coefficient .....	139
2.7.2 Discretizing the variable coefficient .....	140
2.7.3 Computing the coefficient between mesh points .....	141
2.7.4 How a variable coefficient affects the stability .....	142
2.7.5 Neumann condition and a variable coefficient .....	143
2.7.6 Implementation of variable coefficients .....	144
2.7.7 A more general PDE model with variable coefficients ..	145
2.7.8 Generalization: damping .....	146
2.8 Building a general 1D wave equation solver .....	147
2.8.1 User action function as a class .....	147
2.8.2 Pulse propagation in two media .....	150
2.9 Exercises .....	153
2.10 Analysis of the difference equations .....	162
2.10.1 Properties of the solution of the wave equation .....	162
2.10.2 More precise definition of Fourier representations ..	164
2.10.3 Stability .....	166
2.10.4 Numerical dispersion relation .....	169
2.10.5 Extending the analysis to 2D and 3D .....	172
2.11 Finite difference methods for 2D and 3D wave equations ...	176
2.11.1 Multi-dimensional wave equations .....	176
2.11.2 Mesh .....	178
2.11.3 Discretization .....	179
2.12 Implementation .....	181
2.12.1 Scalar computations .....	183
2.12.2 Vectorized computations .....	185
2.12.3 Verification .....	187
2.13 Exercises .....	188
2.14 Applications of wave equations .....	189
2.14.1 Waves on a string .....	190
2.14.2 Waves on a membrane .....	193
2.14.3 Elastic waves in a rod .....	194
2.14.4 The acoustic model for seismic waves .....	194
2.14.5 Sound waves in liquids and gases .....	196
2.14.6 Spherical waves .....	198

---

2.14.7 The linear shallow water equations . . . . .	199
2.14.8 Waves in blood vessels . . . . .	201
2.14.9 Electromagnetic waves . . . . .	204
2.15 Exercises . . . . .	205
<b>3 Diffusion equations . . . . .</b>	<b>219</b>
3.1 An explicit method for the 1D diffusion equation . . . . .	220
3.1.1 The initial-boundary value problem for 1D diffusion . . . . .	220
3.1.2 Forward Euler scheme . . . . .	221
3.1.3 Implementation . . . . .	223
3.1.4 Verification . . . . .	225
3.1.5 Numerical experiments . . . . .	227
3.2 Implicit methods for the 1D diffusion equation . . . . .	234
3.2.1 Backward Euler scheme . . . . .	234
3.2.2 Sparse matrix implementation . . . . .	238
3.2.3 Crank-Nicolson scheme . . . . .	239
3.2.4 The unifying $\theta$ rule . . . . .	241
3.2.5 Experiments . . . . .	242
3.2.6 The Laplace and Poisson equation . . . . .	243
3.3 Analysis of schemes for the diffusion equation . . . . .	244
3.3.1 Properties of the solution . . . . .	245
3.3.2 Analysis of discrete equations . . . . .	248
3.3.3 Analysis of the finite difference schemes . . . . .	248
3.3.4 Analysis of the Forward Euler scheme . . . . .	249
3.3.5 Analysis of the Backward Euler scheme . . . . .	251
3.3.6 Analysis of the Crank-Nicolson scheme . . . . .	252
3.3.7 Summary of accuracy of amplification factors . . . . .	253
3.3.8 Analysis of the 2D diffusion equation . . . . .	254
3.3.9 Explanation of numerical artifacts . . . . .	256
3.4 Exercises . . . . .	258
3.5 Diffusion in heterogeneous media . . . . .	261
3.5.1 Discretization . . . . .	261
3.5.2 Implementation . . . . .	262
3.5.3 Stationary solution . . . . .	263
3.5.4 Piecewise constant medium . . . . .	264
3.5.5 Implementation of diffusion in a piecewise constant medium . . . . .	264

---

3.5.6 Axi-symmetric diffusion . . . . .	266
3.5.7 Spherically-symmetric diffusion . . . . .	270
3.6 Diffusion in 2D . . . . .	271
3.6.1 Discretization . . . . .	271
3.6.2 Numbering of mesh points versus equations and unknowns . . . . .	272
3.6.3 Algorithm for setting up the coefficient matrix . . . . .	277
3.6.4 Implementation with a dense coefficient matrix . . . . .	279
3.6.5 Verification: exact numerical solution . . . . .	282
3.6.6 Verification: convergence rates . . . . .	284
3.6.7 Implementation with a sparse coefficient matrix . . . . .	285
3.6.8 The Jacobi iterative method . . . . .	290
3.6.9 Implementation of the Jacobi method . . . . .	292
3.6.10 Test problem: diffusion of a sine hill . . . . .	294
3.6.11 The relaxed Jacobi method and its relation to the Forward Euler method . . . . .	296
3.6.12 The Gauss-Seidel and SOR methods . . . . .	297
3.6.13 Scalar implementation of the SOR method . . . . .	298
3.6.14 Vectorized implementation of the SOR method . . . . .	299
3.6.15 Direct versus iterative methods . . . . .	303
3.6.16 The Conjugate gradient method . . . . .	306
3.7 Random walk . . . . .	309
3.7.1 Random walk in 1D . . . . .	310
3.7.2 Statistical considerations . . . . .	310
3.7.3 Playing around with some code . . . . .	312
3.7.4 Equivalence with diffusion . . . . .	315
3.7.5 Implementation of multiple walks . . . . .	316
3.7.6 Demonstration of multiple walks . . . . .	323
3.7.7 Ascii visualization of 1D random walk . . . . .	327
3.7.8 Random walk as a stochastic equation . . . . .	328
3.7.9 Random walk in 2D . . . . .	329
3.7.10 Random walk in any number of space dimensions . . . . .	330
3.7.11 Multiple random walks in any number of space dimensions . . . . .	331
3.8 Applications . . . . .	333
3.8.1 Diffusion of a substance . . . . .	333
3.8.2 Heat conduction . . . . .	335
3.8.3 Development of flow between two flat plates . . . . .	338

3.8.4 Flow in a straight tube . . . . .	339
3.8.5 Tribology: thin film fluid flow . . . . .	340
3.8.6 Propagation of electrical signals in the brain . . . . .	341
3.9 Exercises . . . . .	341
<b>4 Advection-dominated equations . . . . .</b>	<b>349</b>
4.1 One-dimensional time-dependent advection equations . . . . .	349
4.1.1 Simplest scheme: forward in time, centered in space . . . . .	350
4.1.2 Analysis of the scheme . . . . .	352
4.1.3 Leapfrog in time, centered differences in space . . . . .	353
4.1.4 Upwind differences in space . . . . .	355
4.1.5 A Crank-Nicolson discretization in time and centered differences in space . . . . .	356
4.1.6 The Lax-Wendroff method . . . . .	356
4.1.7 Analysis of dispersion relations . . . . .	357
4.2 One-dimensional stationary advection-diffusion equation . . . . .	359
4.3 Two-dimensional advection-diffusion equations . . . . .	359
4.4 Applications of advection equations . . . . .	359
<b>5 Staggered mesh discretization . . . . .</b>	<b>361</b>
5.1 Ordinary differential equations . . . . .	361
5.1.1 The Euler-Cromer scheme on a standard mesh . . . . .	361
5.1.2 The Euler-Cromer scheme on a staggered mesh . . . . .	362
5.1.3 Implementation of the scheme on a staggered mesh . . . . .	365
5.1.4 A staggered Euler-Cromer scheme for a generalized model . . . . .	367
5.2 Exercises . . . . .	368
5.3 Partial differential equations . . . . .	369
<b>6 Nonlinear problems . . . . .</b>	<b>371</b>
6.1 Introduction of basic concepts . . . . .	371
6.1.1 Linear versus nonlinear equations . . . . .	371
6.1.2 A simple model problem . . . . .	373
6.1.3 Linearization by explicit time discretization . . . . .	374
6.1.4 Exact solution of nonlinear algebraic equations . . . . .	375
6.1.5 Linearization . . . . .	376

6.1.6	Picard iteration . . . . .	376
6.1.7	Linearization by a geometric mean . . . . .	379
6.1.8	Newton's method . . . . .	380
6.1.9	Relaxation . . . . .	382
6.1.10	Implementation and experiments . . . . .	382
6.1.11	Generalization to a general nonlinear ODE . . . . .	385
6.1.12	Systems of ODEs . . . . .	388
6.2	Systems of nonlinear algebraic equations . . . . .	391
6.2.1	Picard iteration . . . . .	391
6.2.2	Newton's method . . . . .	392
6.2.3	Stopping criteria . . . . .	394
6.2.4	Example: A nonlinear ODE model from epidemiology	395
6.3	Linearization at the differential equation level . . . . .	397
6.3.1	Explicit time integration . . . . .	398
6.3.2	Backward Euler scheme and Picard iteration . . . . .	398
6.3.3	Backward Euler scheme and Newton's method . . . . .	399
6.3.4	Crank-Nicolson discretization . . . . .	402
6.4	Discretization of 1D stationary nonlinear differential equations	403
6.4.1	Finite difference discretization . . . . .	404
6.4.2	Solution of algebraic equations . . . . .	405
6.5	Multi-dimensional PDE problems . . . . .	410
6.5.1	Finite difference discretization . . . . .	410
6.5.2	Continuation methods . . . . .	413
6.6	Exercises . . . . .	414
<b>A</b>	<b>Useful formulas</b> . . . . .	423
A.1	Finite difference operator notation . . . . .	423
A.2	Truncation errors of finite difference approximations . . . . .	424
A.3	Finite differences of exponential functions . . . . .	425
A.4	Finite differences of $t^n$ . . . . .	426
A.4.1	Software . . . . .	427
<b>B</b>	<b>Truncation error analysis</b> . . . . .	429
B.1	Overview of truncation error analysis . . . . .	430
B.1.1	Abstract problem setting . . . . .	430
B.1.2	Error measures . . . . .	430

B.2	Truncation errors in finite difference formulas .....	432
B.2.1	Example: The backward difference for $u'(t)$ .....	432
B.2.2	Example: The forward difference for $u'(t)$ .....	433
B.2.3	Example: The central difference for $u'(t)$ .....	434
B.2.4	Overview of leading-order error terms in finite difference formulas .....	435
B.2.5	Software for computing truncation errors .....	436
B.3	Truncation errors in exponential decay ODE .....	438
B.3.1	Truncation error of the Forward Euler scheme .....	438
B.3.2	Truncation error of the Crank-Nicolson scheme .....	439
B.3.3	Truncation error of the $\theta$ -rule .....	439
B.3.4	Using symbolic software .....	440
B.3.5	Empirical verification of the truncation error .....	441
B.3.6	Increasing the accuracy by adding correction terms ..	445
B.3.7	Extension to variable coefficients .....	449
B.3.8	Exact solutions of the finite difference equations .....	450
B.3.9	Computing truncation errors in nonlinear problems ..	450
B.4	Truncation errors in vibration ODEs .....	451
B.4.1	Linear model without damping .....	451
B.4.2	Model with damping and nonlinearity .....	454
B.4.3	Extension to quadratic damping .....	456
B.4.4	The general model formulated as first-order ODEs ..	457
B.5	Truncation errors in wave equations .....	460
B.5.1	Linear wave equation in 1D .....	460
B.5.2	Finding correction terms .....	461
B.5.3	Extension to variable coefficients .....	462
B.5.4	1D wave equation on a staggered mesh .....	465
B.5.5	Linear wave equation in 2D/3D .....	465
B.6	Truncation errors in diffusion equations .....	466
B.6.1	Linear diffusion equation in 1D .....	466
B.6.2	Linear diffusion equation in 2D/3D .....	468
B.6.3	A nonlinear diffusion equation in 2D .....	468
B.7	Exercises .....	468
C	<b>Software engineering; wave equation model .....</b>	473
C.1	A 1D wave equation simulator .....	473
C.1.1	Mathematical model .....	473

C.1.2 Numerical discretization . . . . .	473
C.1.3 A solver function . . . . .	474
C.2 Saving large arrays in files . . . . .	477
C.2.1 Using <code>savez</code> to store arrays in files . . . . .	478
C.2.2 Using <code>joblib</code> to store arrays in files . . . . .	479
C.2.3 Using a hash to create a file or directory name . . . . .	480
C.3 Software for the 1D wave equation . . . . .	482
C.3.1 Making hash strings from input data . . . . .	483
C.3.2 Avoiding rerunning previously run cases . . . . .	484
C.3.3 Verification . . . . .	484
C.4 Programming the solver with classes . . . . .	485
C.4.1 Class Problem . . . . .	485
C.4.2 Class Mesh . . . . .	485
C.4.3 Class Function . . . . .	488
C.4.4 Class Solver . . . . .	491
C.5 Migrating loops to Cython . . . . .	491
C.5.1 Declaring variables and annotating the code . . . . .	492
C.5.2 Visual inspection of the C translation . . . . .	495
C.5.3 Building the extension module . . . . .	496
C.5.4 Calling the Cython function from Python . . . . .	497
C.6 Migrating loops to Fortran . . . . .	497
C.6.1 The Fortran subroutine . . . . .	498
C.6.2 Building the Fortran module with f2py . . . . .	499
C.6.3 How to avoid array copying . . . . .	501
C.7 Migrating loops to C via Cython . . . . .	503
C.7.1 Translating index pairs to single indices . . . . .	503
C.7.2 The complete C code . . . . .	504
C.7.3 The Cython interface file . . . . .	504
C.7.4 Building the extension module . . . . .	505
C.8 Migrating loops to C via f2py . . . . .	506
C.8.1 Migrating loops to C++ via f2py . . . . .	507
C.9 Exercises . . . . .	508
<b>References</b> . . . . .	511
<b>Index</b> . . . . .	513

# List of Exercises, Problems, and Projects

Problem 1.1: Use linear/quadratic functions for verification .....	59
Exercise 1.2: Show linear growth of the phase with time .....	61
Exercise 1.3: Improve the accuracy by adjusting the frequency...	61
Exercise 1.4: See if adaptive methods improve the phase error ...	61
Exercise 1.5: Use a Taylor polynomial to compute $u^1$ .....	62
Exercise 1.6: Find the minimal resolution of an oscillatory function	62
Exercise 1.7: Visualize the accuracy of finite differences for a cosine function .....	62
Exercise 1.8: Verify convergence rates of the error in energy .....	63
Exercise 1.9: Use linear/quadratic functions for verification .....	63
Exercise 1.10: Use an exact discrete solution for verification .....	63
Exercise 1.11: Use analytical solution for convergence rate tests..	64
Exercise 1.12: Investigate the amplitude errors of many solvers ..	64
Exercise 1.13: Minimize memory usage of a vibration solver .....	64
Exercise 1.14: Implement the solver via classes .....	65
Exercise 1.15: Interpret $[D_tD_t u]^n$ as a forward-backward difference	65
Exercise 1.16: Use a backward difference for the damping term ..	65
Exercise 1.17: Analysis of the Euler-Cromer scheme .....	66
Exercise 1.18: Simulate resonance .....	88
Exercise 1.19: Simulate oscillations of a sliding box .....	88
Exercise 1.20: Simulate a bouncing ball .....	88
Exercise 1.21: Simulate a simple pendulum .....	89
Exercise 1.22: Simulate an elastic pendulum .....	89
Exercise 1.23: Simulate an elastic pendulum with air resistance ..	90
Exercise 2.1: Simulate a standing wave .....	125

Exercise 2.2: Add storage of solution in a user action function . . . . .	126
Exercise 2.3: Use a class for the user action function . . . . .	127
Exercise 2.4: Compare several Courant numbers in one movie . . . . .	127
Project 2.5: Calculus with 1D mesh functions . . . . .	127
Exercise 2.6: Find the analytical solution to a damped wave equation . . . . .	153
Problem 2.7: Explore symmetry boundary conditions . . . . .	153
Exercise 2.8: Send pulse waves through a layered medium . . . . .	154
Exercise 2.9: Explain why numerical noise occurs . . . . .	154
Exercise 2.10: Investigate harmonic averaging in a 1D model . . . . .	155
Problem 2.11: Implement open boundary conditions . . . . .	155
Exercise 2.12: Implement periodic boundary conditions . . . . .	157
Exercise 2.13: Compare discretizations of a Neumann condition . . . . .	158
Exercise 2.14: Verification by a cubic polynomial in space . . . . .	159
Exercise 2.15: Check that a solution fulfills the discrete model . . . . .	188
Project 2.16: Calculus with 2D mesh functions . . . . .	188
Exercise 2.17: Implement Neumann conditions in 2D . . . . .	189
Exercise 2.18: Test the efficiency of compiled loops in 3D . . . . .	189
Exercise 2.19: Simulate waves on a non-homogeneous string . . . . .	205
Exercise 2.20: Simulate damped waves on a string . . . . .	205
Exercise 2.21: Simulate elastic waves in a rod . . . . .	205
Exercise 2.22: Simulate spherical waves . . . . .	206
Problem 2.23: Earthquake-generated tsunami over a subsea hill . . . . .	206
Problem 2.24: Earthquake-generated tsunami over a 3D hill . . . . .	209
Problem 2.25: Investigate Matplotlib for visualization . . . . .	210
Problem 2.26: Investigate visualization packages . . . . .	210
Problem 2.27: Implement loops in compiled languages . . . . .	210
Exercise 2.28: Simulate seismic waves in 2D . . . . .	211
Project 2.29: Model 3D acoustic waves in a room . . . . .	211
Project 2.30: Solve a 1D transport equation . . . . .	212
Problem 2.31: General analytical solution of a 1D damped wave equation . . . . .	216
Problem 2.32: General analytical solution of a 2D damped wave equation . . . . .	218
Exercise 3.1: Explore symmetry in a 1D problem . . . . .	258
Exercise 3.2: Investigate approximation errors from a $u_x = 0$ boundary condition . . . . .	258
Exercise 3.3: Experiment with open boundary conditions in 1D . . . . .	258
Exercise 3.4: Simulate a diffused Gaussian peak in 2D/3D . . . . .	260

Exercise 3.5: Examine stability of a diffusion model with a source term .....	260
Exercise 3.6: Stabilizing the Crank-Nicolson method by Rannacher time stepping .....	341
Project 3.7: Energy estimates for diffusion problems .....	342
Exercise 3.8: Splitting methods and preconditioning .....	345
Exercise 3.9: Oscillating surface temperature of the earth .....	345
Exercise 3.10: Oscillating and pulsating flow in tubes .....	346
Exercise 5.1: Use the forward-backward scheme with quadratic damping .....	368
Problem 6.1: Determine if equations are nonlinear or not .....	414
Exercise 6.2: Derive and investigate a generalized logistic model ..	414
Problem 6.3: Experience the behavior of Newton's method .....	415
Problem 6.4: Compute the Jacobian of a $2 \times 2$ system .....	416
Problem 6.5: Solve nonlinear equations arising from a vibration ODE .....	416
Exercise 6.6: Find the truncation error of arithmetic mean of products .....	417
Problem 6.7: Newton's method for linear problems .....	418
Exercise 6.8: Discretize a 1D problem with a nonlinear coefficient .....	418
Exercise 6.9: Linearize a 1D problem with a nonlinear coefficient .....	418
Problem 6.10: Finite differences for the 1D Bratu problem .....	419
Exercise 6.11: Discretize a nonlinear 1D heat conduction PDE by finite differences .....	420
Exercise 6.12: Differentiate a highly nonlinear term .....	420
Exercise 6.13: Crank-Nicolson for a nonlinear 3D diffusion equation .....	421
Exercise 6.14: Find the sparsity of the Jacobian .....	421
Problem 6.15: Investigate a 1D problem with a continuation method .....	421
Exercise B.1: Truncation error of a weighted mean .....	468
Exercise B.2: Simulate the error of a weighted mean .....	468
Exercise B.3: Verify a truncation error formula .....	468
Exercise B.4: Truncation error of the Backward Euler scheme ...	468
Exercise B.5: Empirical estimation of truncation errors .....	469
Exercise B.6: Correction term for a Backward Euler scheme .....	469
Exercise B.7: Verify the effect of correction terms .....	469
Exercise B.8: Truncation error of the Crank-Nicolson scheme ....	469
Exercise B.9: Truncation error of $u' = f(u, t)$ .....	470
Exercise B.10: Truncation error of $[D_t D_t u]^n$ .....	470
Exercise B.11: Investigate the impact of approximating $u'(0)$ ....	471
Exercise B.12: Investigate the accuracy of a simplified scheme ...	471

Exercise C.1: Make an improved `numpy.savez` function ..... 508

Vibration problems lead to differential equations with solutions that oscillate in time, typically in a damped or undamped sinusoidal fashion. Such solutions put certain demands on the numerical methods compared to other phenomena whose solutions are monotone or very smooth. Both the frequency and amplitude of the oscillations need to be accurately handled by the numerical schemes. Most of the reasoning and specific building blocks introduced in the forthcoming text can be reused to construct sound methods for partial differential equations of wave nature in multiple spatial dimensions.

**hpl 2:** Need to discuss errors also for the damped and nonlinear models. At least the frequency errors must be illustrated here as well and investigated numerically, either in text or exercises.

## 1.1 Finite difference discretization

Many of the numerical challenges faced when computing oscillatory solutions to ODEs and PDEs can be captured by the very simple ODE  $u'' + u = 0$ . This ODE is thus chosen as our starting point for method development, implementation, and analysis.

### 1.1.1 A basic model for vibrations

The simplest model of a vibrating mechanical system has the following form:

$$u'' + \omega^2 u = 0, \quad u(0) = I, \quad u'(0) = 0, \quad t \in (0, T]. \quad (1.1)$$

Here,  $\omega$  and  $I$  are given constants. Section 1.10.1 derives (1.1) from physical principles and explains what the constants mean.

The exact solution of (1.1) is

$$u(t) = I \cos(\omega t). \quad (1.2)$$

That is,  $u$  oscillates with constant amplitude  $I$  and angular frequency  $\omega$ . The corresponding period of oscillations (i.e., the time between two neighboring peaks in the cosine function) is  $P = 2\pi/\omega$ . The number of periods per second is  $f = \omega/(2\pi)$  and measured in the unit Hz. Both  $f$  and  $\omega$  are referred to as frequency, but  $\omega$  is more precisely named *angular frequency*, measured in rad/s.

In vibrating mechanical systems modeled by (1.1),  $u(t)$  very often represents a position or a displacement of a particular point in the system. The derivative  $u'(t)$  then has the interpretation of velocity, and  $u''(t)$  is the associated acceleration. The model (1.1) is not only applicable to vibrating mechanical systems, but also to oscillations in electrical circuits.

### 1.1.2 A centered finite difference scheme

To formulate a finite difference method for the model problem (1.1) we follow the four steps explained in Section 1.1.2 in [4].

**Step 1: Discretizing the domain.** The domain is discretized by introducing a uniformly partitioned time mesh. The points in the mesh are  $t_n = n\Delta t$ ,  $n = 0, 1, \dots, N_t$ , where  $\Delta t = T/N_t$  is the constant length of the time steps. We introduce a mesh function  $u^n$  for  $n = 0, 1, \dots, N_t$ , which approximates the exact solution at the mesh points. The mesh function will be computed from algebraic equations derived from the differential equation problem.

**Step 2: Fulfilling the equation at discrete time points.** The ODE is to be satisfied at each mesh point:

$$u''(t_n) + \omega^2 u(t_n) = 0, \quad n = 1, \dots, N_t. \quad (1.3)$$

**Step 3: Replacing derivatives by finite differences.** The derivative  $u''(t_n)$  is to be replaced by a finite difference approximation. A common second-order accurate approximation to the second-order derivative is

$$u''(t_n) \approx \frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2}. \quad (1.4)$$

Inserting (1.4) in (1.3) yields

$$\frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} = -\omega^2 u^n. \quad (1.5)$$

We also need to replace the derivative in the initial condition by a finite difference. Here we choose a centered difference, whose accuracy is similar to the centered difference we used for  $u''$ :

$$\frac{u^1 - u^{-1}}{2\Delta t} = 0. \quad (1.6)$$

**Step 4: Formulating a recursive algorithm.** To formulate the computational algorithm, we assume that we have already computed  $u^{n-1}$  and  $u^n$  such that  $u^{n+1}$  is the unknown value, which we can readily solve for:

$$u^{n+1} = 2u^n - u^{n-1} - \Delta t^2 \omega^2 u^n. \quad (1.7)$$

The computational algorithm is simply to apply (1.7) successively for  $n = 1, 2, \dots, N_t - 1$ . This numerical scheme sometimes goes under the name Störmer's method or Verlet integration.

**Computing the first step.** We observe that (1.7) cannot be used for  $n = 0$  since the computation of  $u^1$  then involves the undefined value  $u^{-1}$  at  $t = -\Delta t$ . The discretization of the initial condition then comes to our rescue: (1.6) implies  $u^{-1} = u^1$  and this relation can be combined with (1.7) for  $n = 0$  to yield a value for  $u^1$ :

$$u^1 = 2u^0 - u^1 - \Delta t^2 \omega^2 u^0,$$

which reduces to

$$u^1 = u^0 - \frac{1}{2} \Delta t^2 \omega^2 u^0. \quad (1.8)$$

Exercise 1.5 asks you to perform an alternative derivation and also to generalize the initial condition to  $u'(0) = V \neq 0$ .

**The computational algorithm.** The steps for solving (1.1) become

1.  $u^0 = I$
2. compute  $u^1$  from (1.8)
3. for  $n = 1, 2, \dots, N_t - 1$ : compute  $u^{n+1}$  from (1.7)

The algorithm is more precisely expressed directly in Python:

```
t = linspace(0, T, Nt+1) # mesh points in time
dt = t[1] - t[0]          # constant time step
u = zeros(Nt+1)           # solution

u[0] = I
u[1] = u[0] - 0.5*dt**2*w**2*u[0]
for n in range(1, Nt):
    u[n+1] = 2*u[n] - u[n-1] - dt**2*w**2*u[n]
```

### Remark on using w for $\omega$ in computer code

In the code, we use `w` as the symbol for  $\omega$ . The reason is that the authors prefer `w` for readability and comparison with the mathematical  $\omega$  instead of the full word `omega` as variable name.

**Operator notation.** We may write the scheme using a compact difference notation listed in Appendix A.1 (see also Section 1.1.8 in [4]). The difference (1.4) has the operator notation  $[D_tD_t u]^n$  such that we can write:

$$[D_tD_t u + \omega^2 u = 0]^n. \quad (1.9)$$

Note that  $[D_tD_t u]^n$  means applying a central difference with step  $\Delta t/2$  twice:

$$[D_t(D_t u)]^n = \frac{[D_t u]^{n+\frac{1}{2}} - [D_t u]^{n-\frac{1}{2}}}{\Delta t}$$

which is written out as

$$\frac{1}{\Delta t} \left( \frac{u^{n+1} - u^n}{\Delta t} - \frac{u^n - u^{n-1}}{\Delta t} \right) = \frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2}.$$

The discretization of initial conditions can in the operator notation be expressed as

$$[u = I]^0, \quad [D_{2t} u = 0]^0, \quad (1.10)$$

where the operator  $[D_{2t}u]^n$  is defined as

$$[D_{2t}u]^n = \frac{u^{n+1} - u^{n-1}}{2\Delta t}. \quad (1.11)$$

## 1.2 Implementation

### 1.2.1 Making a solver function

The algorithm from the previous section is readily translated to a complete Python function for computing and returning  $u^0, u^1, \dots, u^{N_t}$  and  $t_0, t_1, \dots, t_{N_t}$ , given the input  $I$ ,  $\omega$ ,  $\Delta t$ , and  $T$ :

```
import numpy as np
import matplotlib.pyplot as plt

def solver(I, w, dt, T):
    """
    Solve u'' + w**2*u = 0 for t in (0,T], u(0)=I and u'(0)=0,
    by a central finite difference method with time step dt.
    """
    dt = float(dt)
    Nt = int(round(T/dt))
    u = np.zeros(Nt+1)
    t = np.linspace(0, Nt*dt, Nt+1)

    u[0] = I
    u[1] = u[0] - 0.5*dt**2*w**2*u[0]
    for n in range(1, Nt):
        u[n+1] = 2*u[n] - u[n-1] - dt**2*w**2*u[n]
    return u, t
```

We have imported `numpy` and `matplotlib` under the names `np` and `plt`, respectively, as this is very common in the Python scientific computing community and a good programming habit (since we explicitly see where the different functions come from). An alternative is to do `from numpy import *` and a similar “import all” for `Matplotlib` to avoid the `np` and `plt` prefixes and make the code as close as possible to MATLAB. (See Section 5.1.4 in [4] for a discussion of the two types of import in Python.)

A function for plotting the numerical and the exact solution is also convenient to have:

```
def u_exact(t, I, w):
    return I*np.cos(w*t)

def visualize(u, t, I, w):
    plt.plot(t, u, 'r--o')
```

```
t_fine = np.linspace(0, t[-1], 1001) # very fine mesh for u_e
u_e = u_exact(t_fine, I, w)
plt.hold('on')
plt.plot(t_fine, u_e, 'b-')
plt.legend(['numerical', 'exact'], loc='upper left')
plt.xlabel('t')
plt.ylabel('u')
dt = t[1] - t[0]
plt.title('dt=%g % dt' % dt)
umin = 1.2*u.min(); umax = -umin
plt.axis([t[0], t[-1], umin, umax])
plt.savefig('tmp1.png'); plt.savefig('tmp1.pdf')
```

A corresponding main program calling these functions to simulate a given number of periods (`num_periods`) may take the form

```
I = 1
w = 2*pi
dt = 0.05
num_periods = 5
P = 2*pi/w # one period
T = P*num_periods
u, t = solver(I, w, dt, T)
visualize(u, t, I, w, dt)
```

Adjusting some of the input parameters via the command line can be handy. Here is a code segment using the `ArgumentParser` tool in the `argparse` module to define option value (`-option value`) pairs on the command line:

```
import argparse
parser = argparse.ArgumentParser()
parser.add_argument('--I', type=float, default=1.0)
parser.add_argument('--w', type=float, default=2*pi)
parser.add_argument('--dt', type=float, default=0.05)
parser.add_argument('--num_periods', type=int, default=5)
a = parser.parse_args()
I, w, dt, num_periods = a.I, a.w, a.dt, a.num_periods
```

Such parsing of the command line is explained in more detail in Section 5.2.3 in [4].

A typical execution goes like

---

	Terminal	
--	----------	--

---

Terminal> python vib\_undamped.py --num\_periods 20 --dt 0.1

---

**Computing  $u'$ .** In mechanical vibration applications one is often interested in computing the velocity  $v(t) = u'(t)$  after  $u(t)$  has been computed. This can be done by a central difference,

$$v(t_n) = u'(t_n) \approx v^n = \frac{u^{n+1} - u^{n-1}}{2\Delta t} = [D_{2t}u]^n. \quad (1.12)$$

This formula applies for all inner mesh points,  $n = 1, \dots, N_t - 1$ . For  $n = 0$ ,  $v(0)$  is given by the initial condition on  $u'(0)$ , and for  $n = N_t$  we can use a one-sided, backward difference:

$$v^n = [D_t^- u]^n = \frac{u^n - u^{n-1}}{\Delta t}.$$

Typical (scalar) code is

```
v = np.zeros_like(u) # or v = np.zeros(len(u))
# Use central difference for internal points
for i in range(1, len(u)-1):
    v[i] = (u[i+1] - u[i-1])/(2*dt)
# Use initial condition for u'(0) when i=0
v[0] = 0
# Use backward difference at the final mesh point
v[-1] = (u[-1] - u[-2])/dt
```

Since the loop is slow for large  $N_t$ , we can get rid of the loop by vectorizing the central difference. The above code segment goes as follows in its vectorized version (see Problem 1.2 in [4] for explanation of details):

```
v = np.zeros_like(u)
v[1:-1] = (u[2:] - u[:-2])/(2*dt) # central difference
v[0] = 0 # boundary condition u'(0)
v[-1] = (u[-1] - u[-2])/dt # backward difference
```

## 1.2.2 Verification

**Manual calculation.** The simplest type of verification, which is also instructive for understanding the algorithm, is to compute  $u^1$ ,  $u^2$ , and  $u^3$  with the aid of a calculator and make a function for comparing these results with those from the `solver` function. The `test_three_steps` function in the file `vib_undamped.py` shows the details of how we use the hand calculations to test the code:

```
def test_three_steps():
    from math import pi
    I = 1; w = 2*pi; dt = 0.1; T = 1
    u_by_hand = np.array([1.000000000000000,
                          0.802607911978213,
                          0.288358920740053])
    u, t = solver(I, w, dt, T)
    diff = np.abs(u_by_hand - u[:3]).max()
    tol = 1E-14
```

```
assert diff < tol
```

This function is a proper *test function*, compliant with the pytest and nose testing framework for Python code, because

- the function name begins with `test_`
- the function takes no arguments
- the test is formulated as a boolean condition and executed by `assert`

We shall in this book implement all software verification via such proper test functions, also known as unit testing. See Section 5.3.2 in [4] for more details on how to construct test functions and utilize nose or pytest for automatic execution of tests. Our recommendation is to use pytest and run all test functions in `vib_undamped.py` by

---

Terminal
=====

```
Terminal> py.test -s -v vib_undamped.py
=====
===== test session starts =====
platform linux2 -- Python 2.7.9 -- ...
collected 2 items

vib_undamped.py::test_three_steps PASSED
vib_undamped.py::test_convergence_rates PASSED

===== 2 passed in 0.19 seconds =====
```

---

**Testing very simple polynomial solutions.** Constructing test problems where the exact solution is constant or linear helps initial debugging and verification as one expects any reasonable numerical method to reproduce such solutions to machine precision. Second-order accurate methods will often also reproduce a quadratic solution. Here  $[D_t D_t t^2]^n = 2$ , which is the exact result. A solution  $u = t^2$  leads to  $u'' + \omega^2 u = 2 + (\omega t)^2 \neq 0$ . We must therefore add a source in the equation:  $u'' + \omega^2 u = f$  to allow a solution  $u = t^2$  for  $f = 2 + (\omega t)^2$ . By simple insertion we can show that the mesh function  $u^n = t_n^2$  is also a solution of the discrete equations. Problem 1.1 asks you to carry out all details to show that linear and quadratic solutions are solutions of the discrete equations. Such results are very useful for debugging and verification. You are strongly encouraged to do this problem now!

**Checking convergence rates.** Empirical computation of convergence rates yields a good method for verification. The method and its computational details are explained in detail in Section 3.1.6 in [4]. Readers not familiar with the concept should look up this reference before proceeding.

In the present problem, computing convergence rates means that we must

- perform  $m$  simulations, halving the time steps as:  $\Delta t_i = 2^{-i} \Delta t_0$ ,  $i = 0, \dots, m - 1$ ,
- compute the  $L^2$  norm of the error,  $E_i = \sqrt{\Delta t_i \sum_{n=0}^{N_t-1} (u^n - u_e(t_n))^2}$  in each case,
- estimate the convergence rates  $r_i$  based on two consecutive experiments  $(\Delta t_{i-1}, E_{i-1})$  and  $(\Delta t_i, E_i)$ , assuming  $E_i = C(\Delta t_i)^r$  and  $E_{i-1} = C(\Delta t_{i-1})^r$ . From these equations it follows that  $r = \ln(E_{i-1}/E_i)/\ln(\Delta t_{i-1}/\Delta t_i)$ . Since this  $r$  will vary with  $i$ , we equip it with an index and call it  $r_{i-1}$ , where  $i$  runs from 1 to  $m - 1$ .

The computed rates  $r_0, r_1, \dots, r_{m-2}$  hopefully converge to the number 2 in the present problem, because theory (from Section 1.4) shows that the error of the numerical method we use behaves like  $\Delta t^2$ . The convergence of the sequence of rates  $r_0, r_1, \dots, r_{m-2}$  demands that the time steps  $\Delta t_i$  are sufficiently small for the error model  $E_i = C(\Delta t_i)^r$  to be valid.

All the implementational details of computing the sequence  $r_0, r_1, \dots, r_{m-2}$  appear below.

```
def convergence_rates(m, solver_function, num_periods=8):
    """
    Return m-1 empirical estimates of the convergence rate
    based on m simulations, where the time step is halved
    for each simulation.
    solver_function(I, w, dt, T) solves each problem, where T
    is based on simulation for num_periods periods.
    """
    from math import pi
    w = 0.35; I = 0.3      # just chosen values
    P = 2*pi/w             # period
    dt = P/30               # 30 time step per period 2*pi/w
    T = P*num_periods

    dt_values = []
    E_values = []
    for i in range(m):
        u, t = solver_function(I, w, dt, T)
        u_e = u_exact(t, I, w)
        E = np.sqrt(dt*np.sum((u_e-u)**2))
        dt_values.append(dt)
        E_values.append(E)
        dt = dt/2

    r = [np.log(E_values[i-1]/E_values[i])/
         np.log(dt_values[i-1]/dt_values[i])
         for i in range(1, m, 1)]
    return r
```

The error analysis in Section 1.4 is quite detailed and results in  $r = 2$ , but other methods like truncation error analysis (see Appendix B.4.1) also points to  $r = 2$ , as well as the related quick reasoning that we have used a second-order accurate finite difference approximation  $[D_t D_t u]^n$  to the ODE and a second-order accurate finite difference formula for the initial condition for  $u'$ .

In the present problem, when  $\Delta t_0$  corresponds to 30 time steps per period, the returned  $\mathbf{r}$  list has all its values equal to 2.00 (if rounded to two decimals). This amazingly accurate result means that all  $\Delta t_i$  values are well into the asymptotic regime where the error model  $E_i = C(\Delta t_i)^r$  is valid.

We can now construct a proper test function that computes convergence rates and checks that the final (and usually the best) estimate is sufficiently close to 2. Here, a rough tolerance of 0.1 is enough. This unit test goes like

```
def test_convergence_rates():
    r = convergence_rates(m=5, solver_function=solver, num_periods=8)
    # Accept rate to 1 decimal place
    tol = 0.1
    assert abs(r[-1] - 2.0) < tol
```

The complete code appears in the file `vib_undamped.py`.

### 1.2.3 Scaled model

It is advantageous to use dimensionless variables in simulations, because fewer parameters need to be set. The present problem is made dimensionless by introducing dimensionless variables  $\bar{t} = t/t_c$  and  $\bar{u} = u/u_c$ , where  $t_c$  and  $u_c$  are characteristic scales for  $t$  and  $u$ , respectively. We refer to Section 2.2.1 in [5] for all details about this scaling.

The scaled ODE problem reads

$$\frac{u_c}{t_c^2} \frac{d^2 \bar{u}}{d\bar{t}^2} + u_c \bar{u} = 0, \quad u_c \bar{u}(0) = I, \quad \frac{u_c}{t_c} \frac{d\bar{u}}{d\bar{t}}(0) = 0.$$

A common choice is to take  $t_c$  as one period of the oscillations,  $t_c = 2\pi/w$ , and  $u_c = I$ . This gives the dimensionless model

$$\frac{d^2 \bar{u}}{\bar{t}^2} + 4\pi^2 \bar{u} = 0, \quad \bar{u}(0) = 1, \quad \bar{u}'(0) = 0. \quad (1.13)$$

Observe that there are no physical parameters in (1.13)! We can therefore perform a single numerical simulation  $\bar{u}(\bar{t})$  and afterwards recover any  $u(t; \omega, I)$  by

$$u(t; \omega, I) = u_c \bar{u}(t/t_c) = I \bar{u}(\omega t / (2\pi)).$$

We can easily check this assertion: the solution of the scaled problem is  $\bar{u}(\bar{t}) = \cos(2\pi\bar{t})$ . The formula for  $u$  in terms of  $\bar{u}$  gives  $u = I \cos(\omega t)$ , which is nothing but the solution of the original problem with dimensions.

The scaled model can be run by calling `solver(I=1, w=2*pi, dt, T)`. Each period is now 1 and  $T$  simply counts the number of periods. Choosing `dt` as `1./M` gives  $M$  time steps per period.

## 1.3 Long time simulations

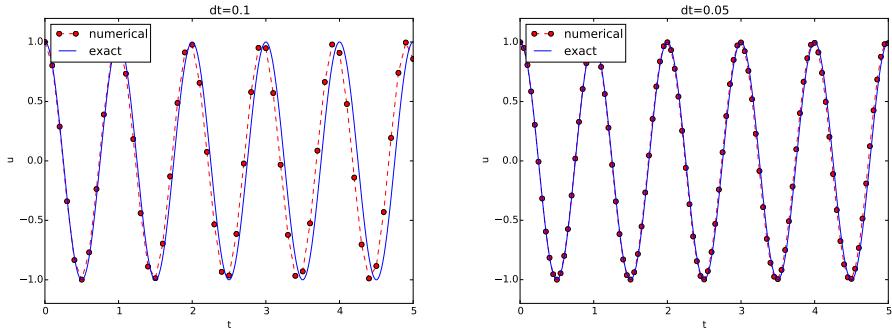
Figure 1.1 shows a comparison of the exact and numerical solution for the scaled model (1.13) with  $\Delta t = 0.1, 0.05$ . From the plot we make the following observations:

- The numerical solution seems to have correct amplitude.
- There is a angular frequency error which is reduced by reducing the time step.
- The total angular frequency error grows with time.

By angular frequency error we mean that the numerical angular frequency differs from the exact  $\omega$ . This is evident by looking at the peaks of the numerical solution: these have incorrect positions compared with the peaks of the exact cosine solution. The effect can be mathematical expressed by writing the numerical solution as  $I \cos \tilde{\omega} t$ , where  $\tilde{\omega}$  is not exactly equal to  $\omega$ . Later, we shall mathematically quantify this numerical angular frequency  $\tilde{\omega}$ .

### 1.3.1 Using a moving plot window

In vibration problems it is often of interest to investigate the system's behavior over long time intervals. Errors in the angular frequency accumulate and become more visible as time grows. We can investigate long time series by introducing a moving plot window that can move along with the  $p$  most recently computed periods of the solution. The `SciTools`



**Fig. 1.1** Effect of halving the time step.

package contains a convenient tool for this: `MovingPlotWindow`. Typing `pydoc scitools.MovingPlotWindow` shows a demo and a description of its use. The function below utilizes the moving plot window and is in fact called by the `main` function the `vib_undamped` module if the number of periods in the simulation exceeds 10.

```
def visualize_front(u, t, I, w, savefig=False, skip_frames=1):
    """
    Visualize u and the exact solution vs t, using a
    moving plot window and continuous drawing of the
    curves as they evolve in time.
    Makes it easy to plot very long time series.
    Plots are saved to files if savefig is True.
    Only each skip_frames-th plot is saved (e.g., if
    skip_frame=10, only each 10th plot is saved to file;
    this is convenient if plot files corresponding to
    different time steps are to be compared).
    """
    import scitools.std as st
    from scitools.MovingPlotWindow import MovingPlotWindow
    from math import pi

    # Remove all old plot files tmp_*.png
    import glob, os
    for filename in glob.glob('tmp_*.png'):
        os.remove(filename)

    P = 2*pi/w # one period
    umin = 1.2*u.min();  umax = -umin
    dt = t[1] - t[0]
    plot_manager = MovingPlotWindow(
        window_width=8*P,
        dt=dt,
        yaxis=[umin, umax],
        mode='continuous drawing')
    frame_counter = 0
```

```

for n in range(1,len(u)):
    if plot_manager.plot(n):
        s = plot_manager.first_index_in_plot
        st.plot(t[s:n+1], u[s:n+1], 'r-1',
                t[s:n+1], I*cos(w*t)[s:n+1], 'b-1',
                title='t=%6.3f' % t[n],
                axis=plot_manager.axis(),
                show=not savefig) # drop window if savefig
    if savefig and n % skip_frames == 0:
        filename = 'tmp_%04d.png' % frame_counter
        st.savefig(filename)
        print 'making plot file', filename, 'at t=%g' % t[n]
        frame_counter += 1
plot_manager.update(n)

```

We run the scaled problem (the default values for the command-line arguments `-I` and `-w` correspond to the scaled problem) for 40 periods with 20 time steps per period:

---

Terminal

---

```
Terminal> python vib_undamped.py --dt 0.05 --num_periods 40
```

---

The moving plot window is invoked, and we can follow the numerical and exact solutions as time progresses. From this demo we see that the angular frequency error is small in the beginning, but it becomes more prominent with time. A new run with  $\Delta t = 0.1$  (i.e., only 10 time steps per period) clearly shows that the phase errors become significant even earlier in the time series, deteriorating the solution further.

### 1.3.2 Making animations

**Producing standard video formats.** The `visualize_front` function stores all the plots in files whose names are numbered: `tmp_0000.png`, `tmp_0001.png`, `tmp_0002.png`, and so on. From these files we may make a movie. The Flash format is popular,

---

Terminal

---

```
Terminal> ffmpeg -r 12 -i tmp_%04d.png -c:v flv movie.flv
```

---

The `ffmpeg` program can be replaced by the `avconv` program in the above command if desired (but at the time of this writing it seems to be more momentum in the `ffmpeg` project). The `-r` option should come first and describes the number of frames per second in the movie. The `-i` option describes the name of the plot files. Other formats can be

generated by changing the video codec and equipping the video file with the right extension:

Format	Codec and filename
Flash	<code>-c:v flv movie.flv</code>
MP4	<code>-c:v libx264 movie.mp4</code>
WebM	<code>-c:v libvpx movie.webm</code>
Ogg	<code>-c:v libtheora movie.ogv</code>

The video file can be played by some video player like `vlc`, `mplayer`, `gxine`, or `totem`, e.g.,

---

Terminal> `vlc movie.webm`

---

A web page can also be used to play the movie. Today's standard is to use the HTML5 `video` tag:

```
<video autoplay loop controls
    width='640' height='365' preload='none'>
<source src='movie.webm' type='video/webm; codecs="vp8, vorbis"'>
</video>
```

Modern browsers do not support all of the video formats. MP4 is needed to successfully play the videos on Apple devices that use the Safari browser. WebM is the preferred format for Chrome, Opera, Firefox, and Internet Explorer v9+. Flash was a popular format, but older browsers that required Flash can play MP4. All browsers that work with Ogg can also work with WebM. This means that to have a video work in all browsers, the video should be available in the MP4 and WebM formats. The proper HTML code reads

```
<video autoplay loop controls
    width='640' height='365' preload='none'>
<source src='movie.mp4' type='video/mp4;
    codecs="avc1.42E01E, mp4a.40.2"'>
<source src='movie.webm' type='video/webm;
    codecs="vp8, vorbis"'>
</video>
```

The MP4 format should appear first to ensure that Apple devices will load the video correctly.

**Caution: number the plot files correctly**

To ensure that the individual plot frames are shown in correct order, it is important to number the files with zero-padded numbers (0000, 0001, 0002, etc.). The printf format %04d specifies an integer in a field of width 4, padded with zeros from the left. A simple Unix wildcard file specification like `tmp_*.png` will then list the frames in the right order. If the numbers in the filenames were not zero-padded, the frame `tmp_11.png` would appear before `tmp_2.png` in the movie.

**Paying PNG files in a web browser.** The `scitools movie` command can create a movie player for a set of PNG files such that a web browser can be used to watch the movie. This interface has the advantage that the speed of the movie can easily be controlled, a feature that scientists often appreciate. The command for creating an HTML with a player for a set of PNG files `tmp_*.png` goes like

---

Terminal

```
Terminal> scitools movie output_file=vib.html fps=4 tmp_*.png
```

---

The `fps` argument controls the speed of the movie (“frames per second”).

To watch the movie, load the video file `vib.html` into some browser, e.g.,

---

Terminal

```
Terminal> google-chrome vib.html # invoke web page
```

---

Clicking on `Start movie` to see the result. Moving this movie to some other place requires moving `vib.html` and all the *PNG files* `tmp_*.png`:

---

Terminal

```
Terminal> mkdir vib_dt0.1
Terminal> mv tmp_*.png vib_dt0.1
Terminal> mv vib.html vib_dt0.1/index.html
```

---

**Making animated GIF files.** The `convert` program from the ImageMagick software suite can be used to produce animated GIF files from a set of PNG files:

---

Terminal

```
Terminal> convert -delay 25 tmp_vib*.png tmp_vib.gif
```

---

The `-delay` option needs an argument of the delay between each frame, measured in 1/100 s, so 4 frames/s here gives 25/100 s delay. Note,

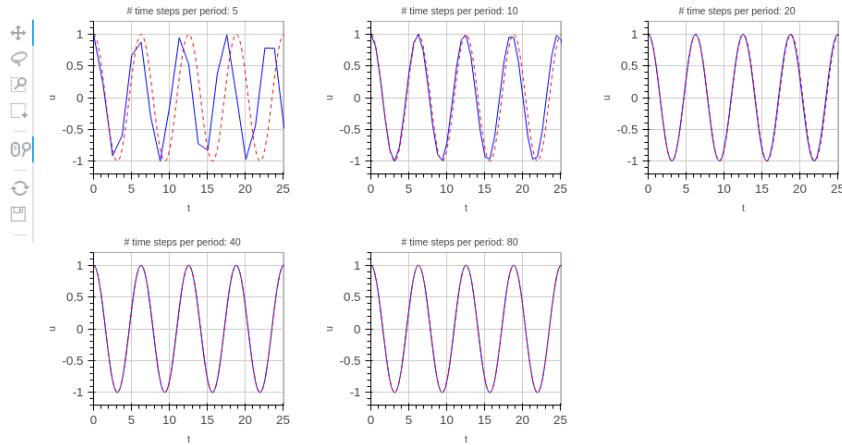
however, that in this particular example with  $\Delta t = 0.05$  and 40 periods, making an animated GIF file out of the large number of PNG files is a very heavy process and not considered feasible. Animated GIFs are best suited for animations with not so many frames and where you want to see each frame and play them slowly.

**hpl 3:** Combine two simulations side by side!

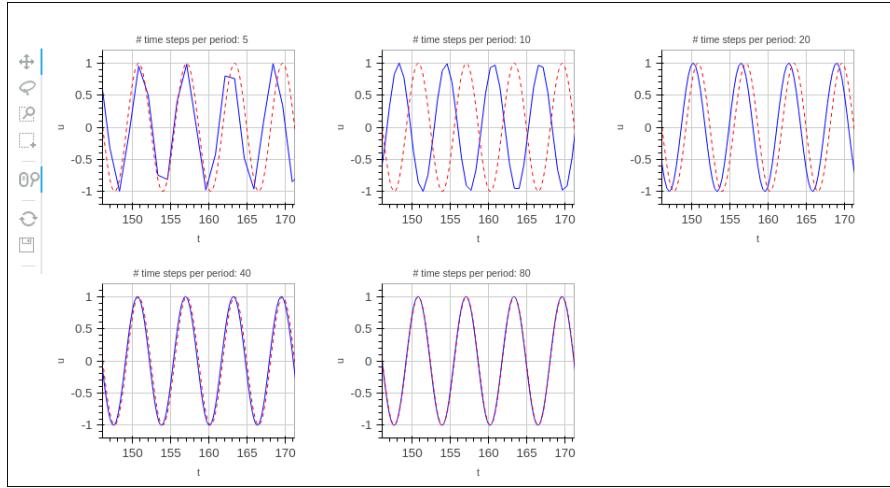
### 1.3.3 Using Bokeh to compare graphs

Instead of a moving plot frame, one can use tools that allows panning by the mouse. For example, we can show four periods of several signals in several plots and then scroll with the mouse through the rest of the simulation *simultaneously* in all the plot windows. The **Bokeh** plotting library offers such tools, but the plots must be displayed in a web browser. The documentation of Bokeh is excellent, so here we just show how the library can be used to compare a set of  $u$  curves corresponding to long time simulations. (By the way, the guidance to correct pronunciation of Bokeh in the [documentation](#) and on [Wikipedia](#) is not directly compatible with a [YouTube video](#)...).

Imagine we have performed experiments for a set of  $\Delta t$  values. We want each curve, together with the exact solution, to appear in a plot, and then arrange all plots in a grid-like fashion:



Furthermore, we want the axes to couple such that if we move into the future in one plot, all the other plots follows (note the displaced  $t$  axes!):



A function for creating a Bokeh plot, given a list of  $u$  arrays and corresponding  $t$  arrays, from different simulations, described compactly in a list of strings `legends`, is implemented below.

```
def bokeh_plot(u, t, legends, I, w, t_range, filename):
    """
    Make plots for u vs t using the Bokeh library.
    u and t are lists (several experiments can be compared).
    legends contain legend strings for the various u,t pairs.
    """
    if not isinstance(u, (list,tuple)):
        u = [u] # wrap in list
    if not isinstance(t, (list,tuple)):
        t = [t] # wrap in list
    if not isinstance(legends, (list,tuple)):
        legends = [legends] # wrap in list

    import bokeh.plotting as plt
    plt.output_file(filename, mode='cdn', title='Comparison')
    # Assume that all t arrays have the same range
    t_fine = np.linspace(0, t[0][-1], 1001) # fine mesh for u_e
    tools = 'pan,wheel_zoom,box_zoom,reset,' \
            'save,box_select,lasso_select'
    u_range = [-1.2*I, 1.2*I]
    font_size = '8pt'
    p = [] # list of plot objects
    # Make the first figure
    p_ = plt.figure(
        width=300, plot_height=250, title=legends[0],
        x_axis_label='t', y_axis_label='u',
        x_range=t_range, y_range=u_range, tools=tools,
        title_text_font_size=font_size)
    p_.xaxis.axis_label_text_font_size=font_size
    p_.yaxis.axis_label_text_font_size=font_size
```

```

p_.line(t[0], u[0], line_color='blue')
# Add exact solution
u_e = u_exact(t_fine, I, w)
p_.line(t_fine, u_e, line_color='red', line_dash='4 4')
p.append(p_)
# Make the rest of the figures and attach their axes to
# the first figure's axes
for i in range(1, len(t)):
    p_ = plt.figure(
        width=300, plot_height=250, title=legends[i],
        x_axis_label='t', y_axis_label='u',
        x_range=p[0].x_range, y_range=p[0].y_range, tools=tools,
        title_text_font_size=font_size)
    p_.xaxis.axis_label_text_font_size = font_size
    p_.yaxis.axis_label_text_font_size = font_size
    p_.line(t[i], u[i], line_color='blue')
    p_.line(t_fine, u_e, line_color='red', line_dash='4 4')
    p.append(p_)

# Arrange all plots in a grid with 3 plots per row
grid = [[]]
for i, p_ in enumerate(p):
    grid[-1].append(p_)
    if (i+1) % 3 == 0:
        # New row
        grid.append([])
plot = plt.gridplot(grid, toolbar_location='left')
plt.save(plot)
plt.show(plot)

```

A particular example using the `bokeh_plot` function appears below.

```

def demo_bokeh():
    """Solve a scaled ODE  $u'' + u = 0$ """
    from math import pi
    w = 1.0          # Scaled problem (frequency)
    P = 2*np.pi/w   # Period
    num_steps_per_period = [5, 10, 20, 40, 80]
    T = 40*P         # Simulation time: 40 periods
    u = []           # List of numerical solutions
    t = []           # List of corresponding meshes
    legends = []
    for n in num_steps_per_period:
        dt = P/n
        u_, t_ = solver(I=1, w=w, dt=dt, T=T)
        u.append(u_)
        t.append(t_)
        legends.append('# time steps per period: %d' % n)
    bokeh_plot(u, t, legends, I=1, w=w, t_range=[0, 4*P],
               filename='tmp.html')

```

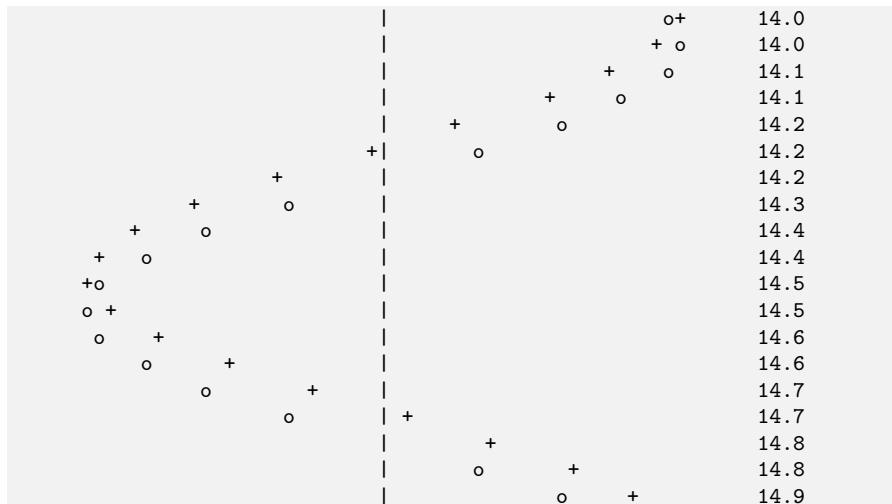
### 1.3.4 Using a line-by-line ascii plotter

Plotting functions vertically, line by line, in the terminal window using ascii characters only is a simple, fast, and convenient visualization technique for long time series. Note that the time axis then is positive downwards on the screen. The tool `scitools.avplotter.Plotter` makes it easy to create such plots:

```
def visualize_front_ascii(u, t, I, w, fps=10):
    """
    Plot u and the exact solution vs t line by line in a
    terminal window (only using ascii characters).
    Makes it easy to plot very long time series.
    """
    from scitools.avplotter import Plotter
    import time
    from math import pi
    P = 2*pi/w
    umin = 1.2*u.min();  umax = -umin

    p = Plotter(ymin=umin, ymax=umax, width=60, symbols='+o')
    for n in range(len(u)):
        print p.plot(t[n], u[n], I*cos(w*t[n])), \
            '%.1f' % (t[n]/P)
        time.sleep(1/float(fps))
```

The call `p.plot` returns a line of text, with the  $t$  axis marked and a symbol + for the first function ( $u$ ) and o for the second function (the exact solution). Here we append to this text a time counter reflecting how many periods the current time point corresponds to. A typical output ( $\omega = 2\pi$ ,  $\Delta t = 0.05$ ) looks like this:



		o	+	14.9
		o+		15.0

### 1.3.5 Empirical analysis of the solution

For oscillating functions like those in Figure 1.1 we may compute the amplitude and frequency (or period) empirically. That is, we run through the discrete solution points  $(t_n, u_n)$  and find all maxima and minima points. The distance between two consecutive maxima (or minima) points can be used as estimate of the local period, while half the difference between the  $u$  value at a maximum and a nearby minimum gives an estimate of the local amplitude.

The local maxima are the points where

$$u^{n-1} < u^n > u^{n+1}, \quad n = 1, \dots, N_t - 1, \quad (1.14)$$

and the local minima are recognized by

$$u^{n-1} > u^n < u^{n+1}, \quad n = 1, \dots, N_t - 1. \quad (1.15)$$

In computer code this becomes

```
def minmax(t, u):
    minima = []
    maxima = []
    for n in range(1, len(u)-1, 1):
        if u[n-1] > u[n] < u[n+1]:
            minima.append((t[n], u[n]))
        if u[n-1] < u[n] > u[n+1]:
            maxima.append((t[n], u[n]))
    return minima, maxima
```

Note that the two returned objects are lists of tuples.

Let  $(t_i, e_i)$ ,  $i = 0, \dots, M - 1$ , be the sequence of all the  $M$  maxima points, where  $t_i$  is the time value and  $e_i$  the corresponding  $u$  value. The local period can be defined as  $p_i = t_{i+1} - t_i$ . With Python syntax this reads

```
def periods(maxima):
    p = [extrema[n][0] - maxima[n-1][0]
         for n in range(1, len(maxima))]
    return np.array(p)
```

The list  $p$  created by a list comprehension is converted to an array since we probably want to compute with it, e.g., find the corresponding frequencies  $2\pi/p$ .

Having the minima and the maxima, the local amplitude can be calculated as the difference between two neighboring minimum and maximum points:

```
def amplitudes(minima, maxima):
    a = [(abs(maxima[n][1] - minima[n][1]))/2.0
          for n in range(min(len(minima), len(maxima)))]
    return np.array(a)
```

The code segments are found in the file `vib_empirical_analysis.py`.

Since  $a[i]$  and  $p[i]$  correspond to the  $i$ -th amplitude estimate and the  $i$ -th period estimate, respectively, it is most convenient to visualize the  $a$  and  $p$  values with the index  $i$  on the horizontal axis. (There is no unique time point associated with either of these estimate since values at two different time points were used in the computations.)

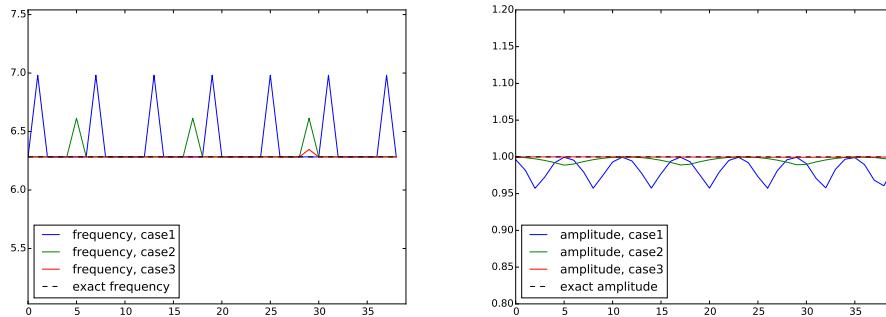
In the analysis of very long time series, it is advantageous to compute and plot  $p$  and  $a$  instead of  $u$  to get an impression of the development of the oscillations. Let us do this for the scaled problem and  $\Delta t = 0.1, 0.05, 0.01$ . A ready-made function

```
plot_empirical_freq_and_amplitude(u, t, I, w)
```

computes the empirical amplitudes and periods, and creates a plot where the amplitudes and angular frequencies are visualized together with the exact amplitude  $I$  and the exact angular frequency  $w$ . We can make a little program for creating the plot:

```
from vib_undamped import solver, plot_empirical_freq_and_amplitude
from math import pi
dt_values = [0.1, 0.05, 0.01]
u_cases = []
t_cases = []
for dt in dt_values:
    # Simulate scaled problem for 40 periods
    u, t = solver(I=1, w=2*pi, dt=dt, T=40)
    u_cases.append(u)
    t_cases.append(t)
plot_empirical_freq_and_amplitude(u_cases, t_cases, I=1, w=2*pi)
```

Figure 1.2 shows the result: we clearly see that lowering  $\Delta t$  improves the angular frequency significantly, while the amplitude seems to be more accurate. The lines with  $\Delta t = 0.01$ , corresponding to 100 steps per period, can hardly be distinguished from the exact values. The next section shows how we can get mathematical insight into why amplitudes are good and frequencies are more inaccurate.



**Fig. 1.2** Empirical amplitude and angular frequency for three cases of time steps.

## 1.4 Analysis of the numerical scheme

### 1.4.1 Deriving a solution of the numerical scheme

After having seen the phase error grow with time in the previous section, we shall now quantify this error through mathematical analysis. The key tool in the analysis will be to establish an exact solution of the discrete equations. The difference equation (1.7) has constant coefficients and is homogeneous. Such equations are known to have solutions on the form  $u^n = CA^n$ , where  $A$  is some number to be determined from the difference equation and  $C$  is found as the initial condition ( $C = I$ ). Recall that  $n$  in  $u^n$  is a superscript labeling the time level, while  $n$  in  $A^n$  is an exponent.

With oscillating functions as solutions, the algebra will be considerably simplified if we seek an  $A$  on the form

$$A = e^{i\tilde{\omega}\Delta t},$$

and solve for the numerical frequency  $\tilde{\omega}$  rather than  $A$ . Note that  $i = \sqrt{-1}$  is the imaginary unit. (Using a complex exponential function gives simpler arithmetics than working with a sine or cosine function.) We have

$$A^n = e^{i\tilde{\omega}\Delta t n} = e^{i\tilde{\omega}t_n} = \cos(\tilde{\omega}t_n) + i \sin(\tilde{\omega}t_n).$$

The physically relevant numerical solution can be taken as the real part of this complex expression.

The calculations go as

$$\begin{aligned}
[D_t D_t u]^n &= \frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} \\
&= I \frac{A^{n+1} - 2A^n + A^{n-1}}{\Delta t^2} \\
&= \frac{I}{\Delta t^2} (e^{i\tilde{\omega}(t_n + \Delta t)} - 2e^{i\tilde{\omega}t_n} + e^{i\tilde{\omega}(t_n - \Delta t)}) \\
&= I e^{i\tilde{\omega}t_n} \frac{1}{\Delta t^2} (e^{i\tilde{\omega}\Delta t} + e^{i\tilde{\omega}(-\Delta t)} - 2) \\
&= I e^{i\tilde{\omega}t_n} \frac{2}{\Delta t^2} (\cosh(i\tilde{\omega}\Delta t) - 1) \\
&= I e^{i\tilde{\omega}t_n} \frac{2}{\Delta t^2} (\cos(\tilde{\omega}\Delta t) - 1) \\
&= -I e^{i\tilde{\omega}t_n} \frac{4}{\Delta t^2} \sin^2\left(\frac{\tilde{\omega}\Delta t}{2}\right)
\end{aligned}$$

The last line follows from the relation  $\cos x - 1 = -2\sin^2(x/2)$  (try `cos(x)-1` in [wolframalpha.com](#) to see the formula).

The scheme (1.7) with  $u^n = I e^{i\tilde{\omega}\Delta t n}$  inserted now gives

$$-I e^{i\tilde{\omega}t_n} \frac{4}{\Delta t^2} \sin^2\left(\frac{\tilde{\omega}\Delta t}{2}\right) + \omega^2 I e^{i\tilde{\omega}t_n} = 0, \quad (1.16)$$

which after dividing by  $I e^{i\tilde{\omega}t_n}$  results in

$$\frac{4}{\Delta t^2} \sin^2\left(\frac{\tilde{\omega}\Delta t}{2}\right) = \omega^2. \quad (1.17)$$

The first step in solving for the unknown  $\tilde{\omega}$  is

$$\sin^2\left(\frac{\tilde{\omega}\Delta t}{2}\right) = \left(\frac{\omega\Delta t}{2}\right)^2.$$

Then, taking the square root, applying the inverse sine function, and multiplying by  $2/\Delta t$ , results in

$$\tilde{\omega} = \pm \frac{2}{\Delta t} \sin^{-1}\left(\frac{\omega\Delta t}{2}\right). \quad (1.18)$$

The first observation of (1.18) tells that there is a phase error since the numerical frequency  $\tilde{\omega}$  never equals the exact frequency  $\omega$ . But how good is the approximation (1.18)? That is, what is the error  $\omega - \tilde{\omega}$  or  $\tilde{\omega}/\omega$ ? Taylor series expansion for small  $\Delta t$  may give an expression that is easier to understand than the complicated function in (1.18):

```
>>> from sympy import *
```

```
>>> dt, w = symbols('dt w')
>>> w_tilde_e = 2/dt*asin(w*dt/2)
>>> w_tilde_series = w_tilde_e.series(dt, 0, 4)
>>> print w_tilde_series
w + dt**2*w**3/24 + O(dt**4)
```

This means that

$$\tilde{\omega} = \omega \left( 1 + \frac{1}{24} \omega^2 \Delta t^2 \right) + \mathcal{O}(\Delta t^4). \quad (1.19)$$

The error in the numerical frequency is of second-order in  $\Delta t$ , and the error vanishes as  $\Delta t \rightarrow 0$ . We see that  $\tilde{\omega} > \omega$  since the term  $\omega^3 \Delta t^2 / 24 > 0$  and this is by far the biggest term in the series expansion for small  $\omega \Delta t$ . A numerical frequency that is too large gives an oscillating curve that oscillates too fast and therefore “lags behind” the exact oscillations, a feature that can be seen in the left plot in Figure 1.1.

Figure 1.3 plots the discrete frequency (1.18) and its approximation (1.19) for  $\omega = 1$  (based on the program `vib_plot_freq.py`). Although  $\tilde{\omega}$  is a function of  $\Delta t$  in (1.19), it is misleading to think of  $\Delta t$  as the important discretization parameter. It is the product  $\omega \Delta t$  that is the key discretization parameter. This quantity reflects the *number of time steps per period* of the oscillations. To see this, we set  $P = N_P \Delta t$ , where  $P$  is the length of a period, and  $N_P$  is the number of time steps during a period. Since  $P$  and  $\omega$  are related by  $P = 2\pi/\omega$ , we get that  $\omega \Delta t = 2\pi/N_P$ , which shows that  $\omega \Delta t$  is directly related to  $N_P$ .

The plot shows that at least  $N_P \sim 25 - 30$  points per period are necessary for reasonable accuracy, but this depends on the length of the simulation ( $T$ ) as the total phase error due to the frequency error grows linearly with time (see Exercise 1.2).

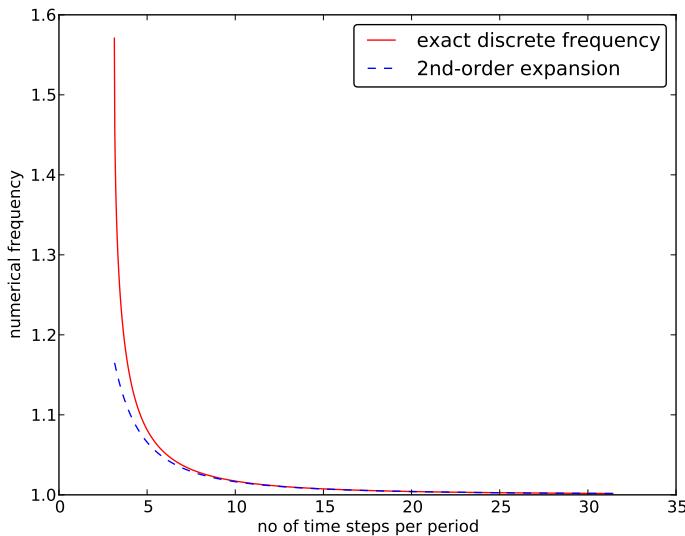
### 1.4.2 Exact discrete solution

Perhaps more important than the  $\tilde{\omega} = \omega + \mathcal{O}(\Delta t^2)$  result found above is the fact that we have an exact discrete solution of the problem:

$$u^n = I \cos(\tilde{\omega} n \Delta t), \quad \tilde{\omega} = \frac{2}{\Delta t} \sin^{-1} \left( \frac{\omega \Delta t}{2} \right). \quad (1.20)$$

We can then compute the error mesh function

$$e^n = u_e(t_n) - u^n = I \cos(\omega n \Delta t) - I \cos(\tilde{\omega} n \Delta t). \quad (1.21)$$



**Fig. 1.3** Exact discrete frequency and its second-order series expansion.

From the formula  $\cos 2x - \cos 2y = -2 \sin(x-y) \sin(x+y)$  we can rewrite  $e^n$  so the expression is easier to interpret:

$$e^n = -2I \sin\left(t \frac{1}{2} (\omega - \tilde{\omega})\right) \sin\left(t \frac{1}{2} (\omega + \tilde{\omega})\right). \quad (1.22)$$

The error mesh function is ideal for verification purposes and you are strongly encouraged to make a test based on (1.20) by doing Exercise 1.10.

### 1.4.3 Convergence

We can use (1.19), (1.21), or (1.22) to show *convergence* of the numerical scheme, i.e.,  $e^n \rightarrow 0$  as  $\Delta t \rightarrow 0$ , which implies that the numerical solution approaches the exact solution as  $\Delta t$  approaches to zero. We have that

$$\lim_{\Delta t \rightarrow 0} \tilde{\omega} = \lim_{\Delta t \rightarrow 0} \frac{2}{\Delta t} \sin^{-1}\left(\frac{\omega \Delta t}{2}\right) = \omega,$$

by L'Hopital's rule. This result could also been computed [WolframAlpha](#), or we could use the limit functionality in `sympy`:

```
>>> import sympy as sym
>>> dt, w = sym.symbols('x w')
```

```
>>> sym.limit((2/dt)*sym.asin(w*dt/2), dt, 0, dir='+')
w
```

Also (1.19) can be used to establish this result that  $\tilde{\omega} \rightarrow \omega$ . It then follows from the expression(s) for  $e^n$  that  $e^n \rightarrow 0$ .

#### 1.4.4 The global error

To achieve more analytical insight into the nature of the global error, we can Taylor expand the error mesh function (1.21). Since  $\tilde{\omega}$  in (1.18) contains  $\Delta t$  in the denominator we use the series expansion for  $\tilde{\omega}$  inside the cosine function. A relevant `sympy` session is

```
>>> from sympy import *
>>> dt, w, t = symbols('dt w t')
>>> w_tilde_e = 2/dt*asin(w*dt/2)
>>> w_tilde_series = w_tilde_e.series(dt, 0, 4)
>>> w_tilde_series
w + dt**2*w**3/24 + O(dt**4)
```

Series expansions in `sympy` have the inconvenient `O()` term that prevents further calculations with the series. We can use the `removeO()` command to get rid of the `O()` term:

```
>>> w_tilde_series = w_tilde_series.removeO()
>>> w_tilde_series
dt**2*w**3/24 + w
```

Using this `w_tilde_series` expression for  $\tilde{\omega}$  in (1.21), dropping  $I$  (which is a common factor), and performing a series expansion of the error yields

```
>>> error = cos(w*t) - cos(w_tilde_series*t)
>>> error.series(dt, 0, 6)
dt**2*t*w**3*sin(t*w)/24 + dt**4*t**2*w**6*cos(t*w)/1152 + O(dt**6)
```

Since we are mainly interested in the leading-order term in such expansions (the term with lowest power in  $\Delta t$ , which goes most slowly to zero), we use the `.as_leading_term(dt)` construction to pick out this term:

```
>>> error.series(dt, 0, 6).as_leading_term(dt)
dt**2*t*w**3*sin(t*w)/24
```

The last result means that the leading order global (true) error at a point  $t$  is proportional to  $\omega^3 t \Delta t^2$ . Considering only the discrete  $t_n$  values for  $t$ ,  $t_n$  is related to  $\Delta t$  through  $t_n = n\Delta t$ . The factor  $\sin(\omega t)$  can at most be 1, so we use this value to bound the leading-order expression to its maximum value

$$e^n = \frac{1}{24} n \omega^3 \Delta t^3.$$

This is the dominating term of the error *at a point*.

We are interested in the accumulated global error, which can be taken as the  $\ell^2$  norm of  $e^n$ . The norm is simply computed by summing contributions from all mesh points:

$$\|e^n\|_{\ell^2}^2 = \Delta t \sum_{n=0}^{N_t} \frac{1}{24^2} n^2 \omega^6 \Delta t^6 = \frac{1}{24^2} \omega^6 \Delta t^7 \sum_{n=0}^{N_t} n^2.$$

The sum  $\sum_{n=0}^{N_t} n^2$  is approximately equal to  $\frac{1}{3} N_t^3$ . Replacing  $N_t$  by  $T/\Delta t$  and taking the square root gives the expression

$$\|e^n\|_{\ell^2} = \frac{1}{24} \sqrt{\frac{T^3}{3}} \omega^3 \Delta t^2.$$

This is our expression for the global (or integrated) error. The main result from this expression is that the global error is proportional to  $\Delta t^2$ .

## 1.4.5 Stability

Looking at (1.20), it appears that the numerical solution has constant and correct amplitude, but an error in the angular frequency. A constant amplitude is not necessarily the case, however! To see this, note that if only  $\Delta t$  is large enough, the magnitude of the argument to  $\sin^{-1}$  in (1.18) may be larger than 1, i.e.,  $\omega \Delta t / 2 > 1$ . In this case,  $\sin^{-1}(\omega \Delta t / 2)$  has a complex value and therefore  $\tilde{\omega}$  becomes complex. Type, for example, `asin(x)` in [wolframalpha.com](#) to see basic properties of  $\sin^{-1}(x)$ .

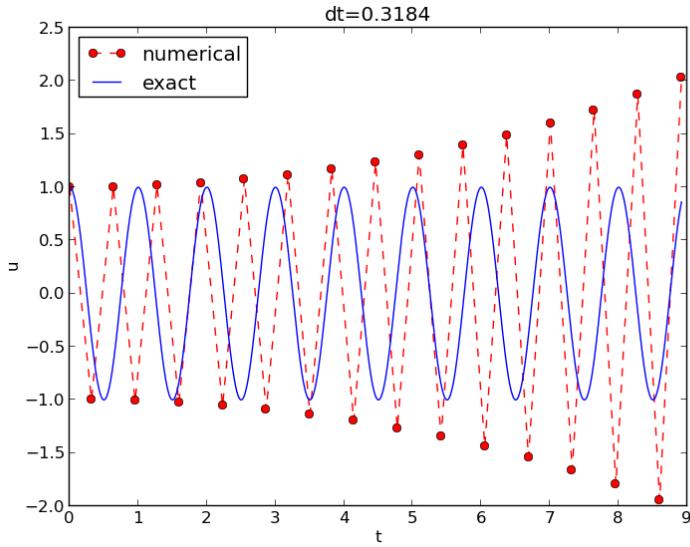
A complex  $\tilde{\omega}$  can be written  $\tilde{\omega} = \tilde{\omega}_r + i \tilde{\omega}_i$ . Since  $\sin^{-1}(x)$  has a *negative* imaginary part for  $x > 1$ ,  $\tilde{\omega}_i < 0$ , which means that  $e^{i\tilde{\omega}t} = e^{-\tilde{\omega}_i t} e^{i\tilde{\omega}_r t}$  will lead to exponential growth in time because  $e^{-\tilde{\omega}_i t}$  with  $\tilde{\omega}_i < 0$  has a positive exponent.

### Stability criterion

We do not tolerate growth in the amplitude since such growth is not present in the exact solution. Therefore, we must impose a *stability criterion* so that the argument in the inverse sine function leads to real and not complex values of  $\tilde{\omega}$ . The stability criterion reads

$$\frac{\omega \Delta t}{2} \leq 1 \quad \Rightarrow \quad \Delta t \leq \frac{2}{\omega}. \quad (1.23)$$

With  $\omega = 2\pi$ ,  $\Delta t > \pi^{-1} = 0.3183098861837907$  will give growing solutions. Figure 1.4 displays what happens when  $\Delta t = 0.3184$ , which is slightly above the critical value:  $\Delta t = \pi^{-1} + 9.01 \cdot 10^{-5}$ .



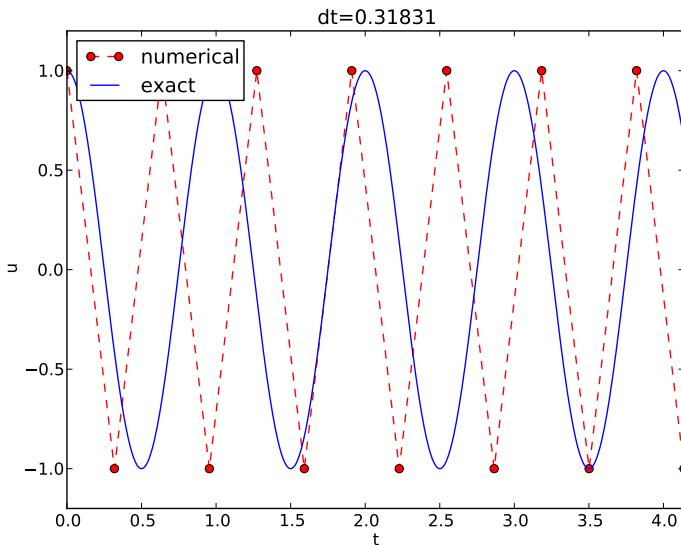
**Fig. 1.4** Growing, unstable solution because of a time step slightly beyond the stability limit.

### 1.4.6 About the accuracy at the stability limit

An interesting question is whether the stability condition  $\Delta t < 2/\omega$  is unfortunate, or more precisely: would it be meaningful to take larger time steps to speed up computations? The answer is a clear no. At the stability limit, we have that  $\sin^{-1} \omega \Delta t / 2 = \sin^{-1} 1 = \pi/2$ , and therefore  $\tilde{\omega} = \pi / \Delta t$ . (Note that the approximate formula (1.19) is very inaccurate for this value of  $\Delta t$  as it predicts  $\tilde{\omega} = 2.34/\pi$ , which is a 25 percent reduction.) The corresponding period of the numerical solution is  $\tilde{P} = 2\pi / \tilde{\omega} = 2\Delta t$ , which means that there is just one time step  $\Delta t$

between a peak (maximum) and a [through](#) (minimum) in the numerical solution. This is the shortest possible wave that can be represented in the mesh! In other words, it is not meaningful to use a larger time step than the stability limit.

Also, the error in angular frequency when  $\Delta t = 2/\omega$  is severe: Figure 1.5 shows a comparison of the numerical and analytical solution with  $\omega = 2\pi$  and  $\Delta t = 2/\omega = \pi^{-1}$ . Already after one period, the numerical solution has a through while the exact solution has a peak (!). The error in frequency when  $\Delta t$  is at the stability limit becomes  $\omega - \tilde{\omega} = \omega(1 - \pi/2) \approx -0.57\omega$ . The corresponding error in the period is  $P - \tilde{P} \approx 0.36P$ . The error after  $m$  periods is then  $0.36mP$ . This error has reached half a period when  $m = 1/(2 \cdot 0.36) \approx 1.38$ , which theoretically confirms the observations in Figure 1.5 that the numerical solution is a through ahead of a peak already after one and a half period. Consequently,  $\Delta t$  should be chosen much less than the stability limit to achieve meaningful numerical computations.



**Fig. 1.5** Numerical solution with  $\Delta t$  exactly at the stability limit.

### Summary

From the accuracy and stability analysis we can draw three important conclusions:

1. The key parameter in the formulas is  $p = \omega\Delta t$ . The period of oscillations is  $P = 2\pi/\omega$ , and the number of time steps per period is  $N_P = P/\Delta t$ . Therefore,  $p = \omega\Delta t = 2\pi/N_P$ , showing that the critical parameter is the number of time steps per period. The smallest possible  $N_P$  is 2, showing that  $p \in (0, \pi]$ .
2. Provided  $p \leq 2$ , the amplitude of the numerical solution is constant.
3. The ratio of the numerical angular frequency and the exact one is  $\tilde{\omega}/\omega \approx 1 + \frac{1}{24}p^2$ . The error  $\frac{1}{24}p^2$  leads to wrongly displaced peaks of the numerical solution, and the error in peak location grows linearly with time (see Exercise 1.2).

## 1.5 Alternative schemes based on 1st-order equations

A standard technique for solving second-order ODEs is to rewrite them as a system of first-order ODEs and then choose a solution strategy from the vast collection of methods for first-order ODE systems. Given the second-order ODE problem

$$u'' + \omega^2 u = 0, \quad u(0) = I, \quad u'(0) = 0,$$

we introduce the auxiliary variable  $v = u'$  and express the ODE problem in terms of first-order derivatives of  $u$  and  $v$ :

$$u' = v, \tag{1.24}$$

$$v' = -\omega^2 u. \tag{1.25}$$

The initial conditions become  $u(0) = I$  and  $v(0) = 0$ .

### 1.5.1 The Forward Euler scheme

A Forward Euler approximation to our  $2 \times 2$  system of ODEs (1.24)-(1.25) becomes

$$[D_t^+ u = v]^n, [D_t^+ v = -\omega^2 u]^n, \quad (1.26)$$

or written out,

$$u^{n+1} = u^n + \Delta t v^n, \quad (1.27)$$

$$v^{n+1} = v^n - \Delta t \omega^2 u^n. \quad (1.28)$$

Let us briefly compare this Forward Euler method with the centered difference scheme for the second-order differential equation. We have from (1.27) and (1.28) applied at levels  $n$  and  $n-1$  that

$$u^{n+1} = u^n + \Delta t v^n = u^n + \Delta t (v^{n-1} - \Delta t \omega^2 u^{n-1}).$$

Since from (1.27)

$$v^{n-1} = \frac{1}{\Delta t} (u^n - u^{n-1}),$$

it follows that

$$u^{n+1} = 2u^n - u^{n-1} - \Delta t^2 \omega^2 u^{n-1},$$

which is very close to the centered difference scheme, but the last term is evaluated at  $t_{n-1}$  instead of  $t_n$ . Rewriting, so that  $\Delta t^2 \omega^2 u^{n-1}$  appears alone on the right-hand side, and then dividing by  $\Delta t^2$ , the new left-hand side is an approximation to  $u''$  at  $t_n$ , while the right-hand side is sampled at  $t_{n-1}$ . All terms should be sampled at the same mesh point, so using  $\omega^2 u^{n-1}$  instead of  $\omega^2 u^n$  points to a kind of mathematical error in the derivation of the scheme. This error turns out to be rather crucial for the accuracy of the Forward Euler method applied to vibration problems (Section 1.5.4 has examples).

The reasoning above does not imply that the Forward Euler scheme is not correct, but more that it is almost equivalent to a second-order accurate scheme for the second-order ODE formulation, and that the error committed has to do with a wrong sampling point.

### 1.5.2 The Backward Euler scheme

A Backward Euler approximation to the ODE system is equally easy to write up in the operator notation:

$$[D_t^- u = v]^{n+1}, \quad (1.29)$$

$$[D_t^- v = -\omega u]^{n+1}. \quad (1.30)$$

This becomes a coupled system for  $u^{n+1}$  and  $v^{n+1}$ :

$$u^{n+1} - \Delta t v^{n+1} = u^n, \quad (1.31)$$

$$v^{n+1} + \Delta t \omega^2 u^{n+1} = v^n. \quad (1.32)$$

We can compare (1.31)-(1.32) with the centered scheme (1.7) for the second-order differential equation. To this end, we eliminate  $v^{n+1}$  in (1.31) using (1.32) solved with respect to  $v^{n+1}$ . Thereafter, we eliminate  $v^n$  using (1.31) solved with respect to  $v^{n+1}$  and also replacing  $n+1$  by  $n$  and  $n$  by  $n-1$ . The resulting equation involving only  $u^{n+1}$ ,  $u^n$ , and  $u^{n-1}$  can be ordered as

$$\frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} = -\omega^2 u^{n+1},$$

which has almost the same form as the centered scheme for the second-order differential equation, but the right-hand side is evaluated at  $u^{n+1}$  and not  $u^n$ . This inconsistent sampling of terms has a dramatic effect on the numerical solution, as we demonstrate in Section 1.5.4.

### 1.5.3 The Crank-Nicolson scheme

The Crank-Nicolson scheme takes this form in the operator notation:

$$[D_t u = \bar{v}^t]^{n+\frac{1}{2}}, \quad (1.33)$$

$$[D_t v = -\omega^2 \bar{u}^t]^{n+\frac{1}{2}}. \quad (1.34)$$

Writing the equations out and rearranging terms, shows that this is also a coupled system of two linear equations at each time level:

$$u^{n+1} - \frac{1}{2}\Delta t v^{n+1} = u^n + \frac{1}{2}\Delta t v^n, \quad (1.35)$$

$$v^{n+1} + \frac{1}{2}\Delta t \omega^2 u^{n+1} = v^n - \frac{1}{2}\Delta t \omega^2 u^n. \quad (1.36)$$

We may compare also this scheme to the centered discretization of the second-order ODE. It turns out that the Crank-Nicolson scheme is equivalent to the discretization

$$\frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} = -\omega^2 \frac{1}{4}(u^{n+1} + 2u^n + u^{n-1}) = -\omega^2 u^n + \mathcal{O}(\Delta t^2). \quad (1.37)$$

That is, the Crank-Nicolson is equivalent to (1.7) for the second-order ODE, apart from an extra term of size  $\Delta t^2$ , but this is an error of the same order as in the finite difference approximation on the left-hand side of the equation anyway. The fact that the Crank-Nicolson scheme is so close to (1.7) makes it a much better method than the Forward or Backward Euler methods for vibration problems, as will be illustrated in Section 1.5.4.

Deriving (1.37) is a bit tricky. We start with rewriting the Crank-Nicolson equations as follows

$$u^{n+1} - u^n = \frac{1}{2}\Delta t(v^{n+1} + v^n), \quad (1.38)$$

$$v^{n+1} = v^n - \frac{1}{2}\Delta t \omega^2(u^{n+1} + u^n), \quad (1.39)$$

and add the latter at the previous time level as well:

$$v^n = v^{n-1} - \frac{1}{2}\Delta t \omega^2(u^n + u^{n-1}) \quad (1.40)$$

We can also rewrite (1.38) at the previous time level as

$$v^n + v^{n-1} = \frac{2}{\Delta t}(u^n - u^{n-1}). \quad (1.41)$$

Inserting (1.39) for  $v^{n+1}$  in (1.38) and (1.40) for  $v^n$  in (1.38) yields after some reordering:

$$u^{n+1} - u^n = \frac{1}{2}(-\frac{1}{2}\Delta t \omega^2(u^{n+1} + 2u^n + u^{n-1}) + v^n + v^{n-1}).$$

Now,  $v^n + v^{n-1}$  can be eliminated by means of (1.41). The result becomes

$$u^{n+1} - 2u^n + u^{n-1} = -\Delta t^2 \omega^2 \frac{1}{4} (u^{n+1} + 2u^n + u^{n-1}). \quad (1.42)$$

It can be shown that

$$\frac{1}{4} (u^{n+1} + 2u^n + u^{n-1}) \approx u^n + \mathcal{O}(\Delta t^2),$$

meaning that (1.42) is an approximation to the centered scheme (1.7) for the second-order ODE where the sampling error in the term  $\Delta t^2 \omega^2 u^n$  is of the same order as the approximation errors in the finite differences, i.e.,  $\mathcal{O}(\Delta t^2)$ . The Crank-Nicolson scheme written as (1.42) therefore has consistent sampling of all terms at the same time point  $t_n$ .

### 1.5.4 Comparison of schemes

We can easily compare methods like the ones above (and many more!) with the aid of the `Odespy` package. Below is a sketch of the code.

```
import odespy
import numpy as np

def f(u, t, w=1):
    # v, u numbering for EulerCromer to work well
    v, u = u # u is array of length 2 holding our [v, u]
    return [-w**2*u, v]

def run_solvers_and_plot(solvers, timesteps_per_period=20,
                        num_periods=1, I=1, w=2*np.pi):
    P = 2*np.pi/w # duration of one period
    dt = P/timesteps_per_period
    Nt = num_periods*timesteps_per_period
    T = Nt*dt
    t_mesh = np.linspace(0, T, Nt+1)

    legends = []
    for solver in solvers:
        solver.set(f_kwarg={‘w’: w})
        solver.set_initial_condition([0, I])
        u, t = solver.solve(t_mesh)
```

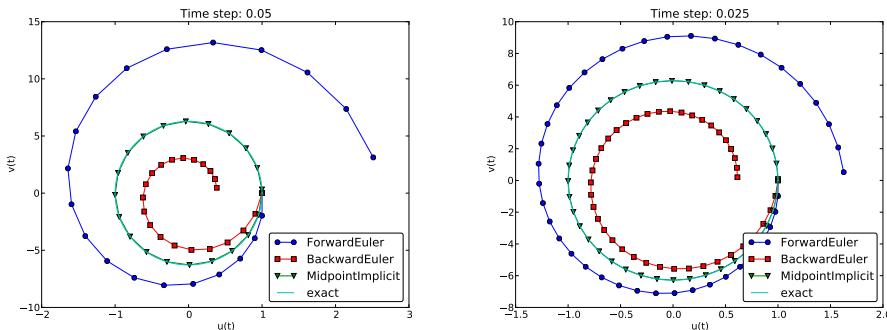
There is quite some more code dealing with plots also, and we refer to the source file `vib_undamped_odespy.py` for details. Observe that keyword arguments in `f(u,t,w=1)` can be supplied through a solver parameter `f_kwarg` (dictionary of additional keyword arguments to `f`).

Specification of the Forward Euler, Backward Euler, and Crank-Nicolson schemes is done like this:

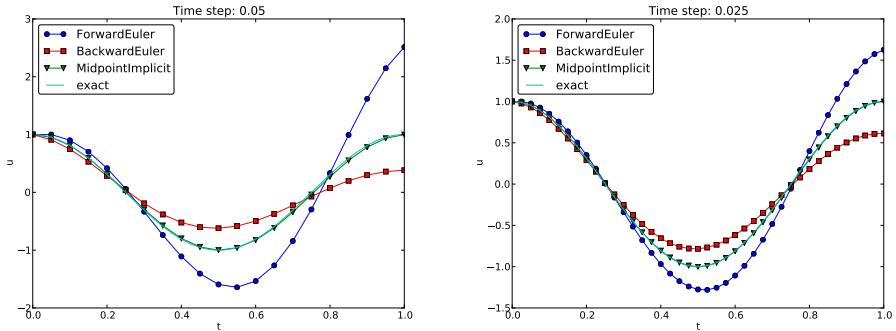
```
solvers = [
    odespy.ForwardEuler(f),
    # Implicit methods must use Newton solver to converge
    odespy.BackwardEuler(f, nonlinear_solver='Newton'),
    odespy.CrankNicolson(f, nonlinear_solver='Newton'),
]
```

The `vib_undamped_odespy.py` program makes two plots of the computed solutions with the various methods in the `solvers` list: one plot with  $u(t)$  versus  $t$ , and one *phase plane plot* where  $v$  is plotted against  $u$ . That is, the phase plane plot is the curve  $(u(t), v(t))$  parameterized by  $t$ . Analytically,  $u = I \cos(\omega t)$  and  $v = u' = -\omega I \sin(\omega t)$ . The exact curve  $(u(t), v(t))$  is therefore an ellipse, which often looks like a circle in a plot if the axes are automatically scaled. The important feature, however, is that the exact curve  $(u(t), v(t))$  is closed and repeats itself for every period. Not all numerical schemes are capable of doing that, meaning that the amplitude instead shrinks or grows with time.

Figure 1.6 show the results. Note that Odespy applies the label `MidpointImplicit` for what we have specified as `CrankNicolson` in the code (`CrankNicolson` is just a synonym for class `MidpointImplicit` in the Odespy code). The Forward Euler scheme in Figure 1.6 has a pronounced spiral curve, pointing to the fact that the amplitude steadily grows, which is also evident in Figure 1.7. The Backward Euler scheme has a similar feature, except that the spiral goes inward and the amplitude is significantly damped. The changing amplitude and the sprial form decreases with decreasing time step. The Crank-Nicolson scheme looks much more accurate. In fact, these plots tell that the Forward and Backward Euler schemes are not suitable for solving our ODEs with oscillating solutions.



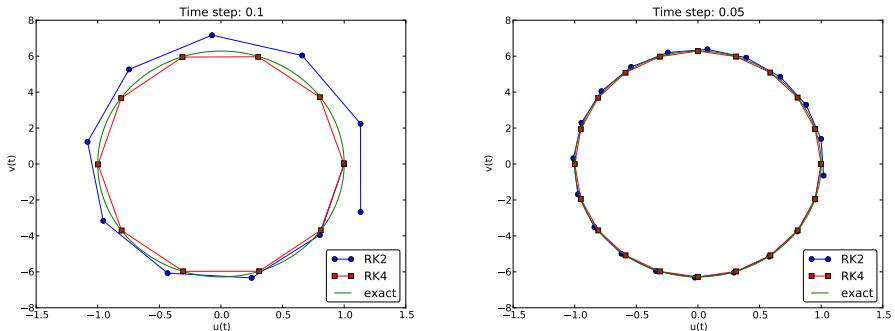
**Fig. 1.6** Comparison of classical schemes in the phase plane for two time step values.



**Fig. 1.7** Comparison of solution curves for classical schemes.

### 1.5.5 Runge-Kutta methods

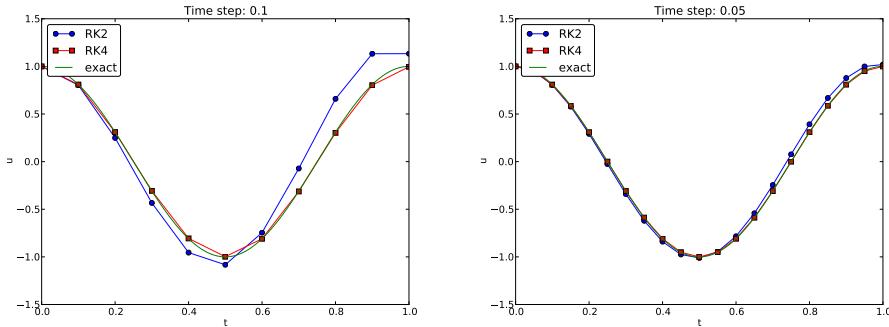
We may run two other popular standard methods for first-order ODEs, the 2nd- and 4th-order Runge-Kutta methods, to see how they perform. Figures 1.8 and 1.9 show the solutions with larger  $\Delta t$  values than what was used in the previous two plots.



**Fig. 1.8** Comparison of Runge-Kutta schemes in the phase plane.

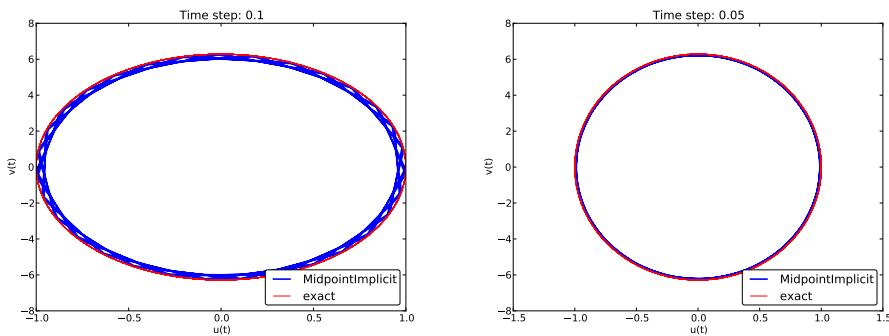
The visual impression is that the 4th-order Runge-Kutta method is very accurate, under all circumstances in these tests, while the 2nd-order scheme suffers from amplitude errors unless the time step is very small.

The corresponding results for the Crank-Nicolson scheme are shown in Figure 1.10. It is clear that the Crank-Nicolson scheme outperforms the 2nd-order Runge-Kutta method. Both schemes have the same order of accuracy  $\mathcal{O}(\Delta t^2)$ , but their differences in the accuracy that matters in a real physical application is very clearly pronounced in this example.



**Fig. 1.9** Comparison of Runge-Kutta schemes.

Exercise 1.12 invites you to investigate how the amplitude is computed by a series of famous methods for first-order ODEs.



**Fig. 1.10** Long-time behavior of the Crank-Nicolson scheme in the phase plane.

### 1.5.6 Analysis of the Forward Euler scheme

We may try to find exact solutions of the discrete equations (1.27)-(1.28) in the Forward Euler method to better understand why this otherwise useful method has so bad performance for vibration ODEs. An “ansatz” for the solution of the discrete equations is

$$\begin{aligned} u^n &= IA^n, \\ v^n &= qIA^n, \end{aligned}$$

where  $q$  and  $A$  are scalars to be determined. We could have used a complex exponential form  $e^{i\tilde{\omega}n\Delta t}$  since we get oscillatory solutions, but the oscillations grow in the Forward Euler method, so the numerical frequency  $\tilde{\omega}$  will be complex anyway (producing an exponentially growing amplitude). Therefore, it is easier to just work with potentially complex  $A$  and  $q$  as introduced above.

The Forward Euler scheme leads to

$$\begin{aligned} A &= 1 + \Delta tq, \\ A &= 1 - \Delta t\omega^2 q^{-1}. \end{aligned}$$

We can easily eliminate  $A$ , get  $q^2 + \omega^2 = 0$ , and solve for

$$q = \pm i\omega,$$

which gives

$$A = 1 \pm \Delta ti\omega.$$

We shall take the real part of  $A^n$  as the solution. The two values of  $A$  are complex conjugates, and the real part of  $A^n$  will be the same for both roots. This is easy to realize if we rewrite the complex numbers in polar form, which is also convenient for further analysis and understanding. The polar form  $re^{i\theta}$  of a complex number  $x + iy$  has  $r = \sqrt{x^2 + y^2}$  and  $\theta = \tan^{-1}(y/x)$ . Hence, the polar form of the two values for  $A$  becomes

$$1 \pm \Delta ti\omega = \sqrt{1 + \omega^2 \Delta t^2} e^{\pm i \tan^{-1}(\omega \Delta t)}.$$

Now it is very easy to compute  $A^n$ :

$$(1 \pm \Delta ti\omega)^n = (1 + \omega^2 \Delta t^2)^{n/2} e^{\pm ni \tan^{-1}(\omega \Delta t)}.$$

Since  $\cos(\theta n) = \cos(-\theta n)$ , the real parts of the two numbers become the same. We therefore continue with the solution that has the plus sign.

The general solution is  $u^n = CA^n$ , where  $C$  is a constant determined from the initial condition:  $u^0 = C = I$ . We have  $u^n = IA^n$  and  $v^n = qIA^n$ . The final solutions are just the real part of the expressions in polar form:

$$u^n = I(1 + \omega^2 \Delta t^2)^{n/2} \cos(n \tan^{-1}(\omega \Delta t)), \quad (1.43)$$

$$v^n = -\omega I(1 + \omega^2 \Delta t^2)^{n/2} \sin(n \tan^{-1}(\omega \Delta t)). \quad (1.44)$$

The expression  $(1 + \omega^2 \Delta t^2)^{n/2}$  causes growth of the amplitude, since a number greater than one is raised to a positive exponent  $n/2$ . We can develop a series expression to better understand the formula for the amplitude. Introducing  $p = \omega \Delta t$  as the key variable and using `sympy` gives

```
>>> from sympy import *
>>> p = symbols('p', real=True)
>>> n = symbols('n', integer=True, positive=True)
>>> amplitude = (1 + p**2)**(n/2)
>>> amplitude.series(p, 0, 4)
1 + n*p**2/2 + O(p**4)
```

The amplitude goes like  $1 + \frac{1}{2}n\omega^2 \Delta t^2$ , clearly growing linearly in time (with  $n$ ).

We can also investigate the error in the angular frequency by a series expansion:

```
>>> n*atan(p).series(p, 0, 4)
n*(p - p**3/3 + O(p**4))
```

This means that the solution for  $u^n$  can be written as

$$u^n = \left(1 + \frac{1}{2}n\omega^2 \Delta t^2 + \mathcal{O}(\Delta t^4)\right) \cos\left(\omega t - \frac{1}{3}\omega t \Delta t^2 + \mathcal{O}(\Delta t^4)\right).$$

The error in the angular frequency is of the same order as in the scheme (1.7) for the second-order ODE, but the error in the amplitude is severe.

## 1.6 Energy considerations

The observations of various methods in the previous section can be better interpreted if we compute a quantity reflecting the total *energy of the system*. It turns out that this quantity,

$$E(t) = \frac{1}{2}(u')^2 + \frac{1}{2}\omega^2 u^2,$$

is *constant* for all  $t$ . Checking that  $E(t)$  really remains constant brings evidence that the numerical computations are sound. It turns out that  $E$  is proportional to the mechanical energy in the system. Conservation of energy is much used to check numerical simulations, so it is well invested time to dive into this subject.

### 1.6.1 Derivation of the energy expression

We start out with multiplying

$$u'' + \omega^2 u = 0,$$

by  $u'$  and integrating from 0 to  $T$ :

$$\int_0^T u'' u' dt + \int_0^T \omega^2 u u' dt = 0.$$

Observing that

$$u'' u' = \frac{d}{dt} \frac{1}{2} (u')^2, \quad u u' = \frac{d}{dt} \frac{1}{2} u^2,$$

we get

$$\int_0^T \left( \frac{d}{dt} \frac{1}{2} (u')^2 + \frac{d}{dt} \frac{1}{2} \omega^2 u^2 \right) dt = E(T) - E(0) = 0,$$

where we have introduced

$$E(t) = \frac{1}{2} (u')^2 + \frac{1}{2} \omega^2 u^2. \quad (1.45)$$

The important result from this derivation is that the total energy is constant:

$$E(t) = E(0).$$

#### **$E(t)$ is closely related to the system's energy**

The quantity  $E(t)$  derived above is physically not the mechanical energy of a vibrating mechanical system, but the energy per unit mass. To see this, we start with Newton's second law  $F = ma$  ( $F$  is the sum of forces,  $m$  is the mass of the system, and  $a$  is the acceleration). The displacement  $u$  is related to  $a$  through  $a = u''$ . With a spring force as the only force we have  $F = -ku$ , where  $k$  is a spring constant measuring the stiffness of the spring. Newton's second law then implies the differential equation

$$-ku = mu'' \Rightarrow mu'' + ku = 0.$$

This equation of motion can be turned into an energy balance equation by finding the work done by each term during a time interval  $[0, T]$ . To this end, we multiply the equation by  $du = u' dt$  and integrate:

$$\int_0^T muu' dt + \int_0^T kuu' dt = 0.$$

The result is

$$\tilde{E}(t) = E_k(t) + E_p(t) = 0,$$

where

$$E_k(t) = \frac{1}{2}mv^2, \quad v = u', \quad (1.46)$$

is the *kinetic energy* of the system, and

$$E_p(t) = \frac{1}{2}ku^2 \quad (1.47)$$

is the *potential energy*. The sum  $\tilde{E}(t)$  is the total mechanical energy. The derivation demonstrates the famous energy principle that, under the right physical circumstances, any change in the kinetic energy is due to a change in potential energy and vice versa. (This principle breaks down when we introduce damping in the system, as we do in Section 1.8.)

The equation  $mu'' + ku = 0$  can be divided by  $m$  and written as  $u'' + \omega^2 u = 0$  for  $\omega = \sqrt{k/m}$ . The energy expression  $E(t) = \frac{1}{2}(u')^2 + \frac{1}{2}\omega^2 u^2$  derived earlier is then  $\tilde{E}(t)/m$ , i.e., mechanical energy per unit mass.

**Energy of the exact solution.** Analytically, we have  $u(t) = I \cos \omega t$ , if  $u(0) = I$  and  $u'(0) = 0$ , so we can easily check the energy evolution and confirm that  $E(t)$  is constant:

$$E(t) = \frac{1}{2}I^2(-\omega \sin \omega t)^2 + \frac{1}{2}\omega^2 I^2 \cos^2 \omega t = \frac{1}{2}\omega^2(\sin^2 \omega t + \cos^2 \omega t) = \frac{1}{2}\omega^2.$$

## 1.6.2 An error measure based on energy

The constant energy is well expressed by its initial value  $E(0)$ , so that the error in mechanical energy can be computed as a mesh function by

$$e_E^n = \frac{1}{2} \left( \frac{u^{n+1} - u^{n-1}}{2\Delta t} \right)^2 + \frac{1}{2} \omega^2 (u^n)^2 - E(0), \quad n = 1, \dots, N_t - 1, \quad (1.48)$$

where

$$E(0) = \frac{1}{2} V^2 + \frac{1}{2} \omega^2 I^2,$$

if  $u(0) = I$  and  $u'(0) = V$ . Note that we have used a centered approximation to  $u'$ :  $u'(t_n) \approx [D_{2t}u]^n$ .

A useful norm of the mesh function  $e_E^n$  for the discrete mechanical energy can be the maximum absolute value of  $e_E^n$ :

$$\|e_E^n\|_{\ell^\infty} = \max_{1 \leq n < N_t} |e_E^n|.$$

Alternatively, we can compute other norms involving integration over all mesh points, but we are often interested in worst case deviation of the energy, and then the maximum value is of particular relevance.

A vectorized Python implementation of  $e_E^n$  takes the form

```
# import numpy as np and compute u, t
dt = t[1]-t[0]
E = 0.5*((u[2:] - u[:-2])/(2*dt))**2 + 0.5*w**2*u[1:-1]**2
E0 = 0.5*V**2 + 0.5*w**2*I**2
e_E = E - E0
e_E_norm = np.abs(e_E).max()
```

The convergence rates of the quantity `e_E_norm` can be used for verification. The value of `e_E_norm` is also useful for comparing schemes through their ability to preserve energy. Below is a table demonstrating the relative error in total energy for various schemes (computed by the `vib_undamped_odespy.py` program). The test problem is  $u'' + 4\pi^2u = 0$  with  $u(0) = 1$  and  $u'(0) = 0$ , so the period is 1 and  $E(t) \approx 4.93$ . We clearly see that the Crank-Nicolson and the Runge-Kutta schemes are superior to the Forward and Backward Euler schemes already after one period.

Method	$T$	$\Delta t$	$\max  e_E^n  / e_E^0$
Forward Euler	1	0.025	$1.678 \cdot 10^0$
Backward Euler	1	0.025	$6.235 \cdot 10^{-1}$
Crank-Nicolson	1	0.025	$1.221 \cdot 10^{-2}$
Runge-Kutta 2nd-order	1	0.025	$6.076 \cdot 10^{-3}$
Runge-Kutta 4th-order	1	0.025	$8.214 \cdot 10^{-3}$

However, after 10 periods, the picture is much more dramatic:

Method	$T$	$\Delta t$	$\max  e_E^n  / e_E^0$
Forward Euler	10	0.025	$1.788 \cdot 10^4$
Backward Euler	10	0.025	$1.000 \cdot 10^0$
Crank-Nicolson	10	0.025	$1.221 \cdot 10^{-2}$
Runge-Kutta 2nd-order	10	0.025	$6.250 \cdot 10^{-2}$
Runge-Kutta 4th-order	10	0.025	$8.288 \cdot 10^{-3}$

The Runge-Kutta and Crank-Nicolson methods hardly change their energy error with  $T$ , while the error in the Forward Euler method grows to huge levels and a relative error of 1 in the Backward Euler method points to  $E(t) \rightarrow 0$  as  $t$  grows large.

Running multiple values of  $\Delta t$ , we can get some insight into the convergence of the energy error:

Method	$T$	$\Delta t$	$\max  e_E^n  / e_E^0$
Forward Euler	10	0.05	$1.120 \cdot 10^8$
Forward Euler	10	0.025	$1.788 \cdot 10^4$
Forward Euler	10	0.0125	$1.374 \cdot 10^2$
Backward Euler	10	0.05	$1.000 \cdot 10^0$
Backward Euler	10	0.025	$1.000 \cdot 10^0$
Backward Euler	10	0.0125	$9.928 \cdot 10^{-1}$
Crank-Nicolson	10	0.05	$4.756 \cdot 10^{-2}$
Crank-Nicolson	10	0.025	$1.221 \cdot 10^{-2}$
Crank-Nicolson	10	0.0125	$3.125 \cdot 10^{-3}$
Runge-Kutta 2nd-order	10	0.05	$6.152 \cdot 10^{-1}$
Runge-Kutta 2nd-order	10	0.025	$6.250 \cdot 10^{-2}$
Runge-Kutta 2nd-order	10	0.0125	$7.631 \cdot 10^{-3}$
Runge-Kutta 4th-order	10	0.05	$3.510 \cdot 10^{-2}$
Runge-Kutta 4th-order	10	0.025	$8.288 \cdot 10^{-3}$
Runge-Kutta 4th-order	10	0.0125	$2.058 \cdot 10^{-3}$

A striking fact from this table is that the error of the Forward Euler method is reduced by the same factor as  $\Delta t$  is reduced by, while the error in the Crank-Nicolson method has a reduction proportional to  $\Delta t^2$  (we cannot say anything for the Backward Euler method). However, for the RK2 method, halving  $\Delta t$  reduces the error by almost a factor of 10 (!), and for the RK4 method the reduction seems proportional to  $\Delta t^2$  only (and the trend is confirmed by running smaller time steps, so for

$\Delta t = 3.9 \cdot 10^{-4}$  the relative error of RK2 is a factor 10 smaller than that of RK4!).

## 1.7 The Euler-Cromer method

While the Runge-Kutta methods and the Crank-Nicolson scheme work well for the vibration equation modeled as a first-order ODE system, both were inferior to the straightforward centered difference scheme for the second-order equation  $u'' + \omega^2 u = 0$ . However, there is a similarly successful scheme available for the first-order system  $u' = v$ ,  $v' = -\omega^2 u$ , to be presented next.

### 1.7.1 Forward-backward discretization

The idea is to apply a Forward Euler discretization to the first equation and a Backward Euler discretization to the second. In operator notation this is stated as

$$[D_t^+ u = v]^n, \quad (1.49)$$

$$[D_t^- v = -\omega^2 u]^{n+1}. \quad (1.50)$$

We can write out the formulas and collect the unknowns on the left-hand side:

$$u^{n+1} = u^n + \Delta t v^n, \quad (1.51)$$

$$v^{n+1} = v^n - \Delta t \omega^2 u^{n+1}. \quad (1.52)$$

We realize that after  $u^{n+1}$  has been computed from (1.51), it may be used directly in (1.52) to compute  $v^{n+1}$ .

In physics, it is more common to update the  $v$  equation first, with a forward difference, and thereafter the  $u$  equation, with a backward difference that applies the most recently computed  $v$  value:

$$v^{n+1} = v^n - \Delta t \omega^2 u^n, \quad (1.53)$$

$$u^{n+1} = u^n + \Delta t v^{n+1}. \quad (1.54)$$

The advantage of ordering the ODEs as in (1.53)-(1.54) becomes evident when considering complicated models. Such models are included if we write our vibration ODE more generally as

$$u'' + g(u, u', t) = 0.$$

We can rewrite this second-order ODE as two first-order ODEs,

$$\begin{aligned} v' &= -g(u, v, t), \\ u' &= v. \end{aligned}$$

This rewrite allows the following scheme to be used:

$$\begin{aligned} v^{n+1} &= v^n - \Delta t g(u^n, v^n, t), \\ u^{n+1} &= u^n + \Delta t v^{n+1}. \end{aligned}$$

We realize that the first update works well with any  $g$  since old values  $u^n$  and  $v^n$  are used. Switching the equations would demand  $u^{n+1}$  and  $v^{n+1}$  values in  $g$  and result in nonlinear algebraic equations to be solved at each time level.

The scheme (1.53)-(1.54) goes under several names: forward-backward scheme, **semi-implicit Euler method**, semi-explicit Euler, symplectic Euler, Newton-Störmer-Verlet, and Euler-Cromer. We shall stick to the latter name. Since both time discretizations are based on first-order difference approximation, one may think that the scheme is only of first-order, but this is not true: the use of a forward and then a backward difference make errors cancel so that the overall error in the scheme is  $\mathcal{O}(\Delta t^2)$ . This is explained below.

How does the Euler-Cromer method preserve the total energy? We may run the example from Section 1.6.2:

Method	$T$	$\Delta t$	$\max  e_E^n  / e_E^0$
Euler-Cromer	10	0.05	$2.530 \cdot 10^{-2}$
Euler-Cromer	10	0.025	$6.206 \cdot 10^{-3}$
Euler-Cromer	10	0.0125	$1.544 \cdot 10^{-3}$

The relative error in the total energy decreases as  $\Delta t^2$ , and the error level is slightly lower than for the Crank-Nicolson and Runge-Kutta methods.

### 1.7.2 Equivalence with the scheme for the second-order ODE

We shall now show that the Euler-Cromer scheme for the system of first-order equations is equivalent to the centered finite difference method for the second-order vibration ODE (!).

We may eliminate the  $v^n$  variable from (1.51)-(1.52) or (1.53)-(1.54). The  $v^{n+1}$  term in (1.53) can be eliminated from (1.54):

$$u^{n+1} = u^n + \Delta t(v^n - \omega^2 \Delta t u^n). \quad (1.55)$$

The  $v^n$  quantity can be expressed by  $u^n$  and  $u^{n-1}$  using (1.54):

$$v^n = \frac{u^n - u^{n-1}}{\Delta t},$$

and when this is inserted in (1.55) we get

$$u^{n+1} = 2u^n - u^{n-1} - \Delta t^2 \omega^2 u^n, \quad (1.56)$$

which is nothing but the centered scheme (1.7)! The two seemingly different numerical methods are mathematically equivalent. Consequently, the previous analysis of (1.7) also applies to the Euler-Cromer method. In particular, the amplitude is constant, given that the stability criterion is fulfilled, but there is always an angular frequency error (1.19). Exercise 1.17 gives guidance on how to derive the exact discrete solution of the two equations in the Euler-Cromer method.

Although the Euler-Cromer scheme and the method (1.7) are equivalent, there could be differences in the way they handle the initial conditions. Let us look into this topic. The initial condition  $u' = 0$  means  $u' = v = 0$ . From (1.53) we get

$$v^1 = v^0 - \Delta t \omega^2 u^0 = \Delta t \omega^2 u^0,$$

and from (1.54) it follows that

$$u^1 = u^0 + \Delta t v^1 = u^0 - \omega^2 \Delta t^2 u^0.$$

When we previously used a centered approximation of  $u'(0) = 0$  combined with the discretization (1.7) of the second-order ODE, we got a slightly different result:  $u^1 = u^0 - \frac{1}{2} \omega^2 \Delta t^2 u^0$ . The difference is  $\frac{1}{2} \omega^2 \Delta t^2 u^0$ , which is of second order in  $\Delta t$ , seemingly consistent with the overall error in the scheme for the differential equation model.

A different view can also be taken. If we approximate  $u'(0) = 0$  by a backward difference,  $(u^0 - u^{-1})/\Delta t = 0$ , we get  $u^{-1} = u^0$ , and when combined with (1.7), it results in  $u^1 = u^0 - \omega^2 \Delta t^2 u^0$ . This means that the Euler-Cromer method based on (1.54)-(1.53) corresponds to using only a first-order approximation to the initial condition in the method from Section 1.1.2.

Correspondingly, using the formulation (1.51)-(1.52) with  $v^n = 0$  leads to  $u^1 = u^0$ , which can be interpreted as using a forward difference approximation for the initial condition  $u'(0) = 0$ . Both Euler-Cromer formulations lead to slightly different values for  $u^1$  compared to the method in Section 1.1.2. The error is  $\frac{1}{2}\omega^2 \Delta t^2 u^0$  and of the same order as the overall scheme.

### 1.7.3 Implementation

The function below, found in `vib_EulerCromer.py`, implements the Euler-Cromer scheme (1.53)-(1.54):

```
import numpy as np

def solver(I, w, dt, T):
    """
    Solve v' = - w**2*u, u'=v for t in (0,T], u(0)=I and v(0)=0,
    by an Euler-Cromer method.
    """
    dt = float(dt)
    Nt = int(round(T/dt))
    u = np.zeros(Nt+1)
    v = np.zeros(Nt+1)
    t = np.linspace(0, Nt*dt, Nt+1)

    v[0] = 0
    u[0] = I
    for n in range(0, Nt):
        v[n+1] = v[n] - dt*w**2*u[n]
        u[n+1] = u[n] + dt*v[n+1]
    return u, v, t
```

Since the Euler-Cromer scheme is equivalent to the finite difference method for the second-order ODE  $u'' + \omega^2 u = 0$  (see Section 1.7.2), the performance of the above `solver` function is the same as for the `solver` function in Section 1.2. The only difference is the formula for the first time step, as discussed above. This deviation in the Euler-Cromer scheme means that the discrete solution listed in Section 1.4.2 is not a solution of the Euler-Cromer scheme!

To verify the implementation of the Euler-Cromer method we can adjust `v[1]` so that the computer-generated values can be compared with the formula (1.20) from in Section 1.4.2. This adjustment is done in an alternative solver function, `solver_ic_fix` in `vib_EulerCromer.py`. Since we now have an exact solution of the discrete equations available, we can write a test function `test_solver` for checking the equality of computed values with the formula (1.20):

```
def test_solver():
    """
    Test solver with fixed initial condition against
    equivalent scheme for the 2nd-order ODE  $u'' + u = 0$ .
    """
    I = 1.2; w = 2.0; T = 5
    dt = 2/w # longest possible time step
    u, v, t = solver_ic_fix(I, w, dt, T)
    from vib_undamped import solver as solver2 # 2nd-order ODE
    u2, t2 = solver2(I, w, dt, T)
    error = np.abs(u - u2).max()
    tol = 1E-14
    assert error < tol
```

Another function, `demo`, visualizes the difference between the Euler-Cromer scheme and the scheme (1.7) for the second-order ODE, arising from the mismatch in the first time level.

The Euler-Cromer method is also available in the `Odespy` package. The important thing to remember is that we must order the unknowns as  $v$  and  $u$ , so the  $\mathbf{u}$  vector at each time level consists of the velocity  $v$  as first component and the displacement  $u$  as second component:

```
# Define ODE
def f(u, t, w=1):
    v, u = u
    return [-w**2*u, v]

# Initialize solver
I = 1
w = 2*np.pi
import odespy
solver = odespy.EulerCromer(f, f_kwarg={‘w’: w})
solver.set_initial_condition([0, I])

# Compute time mesh
P = 2*np.pi/w # duration of one period
dt = P/timesteps_per_period
Nt = num_periods*timesteps_per_period
T = Nt*dt
import numpy as np
t_mesh = np.linspace(0, T, Nt+1)
```

```
# Solve ODE
u, t = solver.solve(t_mesh)
u = u[:,1] # Extract displacement
```

### 1.7.4 The velocity Verlet algorithm

Another very popular algorithm for vibration problems  $u'' + \omega^2 u = 0$  can be derived as follows. First, we step  $u$  forward from  $t_n$  to  $t_{n+1}$  using a three-term Taylor series,

$$u(t_{n+1}) = u(t_n) + u'(t_n)\Delta t + \frac{1}{2}u''(t_n)\Delta t^2.$$

Using  $u' = v$  and  $u'' = -\omega^2 u$ , we get the updating formula

$$u^{n+1} = u^n + v^n \Delta t - \frac{1}{2} \Delta t^2 \omega^2 u^n.$$

Second, the first-order equation for  $v$ ,

$$v' = -\omega^2 u,$$

is discretized by a centered difference in a Crank-Nicolson fashion at  $t_{n+\frac{1}{2}}$ :

$$\frac{v^{n+1} - v^n}{\Delta t} = -\omega^2 \frac{1}{2} (u^n + u^{n+1}).$$

To summarize, we have the scheme

$$u^{n+1} = u^n + v^n \Delta t - \frac{1}{2} \Delta t^2 \omega^2 u^n \quad (1.57)$$

$$v^{n+1} = v^n - \frac{1}{2} \Delta t \omega^2 (u^n + u^{n+1}), \quad (1.58)$$

known as the *velocity Verlet* algorithm. Observe that this scheme is explicit since  $u^{n+1}$  in (1.58) is already computed from (1.57).

The algorithm can be straightforwardly implemented as shown below (the code appears in the file `vib_undamped_velocity_Verlet.py`).

```
from vib_undamped import convergence_rates, main

def solver(I, w, dt, T, return_v=False):
    """
    Solve u'=v, v'=-w**2*u for t in (0,T], u(0)=I and v(0)=0,
```

```

    by the velocity Verlet method with time step dt.
"""

dt = float(dt)
Nt = int(round(T/dt))
u = np.zeros(Nt+1)
v = np.zeros(Nt+1)
t = np.linspace(0, Nt*dt, Nt+1)

u[0] = I
v[0] = 0
for n in range(Nt):
    u[n+1] = u[n] + v[n]*dt - 0.5*dt**2*v**2*u[n]
    v[n+1] = v[n] - 0.5*dt*v**2*(u[n] + u[n+1])
if return_v:
    return u, v, t
else:
    # Return just u and t as in the vib_undamped.py's solver
    return u, t

```

We provide the option that this `solver` function returns the same data as the `solver` function from Section 1.2.1 (if `return_v` is `False`), but alternatively, it may return `v` along with `u` and `t`.

The error in the Taylor series expansion behind (1.57) is  $\mathcal{O}(\Delta t^3)$ , while the error in the central difference for  $v$  is  $\mathcal{O}(\Delta t^2)$ . The overall error is then no better than  $\mathcal{O}(\Delta t^2)$ , which can be verified empirically using the `convergence_rates` function from 1.2.2:

```

>>> import vib_undamped_velocity_Verlet as m
>>> m.convergence_rates(4, solver_function=m.solver)
[2.0036366687367346, 2.0009497328124835, 2.000240105995295]

```

## 1.8 Generalization: damping, nonlinear spring, and external excitation

We shall now generalize the simple model problem from Section 1.1 to include a possibly nonlinear damping term  $f(u')$ , a possibly nonlinear spring (or restoring) force  $s(u)$ , and some external excitation  $F(t)$ :

$$mu'' + f(u') + s(u) = F(t), \quad u(0) = I, \quad u'(0) = V, \quad t \in (0, T]. \quad (1.59)$$

We have also included a possibly nonzero initial value of  $u'(0)$ . The parameters  $m$ ,  $f(u')$ ,  $s(u)$ ,  $F(t)$ ,  $I$ ,  $V$ , and  $T$  are input data.

There are two main types of damping (friction) forces: linear  $f(u') = bu'$ , or quadratic  $f(u') = bu'|u'|$ . Spring systems often feature linear damping,

while air resistance usually gives rise to quadratic damping. Spring forces are often linear:  $s(u) = cu$ , but nonlinear versions are also common, the most famous is the gravity force on a pendulum that acts as a spring with  $s(u) \sim \sin(u)$ .

### 1.8.1 A centered scheme for linear damping

Sampling (1.59) at a mesh point  $t_n$ , replacing  $u''(t_n)$  by  $[D_tD_t u]^n$ , and  $u'(t_n)$  by  $[D_{2t} u]^n$  results in the discretization

$$[mD_tD_t u + f(D_{2t} u) + s(u) = F]^n, \quad (1.60)$$

which written out means

$$m \frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} + f\left(\frac{u^{n+1} - u^{n-1}}{2\Delta t}\right) + s(u^n) = F^n, \quad (1.61)$$

where  $F^n$  as usual means  $F(t)$  evaluated at  $t = t_n$ . Solving (1.61) with respect to the unknown  $u^{n+1}$  gives a problem: the  $u^{n+1}$  inside the  $f$  function makes the equation *nonlinear* unless  $f(u')$  is a linear function,  $f(u') = bu'$ . For now we shall assume that  $f$  is linear in  $u'$ . Then

$$m \frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} + b \frac{u^{n+1} - u^{n-1}}{2\Delta t} + s(u^n) = F^n, \quad (1.62)$$

which gives an explicit formula for  $u$  at each new time level:

$$u^{n+1} = (2mu^n + (\frac{b}{2}\Delta t - m)u^{n-1} + \Delta t^2(F^n - s(u^n)))(m + \frac{b}{2}\Delta t)^{-1}. \quad (1.63)$$

For the first time step we need to discretize  $u'(0) = V$  as  $[D_{2t} u = V]^0$  and combine with (1.63) for  $n = 0$ . The discretized initial condition leads to

$$u^{-1} = u^1 - 2\Delta t V, \quad (1.64)$$

which inserted in (1.63) for  $n = 0$  gives an equation that can be solved for  $u^1$ :

$$u^1 = u^0 + \Delta t V + \frac{\Delta t^2}{2m}(-bV - s(u^0) + F^0). \quad (1.65)$$

### 1.8.2 A centered scheme for quadratic damping

When  $f(u') = bu'|u'|$ , we get a quadratic equation for  $u^{n+1}$  in (1.61). This equation can be straightforwardly solved by the well-known formula for the roots of a quadratic equation. However, we can also avoid the nonlinearity by introducing an approximation with an error of order no higher than what we already have from replacing derivatives with finite differences.

We start with (1.59) and only replace  $u''$  by  $D_tD_t u$ , resulting in

$$[mD_tD_t u + bu'|u'| + s(u) = F]^n. \quad (1.66)$$

Here,  $u'|u'|$  is to be computed at time  $t_n$ . The idea is now to introduce a *geometric mean*, defined by

$$(w^2)^n \approx w^{n-\frac{1}{2}} w^{n+\frac{1}{2}},$$

for some quantity  $w$  depending on time. The error in the geometric mean approximation is  $\mathcal{O}(\Delta t^2)$ , the same as in the approximation  $u'' \approx D_tD_t u$ . With  $w = u'$  it follows that

$$[u'|u'|]^n \approx u'(t_{n+\frac{1}{2}})|u'(t_{n-\frac{1}{2}})|.$$

The next step is to approximate  $u'$  at  $t_{n\pm 1/2}$ , and fortunately a centered difference fits perfectly into the formulas since it involves  $u$  values at the mesh points only. With the approximations

$$u'(t_{n+1/2}) \approx [D_t u]^{n+\frac{1}{2}}, \quad u'(t_{n-1/2}) \approx [D_t u]^{n-\frac{1}{2}}, \quad (1.67)$$

we get

$$[u'|u'|]^n \approx [D_t u]^{n+\frac{1}{2}}|[D_t u]^{n-\frac{1}{2}}| = \frac{u^{n+1} - u^n}{\Delta t} \frac{|u^n - u^{n-1}|}{\Delta t}. \quad (1.68)$$

The counterpart to (1.61) is then

$$m \frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} + b \frac{u^{n+1} - u^n}{\Delta t} \frac{|u^n - u^{n-1}|}{\Delta t} + s(u^n) = F^n, \quad (1.69)$$

which is linear in the unknown  $u^{n+1}$ . Therefore, we can easily solve (1.69) with respect to  $u^{n+1}$  and achieve the explicit updating formula

$$\begin{aligned} u^{n+1} = & \left( m + b|u^n - u^{n-1}| \right)^{-1} \times \\ & \left( 2mu^n - mu^{n-1} + bu^n|u^n - u^{n-1}| + \Delta t^2(F^n - s(u^n)) \right). \end{aligned} \quad (1.70)$$

In the derivation of a special equation for the first time step we run into some trouble: inserting (1.64) in (1.70) for  $n = 0$  results in a complicated nonlinear equation for  $u^1$ . By thinking differently about the problem we can easily get away with the nonlinearity again. We have for  $n = 0$  that  $b[u'|u'|]^0 = bV|V|$ . Using this value in (1.66) gives

$$[mD_tD_t u + bV|V| + s(u) = F]^0. \quad (1.71)$$

Writing this equation out and using (1.64) results in the special equation for the first time step:

$$u^1 = u^0 + \Delta t V + \frac{\Delta t^2}{2m} \left( -bV|V| - s(u^0) + F^0 \right). \quad (1.72)$$

### 1.8.3 A forward-backward discretization of the quadratic damping term

The previous section first proposed to discretize the quadratic damping term  $|u'|u'$  using centered differences:  $[(D_{2t}|D_{2t}u)]^n$ . As this gives rise to a nonlinearity in  $u^{n+1}$ , it was instead proposed to use a geometric mean combined with centered differences. But there are other alternatives. To get rid of the nonlinearity in  $[(D_{2t}|D_{2t}u)]^n$ , one can think differently: apply a backward difference to  $|u'|$ , such that the term involves known values, and apply a forward difference to  $u'$  to make the term linear in the unknown  $u^{n+1}$ . With mathematics,

$$[\beta|u'|u'|]^n \approx \beta|[D_t^- u]^n|[D_t^+ u]^n = \beta \left| \frac{u^n - u^{n-1}}{\Delta t} \right| \frac{u^{n+1} - u^n}{\Delta t}. \quad (1.73)$$

The forward and backward differences have both an error proportional to  $\Delta t$  so one may think the discretization above leads to a first-order scheme. However, by looking at the formulas, we realize that the forward-backward differences in (1.73) result in exactly the same scheme as in (1.69) where we used a geometric mean and centered differences and committed errors of size  $\mathcal{O}(\Delta t^2)$ . Therefore, the forward-backward differences in (1.73)

act in a symmetric way and actually produce a second-order accurate discretization of the quadratic damping term.

### 1.8.4 Implementation

The algorithm arising from the methods in Sections 1.8.1 and 1.8.2 is very similar to the undamped case in Section 1.1.2. The difference is basically a question of different formulas for  $u^1$  and  $u^{n+1}$ . This is actually quite remarkable. The equation (1.59) is normally impossible to solve by pen and paper, but possible for some special choices of  $F$ ,  $s$ , and  $f$ . On the contrary, the complexity of the nonlinear generalized model (1.59) versus the simple undamped model is not a big deal when we solve the problem numerically!

The computational algorithm takes the form

1.  $u^0 = I$
2. compute  $u^1$  from (1.65) if linear damping or (1.72) if quadratic damping
3. for  $n = 1, 2, \dots, N_t - 1$ :
  - a. compute  $u^{n+1}$  from (1.63) if linear damping or (1.70) if quadratic damping

Modifying the `solver` function for the undamped case is fairly easy, the big difference being many more terms and if tests on the type of damping:

```
def solver(I, V, m, b, s, F, dt, T, damping='linear'):
    """
    Solve m*u'' + f(u') + s(u) = F(t) for t in (0,T],
    u(0)=I and u'(0)=V,
    by a central finite difference method with time step dt.
    If damping is 'linear', f(u')=b*u, while if damping is
    'quadratic', f(u')=b*u'*abs(u').
    F(t) and s(u) are Python functions.
    """
    dt = float(dt); b = float(b); m = float(m) # avoid integer div.
    Nt = int(round(T/dt))
    u = np.zeros(Nt+1)
    t = np.linspace(0, Nt*dt, Nt+1)

    u[0] = I
    if damping == 'linear':
        u[1] = u[0] + dt*V + dt**2/(2*m)*(-b*V - s(u[0]) + F(t[0]))
    elif damping == 'quadratic':
        u[1] = u[0] + dt*V + \
            dt**2/(2*m)*(-b*V*abs(V) - s(u[0]) + F(t[0]))
```

```

for n in range(1, Nt):
    if damping == 'linear':
        u[n+1] = (2*m*u[n] + (b*dt/2 - m)*u[n-1] +
                   dt**2*(F(t[n]) - s(u[n])))/(m + b*dt/2)
    elif damping == 'quadratic':
        u[n+1] = (2*m*u[n] - m*u[n-1] + b*u[n]*abs(u[n] - u[n-1]) +
                   dt**2*(F(t[n]) - s(u[n])))/(m + b*abs(u[n] - u[n-1]))
    return u, t

```

The complete code resides in the file `vib.py`.

### 1.8.5 Verification

**Constant solution.** For debugging and initial verification, a constant solution is often very useful. We choose  $u_e(t) = I$ , which implies  $V = 0$ . Inserted in the ODE, we get  $F(t) = s(I)$  for any choice of  $f$ . Since the discrete derivative of a constant vanishes (in particular,  $[D_{2t}I]^n = 0$ ,  $[D_tI]^n = 0$ , and  $[D_tD_tI]^n = 0$ ), the constant solution also fulfills the discrete equations. The constant should therefore be reproduced to machine precision. The function `test_constant` in `vib.py` implements this test.

**hpl 4:** Add verification tests for constant, linear, quadratic. Check how many bugs that are caught by these tests.

**Linear solution.** Now we choose a linear solution:  $u_e = ct + d$ . The initial condition  $u(0) = I$  implies  $d = I$ , and  $u'(0) = V$  forces  $c$  to be  $V$ . Inserting  $u_e = Vt + I$  in the ODE with linear damping results in

$$0 + bV + s(Vt + I) = F(t),$$

while quadratic damping requires the source term

$$0 + b|V|V + s(Vt + I) = F(t).$$

Since the finite difference approximations used to compute  $u'$  all are exact for a linear function, it turns out that the linear  $u_e$  is also a solution of the discrete equations. Exercise 1.9 asks you to carry out all the details.

**Quadratic solution.** Choosing  $u_e = bt^2 + Vt + I$ , with  $b$  arbitrary, fulfills the initial conditions and fits the ODE if  $F$  is adjusted properly. The solution also solves the discrete equations with linear damping. However, this quadratic polynomial in  $t$  does not fulfill the discrete equations in case of quadratic damping, because the geometric mean used in the

approximation of this term introduces an error. Doing Exercise 1.9 will reveal the details. One can fit  $F^n$  in the discrete equations such that the quadratic polynomial is reproduced by the numerical method (to machine precision).

### 1.8.6 Visualization

The functions for visualizations differ significantly from those in the undamped case in the `vib_undamped.py` program because, in the present general case, we do not have an exact solution to include in the plots. Moreover, we have no good estimate of the periods of the oscillations as there will be one period determined by the system parameters, essentially the approximate frequency  $\sqrt{s'(0)/m}$  for linear  $s$  and small damping, and one period dictated by  $F(t)$  in case the excitation is periodic. This is, however, nothing that the program can depend on or make use of. Therefore, the user has to specify  $T$  and the window width to get a plot that moves with the graph and shows the most recent parts of it in long time simulations.

The `vib.py` code contains several functions for analyzing the time series signal and for visualizing the solutions.

### 1.8.7 User interface

The `main` function is changed substantially from the `vib_undamped.py` code, since we need to specify the new data  $c$ ,  $s(u)$ , and  $F(t)$ . In addition, we must set  $T$  and the plot window width (instead of the number of periods we want to simulate as in `vib_undamped.py`). To figure out whether we can use one plot for the whole time series or if we should follow the most recent part of  $u$ , we can use the `plot_emprical_freq_and_amplitude` function's estimate of the number of local maxima. This number is now returned from the function and used in `main` to decide on the visualization technique.

```
def main():
    import argparse
    parser = argparse.ArgumentParser()
    parser.add_argument('--I', type=float, default=1.0)
    parser.add_argument('--V', type=float, default=0.0)
    parser.add_argument('--m', type=float, default=1.0)
    parser.add_argument('--c', type=float, default=0.0)
    parser.add_argument('--s', type=str, default='u')
```

```

parser.add_argument('--F', type=str, default='0')
parser.add_argument('--dt', type=float, default=0.05)
parser.add_argument('--T', type=float, default=140)
parser.add_argument('--damping', type=str, default='linear')
parser.add_argument('--window_width', type=float, default=30)
parser.add_argument('--savefig', action='store_true')
a = parser.parse_args()
from scitools.std import StringFunction
s = StringFunction(a.s, independent_variable='u')
F = StringFunction(a.F, independent_variable='t')
I, V, m, c, dt, T, window_width, savefig, damping = \
    a.I, a.V, a.m, a.c, a.dt, a.T, a.window_width, a.savefig, \
    a.damping

u, t = solver(I, V, m, c, s, F, dt, T)
num_periods = empirical_freq_and_amplitude(u, t)
if num_periods <= 15:
    figure()
    visualize(u, t)
else:
    visualize_front(u, t, window_width, savefig)
show()

```

The program `vib.py` contains the above code snippets and can solve the model problem (1.59). As a demo of `vib.py`, we consider the case  $I = 1$ ,  $V = 0$ ,  $m = 1$ ,  $c = 0.03$ ,  $s(u) = \sin(u)$ ,  $F(t) = 3\cos(4t)$ ,  $\Delta t = 0.05$ , and  $T = 140$ . The relevant command to run is

---

Terminal

```
Terminal> python vib.py --s 'sin(u)' --F '3*cos(4*t)' --c 0.03
```

---

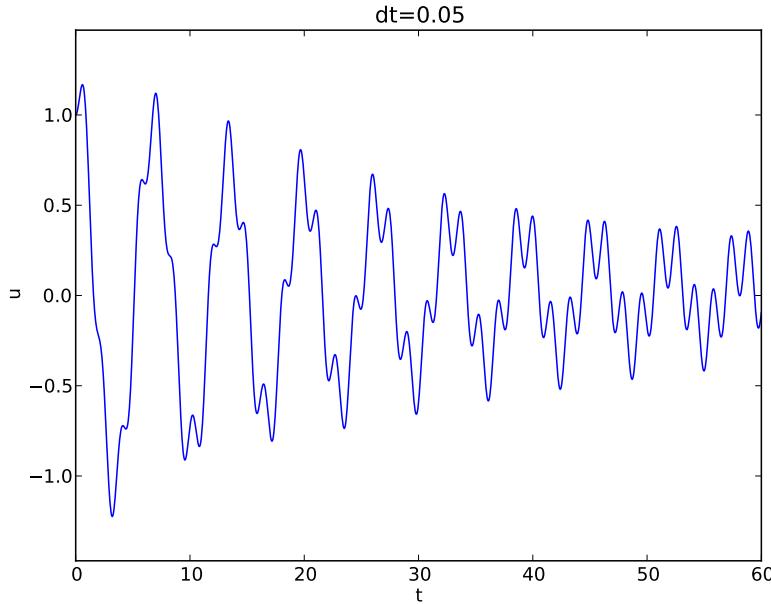
This results in a [moving window](#) following the function on the screen. Figure 1.11 shows a part of the time series.

### 1.8.8 The Euler-Cromer scheme for the generalized model

The ideas of the Euler-Cromer method from Section 1.7 carry over to the generalized model. We write (1.59) as two equations for  $u$  and  $v = u'$ . The first equation is taken as the one with  $v'$  on the left-hand side:

$$v' = \frac{1}{m}(F(t) - s(u) - f(v)), \quad (1.74)$$

$$u' = v. \quad (1.75)$$



**Fig. 1.11** Damped oscillator excited by a sinusoidal function.

The idea is to step (1.74) forward using a standard Forward Euler method, while we update  $u$  from (1.75) with a Backward Euler method, utilizing the recent, computed  $v^{n+1}$  value. In detail,

$$\frac{v^{n+1} - v^n}{\Delta t} = \frac{1}{m}(F(t_n) - s(u^n) - f(v^n)), \quad (1.76)$$

$$\frac{u^{n+1} - u^n}{\Delta t} = v^{n+1}, \quad (1.77)$$

resulting in the explicit scheme

$$v^{n+1} = v^n + \Delta t \frac{1}{m}(F(t_n) - s(u^n) - f(v^n)), \quad (1.78)$$

$$u^{n+1} = u^n + \Delta t v^{n+1}. \quad (1.79)$$

We immediately note one very favorable feature of this scheme: all the nonlinearities in  $s(u)$  and  $f(v)$  are evaluated at a previous time level. This makes the Euler-Cromer method easier to apply and hence much

more convenient than the centered scheme for the second-order ODE (1.59).

The initial conditions are trivially set as

$$v^0 = V, \quad (1.80)$$

$$u^0 = I. \quad (1.81)$$

**hpl 5:** odespy for the generalized problem

## 1.9 Exercises and Problems

### Problem 1.1: Use linear/quadratic functions for verification

Consider the ODE problem

$$u'' + \omega^2 u = f(t), \quad u(0) = I, \quad u'(0) = V, \quad t \in (0, T].$$

- a) Discretize this equation according to  $[D_t D_t u + \omega^2 u = f]^n$  and derive the equation for the first time step ( $u^1$ ).
- b) For verification purposes, we use the method of manufactured solutions (MMS) with the choice of  $u_e(t) = ct + d$ . Find restrictions on  $c$  and  $d$  from the initial conditions. Compute the corresponding source term  $f$ . Show that  $[D_t D_t t]^n = 0$  and use the fact that the  $D_t D_t$  operator is linear,  $[D_t D_t (ct + d)]^n = c[D_t D_t t]^n + [D_t D_t d]^n = 0$ , to show that  $u_e$  is also a perfect solution of the discrete equations.
- c) Use `sympy` to do the symbolic calculations above. Here is a sketch of the program `vib_undamped_verify_mms.py`:

```
import sympy as sym
V, t, I, w, dt = sym.symbols('V t I w dt') # global symbols
f = None # global variable for the source term in the ODE

def ode_source_term(u):
    """Return the terms in the ODE that the source term
    must balance, here u'' + w**2*u.
    u is symbolic Python function of t."""
    return sym.diff(u(t), t, t) + w**2*u(t)

def residual_discrete_eq(u):
    """Return the residual of the discrete eq. with u inserted."""
    R = ...
    return sym.simplify(R)
```

```

def residual_discrete_eq_step1(u):
    """Return the residual of the discrete eq. at the first
    step with u inserted."""
    R = ...
    return sym.simplify(R)

def DtDt(u, dt):
    """Return 2nd-order finite difference for u_tt.
    u is a symbolic Python function of t.
    """
    return ...

def main(u):
    """
    Given some chosen solution u (as a function of t, implemented
    as a Python function), use the method of manufactured solutions
    to compute the source term f, and check if u also solves
    the discrete equations.
    """
    print '*** Testing exact solution: %s ===' % u
    print "Initial conditions u(0)=%s, u'(0)=%s:" % \
        (u(t).subs(t, 0), sym.diff(u(t), t).subs(t, 0))

    # Method of manufactured solution requires fitting f
    global f  # source term in the ODE
    f = sym.simplify(ode_lhs(u))

    # Residual in discrete equations (should be 0)
    print 'residual step1:', residual_discrete_eq_step1(u)
    print 'residual:', residual_discrete_eq(u)

def linear():
    main(lambda t: V*t + I)

if __name__ == '__main__':
    linear()

```

Fill in the various functions such that the calls in the `main` function works.

- d)** The purpose now is to choose a quadratic function  $u_e = bt^2 + ct + d$  as exact solution. Extend the `sympy` code above with a function `quadratic` for fitting `f` and checking if the discrete equations are fulfilled. (The function is very similar to `linear`.)
- e)** Will a polynomial of degree three fulfill the discrete equations?
- f)** Implement a `solver` function for computing the numerical solution of this problem.

**g)** Write a test function for checking that the quadratic solution is computed correctly (to machine precision, but the round-off errors accumulate and increase with  $T$ ) by the `solver` function.

Filename: `vib_undamped_verify_mms`.

### Exercise 1.2: Show linear growth of the phase with time

Consider an exact solution  $I \cos(\omega t)$  and an approximation  $I \cos(\tilde{\omega}t)$ . Define the phase error as the time lag between the peak  $I$  in the exact solution and the corresponding peak in the approximation after  $m$  periods of oscillations. Show that this phase error is linear in  $m$ .

Filename: `vib_phase_error_growth`.

### Exercise 1.3: Improve the accuracy by adjusting the frequency

According to (1.19), the numerical frequency deviates from the exact frequency by a (dominating) amount  $\omega^3 \Delta t^2 / 24 > 0$ . Replace the `w` parameter in the algorithm in the `solver` function in `vib_undamped.py` by `w*(1 - (1./24)*w**2*dt**2)` and test how this adjustment in the numerical algorithm improves the accuracy (use  $\Delta t = 0.1$  and simulate for 80 periods, with and without adjustment of  $\omega$ ).

Filename: `vib_adjust_w`.

### Exercise 1.4: See if adaptive methods improve the phase error

Adaptive methods for solving ODEs aim at adjusting  $\Delta t$  such that the error is within a user-prescribed tolerance. Implement the equation  $u'' + u = 0$  in the `Odespy` software. Use the example from Section 3.2.11 in [4]. Run the scheme with a very low tolerance (say  $10^{-14}$ ) and for a long time, check the number of time points in the solver's mesh (`len(solver.t_all)`), and compare the phase error with that produced by the simple finite difference method from Section 1.1.2 with the same number of (equally spaced) mesh points. The question is whether it pays off to use an adaptive solver or if equally many points with a simple method gives about the same accuracy.

Filename: `vib_undamped_adaptive`.

### Exercise 1.5: Use a Taylor polynomial to compute $u^1$

As an alternative to computing  $u^1$  by (1.8), one can use a Taylor polynomial with three terms:

$$u(t_1) \approx u(0) + u'(0)\Delta t + \frac{1}{2}u''(0)\Delta t^2$$

With  $u'' = -\omega^2 u$  and  $u'(0) = 0$ , show that this method also leads to (1.8). Generalize the condition on  $u'(0)$  to be  $u'(0) = V$  and compute  $u^1$  in this case with both methods.

Filename: `vib_first_step`.

### Exercise 1.6: Find the minimal resolution of an oscillatory function

Sketch the function on a given mesh which has the highest possible frequency. That is, this oscillatory "cos-like" function has its maxima and minima at every two grid points. Find an expression for the frequency of this function, and use the result to find the largest relevant value of  $\omega\Delta t$  when  $\omega$  is the frequency of an oscillating function and  $\Delta t$  is the mesh spacing.

Filename: `vib_largest_wdt`.

### Exercise 1.7: Visualize the accuracy of finite differences for a cosine function

We introduce the error fraction

$$E = \frac{[D_tD_t u]^n}{u''(t_n)}$$

to measure the error in the finite difference approximation  $D_tD_t u$  to  $u''$ . Compute  $E$  for the specific choice of a cosine/sine function of the form  $u = \exp(i\omega t)$  and show that

$$E = \left(\frac{2}{\omega\Delta t}\right)^2 \sin^2\left(\frac{\omega\Delta t}{2}\right).$$

Plot  $E$  as a function of  $p = \omega\Delta t$ . The relevant values of  $p$  are  $[0, \pi]$  (see Exercise 1.6 for why  $p > \pi$  does not make sense). The deviation of the

curve from unity visualizes the error in the approximation. Also expand  $E$  as a Taylor polynomial in  $p$  up to fourth degree (use, e.g., `sympy`).  
Filename: `vib_plot_fd_exp_error`.

### Exercise 1.8: Verify convergence rates of the error in energy

We consider the ODE problem  $u'' + \omega^2 u = 0$ ,  $u(0) = I$ ,  $u'(0) = V$ , for  $t \in (0, T]$ . The total energy of the solution  $E(t) = \frac{1}{2}(u')^2 + \frac{1}{2}\omega^2 u^2$  should stay constant. The error in energy can be computed as explained in Section 1.6.

Make a test function in a file `test_error_conv.py`, where code from `vib_undamped.py` is imported, but the `convergence_rates` and `test_convergence_rates` functions are copied and modified to also incorporate computations of the error in energy and the convergence rate of this error. The expected rate is 2.

Filename: `test_error_conv`.

### Exercise 1.9: Use linear/quadratic functions for verification

This exercise is a generalization of Problem 1.1 to the extended model problem (1.59) where the damping term is either linear or quadratic. Solve the various subproblems and see how the results and problem settings change with the generalized ODE in case of linear or quadratic damping. By modifying the code from Problem 1.1, `sympy` will do most of the work required to analyze the generalized problem.

Filename: `vib_verify_mms`.

### Exercise 1.10: Use an exact discrete solution for verification

Write a test function in a separate file that employs the exact discrete solution (1.20) to verify the implementation of the `solver` function in the file `vib_undamped.py`.

Filename: `test_vib_undamped_exact_discrete_sol`.

### Exercise 1.11: Use analytical solution for convergence rate tests

The purpose of this exercise is to perform convergence tests of the problem (1.59) when  $s(u) = cu$ ,  $F(t) = A \sin \phi t$  and there is no damping. Find the complete analytical solution to the problem in this case (most textbooks on mechanics or ordinary differential equations list the various elements you need to write down the exact solution). Modify the `convergence_rate` function from the `vib_undamped.py` program to perform experiments with the extended model. Verify that the error is of order  $\Delta t^2$ .

Filename: `vib_conv_rate`.

### Exercise 1.12: Investigate the amplitude errors of many solvers

Use the program `vib_undamped_odespy.py` from Section 1.5.4 (utilize the function `amplitudes`) to investigate how well famous methods for 1st-order ODEs can preserve the amplitude of  $u$  in undamped oscillations. Test, for example, the 3rd- and 4th-order Runge-Kutta methods (`RK3`, `RK4`), the Crank-Nicolson method (`CrankNicolson`), the 2nd- and 3rd-order Adams-Bashforth methods (`AdamsBashforth2`, `AdamsBashforth3`), and a 2nd-order Backwards scheme (`Backward2Step`). The relevant governing equations are listed in the beginning of Section 1.5.

Filename: `vib_amplitude_errors`.

### Exercise 1.13: Minimize memory usage of a vibration solver

The program `vib.py` stores the complete solution  $u^0, u^1, \dots, u^{N_t}$  in memory, which is convenient for later plotting. Make a memory minimizing version of this program where only the last three  $u^{n+1}$ ,  $u^n$ , and  $u^{n-1}$  values are stored in memory. Write each computed  $(t_{n+1}, u^{n+1})$  pair to file. Visualize the data in the file (a cool solution is to read one line at a time and plot the  $u$  value using the line-by-line plotter in the `visualize_front_ascii` function - this technique makes it trivial to visualize very long time simulations).

Filename: `vib_memsave`.

### Exercise 1.14: Implement the solver via classes

Reimplement the `vib.py` program using a class `Problem` to hold all the physical parameters of the problem, a class `Solver` to hold the numerical parameters and compute the solution, and a class `Visualizer` to display the solution.

**Hint.** Use the ideas and examples from Section ?? and ?? in [4]. More specifically, make a superclass `Problem` for holding the scalar physical parameters of a problem and let subclasses implement the  $s(u)$  and  $F(t)$  functions as methods. Try to call up as much existing functionality in `vib.py` as possible.

Filename: `vib_class`.

### Exercise 1.15: Interpret $[D_t D_t u]^n$ as a forward-backward difference

Show that the difference  $[D_t D_t u]^n$  is equal to  $[D_t^+ D_t^- u]^n$  and  $D_t^- D_t^+ u]^n$ . That is, instead of applying a centered difference twice one can alternatively apply a mixture forward and backward differences.

Filename: `vib_DtDt_fw_bw`.

### Exercise 1.16: Use a backward difference for the damping term

As an alternative to discretizing the damping terms  $\beta u'$  and  $\beta|u'|u'$  by centered differences, we may apply backward differences:

$$\begin{aligned}[u']^n &\approx [D_t^- u]^n, \\ [|u'|u']^n &\approx [|D_t^- u| D_t^- u]^n = [|D_t^- u|^n |D_t^- u|^n].\end{aligned}$$

The advantage of the backward difference is that the damping term is evaluated using known values  $u^n$  and  $u^{n-1}$  only. Extend the `vib.py` code with a scheme based on using backward differences in the damping terms. Add statements to compare the original approach with centered difference and the new idea launched in this exercise. Perform numerical experiments to investigate how much accuracy that is lost by using the backward differences.

Filename: `vib_gen_bwdamping`.

### Exercise 1.17: Analysis of the Euler-Cromer scheme

The Euler-Cromer scheme for the model problem  $u'' + \omega^2 u = 0$ ,  $u(0) = I$ ,  $u'(0) = 0$ , is given in (1.54)-(1.53). Find the exact discrete solutions of this scheme and show that the solution for  $u^n$  coincides with that found in Section 1.4.

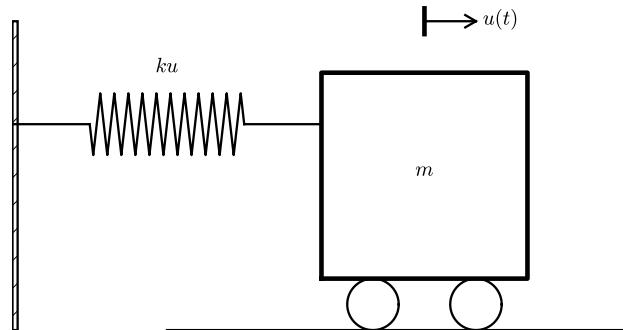
**Hint.** Use an “ansatz”  $u^n = I \exp(i\tilde{\omega}\Delta t n)$  and  $v^n = qu^n$ , where  $\tilde{\omega}$  and  $q$  are unknown parameters. The following formula is handy:

$$e^{i\tilde{\omega}\Delta t} + e^{i\tilde{\omega}(-\Delta t)} - 2 = 2(\cosh(i\tilde{\omega}\Delta t) - 1) = -4 \sin^2\left(\frac{\tilde{\omega}\Delta t}{2}\right).$$

## 1.10 Applications of vibration models

The following text derives some of the most well-known physical problems that lead to second-order ODE models of the type addressed in this book. We consider a simple spring-mass system; thereafter extended with nonlinear spring, damping, and external excitation; a spring-mass system with sliding friction; a simple and a physical (classical) pendulum; and an elastic pendulum.

### 1.10.1 Oscillating mass attached to a spring



**Fig. 1.12** Simple oscillating mass.

The most fundamental mechanical vibration system is depicted in Figure 1.12. A body with mass  $m$  is attached to a spring and can move

horizontally without friction (in the wheels). The position of the body is given by the vector  $\mathbf{r}(t) = u(t)\mathbf{i}$ , where  $\mathbf{i}$  is a unit vector in  $x$  direction. There is only one force acting on the body: a spring force  $\mathbf{F}_s = -k u \mathbf{i}$ , where  $k$  is a constant. The point  $x = 0$ , where  $u = 0$ , must therefore correspond to the body's position where the spring is neither extended nor compressed, so the force vanishes.

The basic physical principle that governs the motion of the body is Newton's second law of motion:  $\mathbf{F} = m\mathbf{a}$ , where  $\mathbf{F}$  is the sum of forces on the body,  $m$  is its mass, and  $\mathbf{a} = \ddot{\mathbf{r}}$  is the acceleration. We use the dot for differentiation with respect to time, which is usual in mechanics. Newton's second law simplifies here to  $-\mathbf{F}_s = m\ddot{\mathbf{u}}\mathbf{i}$ , which translates to

$$-ku = m\ddot{u}.$$

Two initial conditions are needed:  $u(0) = I$ ,  $\dot{u}(0) = V$ . The ODE problem is normally written as

$$m\ddot{u} + ku = 0, \quad u(0) = I, \quad \dot{u}(0) = V. \quad (1.82)$$

It is not uncommon to divide by  $m$  and introduce the frequency  $\omega = \sqrt{k/m}$ :

$$\ddot{u} + \omega^2 u = 0, \quad u(0) = I, \quad \dot{u}(0) = V. \quad (1.83)$$

This is the model problem in the first part of this chapter, with the small difference that we write the time derivative of  $u$  with a dot above, while we used  $u'$  and  $u''$  in previous parts of the book.

Since only one scalar mathematical quantity,  $u(t)$ , describes the complete motion, we say that the mechanical system has one degree of freedom (DOF).

**Scaling.** For numerical simulations it is very convenient to scale (1.83) and thereby get rid of the problem of finding relevant values for all the parameters  $m$ ,  $k$ ,  $I$ , and  $V$ . Since the amplitude of the oscillations are dictated by  $I$  and  $V$  (or more precisely,  $V/\omega$ ), we scale  $u$  by  $I$  (or  $V/\omega$  if  $I = 0$ ):

$$\bar{u} = \frac{u}{I}, \quad \bar{t} = \frac{t}{t_c}.$$

The time scale  $t_c$  is normally chosen as the inverse period  $2\pi/\omega$  or angular frequency  $1/\omega$ , most often as  $t_c = 1/\omega$ . Inserting the dimensionless quantities  $\bar{u}$  and  $\bar{t}$  in (1.83) results in the scaled problem

$$\frac{d^2\bar{u}}{dt^2} + \bar{u} = 0, \quad \bar{u}(0) = 1, \quad \frac{\bar{u}}{t}(0) = \beta = \frac{V}{I\omega},$$

where  $\beta$  is a dimensionless number. Any motion that starts from rest ( $V = 0$ ) is free of parameters in the scaled model!

**The physics.** The typical physics of the system in Figure 1.12 can be described as follows. Initially, we displace the body to some position  $I$ , say at rest ( $V = 0$ ). After releasing the body, the spring, which is extended, will act with a force  $-kI\mathbf{i}$  and pull the body to the left. This force causes an acceleration and therefore increases velocity. The body passes the point  $x = 0$ , where  $u = 0$ , and the spring will then be compressed and act with a force  $kx\mathbf{i}$  against the motion and cause retardation. At some point, the motion stops and the velocity is zero, before the spring force  $kx\mathbf{i}$  accelerates the body in positive direction. The result is that the body accelerates back and forth. As long as there is no friction forces to damp the motion, the oscillations will continue forever.

### 1.10.2 General mechanical vibrating system

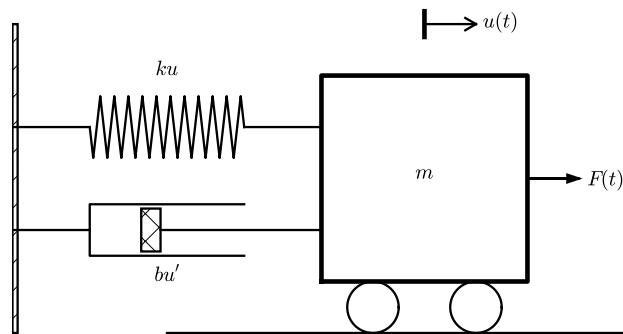


Fig. 1.13 General oscillating system.

The mechanical system in Figure 1.12 can easily be extended to the more general system in Figure 1.13, where the body is attached to a spring and a dashpot, and also subject to an environmental force  $F(t)\mathbf{i}$ . The system has still only one degree of freedom since the body can only move back and forth parallel to the  $x$  axis. The spring force was linear,  $\mathbf{F}_s = -kui\mathbf{i}$ , in Section 1.10.1, but in more general cases it can depend nonlinearly on the position. We therefore set  $\mathbf{F}_s = s(u)\mathbf{i}$ . The dashpot,

which acts as a damper, results in a force  $\mathbf{F}_d$  that depends on the body's velocity  $\dot{\mathbf{u}}$  and that always acts against the motion. The mathematical model of the force is written  $\mathbf{F}_d = f(\dot{\mathbf{u}})\mathbf{i}$ . A positive  $\dot{\mathbf{u}}$  must result in a force acting in the positive  $x$  direction. Finally, we have the external environmental force  $\mathbf{F}_e = F(t)\mathbf{i}$ .

Newton's second law of motion now involves three forces:

$$F(t)\mathbf{i} + f(\dot{\mathbf{u}})\mathbf{i} - s(\mathbf{u})\mathbf{i} = m\ddot{\mathbf{u}}\mathbf{i}.$$

The common mathematical form of the ODE problem is

$$m\ddot{\mathbf{u}} + f(\dot{\mathbf{u}}) + s(\mathbf{u}) = F(t), \quad \mathbf{u}(0) = I, \quad \dot{\mathbf{u}}(0) = V. \quad (1.84)$$

This is the generalized problem treated in the last part of the present chapter, but with prime denoting the derivative instead of the dot.

The most common models for the spring and dashpot are linear:  $f(\dot{\mathbf{u}}) = b\dot{\mathbf{u}}$  with a constant  $b \geq 0$ , and  $s(\mathbf{u}) = k\mathbf{u}$  for a constant  $k$ .

**Scaling.** A specific scaling requires specific choices of  $f$ ,  $s$ , and  $F$ . Suppose we have

$$f(\dot{\mathbf{u}}) = b|\dot{\mathbf{u}}|\dot{\mathbf{u}}, \quad s(\mathbf{u}) = ku, \quad F(t) = A \sin(\phi t).$$

We introduce dimensionless variables as usual,  $\bar{\mathbf{u}} = \mathbf{u}/u_c$  and  $\bar{t} = t/t_c$ . The scale  $u_c$  depends both on the initial conditions and  $F$ , but as time grows, the effect of the initial conditions die out and  $F$  will drive the motion. Inserting  $\bar{\mathbf{u}}$  and  $\bar{t}$  in the ODE gives

$$m \frac{u_c}{t_c^2} \frac{d^2\bar{\mathbf{u}}}{d\bar{t}^2} + b \frac{u_c^2}{t_c^2} \left| \frac{d\bar{\mathbf{u}}}{d\bar{t}} \right| \frac{d\bar{\mathbf{u}}}{d\bar{t}} + ku_c \bar{\mathbf{u}} = A \sin(\phi t_c \bar{t}).$$

We divide by  $u_c/t_c^2$  and demand the coefficients of the  $\bar{\mathbf{u}}$  and the forcing term from  $F(t)$  to have unit coefficients. This leads to the scales

$$t_c = \sqrt{\frac{m}{k}}, \quad u_c = \frac{A}{k}.$$

The scaled ODE becomes

$$\frac{d^2\bar{\mathbf{u}}}{d\bar{t}^2} + 2\beta \left| \frac{d\bar{\mathbf{u}}}{d\bar{t}} \right| \frac{d\bar{\mathbf{u}}}{d\bar{t}} + \bar{\mathbf{u}} = \sin(\gamma \bar{t}), \quad (1.85)$$

where there are two dimensionless numbers:

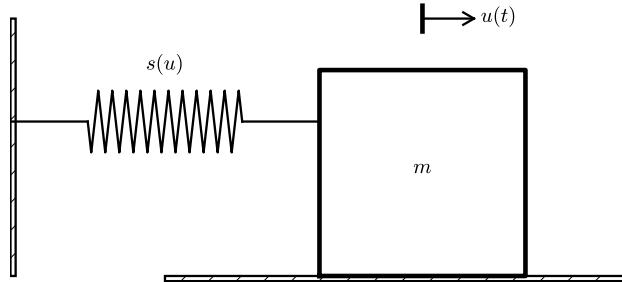
$$\beta = \frac{Ab}{2mk}, \quad \gamma = \phi \sqrt{\frac{m}{k}}.$$

The  $\beta$  number measures the size of the damping term (relative to unity) and is assumed to be small, basically because  $b$  is small. The  $\phi$  number is the ratio of the time scale of free vibrations and the time scale of the forcing. The scaled initial conditions have two other dimensionless numbers as values:

$$\bar{u}(0) = \frac{Ik}{A}, \quad \frac{d\bar{u}}{dt} = \frac{t_c}{u_c} V = \frac{V}{A} \sqrt{mk}.$$

### 1.10.3 A sliding mass attached to a spring

Consider a variant of the oscillating body in Section 1.10.1 and Figure 1.12: the body rests on a flat surface, and there is sliding friction between the body and the surface. Figure 1.14 depicts the problem.



**Fig. 1.14** Sketch of a body sliding on a surface.

The body is attached to a spring with spring force  $-s(u)\mathbf{i}$ . The friction force is proportional to the normal force on the surface,  $-mg\mathbf{j}$ , and given by  $-f(\dot{u})\mathbf{i}$ , where

$$f(\dot{u}) = \begin{cases} -\mu mg, & \dot{u} < 0, \\ \mu mg, & \dot{u} > 0, \\ 0, & \dot{u} = 0 \end{cases}$$

Here,  $\mu$  is a friction coefficient. With the signum function

$$\text{sign}(x) = \begin{cases} -1, & x < 0, \\ 1, & x > 0, \\ 0, & x = 0 \end{cases}$$

we can simply write  $f(\dot{u}) = \mu mg \text{sign}(\dot{u})$  (the sign function is implemented by `numpy.sign`).

The equation of motion becomes

$$m\ddot{u} + \mu mg\text{sign}(\dot{u}) + s(u) = 0, \quad u(0) = I, \quad \dot{u}(0) = V. \quad (1.86)$$

#### 1.10.4 A jumping washing machine

A washing machine is placed on four springs with efficient dampers. If the machine contains just a few clothes, the circular motion of the machine induces a sinusoidal external force and the machine will jump up and down if the frequency of the external force is close to the natural frequency of the machine and its spring-damper system.

**hpl 6:** Not finished. This is a good example on resonance.

#### 1.10.5 Motion of a pendulum

**Simple pendulum.** A classical problem in mechanics is the motion of a pendulum. We first consider a **simple pendulum** (sometimes also called a mathematical pendulum): a small body of mass  $m$  is attached to a massless wire and can oscillate back and forth in the gravity field. Figure 1.15 shows a sketch of the problem.

The motion is governed by Newton's 2nd law, so we need to find expressions for the forces and the acceleration. Three forces on the body are considered: an unknown force  $S$  from the wire, the gravity force  $mg$ , and an air resistance force,  $\frac{1}{2}C_D\rho A|v|v$ , hereafter called the drag force, directed against the velocity of the body. Here,  $C_D$  is a drag coefficient,  $\rho$  is the density of air,  $A$  is the cross section area of the body, and  $v$  is the magnitude of the velocity.

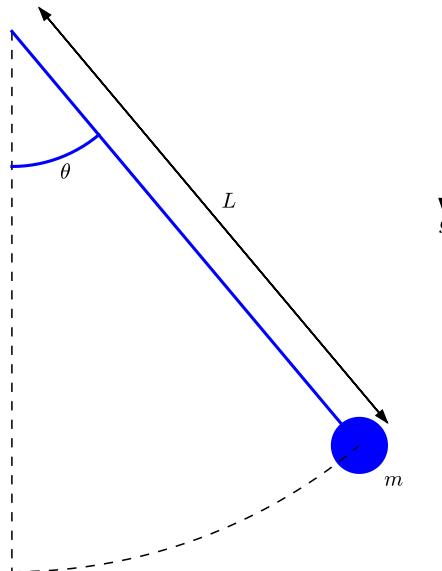
We introduce a coordinate system with polar coordinates and unit vectors  $\mathbf{i}_r$  and  $\mathbf{i}_\theta$  as shown in Figure 1.16. The position of the center of mass of the body is

$$\mathbf{r}(t) = x_0\mathbf{i} + y_0\mathbf{j} + L\mathbf{i}_r,$$

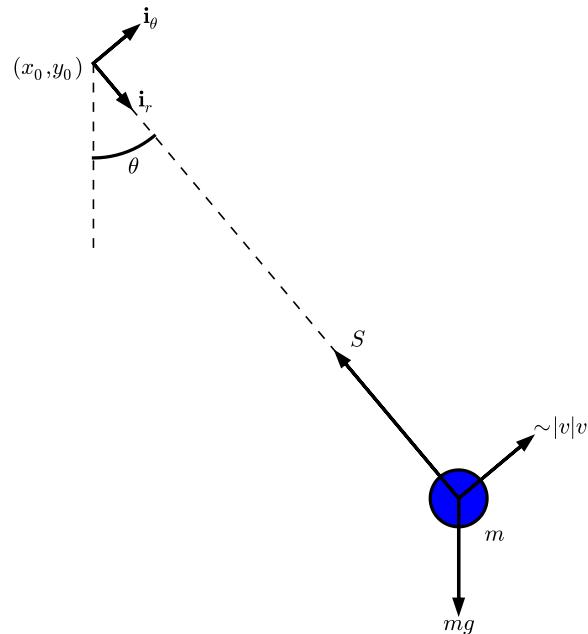
where  $\mathbf{i}$  and  $\mathbf{j}$  are unit vectors in the corresponding Cartesian coordinate system in the  $x$  and  $y$  directions, respectively. We have that  $\mathbf{i}_r = \cos\theta\mathbf{i} + \sin\theta\mathbf{j}$ .

The forces are now expressed as follows.

- Wire force:  $-S\mathbf{i}_r$



**Fig. 1.15** Sketch of a simple pendulum.



**Fig. 1.16** Forces acting on a simple pendulum.

- Gravity force:  $-mg\mathbf{j} = mg(-\sin \theta \mathbf{i}_\theta + \cos \theta \mathbf{i}_r)$
- Drag force:  $-\frac{1}{2}C_D\rho A|v|v \mathbf{i}_\theta$

Since a positive velocity means movement in the direction of  $\mathbf{i}_\theta$ , the drag force must be directed along  $-\mathbf{i}_\theta$  so it works against the motion.

The velocity of the body is found from  $\mathbf{r}$ :

$$\mathbf{v}(t) = \dot{\mathbf{r}}(t) = \frac{d}{d\theta}(x_0\mathbf{i} + y_0\mathbf{j} + L\mathbf{i}_r)\frac{d\theta}{dt} = L\dot{\theta}\mathbf{i}_\theta,$$

since  $\frac{d}{d\theta}\mathbf{i}_r = \mathbf{i}_\theta$ . It follows that  $v = |\mathbf{v}| = L\dot{\theta}$ . The acceleration is

$$\mathbf{a}(t) = \ddot{\mathbf{r}}(r) = \frac{d}{dt}(L\dot{\theta}\mathbf{i}_\theta) = L\ddot{\theta}\mathbf{i}_\theta + L\dot{\theta}\frac{d\mathbf{i}_\theta}{d\theta}\dot{\theta} = L\ddot{\theta}\mathbf{i}_\theta - L\dot{\theta}^2\mathbf{i}_r,$$

since  $\frac{d}{d\theta}\mathbf{i}_\theta = -\mathbf{i}_r$ .

Newton's 2nd law of motion becomes

$$-S\mathbf{i}_r + mg(-\sin\theta\mathbf{i}_\theta + \cos\theta\mathbf{i}_r) - \frac{1}{2}C_D\rho AL^2|\dot{\theta}|\dot{\theta}\mathbf{i}_\theta = mL\ddot{\theta}\mathbf{i}_\theta - L\dot{\theta}^2\mathbf{i}_r,$$

leading to two component equations

$$-S + mg \cos\theta = -L\dot{\theta}^2, \quad (1.87)$$

$$-mg \sin\theta - \frac{1}{2}C_D\rho AL^2|\dot{\theta}|\dot{\theta} = mL\ddot{\theta}. \quad (1.88)$$

From (1.87) we get an expression for  $S = mg \cos\theta + L\dot{\theta}^2$ , and from (1.88) we get a differential equation for the angle  $\theta(t)$ . This latter equation is ordered as

$$m\ddot{\theta} + \frac{1}{2}C_D\rho AL|\dot{\theta}|\dot{\theta} + \frac{mg}{L} \sin\theta = 0. \quad (1.89)$$

Two initial conditions are needed:  $\theta = \Theta$  and  $\dot{\theta} = \Omega$ . Normally, the pendulum motion is started from rest, which means  $\Omega = 0$ .

Equation (1.89) fits the general model used in (1.59) in Section 1.8 if we define  $u = \theta$ ,  $f(u') = \frac{1}{2}C_D\rho AL|\dot{u}|\dot{u}$ ,  $s(u) = L^{-1}mg \sin u$ , and  $F = 0$ . If the body is a sphere with radius  $R$ , we can take  $C_D = 0.4$  and  $A = \pi R^2$ . Exercise 1.21 asks you to scale the equations and carry out specific simulations with this model.

**Physical pendulum.** The motion of a compound or physical pendulum where the wire is a rod with mass, can be modeled very similarly. The governing equation is  $I\mathbf{a} = \mathbf{T}$  where  $I$  is the moment of inertia of the

entire body about the point  $(x_0, y_0)$ , and  $\mathbf{T}$  is the sum of moments of the forces with respect to  $(x_0, y_0)$ . The vector equation reads

$$\mathbf{r} \times (-S\mathbf{i}_r + mg(-\sin \theta \mathbf{i}_\theta + \cos \theta \mathbf{i}_r) - \frac{1}{2} C_D \rho A L^2 |\dot{\theta}| \dot{\theta} \mathbf{i}_\theta) = I(L \ddot{\theta} \mathbf{i}_\theta - L \dot{\theta}^2 \mathbf{i}_r).$$

The component equation in  $\mathbf{i}_\theta$  direction gives the equation of motion for  $\theta(t)$ :

$$I \ddot{\theta} + \frac{1}{2} C_D \rho A L^3 |\dot{\theta}| \dot{\theta} + mgL \sin \theta = 0. \quad (1.90)$$

### 1.10.6 Dynamic free body diagram during pendulum motion

Usually one plots the mathematical quantities as functions of time to visualize the solution of ODE models. Exercise 1.21 asks you to do this for the motion of a pendulum in the previous section. However, sometimes it is more instructive to look at other types of visualizations. For example, we have the pendulum and the free body diagram in Figures 1.15 and 1.16. We may think of these figures as animations in time instead. Especially the free body diagram will show both the motion of the pendulum *and* the size of the forces during the motion. The present section exemplifies how to make such a dynamic body diagram. Two typical snapshots of free body diagrams are displayed below (the drag force is magnified 5 times to become more visual!).



Dynamic physical sketches, coupled to the numerical solution of differential equations, requires a program to produce a sketch for the situation

at each time level. [Pysketcher](#) is such a tool. In fact (and not surprising!) Figures 1.15 and 1.16 were drawn using Pysketcher. The details of the drawings are explained in the [Pysketcher tutorial](#). Here, we outline how this type of sketch can be used to create an animated free body diagram during the motion of a pendulum.

Pysketcher is actually a layer of useful abstractions on top of standard plotting packages. This means that we in fact apply Matplotlib to make the animated free body diagram, but instead of dealing with a wealth of detailed Matplotlib commands, we can express the drawing in terms of more high-level objects, e.g., objects for the wire, angle  $\theta$ , body with mass  $m$ , arrows for forces, etc. When the position of these objects are given through variables, we can just couple those variables to the dynamic solution of our ODE and thereby make a unique drawing for each  $\theta$  value in a simulation.

**Writing the solver.** Let us start with the most familiar part of the current problem: writing the solver function. We use Odespy for this purpose. We also work with dimensionless equations. Since  $\theta$  can be viewed as dimensionless, we only need to introduce a dimensionless time, here taken as  $\bar{t} = t/\sqrt{L/g}$ . The resulting dimensionless mathematical model for  $\theta$ , the dimensionless angular velocity  $\omega$ , the dimensionless wire force  $\bar{S}$ , and the dimensionless drag force  $\bar{D}$  is then

$$\frac{d\omega}{d\bar{t}} = -\alpha|\omega|\omega - \sin \theta, \quad (1.91)$$

$$\frac{d\theta}{d\bar{t}} = \omega, \quad (1.92)$$

$$\bar{S} = \omega^2 + \cos \theta, \quad (1.93)$$

$$\bar{D} = -\alpha|\omega|\omega, \quad (1.94)$$

with

$$\alpha = \frac{C_D \rho \pi R^2 L}{2m}.$$

as a dimensionless parameter expressing the ratio of the drag force and the gravity force.

A suitable function for computing (1.91)-(1.94) is listed below.

```
def simulate(alpha, Theta, dt, T):
    import odespy

    def f(u, t, alpha):
```

```

omega, theta = u
return [-alpha*omega*abs(omega) - sin(theta),
        omega]

import numpy as np
Nt = int(round(T/float(dt)))
t = np.linspace(0, Nt*dt, Nt+1)
solver = odespy.RK4(f, f_args=[alpha])
solver.set_initial_condition([0, Theta])
u, t = solver.solve(
    t, terminate=lambda u, t, n: abs(u[n,1]) < 1E-3)
omega = u[:,0]
theta = u[:,1]
S = omega**2 + np.cos(theta)
drag = -alpha*np.abs(omega)*omega
return t, theta, omega, S, drag

```

**Drawing the free body diagram.** The `sketch` function below applies Pysketcher objects to build a diagram like that in Figure 1.16, except that we have removed the rotation point  $(x_0, y_0)$  and the unit vectors in polar coordinates as these objects are not important for an animated free body diagram.

```

import sys
try:
    from pysketcher import *
except ImportError:
    print 'Pysketcher must be installed from'
    print 'https://github.com/hplgit/pysketcher'
    sys.exit(1)

# Overall dimensions of sketch
H = 15.
W = 17.

drawing_tool.set_coordinate_system(
    xmin=0, xmax=W, ymin=0, ymax=H,
    axis=False)

def sketch(theta, S, mg, drag, t, time_level):
    """
    Draw pendulum sketch with body forces at a time level
    corresponding to time t. The drag force is in
    drag[time_level], the force in the wire is S[time_level],
    the angle is theta[time_level].
    """
    import math
    a = math.degrees(theta[time_level]) # angle in degrees
    L = 0.4*H # Length of pendulum
    P = (W/2, 0.8*H) # Fixed rotation point

    mass_pt = path.geometric_features()['end']

```

```

rod = Line(P, mass_pt)

mass = Circle(center=mass_pt, radius=L/20.)
mass.set_filled_curves(color='blue')
rod_vec = rod.geometric_features()['end'] - \
    rod.geometric_features()['start']
unit_rod_vec = unit_vec(rod_vec)
mass_symbol = Text('m', mass_pt + L/10*unit_rod_vec)

rod_start = rod.geometric_features()['start'] # Point P
vertical = Line(rod_start, rod_start + point(0,-L/3))

def set_dashed_thin_blackline(*objects):
    """Set linestyle of objects to dashed, black, width=1."""
    for obj in objects:
        obj.set_linestyle('dashed')
        obj.set_linecolor('black')
        obj.set_linewidth(1)

set_dashed_thin_blackline(vertical)
set_dashed_thin_blackline(rod)
angle = Arc_wText(r'$\theta$', rod_start, L/6, -90, a,
                  text_spacing=1/30.)

magnitude = 1.2*L/2 # length of a unit force in figure
force = mg[time_level] # constant (scaled eq: about 1)
force *= magnitude
mg_force = Force(mass_pt, mass_pt + force*point(0,-1),
                  ', text_pos='end')
force = S[time_level]
force *= magnitude
rod_force = Force(mass_pt, mass_pt - force*unit_vec(rod_vec),
                  ', text_pos='end',
                  text_spacing=(0.03, 0.01))
force = drag[time_level]
force *= magnitude
air_force = Force(mass_pt, mass_pt -
                  force*unit_vec((rod_vec[1], -rod_vec[0])),
                  ', text_pos='end',
                  text_spacing=(0.04, 0.005))

body_diagram = Composition(
    {'mg': mg_force, 'S': rod_force, 'air': air_force,
     'rod': rod, 'body': mass
     'vertical': vertical, 'theta': angle,})

body_diagram.draw(verbose=0)
drawing_tool.savefig('tmp_%04d.png' % time_level, crop=False)
# (No cropping: otherwise movies will be very strange!)

```

**Making the animated free body diagram.** It now remains to couple the simulate and sketch functions. We first run simulate:

```

from math import pi, radians, degrees
import numpy as np
alpha = 0.4
period = 2*pi    # Use small theta approximation
T = 12*period   # Simulate for 12 periods
dt = period/40  # 40 time steps per period
a = 70          # Initial amplitude in degrees
Theta = radians(a)

t, theta, omega, S, drag = simulate(alpha, Theta, dt, T)

```

The next step is to run through the time levels in the simulation and make a sketch at each level:

```

for time_level, t_ in enumerate(t):
    sketch(theta, S, mg, drag, t_, time_level)

```

The individual sketches are (by the `sketch` function) saved in files with names `tmp_%04d.png`. These can be combined to videos using (e.g.) `ffmpeg`. A complete function `animate` for running the simulation and creating video files is listed below.

```

def animate():
    # Clean up old plot files
    import os, glob
    for filename in glob.glob('tmp_*_.png') + glob.glob('movie.*'):
        os.remove(filename)
    # Solve problem
    from math import pi, radians, degrees
    import numpy as np
    alpha = 0.4
    period = 2*pi    # Use small theta approximation
    T = 12*period   # Simulate for 12 periods
    dt = period/40  # 40 time steps per period
    a = 70          # Initial amplitude in degrees
    Theta = radians(a)

    t, theta, omega, S, drag = simulate(alpha, Theta, dt, T)

    # Visualize drag force 5 times as large
    drag *= 5
    mg = np.ones(S.size)  # Gravity force (needed in sketch)

    # Draw animation
    import time
    for time_level, t_ in enumerate(t):
        sketch(theta, S, mg, drag, t_, time_level)
        time.sleep(0.2)  # Pause between each frame on the screen

    # Make videos
    prog = 'ffmpeg'
    filename = 'tmp_%04d.png'
    fps = 6

```

```

codecs = {'flv': 'flv', 'mp4': 'libx264',
          'webm': 'libvpx', 'ogg': 'libtheora'}
for ext in codecs:
    lib = codecs[ext]
    cmd = '%(prog)s -i %(filename)s -r %(fps)s' % vars()
    cmd += '-vcodec %(lib)s movie.%(ext)s' % vars()
    print(cmd)
    os.system(cmd)

```

### 1.10.7 Motion of an elastic pendulum

Consider a pendulum as in Figure 1.15, but this time the wire is elastic. The length of the wire when it is not stretched is  $L_0$ , while  $L(t)$  is the stretched length at time  $t$  during the motion.

Stretching the elastic wire a distance  $\Delta L$  gives rise to a spring force  $k\Delta L$  in the opposite direction of the stretching. Let  $\mathbf{n}$  be a unit normal vector along the wire from the point  $\mathbf{r}_0 = (x_0, y_0)$  and in the direction of  $\mathbf{i}_\theta$ , see Figure 1.16 for definition of  $(x_0, y_0)$  and  $\mathbf{i}_\theta$ . Obviously, we have  $\mathbf{n} = \mathbf{i}_\theta$ , but in this modeling of an elastic pendulum we do not need polar coordinates. Instead, it is more straightforward to develop the equation in Cartesian coordinates.

A mathematical expression for  $\mathbf{n}$  is

$$\mathbf{n} = \frac{\mathbf{r} - \mathbf{r}_0}{L(t)},$$

where  $L(t) = \|\mathbf{r} - \mathbf{r}_0\|$  is the current length of the elastic wire. The position vector  $\mathbf{r}$  in Cartesian coordinates reads  $\mathbf{r}(t) = x(t)\mathbf{i} + y(t)\mathbf{j}$ , where  $\mathbf{i}$  and  $\mathbf{j}$  are unit vectors in the  $x$  and  $y$  directions, respectively. It is convenient to introduce the Cartesian components  $n_x$  and  $n_y$  of the normal vector:

$$\mathbf{n} = \frac{\mathbf{r} - \mathbf{r}_0}{L(t)} = \frac{x(t) - x_0}{L(t)}\mathbf{i} + \frac{y(t) - y_0}{L(t)}\mathbf{j} = n_x\mathbf{i} + n_y\mathbf{j}.$$

The stretch  $\Delta L$  in the wire is

$$\Delta t = L(t) - L_0.$$

The force in the wire is then  $-S\mathbf{n} = -k\Delta L\mathbf{n}$ .

The other forces are the gravity and the air resistance, just as in Figure 1.16. The main difference is that we have a *model* for  $S$  in terms

of the motion (as soon as we have expressed  $\Delta L$  by  $\mathbf{r}$ ). For simplicity, we drop the air resistance term (but Exercise 1.23 asks you to include it).

Newton's second law of motion applied to the body now results in

$$m\ddot{\mathbf{r}} = -k(L - L_0)\mathbf{n} - mg\mathbf{j} \quad (1.95)$$

The two components of (1.95) are

$$\ddot{x} = -\frac{k}{m}(L - L_0)n_x, \quad (1.96)$$

(1.97)

$$\ddot{y} = -\frac{k}{m}(L - L_0)n_y - g. \quad (1.98)$$

**Remarks about an elastic vs a non-elastic pendulum.** Note that the derivation of the ODEs for an elastic pendulum is more straightforward than for a classical, non-elastic pendulum, since we avoid the details with polar coordinates, but instead work with Newton's second law directly in Cartesian coordinates. The reason why we can do this is that the elastic pendulum undergoes a general two-dimensional motion where all the forces are known or expressed as functions of  $x(t)$  and  $y(t)$ , such that we get two ordinary differential equations. The motion of the non-elastic pendulum, on the other hand, is constrained: the body has to move along a circular path, and the force  $S$  in the wire is unknown.

The non-elastic pendulum therefore leads to a *differential-algebraic* equation, i.e., ODEs for  $x(t)$  and  $y(t)$  combined with an extra constraint  $(x - x_0)^2 + (y - y_0)^2 = L^2$  ensuring that the motion takes place along a circular path. The extra constraint (equation) is compensated by an extra unknown force  $-S\mathbf{n}$ . Differential-algebraic equations are normally hard to solve, especially with pen and paper. Fortunately, for the non-elastic pendulum we can do a trick: in polar coordinates the unknown force  $S$  appears only in the radial component of Newton's second law, while the unknown degree of freedom for describing the motion, the angle  $\theta(t)$ , is completely governed by the azimuthal component. This allows us to decouple the unknowns  $S$  and  $\theta$ . But this is a kind of trick and not a widely applicable method. With an elastic pendulum we use straightforward reasoning with Newton's 2nd law and arrive at a standard ODE problem that (after scaling) is easy solve on a computer.

**Initial conditions.** What is the initial position of the body? We imagine that first the pendulum hangs in equilibrium in its vertical position, and

then it is displaced an angle  $\Theta$ . The equilibrium position is governed by the ODEs with the accelerations set to zero. The  $x$  component leads to  $x(t) = x_0$ , while the  $y$  component gives

$$0 = -\frac{k}{m}(L - L_0)n_y - g = \frac{k}{m}(L(0) - L_0) - g \quad \Rightarrow \quad L(0) = L_0 + mg/k,$$

since  $n_y = -1$  in this position. The corresponding  $y$  value is then from  $n_y = -1$ :

$$y(t) = y_0 - L(0) = y_0 - (L_0 + mg/k).$$

Let us now choose  $(x_0, y_0)$  such that the body is at the origin in the equilibrium position:

$$x_0 = 0, \quad y_0 = L_0 + mg/k.$$

Displacing the body an angle  $\Theta$  to the right leads to the initial position

$$x(0) = (L_0 + mg/k) \sin \Theta, \quad y(0) = (L_0 + mg/k)(1 - \cos \Theta).$$

The initial velocities can be set to zero:  $x'(0) = y'(0) = 0$ .

**The complete ODE problem.** We can summarize all the equations as follows:

$$\begin{aligned} \ddot{x} &= -\frac{k}{m}(L - L_0)n_x, \\ \ddot{y} &= -\frac{k}{m}(L - L_0)n_y - g, \\ L &= \sqrt{(x - x_0)^2 + (y - y_0)^2}, \\ n_x &= \frac{x - x_0}{L}, \\ n_y &= \frac{y - y_0}{L}, \\ x(0) &= (L_0 + mg/k) \sin \Theta, \\ x'(0) &= 0, \\ y(0) &= (L_0 + mg/k)(1 - \cos \Theta), \\ y'(0) &= 0. \end{aligned}$$

We insert  $n_x$  and  $n_y$  in the ODEs:

$$\ddot{x} = -\frac{k}{m} \left(1 - \frac{L_0}{L}\right) (x - x_0), \quad (1.99)$$

$$\ddot{y} = -\frac{k}{m} \left(1 - \frac{L_0}{L}\right) (y - y_0) - g, \quad (1.100)$$

$$L = \sqrt{(x - x_0)^2 + (y - y_0)^2}, \quad (1.101)$$

$$x(0) = (L_0 + mg/k) \sin \Theta, \quad (1.102)$$

$$x'(0) = 0, \quad (1.103)$$

$$y(0) = (L_0 + mg/k)(1 - \cos \Theta), \quad (1.104)$$

$$y'(0) = 0. \quad (1.105)$$

**Scaling.** The elastic pendulum model can be used to study both an elastic pendulum and a classic, non-elastic pendulum. The latter problem is obtained by letting  $k \rightarrow \infty$ . Unfortunately, a serious problem with the ODEs (1.99)-(1.100) is that for large  $k$ , we have a very large factor  $k/m$  multiplied by a very small number  $1 - L_0/L$ , since for large  $k$ ,  $L \approx L_0$  (very small deformations of the wire). The product is subject to significant round-off errors for many relevant physical values of the parameters. To circumvent the problem, we introduce a scaling. This will also remove physical parameters from the problem such that we end up with only one dimensionless parameter, closely related to the elasticity of the wire. Simulations can then be done by setting just this dimensionless parameter.

The characteristic length can be taken such that in equilibrium, the scaled length is unity, i.e., the characteristic length is  $L_0 + mg/k$ :

$$\bar{x} = \frac{x}{L_0 + mg/k}, \quad \bar{y} = \frac{y}{L_0 + mg/k}.$$

We must then also work with the scaled length  $\bar{L} = L/(L_0 + mg/k)$ .

Introducing  $\bar{t} = t/t_c$ , where  $t_c$  is a characteristic time we have to decide upon later, one gets

$$\begin{aligned}\frac{d^2\bar{x}}{dt^2} &= -t_c^2 \frac{k}{m} \left(1 - \frac{L_0}{L_0 + mg/k} \frac{1}{\bar{L}}\right) \bar{x}, \\ \frac{d^2\bar{y}}{dt^2} &= -t_c^2 \frac{k}{m} \left(1 - \frac{L_0}{L_0 + mg/k} \frac{1}{\bar{L}}\right) (\bar{y} - 1) - t_c^2 \frac{g}{L_0 + mg/k}, \\ \bar{L} &= \sqrt{\bar{x}^2 + (\bar{y} - 1)^2}, \\ \bar{x}(0) &= \sin \Theta, \\ \bar{x}'(0) &= 0, \\ \bar{y}(0) &= 1 - \cos \Theta, \\ \bar{y}'(0) &= 0.\end{aligned}$$

For a non-elastic pendulum with small angles, we know that the frequency of the oscillations are  $\omega = \sqrt{L/g}$ . It is therefore natural to choose a similar expression here, either the length in the equilibrium position,

$$t_c^2 = \frac{L_0 + mg/k}{g}.$$

or simply the unstretched length,

$$t_c^2 = \frac{L_0}{g}.$$

These quantities are not very different (since the elastic model is valid only for quite small elongations), so we take the latter as it is the simplest one.

The ODEs become

$$\begin{aligned}\frac{d^2\bar{x}}{dt^2} &= -\frac{L_0 k}{mg} \left(1 - \frac{L_0}{L_0 + mg/k} \frac{1}{\bar{L}}\right) \bar{x}, \\ \frac{d^2\bar{y}}{dt^2} &= -\frac{L_0 k}{mg} \left(1 - \frac{L_0}{L_0 + mg/k} \frac{1}{\bar{L}}\right) (\bar{y} - 1) - \frac{L_0}{L_0 + mg/k}, \\ \bar{L} &= \sqrt{\bar{x}^2 + (\bar{y} - 1)^2}.\end{aligned}$$

We can now identify a dimensionless number

$$\beta = \frac{L_0}{L_0 + mg/k} = \frac{1}{1 + \frac{mg}{L_0 k}},$$

which is the ratio of the unstretched length and the stretched length in equilibrium. The non-elastic pendulum will have  $\beta = 1$  ( $k \rightarrow \infty$ ). With  $\beta$  the ODEs read

$$\frac{d^2\bar{x}}{d\bar{t}^2} = -\frac{\beta}{1-\beta} \left(1 - \frac{\beta}{\bar{L}}\right) \bar{x}, \quad (1.106)$$

$$\frac{d^2\bar{y}}{d\bar{t}^2} = -\frac{\beta}{1-\beta} \left(1 - \frac{\beta}{\bar{L}}\right) (\bar{y} - 1) - \beta, \quad (1.107)$$

$$\bar{L} = \sqrt{\bar{x}^2 + (\bar{y} - 1)^2}, \quad (1.108)$$

$$\bar{x}(0) = (1 + \epsilon) \sin \Theta, \quad (1.109)$$

$$\frac{d\bar{x}}{d\bar{t}}(0) = 0, \quad (1.110)$$

$$\bar{y}(0) = 1 - (1 + \epsilon) \cos \Theta, \quad (1.111)$$

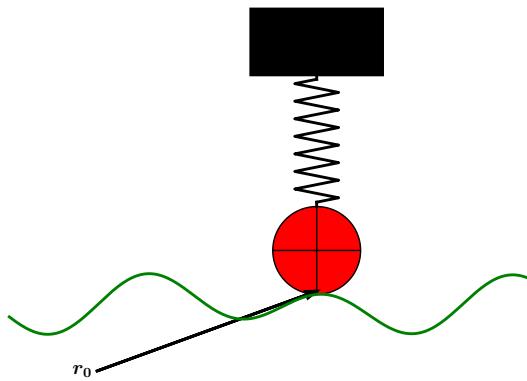
$$\frac{d\bar{y}}{d\bar{t}}(0) = 0, \quad (1.112)$$

We have here added a parameter  $\epsilon$ , which is an additional downward stretch of the wire at  $t = 0$ . This parameter makes it possible to do a desired test: vertical oscillations of the pendulum. Without  $\epsilon$ , starting the motion from  $(0, 0)$  with zero velocity will result in  $x = y = 0$  for all times (also a good test!), but with an initial stretch so the body's position is  $(0, \epsilon)$ , we will have oscillatory vertical motion with amplitude  $\epsilon$  (see Exercise 1.22).

**Remark on the non-elastic limit.** We immediately see that as  $k \rightarrow \infty$  (i.e., we obtain a non-elastic pendulum),  $\beta \rightarrow 1$ ,  $\bar{L} \rightarrow 1$ , and we have very small values  $1 - \beta \bar{L}^{-1}$  divided by very small values  $1 - \beta$  in the ODEs. However, it turns out that we can set  $\beta$  very close to one and obtain a path of the body that within the visual accuracy of a plot does not show any elastic oscillations. (Should the division of very small values become a problem, one can study the limit by L'Hospital's rule:

$$\lim_{\beta \rightarrow 1} \frac{1 - \beta \bar{L}^{-1}}{1 - \beta} = \frac{1}{\bar{L}},$$

and use the limit  $\bar{L}^{-1}$  in the ODEs for  $\beta$  values very close to 1.)



**Fig. 1.17** Sketch of one-wheel vehicle on a bumpy road.

### 1.10.8 Vehicle on a bumpy road

We consider a very simplistic vehicle, on one wheel, rolling along a bumpy road. The oscillatory nature of the road will induce an external forcing on the spring system in the vehicle and cause vibrations. Figure 1.17 outlines the situation.

To derive the equation that governs the motion, we must first establish the position vector of the black mass at the top of the spring. Suppose the spring has length  $L$  without any elongation or compression, suppose the radius of the wheel is  $R$ , and suppose the height of the black mass at the top is  $H$ . With the aid of the  $\mathbf{r}_0$  vector in Figure 1.17, the position  $\mathbf{r}$  of the center point of the mass is

$$\mathbf{r} = \mathbf{r}_0 + 2R\mathbf{j} + L\mathbf{j} + u\mathbf{j} + \frac{1}{2}H\mathbf{j}, \quad (1.113)$$

where  $u$  is the elongation or compression in the spring according to the (unknown and to be computed) vertical displacement  $u$  relative to the road. If the vehicle travels with constant horizontal velocity  $v$  and  $h(x)$  is the shape of the road, then the vector  $\mathbf{r}_0$  is

$$\mathbf{r}_0 = vt\mathbf{i} + h(vt)\mathbf{j},$$

if the motion starts from  $x = 0$  at time  $t = 0$ .

The forces on the mass is the gravity, the spring force, and an optional damping force that is proportional to the vertical velocity  $\dot{u}$ . Newton's second law of motion then tells that

$$m\ddot{\mathbf{r}} = -mg\mathbf{j} - s(u) - b\dot{u}\mathbf{j}.$$

This leads to

$$m\ddot{u} = -s(u) - b\dot{u} - mg - mh''(vt)v^2$$

To simplify a little bit, we omit the gravity force  $mg$  in comparison with the other terms. Introducing  $u'$  for  $\dot{u}$  then gives a standard damped, vibration equation with external forcing:

$$mu'' + bu' + s(u) = -mh''(vt)v^2. \quad (1.114)$$

Since the road is normally known just as a set of array values,  $h''$  must be computed by finite differences. Let  $\Delta x$  be the spacing between measured values  $h_i = h(i\Delta x)$  on the road. The discrete second-order derivative  $h''$  reads

$$q_i = \frac{h_{i-1} - 2h_i + h_{i+1}}{\Delta x^2}, \quad i = 1, \dots, N_x - 1.$$

We may for maximum simplicity set the end points as  $q_0 = q_1$  and  $q_{N_x} = q_{N_x-1}$ . The term  $-mh''(vt)v^2$  corresponds to a force with discrete time values

$$F^n = -mq_nv^2, \quad \Delta t = v^{-1}\Delta x.$$

This force can be directly used in a numerical model

$$[mD_tD_t u + bD_{2t}u + s(u) = F]^n.$$

Software for computing  $u$  and also making an animated sketch of the motion like we did in Section 1.10.6 is found in a separate project on the web: <https://github.com/hplgit/bumpy>. You may start looking at the tutorial.

### 1.10.9 Bouncing ball

A bouncing ball is a ball in free vertically fall until it impacts the ground, but during the impact, some kinetic energy is lost, and a new motion upwards with reduced velocity starts. After the motion is retarded, a new free fall starts, and the process is repeated. At some point the velocity close to the ground is so small that the ball is considered to be finally at rest.

The motion of the ball falling in air is governed by Newton's second law  $F = ma$ , where  $a$  is the acceleration of the body,  $m$  is the mass, and

$F$  is the sum of all forces. Here, we neglect the air resistance so that gravity  $-mg$  is the only force. The height of the ball is denoted by  $h$  and  $v$  is the velocity. The relations between  $h$ ,  $v$ , and  $a$ ,

$$h'(t) = v(t), \quad v'(t) = a(t),$$

combined with Newton's second law gives the ODE model

$$h''(t) = -g, \quad (1.115)$$

or expressed alternatively as a system of first-order equations:

$$v'(t) = -g, \quad (1.116)$$

$$h'(t) = v(t). \quad (1.117)$$

These equations govern the motion as long as the ball is away from the ground by a small distance  $\epsilon_h > 0$ . When  $h < \epsilon_h$ , we have two cases.

1. The ball impacts the ground, recognized by a sufficiently large negative velocity ( $v < -\epsilon_v$ ). The velocity then changes sign and is reduced by a factor  $C_R$ , known as the **coefficient of restitution**. For plotting purposes, one may set  $h = 0$ .
2. The motion stops, recognized by a sufficiently small velocity ( $|v| < \epsilon_v$ ) close to the ground.

### 1.10.10 Electric circuits

Although the term "mechanical vibrations" is used in the present book, we must mention that the same type of equations arise when modeling electric circuits. The current  $I(t)$  in a circuit with an inductor with inductance  $L$ , a capacitor with capacitance  $C$ , and overall resistance  $R$ , is governed by

$$\ddot{I} + \frac{R}{L}\dot{I} + \frac{1}{LC}I = \dot{V}(t), \quad (1.118)$$

where  $V(t)$  is the voltage source powering the circuit. This equation has the same form as the general model considered in Section refrefvib:modell2 if we set  $u = I$ ,  $f(u' = bu'$  and define  $b = R/L$ ,  $s(u) = L^{-1}C^{-1}u$ , and  $F(t) = \dot{V}(t)$ .

## 1.11 Exercises

### Exercise 1.18: Simulate resonance

We consider the scaled ODE model (1.85) from Section 1.10.2. After scaling, the amplitude of  $u$  will have a size about unity as time grows and the effect of the initial conditions die out due to damping. However, as  $\gamma \rightarrow 1$ , the amplitude of  $u$  increases, especially if  $\beta$  is small. This effect is called *resonance*. The purpose of this exercise is to explore resonance.

- a)** Figure out how the `solver` function in `vib.py` can be called for the scaled ODE (1.85).
- b)** Run  $\gamma = 5, 1.5, 1.1, 1$  for  $\beta = 0.005, 0.05, 0.2$ . For each  $\beta$  value, present an image with plots of  $u(t)$  for the four  $\gamma$  values.

Filename: `resonance`.

### Exercise 1.19: Simulate oscillations of a sliding box

Consider a sliding box on a flat surface as modeled in Section 1.10.3. As spring force we choose the nonlinear formula

$$s(u) = \frac{k}{\alpha} \tanh(\alpha u) = ku + \frac{1}{3}\alpha^2 ku^3 + \frac{2}{15}\alpha^4 ku^5 + \mathcal{O}(u^6).$$

- a)** Plot  $g(u) = \alpha^{-1} \tanh(\alpha u)$  for various values of  $\alpha$ . Assume  $u \in [-1, 1]$ .
- b)** Scale the equations using  $I$  as scale for  $u$  and  $\sqrt{m/k}$  as time scale.
- c)** Implement the scaled model in b). Run it for some values of the dimensionless parameters.

Filename: `sliding_box`.

### Exercise 1.20: Simulate a bouncing ball

Section 1.10.9 presents a model for a bouncing ball. Choose one of the two ODE formulation, (1.115) or (1.116)-(1.117), and simulate the motion of a bouncing ball. Plot  $h(t)$ . Think about how to plot  $v(t)$ .

**Hint.** A naive implementation may get stuck in repeated impacts for large time step sizes. To avoid this situation, one can introduce a state variable that holds the mode of the motion: free fall, impact, or rest. Two consecutive impacts imply that the motion has stopped.

Filename: `bouncing_ball`.

### Exercise 1.21: Simulate a simple pendulum

Simulation of simple pendulum can be carried out by using the mathematical model derived in Section 1.10.5 and calling up functionality in the `vib.py` file (i.e., solve the second-order ODE by centered finite differences).

**a)** Scale the model. Set up the dimensionless governing equation for  $\theta$  and expressions for dimensionless drag and wire forces.

**b)** Write a function for computing  $\theta$  and the dimensionless drag force and the force in the wire, using the `solver` function in the `vib.py` file. Plot these three quantities below each other (in subplots) so the graphs can be compared. Run two cases, first one in the limit of  $\Theta$  small and no drag, and then a second one with  $\Theta = 40$  degrees and  $\alpha = 0.8$ .

Filename: `simple_pendulum`.

### Exercise 1.22: Simulate an elastic pendulum

Section 1.10.7 describes a model for an elastic pendulum, resulting in a system of two ODEs. The purpose of this exercise is to implement the scaled model, test the software, and generalize the model.

**a)** Write a function `simulate` that can simulate an elastic pendulum using the scaled model. The function should have the following arguments:

```
def simulate(
    beta=0.9,                      # dimensionless parameter
    Theta=30,                       # initial angle in degrees
    epsilon=0,                        # initial stretch of wire
    num_periods=6,                   # simulate for num_periods
    time_steps_per_period=60,        # time step resolution
    plot=True,                        # make plots or not
):
```

To set the total simulation time and the time step, we use our knowledge of the scaled, classical, non-elastic pendulum:  $u'' + u = 0$ , with solution  $u = \Theta \cos \bar{t}$ . The period of these oscillations is  $P = 2\pi$  and the frequency is unity. The time for simulation is taken as `num_periods` times  $P$ . The time step is set as  $P$  divided by `time_steps_per_period`.

The `simulate` function should return the arrays of  $x$ ,  $y$ ,  $\theta$ , and  $t$ , where  $\theta = \tan^{-1}(x/(1 - y))$  is the angular displacement of the elastic pendulum corresponding to the position  $(x, y)$ .

If `plot` is `True`, make a plot of  $\bar{y}(\bar{t})$  versus  $\bar{x}(\bar{t})$ , i.e., the physical motion of the mass at  $(\bar{x}, \bar{y})$ . Use the equal aspect ratio on the axis such

that we get a physically correct picture of the motion. Also make a plot of  $\theta(\bar{t})$ , where  $\theta$  is measured in degrees. If  $\Theta < 10$  degrees, add a plot that compares the solutions of the scaled, classical, non-elastic pendulum and the elastic pendulum ( $\theta(t)$ ).

Although the mathematics here employs a bar over scaled quantities, the code should feature plain names  $x$  for  $\bar{x}$ ,  $y$  for  $\bar{y}$ , and  $t$  for  $\bar{t}$  (rather than  $x\_bar$ , etc.). These variable names make the code easier to read and compare with the mathematics.

**Hint 1.** Equal aspect ratio is set by `plt.gca().set_aspect('equal')` in Matplotlib (`import matplotlib.pyplot as plt`) and by `plot(..., daspect=[1,1,1], daspectmode='equal')` in SciTools (`import scitools.std as plt`).

**Hint 2.** If you want to use Odespy to solve the equations, order the ODEs like  $\dot{\bar{x}}, \ddot{\bar{x}}, \dot{\bar{y}}, \ddot{\bar{y}}$  such that the Euler-Cromer scheme can (also) be used (`odespy.EulerCromer`).

**b)** Write a test function for testing that  $\Theta = 0$  and  $\epsilon = 0$  gives  $x = y = 0$  for all times.

**c)** Write another test function for checking that the pure vertical motion of the elastic pendulum is correct. Start with simplifying the ODEs for pure vertical motion and show that  $\bar{y}(\bar{t})$  fulfills a vibration equation with frequency  $\sqrt{\beta/(1 - \beta)}$ . Set up the exact solution.

Write a test function that uses this special case to verify the `simulate` function. There will be numerical approximation errors present in the results from `simulate` so you have to believe in correct results and set a (low) tolerance that corresponds to the computed maximum error. Use a small  $\Delta t$  to obtain a small numerical approximation error.

**d)** Make a function `demo(beta, Theta)` for simulating an elastic pendulum with a given  $\beta$  parameter and initial angle  $\Theta$ . Use 600 time steps per period to get every accurate results, and simulate for 3 periods.

Filename: `elastic_pendulum`.

### Exercise 1.23: Simulate an elastic pendulum with air resistance

This is a continuation Exercise 1.23. Air resistance on the body with mass  $m$  can be modeled by the force  $-\frac{1}{2}\varrho C_D A |\mathbf{v}| \mathbf{v}$ , where  $C_D$  is a drag coefficient (0.2 for a sphere),  $\varrho$  is the density of air ( $1.2 \text{ kg m}^{-3}$ ),  $A$  is

the cross section area ( $A = \pi R^2$  for a sphere, where  $R$  is the radius), and  $v$  is the velocity of the body. Include air resistance in the original model, scale the model, write a function `simulate_drag` that is a copy of the `simulate` function from Exercise 1.23, but with the new ODEs included, and show plots of how air resistance influences the motion.

Filename: `elastic_pendulum_drag`.

**Remarks.** Test functions are challenging to construct for the problem with air resistance. You can reuse the tests from Exercise 1.23 for `simulate_drag`, but these tests does not verify the new terms arising from air resistance.



A very wide range of physical processes lead to wave motion, where signals are propagated through a medium in space and time, normally with little or no permanent movement of the medium itself. The shape of the signals may undergo changes as they travel through matter, but usually not so much that the signals cannot be recognized at some later point in space and time. Many types of wave motion can be described by the equation  $u_{tt} = \nabla \cdot (c^2 \nabla u) + f$ , which we will solve in the forthcoming text by finite difference methods.

## 2.1 Simulation of waves on a string

We begin our study of wave equations by simulating one-dimensional waves on a string, say on a guitar or violin. Let the string in the deformed state coincide with the interval  $[0, L]$  on the  $x$  axis, and let  $u(x, t)$  be the displacement at time  $t$  in the  $y$  direction of a point initially at  $x$ . The displacement function  $u$  is governed by the mathematical model

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2}, \quad x \in (0, L), \quad t \in (0, T] \quad (2.1)$$

$$u(x, 0) = I(x), \quad x \in [0, L] \quad (2.2)$$

$$\frac{\partial}{\partial t} u(x, 0) = 0, \quad x \in [0, L] \quad (2.3)$$

$$u(0, t) = 0, \quad t \in (0, T] \quad (2.4)$$

$$u(L, t) = 0, \quad t \in (0, T] \quad (2.5)$$

The constant  $c$  and the function  $I(x)$  must be prescribed.

Equation (2.1) is known as the one-dimensional *wave equation*. Since this PDE contains a second-order derivative in time, we need *two initial conditions*. The condition (2.2) specifies the initial shape of the string,  $I(x)$ , and (2.3) expresses that the initial velocity of the string is zero. In addition, PDEs need *boundary conditions*, given here as (2.4) and (2.5). These two conditions specify that the string is fixed at the ends, i.e., that the displacement  $u$  is zero.

The solution  $u(x, t)$  varies in space and time and describes waves that move with velocity  $c$  to the left and right.

Sometimes we will use a more compact notation for the partial derivatives to save space:

$$u_t = \frac{\partial u}{\partial t}, \quad u_{tt} = \frac{\partial^2 u}{\partial t^2}, \quad (2.6)$$

and similar expressions for derivatives with respect to other variables. Then the wave equation can be written compactly as  $u_{tt} = c^2 u_{xx}$ .

The PDE problem (2.1)-(2.5) will now be discretized in space and time by a finite difference method.

### 2.1.1 Discretizing the domain

The temporal domain  $[0, T]$  is represented by a finite number of mesh points

$$0 = t_0 < t_1 < t_2 < \cdots < t_{N_t-1} < t_{N_t} = T. \quad (2.7)$$

Similarly, the spatial domain  $[0, L]$  is replaced by a set of mesh points

$$0 = x_0 < x_1 < x_2 < \cdots < x_{N_x-1} < x_{N_x} = L. \quad (2.8)$$

One may view the mesh as two-dimensional in the  $x, t$  plane, consisting of points  $(x_i, t_n)$ , with  $i = 0, \dots, N_x$  and  $n = 0, \dots, N_t$ .

**Uniform meshes.** For uniformly distributed mesh points we can introduce the constant mesh spacings  $\Delta t$  and  $\Delta x$ . We have that

$$x_i = i\Delta x, \quad i = 0, \dots, N_x, \quad t_n = n\Delta t, \quad n = 0, \dots, N_t. \quad (2.9)$$

We also have that  $\Delta x = x_i - x_{i-1}$ ,  $i = 1, \dots, N_x$ , and  $\Delta t = t_n - t_{n-1}$ ,  $n = 1, \dots, N_t$ . Figure 2.1 displays a mesh in the  $x, t$  plane with  $N_t = 5$ ,  $N_x = 5$ , and constant mesh spacings.

### 2.1.2 The discrete solution

The solution  $u(x, t)$  is sought at the mesh points. We introduce the mesh function  $u_i^n$ , which approximates the exact solution at the mesh point  $(x_i, t_n)$  for  $i = 0, \dots, N_x$  and  $n = 0, \dots, N_t$ . Using the finite difference method, we shall develop algebraic equations for computing the mesh function.

### 2.1.3 Fulfilling the equation at the mesh points

In the finite difference method, we relax the condition that (2.1) holds at all points in the space-time domain  $(0, L) \times (0, T]$  to the requirement that the PDE is fulfilled at the *interior* mesh points only:

$$\frac{\partial^2}{\partial t^2} u(x_i, t_n) = c^2 \frac{\partial^2}{\partial x^2} u(x_i, t_n), \quad (2.10)$$

for  $i = 1, \dots, N_x - 1$  and  $n = 1, \dots, N_t - 1$ . For  $n = 0$  we have the initial conditions  $u = I(x)$  and  $u_t = 0$ , and at the boundaries  $i = 0, N_x$  we have the boundary condition  $u = 0$ .

### 2.1.4 Replacing derivatives by finite differences

The second-order derivatives can be replaced by central differences. The most widely used difference approximation of the second-order derivative is

$$\frac{\partial^2}{\partial t^2} u(x_i, t_n) \approx \frac{u_i^{n+1} - 2u_i^n + u_i^{n-1}}{\Delta t^2}.$$

It is convenient to introduce the finite difference operator notation

$$[D_t D_t u]_i^n = \frac{u_i^{n+1} - 2u_i^n + u_i^{n-1}}{\Delta t^2}.$$

A similar approximation of the second-order derivative in the  $x$  direction reads

$$\frac{\partial^2}{\partial x^2} u(x_i, t_n) \approx \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2} = [D_x D_x u]_i^n.$$

**Algebraic version of the PDE.** We can now replace the derivatives in (2.10) and get

$$\frac{u_i^{n+1} - 2u_i^n + u_i^{n-1}}{\Delta t^2} = c^2 \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2}, \quad (2.11)$$

or written more compactly using the operator notation:

$$[D_t D_t u = c^2 D_x D_x u]_i^n. \quad (2.12)$$

**Interpretation of the equation as a stencil.** A typical feature of (2.11) is that it involves  $u$  values from neighboring points only:  $u_i^{n+1}$ ,  $u_{i\pm 1}^n$ ,  $u_i^n$ , and  $u_i^{n-1}$ . The circles in Figure 2.1 illustrate such neighboring mesh points that contribute to an algebraic equation. In this particular case, we have sampled the PDE at the point  $(2, 2)$  and constructed (2.11), which then involves a coupling of  $u_2^1$ ,  $u_1^2$ ,  $u_2^2$ ,  $u_3^2$ , and  $u_2^3$ . The term *stencil* is often used about the algebraic equation at a mesh point, and the geometry of a typical stencil is illustrated in Figure 2.1. One also often refers to the algebraic equations as *discrete equations*, *(finite) difference equations* or a *finite difference scheme*.

**Algebraic version of the initial conditions.** We also need to replace the derivative in the initial condition (2.3) by a finite difference approximation. A centered difference of the type

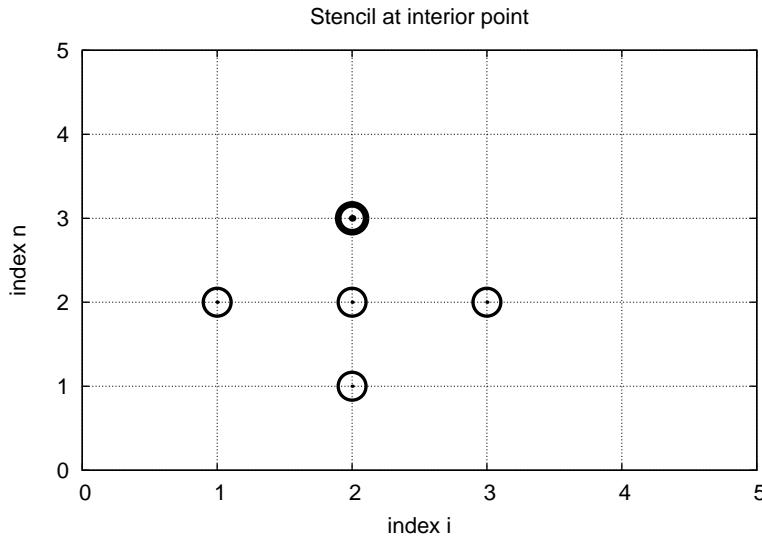
$$\frac{\partial}{\partial t} u(x_i, t_n) \approx \frac{u_i^1 - u_i^{-1}}{2\Delta t} = [D_{2t} u]_i^0,$$

seems appropriate. In operator notation the initial condition is written as

$$[D_{2t} u]_i^n = 0, \quad n = 0.$$

Writing out this equation and ordering the terms give

$$u_i^{-1} = u_i^1, \quad i = 0, \dots, N_x. \quad (2.13)$$



**Fig. 2.1** Mesh in space and time. The circles show points connected in a finite difference equation.

The other initial condition can be computed by

$$u_i^0 = I(x_i), \quad i = 0, \dots, N_x.$$

### 2.1.5 Formulating a recursive algorithm

We assume that  $u_i^n$  and  $u_i^{n-1}$  are already computed for  $i = 0, \dots, N_x$ . The only unknown quantity in (2.11) is therefore  $u_i^{n+1}$ , which we can solve for:

$$u_i^{n+1} = -u_i^{n-1} + 2u_i^n + C^2 (u_{i+1}^n - 2u_i^n + u_{i-1}^n), \quad (2.14)$$

where we have introduced the parameter

$$C = c \frac{\Delta t}{\Delta x}, \quad (2.15)$$

known as the *Courant number*.

**$C$  is the key parameter in the discrete wave equation**

We see that the discrete version of the PDE features only one parameter,  $C$ , which is therefore the key parameter that governs the quality of the numerical solution (see Section 2.10 for details). Both the primary physical parameter  $c$  and the numerical parameters  $\Delta x$  and  $\Delta t$  are lumped together in  $C$ . Note that  $C$  is a dimensionless parameter.

Given that  $u_i^{n-1}$  and  $u_i^n$  are computed for  $i = 0, \dots, N_x$ , we find new values at the next time level by applying the formula (2.14) for  $i = 1, \dots, N_x - 1$ . Figure 2.1 illustrates the points that are used to compute  $u_2^3$ . For the boundary points,  $i = 0$  and  $i = N_x$ , we apply the boundary conditions  $u_i^{n+1} = 0$ .

A problem with (2.14) arises when  $n = 0$  since the formula for  $u_i^1$  involves  $u_i^{-1}$ , which is an undefined quantity outside the time mesh (and the time domain). However, we can use the initial condition (2.13) in combination with (2.14) when  $n = 0$  to eliminate  $u_i^{-1}$  and arrive at a special formula for  $u_i^1$ :

$$u_i^1 = u_i^0 - \frac{1}{2}C^2(u_{i+1}^0 - 2u_i^0 + u_{i-1}^0). \quad (2.16)$$

Figure 2.2 illustrates how (2.16) connects four instead of five points:  $u_2^1$ ,  $u_1^0$ ,  $u_2^0$ , and  $u_3^0$ .

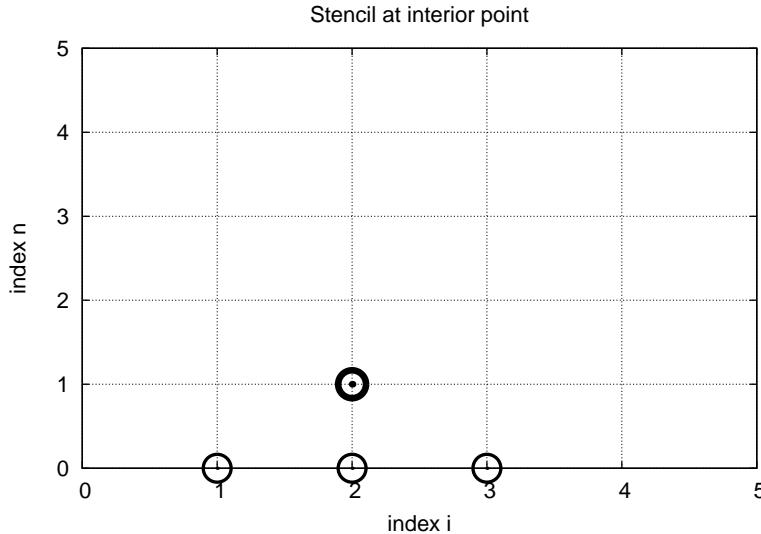
We can now summarize the computational algorithm:

1. Compute  $u_i^0 = I(x_i)$  for  $i = 0, \dots, N_x$
2. Compute  $u_i^1$  by (2.16) and set  $u_i^1 = 0$  for the boundary points  $i = 0$  and  $i = N_x$ , for  $n = 1, 2, \dots, N - 1$ ,
3. For each time level  $n = 1, 2, \dots, N_t - 1$ 
  - a. apply (2.14) to find  $u_i^{n+1}$  for  $i = 1, \dots, N_x - 1$
  - b. set  $u_i^{n+1} = 0$  for the boundary points  $i = 0, i = N_x$ .

The algorithm essentially consists of moving a finite difference stencil through all the mesh points, which can be seen as an animation in a [web page](#) or a [movie file](#).

### 2.1.6 Sketch of an implementation

In a Python implementation of this algorithm, we use the array elements  $u[i]$  to store  $u_i^{n+1}$ ,  $u\_1[i]$  to store  $u_i^n$ , and  $u\_2[i]$  to store  $u_i^{n-1}$ . Our



**Fig. 2.2** Modified stencil for the first time step.

naming convention is to use  $u$  for the unknown new spatial field to be computed,  $u_{-1}$  as the solution one time step back in time,  $u_{-2}$  as the solution two time steps back in time and so forth.

The algorithm only involves the three most recent time levels, so we need only three arrays for  $u_i^{n+1}$ ,  $u_i^n$ , and  $u_i^{n-1}$ ,  $i = 0, \dots, N_x$ . Storing all the solutions in a two-dimensional array of size  $(N_x + 1) \times (N_t + 1)$  would be possible in this simple one-dimensional PDE problem, but is normally out of the question in three-dimensional (3D) and large two-dimensional (2D) problems. We shall therefore, in all our PDE solving programs, have the unknown in memory at as few time levels as possible.

The following Python snippet realizes the steps in the computational algorithm.

```
# Given mesh points as arrays x and t (x[i], t[n])
dx = x[1] - x[0]
dt = t[1] - t[0]
C = c*dt/dx           # Courant number
Nt = len(t)-1
C2 = C**2             # Help variable in the scheme

# Set initial condition u(x,0) = I(x)
for i in range(0, Nx+1):
    u_1[i] = I(x[i])

# Apply special formula for first step, incorporating du/dt=0
```

```

for i in range(1, Nx):
    u[i] = u_1[i] - 0.5*C**2(u_1[i+1] - 2*u_1[i] + u_1[i-1])
u[0] = 0; u[Nx] = 0 # Enforce boundary conditions

# Switch variables before next step
u_2[:,], u_1[:,] = u_1, u

for n in range(1, Nt):
    # Update all inner mesh points at time t[n+1]
    for i in range(1, Nx):
        u[i] = 2u_1[i] - u_2[i] - \
            C**2(u_1[i+1] - 2*u_1[i] + u_1[i-1])

    # Insert boundary conditions
    u[0] = 0; u[Nx] = 0

    # Switch variables before next step
    u_2[:,], u_1[:,] = u_1, u

```

## 2.2 Verification

Before implementing the algorithm, it is convenient to add a source term to the PDE (2.1), since that gives us more freedom in finding test problems for verification. Physically, a source term acts as a generator for waves in the interior of the domain.

### 2.2.1 A slightly generalized model problem

We now address the following extended initial-boundary value problem for one-dimensional wave phenomena:

$$u_{tt} = c^2 u_{xx} + f(x, t), \quad x \in (0, L), t \in (0, T] \quad (2.17)$$

$$u(x, 0) = I(x), \quad x \in [0, L] \quad (2.18)$$

$$u_t(x, 0) = V(x), \quad x \in [0, L] \quad (2.19)$$

$$u(0, t) = 0, \quad t > 0 \quad (2.20)$$

$$u(L, t) = 0, \quad t > 0 \quad (2.21)$$

Sampling the PDE at  $(x_i, t_n)$  and using the same finite difference approximations as above, yields

$$[D_t D_t u = c^2 D_x D_x u + f]_i^n. \quad (2.22)$$

Writing this out and solving for the unknown  $u_i^{n+1}$  results in

$$u_i^{n+1} = -u_i^{n-1} + 2u_i^n + C^2(u_{i+1}^n - 2u_i^n + u_{i-1}^n) + \Delta t^2 f_i^n. \quad (2.23)$$

The equation for the first time step must be rederived. The discretization of the initial condition  $u_t = V(x)$  at  $t = 0$  becomes

$$[D_{2t} u = V]_i^0 \Rightarrow u_i^{-1} = u_i^1 - 2\Delta t V_i,$$

which, when inserted in (2.23) for  $n = 0$ , gives the special formula

$$u_i^1 = u_i^0 - \Delta t V_i + \frac{1}{2} C^2 (u_{i+1}^0 - 2u_i^0 + u_{i-1}^0) + \frac{1}{2} \Delta t^2 f_i^0. \quad (2.24)$$

## 2.2.2 Using an analytical solution of physical significance

Many wave problems feature sinusoidal oscillations in time and space. For example, the original PDE problem (2.1)-(2.5) allows an exact solution

$$u_e(x, t) = A \sin\left(\frac{\pi}{L}x\right) \cos\left(\frac{\pi}{L}ct\right). \quad (2.25)$$

This  $u_e$  fulfills the PDE with  $f = 0$ , boundary conditions  $u_e(0, t) = u_e(L, 0) = 0$ , as well as initial conditions  $I(x) = A \sin\left(\frac{\pi}{L}x\right)$  and  $V = 0$ .

### How to use exact solutions for verification

It is common to use such exact solutions of physical interest to verify implementations. However, the numerical solution  $u_i^n$  will only be an approximation to  $u_e(x_i, t_n)$ . We have no knowledge of the precise size of the error in this approximation, and therefore we can never know if discrepancies between  $u_i^n$  and  $u_e(x_i, t_n)$  are caused by mathematical approximations or programming errors. In particular, if a plot of the computed solution  $u_i^n$  and the exact one (2.25) look similar, many are tempted to claim that the implementation works. However, even if color plots look nice and the accuracy is “deemed good”, there can still be serious programming errors present!

The only way to use exact physical solutions like (2.25) for serious and thorough verification is to run a series of finer and finer meshes, measure the integrated error in each mesh, and from this information estimate the empirical convergence rate of the method.

An introduction to the computing of convergence rates is given in Section 3.1.6 in [4]. There is also a detailed example on computing convergence rates in Section 1.2.2.

In the present problem, one expects the method to have a convergence rate of 2 (see Section 2.10), so if the computed rates are close to 2 on a sufficiently fine mesh, we have good evidence that the implementation is free of programming mistakes.

### 2.2.3 Manufactured solution

One problem with the exact solution (2.25) is that it requires a simplification ( $V = 0, f = 0$ ) of the implemented problem (2.17)-(2.21). An advantage of using a *manufactured solution* is that we can test all terms in the PDE problem. The idea of this approach is to set up some chosen solution and fit the source term, boundary conditions, and initial conditions to be compatible with the chosen solution. Given that our boundary conditions in the implementation are  $u(0, t) = u(L, t) = 0$ , we must choose a solution that fulfills these conditions. One example is

$$u_e(x, t) = x(L - x) \sin t.$$

Inserted in the PDE  $u_{tt} = c^2 u_{xx} + f$  we get

$$-x(L - x) \sin t = -c^2 2 \sin t + f \quad \Rightarrow f = (2c^2 - x(L - x)) \sin t.$$

The initial conditions become

$$\begin{aligned} u(x, 0) &= I(x) = 0, \\ u_t(x, 0) &= V(x) = x(L - x). \end{aligned}$$

To verify the code, we compute the convergence rates in a series of simulations, letting each simulation use a finer mesh than the previous one. Such empirical estimation of convergence rates tests relies on an

assumption that some measure  $E$  of the numerical error is related to the discretization parameters through

$$E = C_t \Delta t^r + C_x \Delta x^p,$$

where  $C_t$ ,  $C_x$ ,  $r$ , and  $p$  are constants. The constants  $r$  and  $p$  are known as the *convergence rates* in time and space, respectively. From the accuracy in the finite difference approximations, we expect  $r = p = 2$ , since the error terms are of order  $\Delta t^2$  and  $\Delta x^2$ . This is confirmed by truncation error analysis and other types of analysis.

By using an exact solution of the PDE problem, we will next compute the error measure  $E$  on a sequence of refined meshes and see if the rates  $r = p = 2$  are obtained. We will not be concerned with estimating the constants  $C_t$  and  $C_x$ .

It is advantageous to introduce a single discretization parameter  $h = \Delta t = \hat{c} \Delta x$  for some constant  $\hat{c}$ . Since  $\Delta t$  and  $\Delta x$  are related through the Courant number,  $\Delta t = C \Delta x / c$ , we set  $h = \Delta t$ , and then  $\Delta x = hc/C$ . Now the expression for the error measure is greatly simplified:

$$E = C_t \Delta t^r + C_x \Delta x^r = C_t h^r + C_x \left( \frac{c}{C} \right)^r h^r = D h^r, \quad D = C_t + C_x \left( \frac{c}{C} \right)^r.$$

We choose an initial discretization parameter  $h_0$  and run experiments with decreasing  $h$ :  $h_i = 2^{-i} h_0$ ,  $i = 1, 2, \dots, m$ . Halving  $h$  in each experiment is not necessary, but it is a common choice. For each experiment we must record  $E$  and  $h$ . A standard choice of error measure is the  $\ell^2$  or  $\ell^\infty$  norm of the error mesh function  $e_i^n$ :

$$E = \|e_i^n\|_{\ell^2} = \left( \Delta t \Delta x \sum_{n=0}^{N_t} \sum_{i=0}^{N_x} (e_i^n)^2 \right)^{\frac{1}{2}}, \quad e_i^n = u_e(x_i, t_n) - u_i^n, \quad (2.26)$$

$$E = \|e_i^n\|_{\ell^\infty} = \max_{i,n} |e_i^n|. \quad (2.27)$$

In Python, one can compute  $\sum_i (e_i^n)^2$  at each time step and accumulate the value in some sum variable, say `e2_sum`. At the final time step one can do `sqrt(dt*dx*e2_sum)`. For the  $\ell^\infty$  norm one must compare the maximum error at a time level (`e.max()`) with the global maximum over the time domain: `e_max = max(e_max, e.max())`.

An alternative error measure is to use a spatial norm at one time step only, e.g., the end time  $T$  ( $n = N_t$ ):

$$E = \|e_i^n\|_{\ell^2} = \left( \Delta x \sum_{i=0}^{N_x} (e_i^n)^2 \right)^{\frac{1}{2}}, \quad e_i^n = u_e(x_i, t_n) - u_i^n, \quad (2.28)$$

$$E = \|e_i^n\|_{\ell^\infty} = \max_{0 \leq i \leq N_x} |e_i^n|. \quad (2.29)$$

The important issue is that our error measure  $E$  must be one number that represents the error in the simulation.

Let  $E_i$  be the error measure in experiment (mesh) number  $i$  and let  $h_i$  be the corresponding discretization parameter ( $h$ ). With the error model  $E_i = Dh_i^r$ , we can estimate  $r$  by comparing two consecutive experiments:

$$\begin{aligned} E_{i+1} &= Dh_{i+1}^r, \\ E_i &= Dh_i^r. \end{aligned}$$

Dividing the two equations eliminates the (uninteresting) constant  $D$ . Thereafter, solving for  $r$  yields

$$r = \frac{\ln E_{i+1}/E_i}{\ln h_{i+1}/h_i}.$$

Since  $r$  depends on  $i$ , i.e., which simulations we compare, we add an index to  $r$ :  $r_i$ , where  $i = 0, \dots, m-2$ , if we have  $m$  experiments:  $(h_0, E_0), \dots, (h_{m-1}, E_{m-1})$ .

In our present discretization of the wave equation we expect  $r = 2$ , and hence the  $r_i$  values should converge to 2 as  $i$  increases.

## 2.2.4 Constructing an exact solution of the discrete equations

With a manufactured or known analytical solution, as outlined above, we can estimate convergence rates and see if they have the correct asymptotic behavior. Experience shows that this is a quite good verification technique in that many common bugs will destroy the convergence rates. A significantly better test though, would be to check that the numerical solution is exactly what it should be. This will in general require exact knowledge of the numerical error, which we do not normally have (although we in Section 2.10 establish such knowledge in simple cases). However, it is possible to look for solutions where we can show that

the numerical error vanishes, i.e., the solution of the original continuous PDE problem is also a solution of the discrete equations. This property often arises if the exact solution of the PDE is a lower-order polynomial. (Truncation error analysis leads to error measures that involve derivatives of the exact solution. In the present problem, the truncation error involves 4th-order derivatives of  $u$  in space and time. Choosing  $u$  as a polynomial of degree three or less will therefore lead to vanishing error.)

We shall now illustrate the construction of an exact solution to both the PDE itself and the discrete equations. Our chosen manufactured solution is quadratic in space and linear in time. More specifically, we set

$$u_e(x, t) = x(L - x)(1 + \frac{1}{2}t), \quad (2.30)$$

which by insertion in the PDE leads to  $f(x, t) = 2(1 + t)c^2$ . This  $u_e$  fulfills the boundary conditions  $u = 0$  and demands  $I(x) = x(L - x)$  and  $V(x) = \frac{1}{2}x(L - x)$ .

To realize that the chosen  $u_e$  is also an exact solution of the discrete equations, we first remind ourselves that  $t_n = n\Delta t$  before we establish that

$$[D_tD_tt^2]^n = \frac{t_{n+1}^2 - 2t_n^2 + t_{n-1}^2}{\Delta t^2} = (n+1)^2 - 2n^2 + (n-1)^2 = 2, \quad (2.31)$$

$$[D_tD_tt]^n = \frac{t_{n+1} - 2t_n + t_{n-1}}{\Delta t^2} = \frac{((n+1) - 2n + (n-1))\Delta t}{\Delta t^2} = 0. \quad (2.32)$$

Hence,

$$[D_tD_tu_e]_i^n = x_i(L - x_i)[D_tD_t(1 + \frac{1}{2}t)]^n = x_i(L - x_i)\frac{1}{2}[D_tD_tt]^n = 0.$$

Similarly, we get that

$$\begin{aligned} [D_xD_xu_e]_i^n &= (1 + \frac{1}{2}t_n)[D_xD_x(xL - x^2)]_i \\ &= (1 + \frac{1}{2}t_n)[LD_xD_xx - D_xD_xx^2]_i \\ &= -2(1 + \frac{1}{2}t_n). \end{aligned}$$

Now,  $f_i^n = 2(1 + \frac{1}{2}t_n)c^2$ , which results in

$$[D_t D_t u_e - c^2 D_x D_x u_e - f]_i^n = 0 + c^2 2(1 + \frac{1}{2}t_n) + 2(1 + \frac{1}{2}t_n)c^2 = 0.$$

Moreover,  $u_e(x_i, 0) = I(x_i)$ ,  $\partial u_e / \partial t = V(x_i)$  at  $t = 0$ , and  $u_e(x_0, t) = u_e(x_{N_x}, t) = 0$ . Also the modified scheme for the first time step is fulfilled by  $u_e(x_i, t_n)$ .

Therefore, the exact solution  $u_e(x, t) = x(L - x)(1 + t/2)$  of the PDE problem is also an exact solution of the discrete problem. We can use this result to check that the computed  $u_i^n$  values from an implementation equals  $u_e(x_i, t_n)$  within machine precision, *regardless of the mesh spacings  $\Delta x$  and  $\Delta t$* ! Nevertheless, there might be stability restrictions on  $\Delta x$  and  $\Delta t$ , so the test can only be run for a mesh that is compatible with the stability criterion (which in the present case is  $C \leq 1$ , to be derived later).

### Notice

A product of quadratic or linear expressions in the various independent variables, as shown above, will often fulfill both the PDE problem and the discrete equations, and can therefore be very useful solutions for verifying implementations.

However, for 1D wave equations of the type  $u_{tt} = c^2 u_{xx}$  we shall see that there is always another much more powerful way of generating exact solutions (which consists in just setting  $C = 1$  (!), as shown in Section 2.10).

## 2.3 Implementation

This section presents the complete computational algorithm, its implementation in Python code, animation of the solution, and verification of the implementation.

A real implementation of the basic computational algorithm from Sections 2.1.5 and 2.1.6 can be encapsulated in a function, taking all the input data for the problem as arguments. The physical input data consists of  $c$ ,  $I(x)$ ,  $V(x)$ ,  $f(x, t)$ ,  $L$ , and  $T$ . The numerical input is the mesh parameters  $\Delta t$  and  $\Delta x$ .

Instead of specifying  $\Delta t$  and  $\Delta x$ , we can specify one of them and the Courant number  $C$  instead, since having explicit control of the Courant number is convenient when investigating the numerical method. Many find it natural to prescribe the resolution of the spatial grid and set  $N_x$ . The solver function can then compute  $\Delta t = CL/(cN_x)$ . However, for comparing  $u(x, t)$  curves (as functions of  $x$ ) for various Courant numbers it is more convenient to keep  $\Delta t$  fixed for all  $C$  and let  $\Delta x$  vary according to  $\Delta x = c\Delta t/C$ . With  $\Delta t$  fixed, all frames correspond to the same time  $t$ , and this simplifies animations that compare simulations with different mesh resolutions. Plotting functions of  $x$  with different spatial resolution is trivial, so it is easier to let  $\Delta x$  vary in the simulations than  $\Delta t$ .

### 2.3.1 Callback function for user-specific actions

The solution at all spatial points at a new time level is stored in an array  $u$  of length  $N_x + 1$ . We need to decide what do to with this solution, e.g., visualize the curve, analyze the values, or write the array to file for later use. The decision about what to do is left to the user in the form of a user-supplied function

```
user_action(u, x, t, n)
```

where  $u$  is the solution at the spatial points  $x$  at time  $t[n]$ . The `user_action` function is called from the solver at each time level  $n$ .

If the user wants to plot the solution or store the solution at a time point, she needs to write such a function and take appropriate actions inside it. We will show examples on many such `user_action` functions.

Since the solver function makes calls back to the user's code via such a function, this type of function is called a *callback function*. When writing general software, like our solver function, which also needs to carry out special problem-dependent actions (like visualization), it is a common technique to leave those actions to user-supplied callback functions.

### 2.3.2 The solver function

A first attempt at a solver function is listed below.

```
import numpy as np

def solver(I, V, f, c, L, dt, C, T, user_action=None):
    """Solve u_tt=c^2*u_xx + f on (0,L)x(0,T]."""


```

```

Nt = int(round(T/dt))
t = np.linspace(0, Nt*dt, Nt+1)    # Mesh points in time
dx = dt*c/float(C)
Nx = int(round(L/dx))
x = np.linspace(0, L, Nx+1)          # Mesh points in space
C2 = C**2                           # Help variable in the scheme
# Make sure dx and dt are compatible with x and t
dx = x[1] - x[0]
dt = t[1] - t[0]

if f is None or f == 0 :
    f = lambda x, t: 0
if V is None or V == 0:
    V = lambda x: 0

u    = np.zeros(Nx+1)    # Solution array at new time level
u_1 = np.zeros(Nx+1)    # Solution at 1 time level back
u_2 = np.zeros(Nx+1)    # Solution at 2 time levels back

import time; t0 = time.clock() # for measuring CPU time

# Load initial condition into u_1
for i in range(0,Nx+1):
    u_1[i] = I(x[i])

if user_action is not None:
    user_action(u_1, x, t, 0)

# Special formula for first time step
n = 0
for i in range(1, Nx):
    u[i] = u_1[i] + dt*V(x[i]) + \
        0.5*C2*(u_1[i-1] - 2*u_1[i] + u_1[i+1]) + \
        0.5*dt**2*f(x[i], t[n])
u[0] = 0; u[Nx] = 0

if user_action is not None:
    user_action(u, x, t, 1)

# Switch variables before next step
u_2[:] = u_1; u_1[:] = u

for n in range(1, Nt):
    # Update all inner points at time t[n+1]
    for i in range(1, Nx):
        u[i] = - u_2[i] + 2*u_1[i] + \
            C2*(u_1[i-1] - 2*u_1[i] + u_1[i+1]) + \
            dt**2*f(x[i], t[n])

    # Insert boundary conditions
    u[0] = 0; u[Nx] = 0
    if user_action is not None:
        if user_action(u, x, t, n+1):
            break

```

```

# Switch variables before next step
u_2[:] = u_1;  u_1[:] = u

cpu_time = time.clock() - t0
return u, x, t, cpu_time

```

A remark on the computation of  $\text{dx}$  and  $\text{dt}$  is probably necessary. Although we give  $\text{dt}$  and compute  $\text{dx}$  via  $C$  and  $c$ , the resulting  $t$  and  $x$  meshes do not necessarily correspond exactly to these values because of rounding errors. To explicitly ensure that  $\text{dx}$  and  $\text{dt}$  correspond to the cell sizes in  $x$  and  $t$ , we recompute the values.

### 2.3.3 Verification: exact quadratic solution

We use the test problem derived in Section 2.2.1 for verification. Below is a unit test based on this test problem and realized as a proper *test function* compatible with the unit test frameworks nose or pytest.

```

def test_quadratic():
    """Check that  $u(x,t)=x(L-x)(1+t/2)$  is exactly reproduced."""

    def u_exact(x, t):
        return x*(L-x)*(1 + 0.5*t)

    def I(x):
        return u_exact(x, 0)

    def V(x):
        return 0.5*u_exact(x, 0)

    def f(x, t):
        return 2*(1 + 0.5*t)*c**2

    L = 2.5
    c = 1.5
    C = 0.75
    Nx = 6 # Very coarse mesh for this exact test
    dt = C*(L/Nx)/c
    T = 18

    def assert_no_error(u, x, t, n):
        u_e = u_exact(x, t[n])
        diff = np.abs(u - u_e).max()
        tol = 1E-13
        assert diff < tol

    solver(I, V, f, c, L, dt, C, T,
           user_action=assert_no_error)

```

When this function resides in the file `wave1D_u0.py`, one can run `pytest` to check that all test functions with names `test_*`() in this file work:

---

Terminal>	Terminal
-----------	----------

```
Terminal> py.test -s -v wave1D_u0.py
```

---

### 2.3.4 Visualization: animating the solution

Now that we have verified the implementation it is time to do a real computation where we also display the evolution of the waves on the screen. Since the `solver` function knows nothing about what type of visualizations we may want, it calls the callback function `user_action(u, x, t, n)`. We must therefore write this function and find the proper statements for plotting the solution.

**Function for administering the simulation.** The following `viz` function

1. defines a `user_action` callback function for plotting the solution at each time level,
2. calls the `solver` function, and
3. combines all the plots (in files) to video in different formats.

```
def viz(
    I, V, f, c, L, dt, C, T,    # PDE parameters
    umin, umax,                  # Interval for u in plots
    animate=True,                # Simulation with animation?
    tool='matplotlib',           # 'matplotlib' or 'scitoools'
    solver_function=solver,      # Function with numerical algorithm
):
    """Run solver and visualize u at each time level."""

    def plot_u_st(u, x, t, n):
        """user_action function for solver."""
        plt.plot(x, u, 'r-',
                  xlabel='x', ylabel='u',
                  axis=[0, L, umin, umax],
                  title='t=%f' % t[n], show=True)
        # Let the initial condition stay on the screen for 2
        # seconds, else insert a pause of 0.2 s between each plot
        time.sleep(2) if t[n] == 0 else time.sleep(0.2)
        plt.savefig('frame_%04d.png' % n)  # for movie making

    class PlotMatplotlib:
        def __call__(self, u, x, t, n):
            """user_action function for solver."""
            if n == 0:
                plt.ion()
```

```

        self.lines = plt.plot(x, u, 'r-')
        plt.xlabel('x'); plt.ylabel('u')
        plt.axis([0, L, umin, umax])
        plt.legend(['t=%f' % t[n]], loc='lower left')
    else:
        self.lines[0].set_ydata(u)
        plt.legend(['t=%f' % t[n]], loc='lower left')
        plt.draw()
        time.sleep(2) if t[n] == 0 else time.sleep(0.2)
        plt.savefig('tmp_%04d.png' % n) # for movie making

if tool == 'matplotlib':
    import matplotlib.pyplot as plt
    plot_u = PlotMatplotlib()
elif tool == 'scitools':
    import scitools.std as plt # scitools.easyviz interface
    plot_u = plot_u_st
    import time, glob, os

# Clean up old movie frames
for filename in glob.glob('tmp_*.png'):
    os.remove(filename)

# Call solver and do the simulation
user_action = plot_u if animate else None
u, x, t, cpu = solver_function(
    I, V, f, c, L, dt, C, T, user_action)

# Make video files
fps = 4 # frames per second
codec2ext = dict(flv='flv', libx264='mp4', libvpx='webm',
                 libtheora='ogg') # video formats
filespec = 'tmp_%04d.png'
movie_program = 'ffmpeg' # or 'avconv'
for codec in codec2ext:
    ext = codec2ext[codec]
    cmd = '%(movie_program)s -r %(fps)d -i %(filespec)s '\
          '-vcodec %(codec)s movie.%s' % vars()
    os.system(cmd)

if tool == 'scitools':
    # Make an HTML play for showing the animation in a browser
    plt.movie('tmp_*.png', encoder='html', fps=fps,
              output_file='movie.html')
return cpu

```

**Dissection of the code.** The `viz` function can either use SciTools or Matplotlib for visualizing the solution. The `user_action` function based on SciTools is called `plot_u_st`, while the `user_action` function based on Matplotlib is a bit more complicated as it is realized as a class and needs statements that differ from those for making static plots. SciTools can utilize both Matplotlib and Gnuplot (and many other

plotting programs) for doing the graphics, but Gnuplot is a relevant choice for large  $N_x$  or in two-dimensional problems as Gnuplot is significantly faster than Matplotlib for screen animations.

A function inside another function, like `plot_u_st` in the above code segment, has access to *and remembers* all the local variables in the surrounding code inside the `viz` function (!). This is known in computer science as a *closure* and is very convenient to program with. For example, the `plt` and `time` modules defined outside `plot_u` are accessible for `plot_u_st` when the function is called (as `user_action`) in the `solver` function. Some may think, however, that a class instead of a closure is a cleaner and easier-to-understand implementation of the user action function, see Section 2.8.

The `plot_u_st` function just makes a standard SciTools `plot` command for plotting `u` as a function of `x` at time `t[n]`. To achieve a smooth animation, the `plot` command should take keyword arguments instead of being broken into separate calls to `xlabel`, `ylabel`, `axis`, `time`, and `show`. Several `plot` calls will automatically cause an animation on the screen. In addition, we want to save each frame in the animation to file. We then need a filename where the frame number is padded with zeros, here `tmp_0000.png`, `tmp_0001.png`, and so on. The proper `printf` construction is then `tmp_%04d.png`. Section 1.3.2 contains more basic information on making animations.

The solver is called with an argument `plot_u` as `user_function`. If the user chooses to use SciTools, `plot_u` is the `plot_u_st` callback function, but for Matplotlib it is an instance of the class `PlotMatplotlib`. Also this class makes use of variables defined in the `viz` function: `plt` and `time`. With Matplotlib, one has to make the first plot the standard way, and then update the `y` data in the plot at every time level. The update requires active use of the returned value from `plt.plot` in the first plot. This value would need to be stored in a local variable if we were to use a closure for the `user_action` function when doing the animation with Matplotlib. It is much easier to store the variable as a class attribute `self.lines`. Since the class is essentially a function, we implement the function as the special method `__call__` such that the instance `plot_u(u, x, t, n)` can be called as a standard callback function from `solver`.

**Making movie files.** From the `frame_*.png` files containing the frames in the animation we can make video files. Section 1.3.2 presents basic information on how to use the `ffmpeg` (or `avconv`) program for producing video files in different modern formats: Flash, MP4, Webm, and Ogg.

The `viz` function creates an `ffmpeg` or `avconv` command with the proper arguments for each of the formats Flash, MP4, WebM, and Ogg. The task is greatly simplified by having a `codec2ext` dictionary for mapping video codec names to filename extensions. As mentioned in Section 1.3.2, only two formats are actually needed to ensure that all browsers can successfully play the video: MP4 and WebM.

Some animations having a large number of plot files may not be properly combined into a video using `ffmpeg` or `avconv`. A method that always works is to play the PNG files as an animation in a browser using JavaScript code in an HTML file. The SciTools package has a function `movie` (or a stand-alone command `scitools movie`) for creating such an HTML player. The `plt.movie` call in the `viz` function shows how the function is used. The file `movie.html` can be loaded into a browser and features a user interface where the speed of the animation can be controlled. Note that the movie in this case consists of the `movie.html` file and all the frame files `tmp_*.png`.

**Skipping frames for animation speed.** Sometimes the time step is small and  $T$  is large, leading to an inconveniently large number of plot files and a slow animation on the screen. The solution to such a problem is to decide on a total number of frames in the animation, `num_frames`, and plot the solution only for every `skip_frame` frames. For example, setting `skip_frame=5` leads to plots of every 5 frames. The default value `skip_frame=1` plots every frame. The total number of time levels (i.e., maximum possible number of frames) is the length of `t`, `t.size` (or `len(t)`), so if we want `num_frames` frames in the animation, we need to plot every `t.size/num_frames` frames:

```
skip_frame = int(t.size/float(num_frames))
if n % skip_frame == 0 or n == t.size-1:
    st.plot(x, u, 'r-', ...)
```

The initial condition ( $n=0$ ) included by `n % skip_frame == 0`, as well as every `skip_frame`-th frame. As `n % skip_frame == 0` will very seldom be true for the very final frame, we must also check if `n == t.size-1` to get the final frame included.

A simple choice of numbers may illustrate the formulas: say we have 801 frames in total (`t.size`) and we allow only 60 frames to be plotted. Then we need to plot every  $801/60$  frame, which with integer division yields 13 as `every`. Using the mod function, `n % every`, this operation is zero every time `n` can be divided by 13 without a remainder. That is,

the `if` test is true when `n` equals 0, 13, 26, 39, ..., 780, 801. The associated code is included in the `plot_u` function in the file `wave1D_u0v.py`.

### 2.3.5 Running a case

The first demo of our 1D wave equation solver concerns vibrations of a string that is initially deformed to a triangular shape, like when picking a guitar string:

$$I(x) = \begin{cases} ax/x_0, & x < x_0, \\ a(L-x)/(L-x_0), & \text{otherwise} \end{cases} \quad (2.33)$$

We choose  $L = 75$  cm,  $x_0 = 0.8L$ ,  $a = 5$  mm, and a time frequency  $\nu = 440$  Hz. The relation between the wave speed  $c$  and  $\nu$  is  $c = \nu\lambda$ , where  $\lambda$  is the wavelength, taken as  $2L$  because the longest wave on the string forms half a wavelength. There is no external force, so  $f = 0$ , and the string is at rest initially so that  $V = 0$ .

Regarding numerical parameters, we need to specify a  $\Delta t$ . Sometimes it is more natural to think of a spatial resolution instead of a time step. A natural semi-coarse spatial resolution in the present problem is  $N_x = 50$ . We can then choose the associated  $\Delta t$  (as required by the `viz` and `solver` functions) as the stability limit:  $\Delta t = L/(N_x c)$ . This is the  $\Delta t$  to be specified, but notice that if  $C < 1$ , the actual  $\Delta x$  computed in `solver` gets larger than  $L/N_x$ :  $\Delta x = c\Delta t/C = L/(N_x C)$ . (The reason is that we fix  $\Delta t$  and adjust  $\Delta x$ , so if  $C$  gets smaller, the code implements this effect in terms of a larger  $\Delta x$ .)

A function for setting the physical and numerical parameters and calling `viz` in this application goes as follows:

```
def guitar(C):
    """Triangular wave (pulled guitar string)."""
    L = 0.75
    x0 = 0.8*L
    a = 0.005
    freq = 440
    wavelength = 2*L
    c = freq*wavelength
    omega = 2*pi*freq
    num_periods = 1
    T = 2*pi/omega*num_periods
    # Choose dt the same as the stability limit for Nx=50
    dt = L/50./c

    def I(x):
```

```

    return a*x/x0 if x < x0 else a/(L-x0)*(L-x)

umin = -1.2*a; umax = -umin
cpu = viz(I, 0, 0, c, L, dt, C, T, umin, umax,
animate=True, tool='scitools')

```

The associated program has the name `wave1D_u0.py`. Run the program and watch the movie of the vibrating string.

**hpl 7:** Must recompute these movies as  $\Delta t$  is different when  $C < 1$ .

### 2.3.6 Working with a scaled PDE model

Depending on the model, it may be a substantial job to establish consistent and relevant physical parameter values for a case. The guitar string example illustrates the point. However, by *scaling* the mathematical problem we can often reduce the need to estimate physical parameters dramatically. The scaling technique consists of introducing new independent and dependent variables, with the aim that the absolute values of these lie in  $[0, 1]$ . We introduce the dimensionless variables (details are found in Section 3.1.1 in [5]

$$\bar{x} = \frac{x}{L}, \quad \bar{t} = \frac{c}{L}t, \quad \bar{u} = \frac{u}{a}.$$

Here,  $L$  is a typical length scale, e.g., the length of the domain, and  $a$  is a typical size of  $u$ , e.g., determined from the initial condition:  $a = \max_x |I(x)|$ .

We get by the chain rule that

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial t} \left( \frac{\bar{u}}{a} \right) \frac{d\bar{t}}{dt} = \frac{\partial \bar{u}}{\partial \bar{t}} \frac{1}{a} \frac{c}{L} = \frac{ac}{L} \frac{\partial \bar{u}}{\partial \bar{t}}.$$

Similarly,

$$\frac{\partial u}{\partial x} = \frac{a}{L} \frac{\partial \bar{u}}{\partial \bar{x}}.$$

Inserting the dimensionless variables in the PDE gives, in case  $f = 0$ ,

$$\frac{a^2 c^2}{L^2} \frac{\partial^2 \bar{u}}{\partial \bar{t}^2} = \frac{a^2 c^2}{L^2} \frac{\partial^2 \bar{u}}{\partial \bar{x}^2}.$$

Dropping the bars, we arrive at the scaled PDE

$$\frac{\partial^2 u}{\partial t^2} = \frac{\partial^2 u}{\partial x^2}, \quad x \in (0, 1), \quad t \in (0, cT/L), \quad (2.34)$$

which has no parameter  $c^2$  anymore. The initial conditions are scaled as

$$a\bar{u}(\bar{x}, 0) = I(L\bar{x})$$

and

$$\frac{a}{L/c} \frac{\partial \bar{u}}{\partial \bar{t}}(\bar{x}, 0) = V(L\bar{x}),$$

resulting in

$$\bar{u}(\bar{x}, 0) = \frac{I(L\bar{x})}{\max_x |I(x)|}, \quad \frac{\partial \bar{u}}{\partial \bar{t}}(\bar{x}, 0) = \frac{L}{ac} V(L\bar{x}).$$

In the common case  $V = 0$  we see that there are no physical parameters to be estimated in the PDE model!

If we have a program implemented for the physical wave equation with dimensions, we can obtain the dimensionless, scaled version by setting  $c = 1$ . The initial condition of a guitar string, given in (2.33), gets its scaled form by choosing  $a = 1$ ,  $L = 1$ , and  $x_0 \in [0, 1]$ . This means that we only need to decide on the  $x_0$  value as a fraction of unity, because the scaled problem corresponds to setting all other parameters to unity. In the code we can just set `a=c=L=1`, `x0=0.8`, and there is no need to calculate with wavelengths and frequencies to estimate  $c$ !

The only non-trivial parameter to estimate in the scaled problem is the final end time of the simulation, or more precisely, how it relates to periods in periodic solutions in time, since we often want to express the end time as a certain number of periods. The period in the dimensionless problem is 2, so the end time can be set to the desired number of periods times 2.

Why the dimensionless period is 2 can be explained by the following reasoning. Suppose that  $u$  behaves as  $\cos(\omega t)$  in time in the original problem with dimensions. The corresponding period is then  $P = 2\pi/\omega$ , but we need to estimate  $\omega$ . A typical solution of the wave equation is  $u(x, t) = A \cos(kx) \cos(\omega t)$ , where  $A$  is an amplitude and  $k$  is related to the wave length  $\lambda$  in space:  $\lambda = 2\pi/k$ . Both  $\lambda$  and  $A$  will be given by the initial condition  $I(x)$ . Inserting this  $u(x, t)$  in the PDE yields  $-\omega^2 = -c^2 k^2$ , i.e.,  $\omega = kc$ . The period is therefore  $P = 2\pi/(kc)$ . If the boundary conditions are  $u(0, t) = u(0, L)$ , we need to have  $kL = n\pi$  for integer  $n$ . The period becomes  $P = 2L/nc$ . The longest period is  $P = 2L/c$ . The dimensionless period  $\tilde{P}$  is obtained by dividing  $P$  by

the time scale  $L/c$ , which results in  $\tilde{P} = 2$ . Shorter waves in the initial condition will have a dimensionless shorter period  $\tilde{P} = 2/n$  ( $n > 1$ ).

## 2.4 Vectorization

The computational algorithm for solving the wave equation visits one mesh point at a time and evaluates a formula for the new value  $u_i^{n+1}$  at that point. Technically, this is implemented by a loop over array elements in a program. Such loops may run slowly in Python (and similar interpreted languages such as R and MATLAB). One technique for speeding up loops is to perform operations on entire arrays instead of working with one element at a time. This is referred to as *vectorization*, *vector computing*, or *array computing*. Operations on whole arrays are possible if the computations involving each element is independent of each other and therefore can, at least in principle, be performed simultaneously. That is, vectorization not only speeds up the code on serial computers, but also makes it easy to exploit parallel computing. Actually, there are Python tools like [Numba](#) that can automatically turn vectorized code into parallel code.

### 2.4.1 Operations on slices of arrays

Efficient computing with `numpy` arrays demands that we avoid loops and compute with entire arrays at once (or at least large portions of them). Consider this calculation of differences  $d_i = u_{i+1} - u_i$ :

```
n = u.size
for i in range(0, n-1):
    d[i] = u[i+1] - u[i]
```

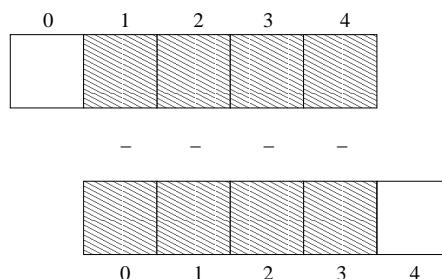
All the differences here are independent of each other. The computation of `d` can therefore alternatively be done by subtracting the array  $(u_0, u_1, \dots, u_{n-1})$  from the array where the elements are shifted one index upwards:  $(u_1, u_2, \dots, u_n)$ , see Figure 2.3. The former subset of the array can be expressed by `u[0:n-1]`, `u[0:-1]`, or just `u[:-1]`, meaning from index 0 up to, but not including, the last element (-1). The latter subset is obtained by `u[1:n]` or `u[1:]`, meaning from index 1 and the rest of the array. The computation of `d` can now be done without an explicit Python loop:

```
d = u[1:] - u[:-1]
```

or with explicit limits if desired:

```
d = u[1:n] - u[0:n-1]
```

Indices with a colon, going from an index to (but not including) another index are called *slices*. With numpy arrays, the computations are still done by loops, but in efficient, compiled, highly optimized C or Fortran code. Such loops are sometimes referred to as *vectorized loops*. Such loops can also easily be distributed among many processors on parallel computers. We say that the *scalar code* above, working on an element (a scalar) at a time, has been replaced by an equivalent *vectorized code*. The process of vectorizing code is called *vectorization*.



**Fig. 2.3** Illustration of subtracting two slices of two arrays.

### Test your understanding

Newcomers to vectorization are encouraged to choose a small array `u`, say with five elements, and simulate with pen and paper both the loop version and the vectorized version above.

Finite difference schemes basically contain differences between array elements with shifted indices. As an example, consider the updating formula

```
for i in range(1, n-1):
    u2[i] = u[i-1] - 2*u[i] + u[i+1]
```

The vectorization consists of replacing the loop by arithmetics on slices of arrays of length `n-2`:

```
u2 = u[:-2] - 2*u[1:-1] + u[2:]
u2 = u[0:n-2] - 2*u[1:n-1] + u[2:n] # alternative
```

Note that the length of  $u_2$  becomes  $n-2$ . If  $u_2$  is already an array of length  $n$  and we want to use the formula to update all the “inner” elements of  $u_2$ , as we will when solving a 1D wave equation, we can write

```
u2[1:-1] = u[:-2] - 2*u[1:-1] + u[2:]
u2[1:n-1] = u[0:n-2] - 2*u[1:n-1] + u[2:n] # alternative
```

The first expression’s right-hand side is realized by the following steps, involving temporary arrays with intermediate results, since each array operation can only involve one or two arrays. The `numpy` package performs (behind the scenes) the first line above in four steps:

```
temp1 = 2*u[1:-1]
temp2 = u[:-2] - temp1
temp3 = temp2 + u[2:]
u2[1:-1] = temp3
```

We need three temporary arrays, but a user does not need to worry about such temporary arrays.

### Common mistakes with array slices

Array expressions with slices demand that the slices have the same shape. It is easy to make a mistake in, e.g.,

```
u2[1:n-1] = u[0:n-2] - 2*u[1:n-1] + u[2:n]
```

and write

```
u2[1:n-1] = u[0:n-2] - 2*u[1:n-1] + u[1:n]
```

Now  $u[1:n]$  has wrong length ( $n-1$ ) compared to the other array slices, causing a `ValueError` and the message `could not broadcast input array from shape 103 into shape 104` (if  $n$  is 105). When such errors occur one must closely examine all the slices. Usually, it is easier to get upper limits of slices right when they use  $-1$  or  $-2$  or empty limit rather than expressions involving the length.

Another common mistake is to forget the slice in the array on the left-hand side,

```
u2 = u[0:n-2] - 2*u[1:n-1] + u[1:n]
```

This is really crucial: now  $u_2$  becomes a *new* array of length  $n-2$ , which is the wrong length as we have no entries for the boundary values. We meant to insert the right-hand side array *into* the original

u2 array for the entries that correspond to the internal points in the mesh (1:n-1 or 1:-1).

Vectorization may also work nicely with functions. To illustrate, we may extend the previous example as follows:

```
def f(x):
    return x**2 + 1

for i in range(1, n-1):
    u2[i] = u[i-1] - 2*u[i] + u[i+1] + f(x[i])
```

Assuming u2, u, and x all have length n, the vectorized version becomes

```
u2[1:-1] = u[:-2] - 2*u[1:-1] + u[2:] + f(x[1:-1])
```

Obviously, f must be able to take an array as argument for  $f(x[1:-1])$  to make sense.

#### 2.4.2 Finite difference schemes expressed as slices

We now have the necessary tools to vectorize the wave equation algorithm as described mathematically in Section 2.1.5 and through code in Section 2.3.2. There are three loops: one for the initial condition, one for the first time step, and finally the loop that is repeated for all subsequent time levels. Since only the latter is repeated a potentially large number of times, we limit our vectorization efforts to this loop:

```
for i in range(1, Nx):
    u[i] = 2*u_1[i] - u_2[i] + \
            C2*(u_1[i-1] - 2*u_1[i] + u_1[i+1])
```

The vectorized version becomes

```
u[1:-1] = - u_2[1:-1] + 2*u_1[1:-1] + \
           C2*(u_1[:-2] - 2*u_1[1:-1] + u_1[2:])
```

or

```
u[1:Nx] = 2*u_1[1:Nx] - u_2[1:Nx] + \
           C2*(u_1[0:Nx-1] - 2*u_1[1:Nx] + u_1[2:Nx+1])
```

The program `wave1D_u0v.py` contains a new version of the function `solver` where both the scalar and the vectorized loops are included (the argument `version` is set to `scalar` or `vectorized`, respectively).

### 2.4.3 Verification

We may reuse the quadratic solution  $u_e(x, t) = x(L - x)(1 + \frac{1}{2}t)$  for verifying also the vectorized code. A test function can now verify both the scalar and the vectorized version. Moreover, we may use a `user_action` function that compares the computed and exact solution at each time level and performs a test:

```
def test_quadratic():
    """
    Check the scalar and vectorized versions for
    a quadratic u(x,t)=x(L-x)(1+t/2) that is exactly reproduced.
    """

    # The following function must work for x as array or scalar
    u_exact = lambda x, t: x*(L - x)*(1 + 0.5*t)
    I = lambda x: u_exact(x, 0)
    V = lambda x: 0.5*u_exact(x, 0)
    # f is a scalar (zeros_like(x) works for scalar x too)
    f = lambda x, t: np.zeros_like(x) + 2*c**2*(1 + 0.5*t)

    L = 2.5
    c = 1.5
    C = 0.75
    Nx = 3 # Very coarse mesh for this exact test
    dt = C*(L/Nx)/c
    T = 18

    def assert_no_error(u, x, t, n):
        u_e = u_exact(x, t[n])
        tol = 1E-13
        diff = np.abs(u - u_e).max()
        assert diff < tol

        solver(I, V, f, c, L, dt, C, T,
               user_action=assert_no_error, version='scalar')
        solver(I, V, f, c, L, dt, C, T,
               user_action=assert_no_error, version='vectorized')
```

#### Lambda functions

The code segment above demonstrates how to achieve very compact code, without degraded readability, by use of lambda functions for the various input parameters that require a Python function. In essence,

```
f = lambda x, t: L*(x-t)**2
```

is equivalent to

```
def f(x, t):
    return L(x-t)**2
```

Note that lambda functions can just contain a single expression and no statements.

One advantage with lambda functions is that they can be used directly in calls:

```
solver(I=lambda x: sin(pi*x/L), V=0, f=0, ...)
```

#### 2.4.4 Efficiency measurements

The `wave1D_u0v.py` contains our new `solver` function with both scalar and vectorized code. For comparing the efficiency of scalar versus vectorized code, we need a `viz` function as discussed in Section 2.3.4. All of this `viz` function can be reused, except the call to `solver_function`. This call lacks the parameter `version`, which we want to set to `vectorized` and `scalar` for our efficiency measurements.

One solution is to copy the `viz` code from `wave1D_u0` into `wave1D_u0v.py` and add a `version` argument to the `solver_function` call. Taking into account how much animation code we then duplicate, this is not a good idea. Alternatively, introducing the `version` argument in `wave1D_u0.viz`, so that this function can be imported into `wave1D_u0v.py`, is not a good solution either, since `version` has no meaning in that file. We need better ideas!

**Solution 1.** Calling `viz` in `wave1D_u0` with `solver_function` as our new solver in `wave1D_u0v` works fine, since this solver has `version='vectorized'` as default value. The problem arises when we want to test `version='vectorized'`. The simplest solution is then to use `wave1D_u0.solver` instead. We make a new `viz` function in `wave1D_u0v.py` that has a `version` argument and that just calls `wave1D_u0.viz`:

```
def viz(
    I, V, f, c, L, dt, C, T, # PDE parameters
    umin, umax,             # Interval for u in plots
    animate=True,            # Simulation with animation?
    tool='matplotlib',       # 'matplotlib' or 'scitools'
    solver_function=solver, # Function with numerical algorithm
    version='vectorized',   # 'scalar' or 'vectorized'
):
```

```

import wave1D_u0
if version == 'vectorized':
    # Reuse viz from wave1D_u0, but with the present
    # modules' new vectorized solver (which has
    # version='vectorized' as default argument;
    # wave1D_u0.viz does not feature this argument)
    cpu = wave1D_u0.viz(
        I, V, f, c, L, dt, C, T, umin, umax,
        animate, tool, solver_function=solver)
elif version == 'scalar':
    # Call wave1D_u0.viz with a solver with
    # scalar code and use wave1D_u0.solver.
    cpu = wave1D_u0.viz(
        I, V, f, c, L, dt, C, T, umin, umax,
        animate, tool,
        solver_function=wave1D_u0.solver)

```

**Solution 2.** There is a more advanced and fancier solution featuring a very useful trick: we can make a new function that will always call `wave1D_u0v.solver` with `version='scalar'`. The `functools.partial` function from standard Python takes a function `func` as argument and a series of positional and keyword arguments and returns a new function that will call `func` with the supplied arguments, while the user can control all the other arguments in `func`. Consider a trivial example,

```

def f(a, b, c=2):
    return a + b + c

```

We want to ensure that `f` is always called with `c=3`, i.e., `f` has only two “free” arguments `a` and `b`. This functionality is obtained by

```

import functools
f2 = functools.partial(f, c=3)

print f2(1, 2) # results in 1+2+3=6

```

Now `f2` calls `f` with whatever the user supplies as `a` and `b`, but `c` is always 3.

Back to our `viz` code, we can do

```

import functools
# Call scalar with version fixed to 'scalar'
scalar_solver = functools.partial(scalar, version='scalar')
cpu = wave1D_u0.viz(
    I, V, f, c, L, dt, C, T, umin, umax,
    animate, tool, solver_function=scalar_solver)

```

The new `scalar_solver` takes the same arguments as `wave1D_u0.scalar` and calls `wave1D_u0v.scalar`, but always supplies the extra argument `version='scalar'`. When sending this `solver_function` to

`wave1D_u0.viz`, the latter will call `wave1D_u0v.solver` with all the `I`, `V`, `f`, etc., arguments we supply, plus `version='scalar'`.

**Efficiency experiments.** We now have a `viz` function that can call our solver function both in scalar and vectorized mode. The function `run_efficiency_experiments` in `wave1D_u0v.py` performs a set of experiments and reports the CPU time spent in the scalar and vectorized solver for the previous string vibration example with spatial mesh resolutions  $N_x = 50, 100, 200, 400, 800$ . Running this function reveals that the vectorized code runs substantially faster: the vectorized code runs approximately  $N_x/10$  times as fast as the scalar code!

#### 2.4.5 Remark on the updating of arrays

At the end of each time step we need to update the `u_2` and `u_1` arrays such that they have the right content for the next time step:

```
u_2[:] = u_1
u_1[:] = u
```

The order here is important! (Updating `u_1` first, makes `u_2` equal to `u`, which is wrong.)

The assignment `u_1[:] = u` copies the content of the `u` array into the elements of the `u_1` array. Such copying takes time, but that time is negligible compared to the time needed for computing `u` from the finite difference formula, even when the formula has a vectorized implementation. However, efficiency of program code is a key topic when solving PDEs numerically (particularly when there are two or three space dimensions), so it must be mentioned that there exists a much more efficient way of making the arrays `u_2` and `u_1` ready for the next time step. The idea is based on *switching references* and explained as follows.

A Python variable is actually a reference to some object (C programmers may think of pointers). Instead of copying data, we can let `u_2` refer to the `u_1` object and `u_1` refer to the `u` object. This is a very efficient operation (like switching pointers in C). A naive implementation like

```
u_2 = u_1
u_1 = u
```

will fail, however, because now `u_2` refers to the `u_1` object, but then the name `u_1` refers to `u`, so that this `u` object has two references, `u_1` and `u`, while our third array, originally referred to by `u_2`, has no more references and is lost. This means that the variables `u`, `u_1`, and `u_2` refer

to two arrays and not three. Consequently, the computations at the next time level will be messed up since updating the elements in  $\mathbf{u}$  will imply updating the elements in  $\mathbf{u}_1$  too so the solution at the previous time step, which is crucial in our formulas, is destroyed.

While  $\mathbf{u}_2 = \mathbf{u}_1$  is fine,  $\mathbf{u}_1 = \mathbf{u}$  is problematic, so the solution to this problem is to ensure that  $\mathbf{u}$  points to the  $\mathbf{u}_2$  array. This is mathematically wrong, but new correct values will be filled into  $\mathbf{u}$  at the next time step and make it right.

The correct switch of references is

```
tmp = u_2
u_2 = u_1
u_1 = u
u = tmp
```

We can get rid of the temporary reference  $\mathbf{tmp}$  by writing

```
u_2, u_1, u = u_1, u, u_2
```

This switching of references for updating our arrays will be used in later implementations.

#### Caution:

The update  $\mathbf{u}_2, \mathbf{u}_1, \mathbf{u} = \mathbf{u}_1, \mathbf{u}, \mathbf{u}_2$  leaves wrong content in  $\mathbf{u}$  at the final time step. This means that if we return  $\mathbf{u}$ , as we do in the example codes here, we actually return  $\mathbf{u}_2$ , which is obviously wrong. It is therefore important to adjust the content of  $\mathbf{u}$  to  $\mathbf{u} = \mathbf{u}_1$  before returning  $\mathbf{u}$ .

## 2.5 Exercises

### Exercise 2.1: Simulate a standing wave

The purpose of this exercise is to simulate standing waves on  $[0, L]$  and illustrate the error in the simulation. Standing waves arise from an initial condition

$$u(x, 0) = A \sin\left(\frac{\pi}{L} mx\right),$$

where  $m$  is an integer and  $A$  is a freely chosen amplitude. The corresponding exact solution can be computed and reads

$$u_e(x, t) = A \sin\left(\frac{\pi}{L}mx\right) \cos\left(\frac{\pi}{L}mct\right).$$

- a)** Explain that for a function  $\sin kx \cos \omega t$  the wave length in space is  $\lambda = 2\pi/k$  and the period in time is  $P = 2\pi/\omega$ . Use these expressions to find the wave length in space and period in time of  $u_e$  above.
- b)** Import the `solver` function from `wave1D_u0.py` into a new file where the `viz` function is reimplemented such that it plots either the numerical *and* the exact solution, *or* the error.
- c)** Make animations where you illustrate how the error  $e_i^n = u_e(x_i, t_n) - u_i^n$  develops and increases in time. Also make animations of  $u$  and  $u_e$  simultaneously.

**Hint 1.** Quite long time simulations are needed in order to display significant discrepancies between the numerical and exact solution.

**Hint 2.** A possible set of parameters is  $L = 12$ ,  $m = 9$ ,  $c = 2$ ,  $A = 1$ ,  $N_x = 80$ ,  $C = 0.8$ . The error mesh function  $e^n$  can be simulated for 10 periods, while 20-30 periods are needed to show significant differences between the curves for the numerical and exact solution.

Filename: `wave_standing`.

**Remarks.** The important parameters for numerical quality are  $C$  and  $k\Delta x$ , where  $C = c\Delta t/\Delta x$  is the Courant number and  $k$  is defined above ( $k\Delta x$  is proportional to how many mesh points we have per wave length in space, see Section 2.10.4 for explanation).

### Exercise 2.2: Add storage of solution in a user action function

Extend the `plot_u` function in the file `wave1D_u0.py` to also store the solutions `u` in a list. To this end, declare `all_u` as an empty list in the `viz` function, outside `plot_u`, and perform an append operation inside the `plot_u` function. Note that a function, like `plot_u`, inside another function, like `viz`, remembers all local variables in `viz` function, including `all_u`, even when `plot_u` is called (as `user_action`) in the `solver` function. Test both `all_u.append(u)` and `all_u.append(u.copy())`. Why does one of these constructions fail to store the solution correctly? Let the `viz` function return the `all_u` list converted to a two-dimensional `numpy` array.

Filename: `wave1D_u0_s_store`.

### Exercise 2.3: Use a class for the user action function

Redo Exercise 2.2 using a class for the user action function. Let the `all_u` list be an attribute in this class and implement the user action function as a method (the special method `__call__` is a natural choice). The class versions avoid that the user action function depends on parameters defined outside the function (such as `all_u` in Exercise 2.2).

Filename: `wave1D_u0_s2c.py`.

### Exercise 2.4: Compare several Courant numbers in one movie

The goal of this exercise is to make movies where several curves, corresponding to different Courant numbers, are visualized. Write a program that resembles `wave1D_u0_s2c.py` in Exercise 2.3, but with a `viz` function that can take a list of  $C$  values as argument and create a movie with solutions corresponding to the given  $C$  values. The `plot_u` function must be changed to store the solution in an array (see Exercise 2.2 or 2.3 for details), `solver` must be computed for each value of the Courant number, and finally one must run through each time step and plot all the spatial solution curves in one figure and store it in a file.

The challenge in such a visualization is to ensure that the curves in one plot correspond to the same time point. The easiest remedy is to keep the time resolution constant and change the space resolution to change the Courant number. Note that each spatial grid is needed for the final plotting, so it is an option to store those grids too.

Filename: `wave_numerics_comparison.py`.

### Project 2.5: Calculus with 1D mesh functions

This project explores integration and differentiation of mesh functions, both with scalar and vectorized implementations. We are given a mesh function  $f_i$  on a spatial one-dimensional mesh  $x_i = i\Delta x$ ,  $i = 0, \dots, N_x$ , over the interval  $[a, b]$ .

- Define the discrete derivative of  $f_i$  by using centered differences at internal mesh points and one-sided differences at the end points. Implement a scalar version of the computation in a Python function and write an associated unit test for the linear case  $f(x) = 4x - 2.5$  where the discrete derivative should be exact.

- b)** Vectorize the implementation of the discrete derivative. Extend the unit test to check the validity of the implementation.
- c)** To compute the discrete integral  $F_i$  of  $f_i$ , we assume that the mesh function  $f_i$  varies linearly between the mesh points. Let  $f(x)$  be such a linear interpolant of  $f_i$ . We then have

$$F_i = \int_{x_0}^{x_i} f(x) dx .$$

The exact integral of a piecewise linear function  $f(x)$  is given by the Trapezoidal rule. Show that if  $F_i$  is already computed, we can find  $F_{i+1}$  from

$$F_{i+1} = F_i + \frac{1}{2}(f_i + f_{i+1})\Delta x .$$

Make a function for the scalar implementation of the discrete integral as a mesh function. That is, the function should return  $F_i$  for  $i = 0, \dots, N_x$ . For a unit test one can use the fact that the above defined discrete integral of a linear function (say  $f(x) = 4x - 2.5$ ) is exact.

- d)** Vectorize the implementation of the discrete integral. Extend the unit test to check the validity of the implementation.

**Hint.** Interpret the recursive formula for  $F_{i+1}$  as a sum. Make an array with each element of the sum and use the "cumsum" (`numpy.cumsum`) operation to compute the accumulative sum: `numpy.cumsum([1, 3, 5])` is `[1, 4, 9]`.

- e)** Create a class `MeshCalculus` that can integrate and differentiate mesh functions. The class can just define some methods that call the previously implemented Python functions. Here is an example on the usage:

```
import numpy as np
calc = MeshCalculus(vectorized=True)
x = np.linspace(0, 1, 11)          # mesh
f = np.exp(x)                     # mesh function
df = calc.differentiate(f, x)     # discrete derivative
F = calc.integrate(f, x)          # discrete anti-derivative
```

Filename: `mesh_calculus_1D`.

## 2.6 Generalization: reflecting boundaries

The boundary condition  $u = 0$  in a wave equation reflects the wave, but  $u$  changes sign at the boundary, while the condition  $u_x = 0$  reflects the wave as a mirror and preserves the sign, see a [web page](#) or a [movie file](#) for demonstration.

Our next task is to explain how to implement the boundary condition  $u_x = 0$ , which is more complicated to express numerically and also to implement than a given value of  $u$ . We shall present two methods for implementing  $u_x = 0$  in a finite difference scheme, one based on deriving a modified stencil at the boundary, and another one based on extending the mesh with ghost cells and ghost points.

### 2.6.1 Neumann boundary condition

When a wave hits a boundary and is to be reflected back, one applies the condition

$$\frac{\partial u}{\partial n} \equiv \mathbf{n} \cdot \nabla u = 0. \quad (2.35)$$

The derivative  $\partial/\partial n$  is in the outward normal direction from a general boundary. For a 1D domain  $[0, L]$ , we have that

$$\left. \frac{\partial}{\partial n} \right|_{x=L} = \left. \frac{\partial}{\partial x} \right|_{x=L}, \quad \left. \frac{\partial}{\partial n} \right|_{x=0} = - \left. \frac{\partial}{\partial x} \right|_{x=0}.$$

#### Boundary condition terminology

Boundary conditions that specify the value of  $\partial u / \partial n$ , or shorter  $u_n$ , are known as **Neumann conditions**, while **Dirichlet conditions** refer to specifications of  $u$ . When the values are zero ( $\partial u / \partial n = 0$  or  $u = 0$ ) we speak about *homogeneous* Neumann or Dirichlet conditions.

### 2.6.2 Discretization of derivatives at the boundary

How can we incorporate the condition (2.35) in the finite difference scheme? Since we have used central differences in all the other approxi-

mations to derivatives in the scheme, it is tempting to implement (2.35) at  $x = 0$  and  $t = t_n$  by the difference

$$[D_{2x}u]_0^n = \frac{u_{-1}^n - u_1^n}{2\Delta x} = 0. \quad (2.36)$$

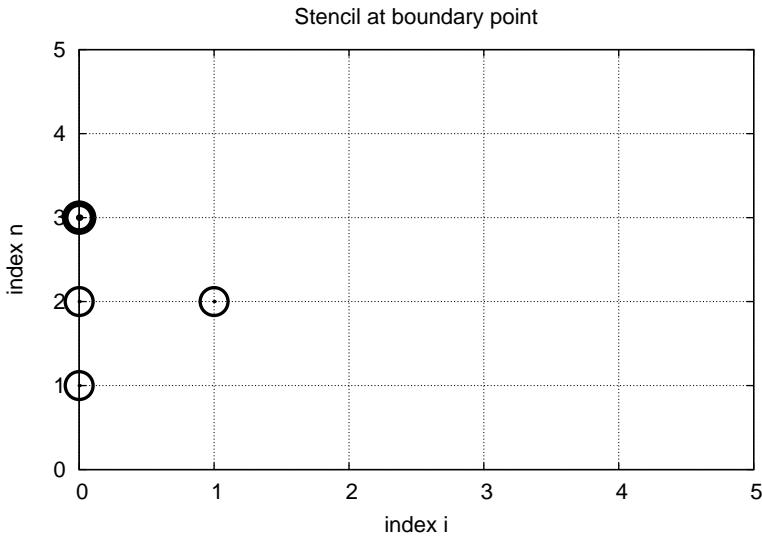
The problem is that  $u_{-1}^n$  is not a  $u$  value that is being computed since the point is outside the mesh. However, if we combine (2.36) with the scheme for  $i = 0$ ,

$$u_i^{n+1} = -u_i^{n-1} + 2u_i^n + C^2 (u_{i+1}^n - 2u_i^n + u_{i-1}^n), \quad (2.37)$$

we can eliminate the fictitious value  $u_{-1}^n$ . We see that  $u_{-1}^n = u_1^n$  from (2.36), which can be used in (2.37) to arrive at a modified scheme for the boundary point  $u_0^{n+1}$ :

$$u_i^{n+1} = -u_i^{n-1} + 2u_i^n + 2C^2 (u_{i+1}^n - u_i^n), \quad i = 0. \quad (2.38)$$

Figure 2.4 visualizes this equation for computing  $u_0^3$  in terms of  $u_0^2$ ,  $u_0^1$ , and  $u_1^2$ .



**Fig. 2.4** Modified stencil at a boundary with a Neumann condition.

Similarly, (2.35) applied at  $x = L$  is discretized by a central difference

$$\frac{u_{N_x+1}^n - u_{N_x-1}^n}{2\Delta x} = 0. \quad (2.39)$$

Combined with the scheme for  $i = N_x$  we get a modified scheme for the boundary value  $u_{N_x}^{n+1}$ :

$$u_i^{n+1} = -u_i^{n-1} + 2u_i^n + 2C^2 (u_{i-1}^n - u_i^n), \quad i = N_x. \quad (2.40)$$

The modification of the scheme at the boundary is also required for the special formula for the first time step. How the stencil moves through the mesh and is modified at the boundary can be illustrated by an animation in a web page or a movie file.

### 2.6.3 Implementation of Neumann conditions

We have seen in the preceding section that the special formulas for the boundary points arise from replacing  $u_{i-1}^n$  by  $u_{i+1}^n$  when computing  $u_i^{n+1}$  from the stencil formula for  $i = 0$ . Similarly, we replace  $u_{i+1}^n$  by  $u_{i-1}^n$  in the stencil formula for  $i = N_x$ . This observation can conveniently be used in the coding: we just work with the general stencil formula, but write the code such that it is easy to replace  $u[i-1]$  by  $u[i+1]$  and vice versa. This is achieved by having the indices  $i+1$  and  $i-1$  as variables  $ip1$  ( $i$  plus 1) and  $im1$  ( $i$  minus 1), respectively. At the boundary we can easily define  $im1=i+1$  while we use  $im1=i-1$  in the internal parts of the mesh. Here are the details of the implementation (note that the updating formula for  $u[i]$  is the general stencil formula):

```
i = 0
ip1 = i+1
im1 = ip1 # i-1 -> i+1
u[i] = u_1[i] + C2*(u_1[im1] - 2*u_1[i] + u_1[ip1])

i = Nx
im1 = i-1
ip1 = im1 # i+1 -> i-1
u[i] = u_1[i] + C2*(u_1[im1] - 2*u_1[i] + u_1[ip1])
```

We can in fact create one loop over both the internal and boundary points and use only one updating formula:

```
for i in range(0, Nx+1):
    ip1 = i+1 if i < Nx else i-1
    im1 = i-1 if i > 0 else i+1
    u[i] = u_1[i] + C2*(u_1[im1] - 2*u_1[i] + u_1[ip1])
```

The program `wave1D_n0.py` contains a complete implementation of the 1D wave equation with boundary conditions  $u_x = 0$  at  $x = 0$  and  $x = L$ .

It would be nice to modify the `test_quadratic` test case from the `wave1D_u0.py` with Dirichlet conditions, described in Section 2.4.3. However, the Neumann conditions require the polynomial variation in the  $x$  direction to be of third degree, which causes challenging problems when designing a test where the numerical solution is known exactly. Exercise 2.14 outlines ideas and code for this purpose. The only test in `wave1D_n0.py` is to start with a plug wave at rest and see that the initial condition is reached again perfectly after one period of motion, but such a test requires  $C = 1$  (so the numerical solution coincides with the exact solution of the PDE, see Section 2.10.4).

## 2.6.4 Index set notation

To improve our mathematical writing and our implementations, it is wise to introduce a special notation for index sets. This means that we write  $x_i$ ,  $i \in \mathcal{I}_x$ , instead of  $i = 0, \dots, N_x$ . Obviously,  $\mathcal{I}_x$  must be the index set  $\mathcal{I}_x = \{0, \dots, N_x\}$ , but it is often advantageous to have a symbol for this set rather than specifying all its elements (all the time, as we have done up to now). This new notation saves writing and makes specifications of algorithms and their implementation of computer code simpler.

The first index in the set will be denoted  $\mathcal{I}_x^0$  and the last  $\mathcal{I}_x^{-1}$ . When we need to skip the first element of the set, we use  $\mathcal{I}_x^+$  for the remaining subset  $\mathcal{I}_x^+ = \{1, \dots, N_x\}$ . Similarly, if the last element is to be dropped, we write  $\mathcal{I}_x^- = \{0, \dots, N_x - 1\}$  for the remaining indices. All the indices corresponding to inner grid points are specified by  $\mathcal{I}_x^i = \{1, \dots, N_x - 1\}$ . For the time domain we find it natural to explicitly use 0 as the first index, so we will usually write  $n = 0$  and  $t_0$  rather than  $n = \mathcal{I}_t^0$ . We also avoid notation like  $x_{\mathcal{I}_x^{-1}}$  and will instead use  $x_i$ ,  $i = \mathcal{I}_x^{-1}$ .

The Python code associated with index sets applies the following conventions:

Notation	Python
$\mathcal{I}_x$	<code>Ix</code>
$\mathcal{I}_x^0$	<code>Ix[0]</code>
$\mathcal{I}_x^{-1}$	<code>Ix[-1]</code>
$\mathcal{I}_x^-$	<code>Ix[:-1]</code>
$\mathcal{I}_x^+$	<code>Ix[1:]</code>
$\mathcal{I}_x^i$	<code>Ix[1:-1]</code>

### Why index sets are useful

An important feature of the index set notation is that it keeps our formulas and code independent of how we count mesh points. For example, the notation  $i \in \mathcal{I}_x$  or  $i = \mathcal{I}_x^0$  remains the same whether  $\mathcal{I}_x$  is defined as above or as starting at 1, i.e.,  $\mathcal{I}_x = \{1, \dots, Q\}$ . Similarly, we can in the code define `Ix=range(Nx+1)` or `Ix=range(1,Q)`, and expressions like `Ix[0]` and `Ix[1:-1]` remain correct. One application where the index set notation is convenient is conversion of code from a language where arrays has base index 0 (e.g., Python and C) to languages where the base index is 1 (e.g., MATLAB and Fortran). Another important application is implementation of Neumann conditions via ghost points (see next section).

For the current problem setting in the  $x, t$  plane, we work with the index sets

$$\mathcal{I}_x = \{0, \dots, N_x\}, \quad \mathcal{I}_t = \{0, \dots, N_t\}, \quad (2.41)$$

defined in Python as

```
Ix = range(0, Nx+1)
It = range(0, Nt+1)
```

A finite difference scheme can with the index set notation be specified as

$$\begin{aligned} u_i^{n+1} &= u_i^n - \frac{1}{2}C^2(u_{i+1}^n - 2u_i^n + u_{i-1}^n), \quad , i \in \mathcal{I}_x^i, \quad n = 0, \\ u_i^{n+1} &= -u_i^{n-1} + 2u_i^n + C^2(u_{i+1}^n - 2u_i^n + u_{i-1}^n), \quad i \in \mathcal{I}_x^i, \quad n \in \mathcal{I}_t, \\ u_i^{n+1} &= 0, \quad i = \mathcal{I}_x^0, \quad n \in \mathcal{I}_t^-, \\ u_i^{n+1} &= 0, \quad i = \mathcal{I}_x^{-1}, \quad n \in \mathcal{I}_t^-. \end{aligned}$$

The corresponding implementation becomes

```
# Initial condition
for i in Ix[1:-1]:
    u[i] = u_1[i] - 0.5*C2*(u_1[i-1] - 2*u_1[i] + u_1[i+1])

# Time loop
for n in It[1:-1]:
    # Compute internal points
    for i in Ix[1:-1]:
        u[i] = -u_2[i] + 2*u_1[i] + \
```

```
C2*(u_1[i-1] - 2*u_1[i] + u_1[i+1])
# Compute boundary conditions
i = Ix[0]; u[i] = 0
i = Ix[-1]; u[i] = 0
```

### Notice

The program `wave1D_dn.py` applies the index set notation and solves the 1D wave equation  $u_{tt} = c^2 u_{xx} + f(x, t)$  with quite general boundary and initial conditions:

- $x = 0$ :  $u = U_0(t)$  or  $u_x = 0$
- $x = L$ :  $u = U_L(t)$  or  $u_x = 0$
- $t = 0$ :  $u = I(x)$
- $t = 0$ :  $u_t = I_t(x)$

The program combines Dirichlet and Neumann conditions, scalar and vectorized implementation of schemes, and the index notation into one piece of code. A lot of test examples are also included in the program:

- A rectangular plug-shaped initial condition. (For  $C = 1$  the solution will be a rectangle that jumps one cell per time step, making the case well suited for verification.)
- A Gaussian function as initial condition.
- A triangular profile as initial condition, which resembles the typical initial shape of a guitar string.
- A sinusoidal variation of  $u$  at  $x = 0$  and either  $u = 0$  or  $u_x = 0$  at  $x = L$ .
- An exact analytical solution  $u(x, t) = \cos(m\pi t/L) \sin(\frac{1}{2}m\pi x/L)$ , which can be used for convergence rate tests.

**hpl 8:** Should include some experiments here or make exercises. Qualitative behavior of the wave equation can be exemplified.

### 2.6.5 Verifying the implementation of Neumann conditions

How can we test that the Neumann conditions are correctly implemented? The `solver` function in the `wave1D_dn.py` program described in the box

above accepts Dirichlet or Neumann conditions at  $x = 0$  and  $x = L$ . It is tempting to apply a quadratic solution as described in Sections 2.2.1 and 2.3.3, but it turns out that this solution is no longer an exact solution of the discrete equations if a Neumann condition is implemented on the boundary. A linear solution does not help since we only have homogeneous Neumann conditions in `wave1D_dn.py`, and we are consequently left with testing just a constant solution:  $u = \text{const}$ .

```
def test_constant():
    """
    Check the scalar and vectorized versions for
    a constant  $u(x,t)$ . We simulate in  $[0, L]$  and apply
    Neumann and Dirichlet conditions at both ends.
    """

    u_const = 0.45
    u_exact = lambda x, t: u_const
    I = lambda x: u_exact(x, 0)
    V = lambda x: 0
    f = lambda x, t: 0

    def assert_no_error(u, x, t, n):
        u_e = u_exact(x, t[n])
        diff = np.abs(u - u_e).max()
        msg = 'diff=%E, t_%d=%g' % (diff, n, t[n])
        tol = 1E-13
        assert diff < tol, msg

    for U_0 in (None, lambda t: u_const):
        for U_L in (None, lambda t: u_const):
            L = 2.5
            c = 1.5
            C = 0.75
            Nx = 3 # Very coarse mesh for this exact test
            dt = C*(L/Nx)/c
            T = 18 # long time integration

            solver(I, V, f, c, U_0, U_L, L, dt, C, T,
                   user_action=assert_no_error,
                   version='scalar')
            solver(I, V, f, c, U_0, U_L, L, dt, C, T,
                   user_action=assert_no_error,
                   version='vectorized')
    print U_0, U_L
```

The quadratic solution is very useful for testing though, but it requires Dirichlet conditions at both ends.

Another test may utilize the fact that the approximation error vanishes when the Courant number is unity. We can, for example, start with a plug profile as initial condition, let this wave split into two plug waves, one in each direction, and check that the two plug waves come back and

form the initial condition again after “one period” of the solution process. Neumann conditions can be applied at both ends. A proper test function reads

```
def test_plug():
    """Check that an initial plug is correct back after one period."""
    L = 1.0
    c = 0.5
    dt = (L/10)/c # Nx=10
    I = lambda x: 0 if abs(x-L/2.0) > 0.1 else 1

    u_s, x, t, cpu = solver(
        I=I,
        V=None, f=None, c=0.5, U_0=None, U_L=None, L=L,
        dt=dt, C=1, T=4, user_action=None, version='scalar')
    u_v, x, t, cpu = solver(
        I=I,
        V=None, f=None, c=0.5, U_0=None, U_L=None, L=L,
        dt=dt, C=1, T=4, user_action=None, version='vectorized')
    tol = 1E-13
    diff = abs(u_s - u_v).max()
    assert diff < tol
    u_0 = np.array([I(x_) for x_ in x])
    diff = np.abs(u_s - u_0).max()
    assert diff < tol
```

Other tests must rely on an unknown approximation error, so effectively we are left with tests on the convergence rate.

## 2.6.6 Alternative implementation via ghost cells

**Idea.** Instead of modifying the scheme at the boundary, we can introduce extra points outside the domain such that the fictitious values  $u_{-1}^n$  and  $u_{N_x+1}^n$  are defined in the mesh. Adding the intervals  $[-\Delta x, 0]$  and  $[L, L + \Delta x]$ , often referred to as *ghost cells*, to the mesh gives us all the needed mesh points, corresponding to  $i = -1, 0, \dots, N_x, N_x + 1$ . The extra points  $i = -1$  and  $i = N_x + 1$  are known as *ghost points*, and values at these points,  $u_{-1}^n$  and  $u_{N_x+1}^n$ , are called *ghost values*.

The important idea is to ensure that we always have

$$u_{-1}^n = u_1^n \text{ and } u_{N_x+1}^n = u_{N_x-1}^n,$$

because then the application of the standard scheme at a boundary point  $i = 0$  or  $i = N_x$  will be correct and guarantee that the solution is compatible with the boundary condition  $u_x = 0$ .

**Implementation.** The `u` array now needs extra elements corresponding to the ghost points. Two new point values are needed:

```
u = zeros(Nx+3)
```

The arrays `u_1` and `u_2` must be defined accordingly.

Unfortunately, a major indexing problem arises with ghost cells. The reason is that Python indices *must* start at 0 and `u[-1]` will always mean the last element in `u`. This fact gives, apparently, a mismatch between the mathematical indices  $i = -1, 0, \dots, N_x + 1$  and the Python indices running over `u`:  $0, \dots, Nx+2$ . One remedy is to change the mathematical indexing of  $i$  in the scheme and write

$$u_i^{n+1} = \dots, \quad i = 1, \dots, N_x + 1,$$

instead of  $i = 0, \dots, N_x$  as we have previously used. The ghost points now correspond to  $i = 0$  and  $i = N_x + 1$ . A better solution is to use the ideas of Section 2.6.4: we hide the specific index value in an index set and operate with inner and boundary points using the index set notation.

To this end, we define `u` with proper length and `Ix` to be the corresponding indices for the real physical mesh points ( $1, 2, \dots, N_x + 1$ ):

```
u = zeros(Nx+3)
Ix = range(1, u.shape[0]-1)
```

That is, the boundary points have indices `Ix[0]` and `Ix[-1]` (as before). We first update the solution at all physical mesh points (i.e., interior points in the mesh):

```
for i in Ix:
    u[i] = -u_2[i] + 2*u_1[i] + \
            C2*(u_1[i-1] - 2*u_1[i] + u_1[i+1])
```

The indexing becomes a bit more complicated when we call functions like `V(x)` and `f(x, t)`, as we must remember that the appropriate  $x$  coordinate is given as `x[i-Ix[0]]`:

```
for i in Ix:
    u[i] = u_1[i] + dt*V(x[i-Ix[0]]) + \
            0.5*C2*(u_1[i-1] - 2*u_1[i] + u_1[i+1]) + \
            0.5*dt2*f(x[i-Ix[0]], t[0])
```

It remains to update the solution at ghost points, i.e., `u[0]` and `u[-1]` (or `u[Nx+2]`). For a boundary condition  $u_x = 0$ , the ghost value must equal the value at the associated inner mesh point. Computer code makes this statement precise:

```
i = Ix[0] # x=0 boundary
```

```

u[i-1] = u[i+1]
i = Ix[-1]           # x=L boundary
u[i+1] = u[i-1]

```

The physical solution to be plotted is now in  $u[1:-1]$ , or equivalently  $u[Ix[0]:Ix[-1]+1]$ , so this slice is the quantity to be returned from a solver function. A complete implementation appears in the program `wave1D_n0_ghost.py`.

### Warning

We have to be careful with how the spatial and temporal mesh points are stored. Say we let  $x$  be the physical mesh points,

```
x = linspace(0, L, Nx+1)
```

"Standard coding" of the initial condition,

```

for i in Ix:
    u_1[i] = I(x[i])

```

becomes wrong, since  $u_1$  and  $x$  have different lengths and the index  $i$  corresponds to two different mesh points. In fact,  $x[i]$  corresponds to  $u[1+i]$ . A correct implementation is

```

for i in Ix:
    u_1[i] = I(x[i-Ix[0]])

```

Similarly, a source term usually coded as  $f(x[i], t[n])$  is incorrect if  $x$  is defined to be the physical points, so  $x[i]$  must be replaced by  $x[i-Ix[0]]$ .

An alternative remedy is to let  $x$  also cover the ghost points such that  $u[i]$  is the value at  $x[i]$ .

The ghost cell is only added to the boundary where we have a Neumann condition. Suppose we have a Dirichlet condition at  $x = L$  and a homogeneous Neumann condition at  $x = 0$ . One ghost cell  $[-\Delta x, 0]$  is added to the mesh, so the index set for the physical points becomes  $\{1, \dots, N_x + 1\}$ . A relevant implementation is

```

u = zeros(Nx+2)
Ix = range(1, u.shape[0])
...
for i in Ix[:-1]:
    u[i] = - u_2[i] + 2*u_1[i] + \
            C2*(u_1[i-1] - 2*u_1[i] + u_1[i+1]) + \
            dt2*f(x[i-Ix[0]], t[n])

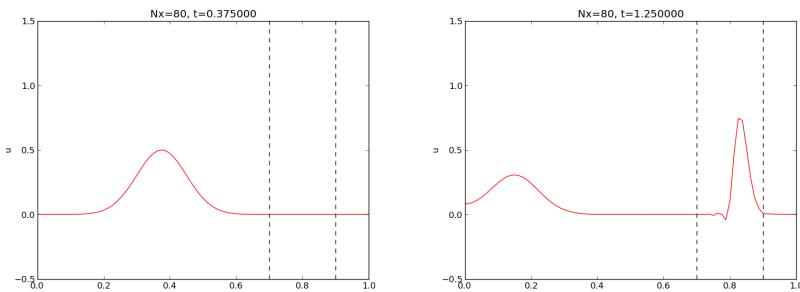
```

```
i = Ix[-1]
u[i] = U_0      # set Dirichlet value
i = Ix[0]
u[i-1] = u[i+1] # update ghost value
```

The physical solution to be plotted is now in  $u[1:]$  or (as always)  $u[Ix[0]:Ix[-1]+1]$ .

## 2.7 Generalization: variable wave velocity

Our next generalization of the 1D wave equation (2.1) or (2.17) is to allow for a variable wave velocity  $c$ :  $c = c(x)$ , usually motivated by wave motion in a domain composed of different physical media. When the media differ in physical properties like density or porosity, the wave velocity  $c$  is affected and will depend on the position in space. Figure 2.5 shows a wave propagating in one medium  $[0, 0.7] \cup [0.9, 1]$  with wave velocity  $c_1$  (left) before it enters a second medium  $(0.7, 0.9)$  with wave velocity  $c_2$  (right). When the wave passes the boundary where  $c$  jumps from  $c_1$  to  $c_2$ , a part of the wave is reflected back into the first medium (the *reflected wave*), while one part is transmitted through the second medium (the *transmitted wave*).



**Fig. 2.5** Left: wave entering another medium; right: transmitted and reflected wave.

### 2.7.1 The model PDE with a variable coefficient

Instead of working with the squared quantity  $c^2(x)$ , we shall for notational convenience introduce  $q(x) = c^2(x)$ . A 1D wave equation with variable wave velocity often takes the form

$$\frac{\partial^2 u}{\partial t^2} = \frac{\partial}{\partial x} \left( q(x) \frac{\partial u}{\partial x} \right) + f(x, t). \quad (2.42)$$

This is the most frequent form of a wave equation with variable wave velocity, but other forms also appear, see Section 2.14.1 and equation (2.125).

As usual, we sample (2.42) at a mesh point,

$$\frac{\partial^2}{\partial t^2} u(x_i, t_n) = \frac{\partial}{\partial x} \left( q(x_i) \frac{\partial}{\partial x} u(x_i, t_n) \right) + f(x_i, t_n),$$

where the only new term to discretize is

$$\frac{\partial}{\partial x} \left( q(x_i) \frac{\partial}{\partial x} u(x_i, t_n) \right) = \left[ \frac{\partial}{\partial x} \left( q(x) \frac{\partial u}{\partial x} \right) \right]_i^n.$$

## 2.7.2 Discretizing the variable coefficient

The principal idea is to first discretize the outer derivative. Define

$$\phi = q(x) \frac{\partial u}{\partial x},$$

and use a centered derivative around  $x = x_i$  for the derivative of  $\phi$ :

$$\left[ \frac{\partial \phi}{\partial x} \right]_i^n \approx \frac{\phi_{i+\frac{1}{2}} - \phi_{i-\frac{1}{2}}}{\Delta x} = [D_x \phi]_i^n.$$

Then discretize

$$\phi_{i+\frac{1}{2}} = q_{i+\frac{1}{2}} \left[ \frac{\partial u}{\partial x} \right]_{i+\frac{1}{2}}^n \approx q_{i+\frac{1}{2}} \frac{u_{i+1}^n - u_i^n}{\Delta x} = [q D_x u]_{i+\frac{1}{2}}^n.$$

Similarly,

$$\phi_{i-\frac{1}{2}} = q_{i-\frac{1}{2}} \left[ \frac{\partial u}{\partial x} \right]_{i-\frac{1}{2}}^n \approx q_{i-\frac{1}{2}} \frac{u_i^n - u_{i-1}^n}{\Delta x} = [q D_x u]_{i-\frac{1}{2}}^n.$$

These intermediate results are now combined to

$$\left[ \frac{\partial}{\partial x} \left( q(x) \frac{\partial u}{\partial x} \right) \right]_i^n \approx \frac{1}{\Delta x^2} \left( q_{i+\frac{1}{2}} (u_{i+1}^n - u_i^n) - q_{i-\frac{1}{2}} (u_i^n - u_{i-1}^n) \right). \quad (2.43)$$

With operator notation we can write the discretization as

$$\left[ \frac{\partial}{\partial x} \left( q(x) \frac{\partial u}{\partial x} \right) \right]_i^n \approx [D_x q D_x u]_i^n. \quad (2.44)$$

### Do not use the chain rule on the spatial derivative term

Many are tempted to use the chain rule on the term  $\frac{\partial}{\partial x} \left( q(x) \frac{\partial u}{\partial x} \right)$ , but this is not a good idea when discretizing such a term.

The term with a variable coefficient expresses the net flux  $qu_x$  into a small volume (i.e., interval in 1D):

$$\frac{\partial}{\partial x} \left( q(x) \frac{\partial u}{\partial x} \right) \approx \frac{1}{\Delta x} (q(x + \Delta x)u_x(x + \Delta x) - q(x)u_x(x)).$$

Our discretization reflects this principle directly:  $qu_x$  at the right end of the cell minus  $qu_x$  at the left end, because this follows from the formula (2.43) or  $[D_x(qD_xu)]_i^n$ .

When using the chain rule, we get two terms  $qu_{xx} + q_xu_x$ . The typical discretization is

$$D_x q D_x u + D_{2x} q D_{2x} u]_i^n, \quad (2.45)$$

Writing this out shows that it is different from  $[D_x(qD_xu)]_i^n$  and lacks the physical interpretation of net flux into a cell. With a smooth and slowly varying  $q(x)$  the differences between the two discretizations are not substantial. However, when  $q$  exhibits (potentially large) jumps,  $[D_x(qD_xu)]_i^n$  with harmonic averaging of  $q$  yields a better solution than arithmetic averaging or (2.45). In the literature, the discretization  $[D_x(qD_xu)]_i^n$  totally dominates and very few mention the possibility of (2.45).

### 2.7.3 Computing the coefficient between mesh points

If  $q$  is a known function of  $x$ , we can easily evaluate  $q_{i+\frac{1}{2}}$  simply as  $q(x_{i+\frac{1}{2}})$  with  $x_{i+\frac{1}{2}} = x_i + \frac{1}{2}\Delta x$ . However, in many cases  $q$ , and hence  $q_{i+\frac{1}{2}}$ , is only known as a discrete function, often at the mesh points  $x_i$ . Evaluating  $q$  between two mesh points  $x_i$  and  $x_{i+1}$  can then be done by averaging in three ways:

$$q_{i+\frac{1}{2}} \approx \frac{1}{2} (q_i + q_{i+1}) = [\bar{q}^x]_i \quad (\text{arithmetic mean}) \quad (2.46)$$

$$q_{i+\frac{1}{2}} \approx 2 \left( \frac{1}{q_i} + \frac{1}{q_{i+1}} \right)^{-1} \quad (\text{harmonic mean}) \quad (2.47)$$

$$q_{i+\frac{1}{2}} \approx (q_i q_{i+1})^{1/2} \quad (\text{geometric mean}) \quad (2.48)$$

The arithmetic mean in (2.46) is by far the most commonly used averaging technique and is well suited for smooth  $q(x)$  functions. The harmonic mean is often preferred when  $q(x)$  exhibits large jumps (which is typical for geological media). The geometric mean is less used, but popular in discretizations to linearize quadratic nonlinearities (see Section 1.8.2 for an example).

With the operator notation from (2.46) we can specify the discretization of the complete variable-coefficient wave equation in a compact way:

$$[D_t D_t u = D_x \bar{q}^x D_x u + f]_i^n. \quad (2.49)$$

From this notation we immediately see what kind of differences that each term is approximated with. The notation  $\bar{q}^x$  also specifies that the variable coefficient is approximated by an arithmetic mean, the definition being  $[\bar{q}^x]_{i+\frac{1}{2}} = (q_i + q_{i+1})/2$ . With the notation  $[D_x q D_x u]_i^n$ , we specify that  $q$  is evaluated directly, as a function, between the mesh points:  $q(x_{i-\frac{1}{2}})$  and  $q(x_{i+\frac{1}{2}})$ .

Before any implementation, it remains to solve (2.49) with respect to  $u_i^{n+1}$ :

$$\begin{aligned} u_i^{n+1} = & -u_i^{n-1} + 2u_i^n + \\ & \left( \frac{\Delta t}{\Delta x} \right)^2 \left( \frac{1}{2}(q_i + q_{i+1})(u_{i+1}^n - u_i^n) - \frac{1}{2}(q_i + q_{i-1})(u_i^n - u_{i-1}^n) \right) + \\ & \Delta t^2 f_i^n. \end{aligned} \quad (2.50)$$

## 2.7.4 How a variable coefficient affects the stability

The stability criterion derived in Section 2.10.3 reads  $\Delta t \leq \Delta x/c$ . If  $c = c(x)$ , the criterion will depend on the spatial location. We must therefore choose a  $\Delta t$  that is small enough such that no mesh cell has  $\Delta x/c(x) > \Delta t$ . That is, we must use the largest  $c$  value in the criterion:

$$\Delta t \leq \beta \frac{\Delta x}{\max_{x \in [0, L]} c(x)}. \quad (2.51)$$

The parameter  $\beta$  is included as a safety factor: in some problems with a significantly varying  $c$  it turns out that one must choose  $\beta < 1$  to have stable solutions ( $\beta = 0.9$  may act as an all-round value).

A different strategy to handle the stability criterion with variable wave velocity is to use a spatially varying  $\Delta t$ . While the idea is mathematically attractive at first sight, the implementation quickly becomes very complicated, so we stick to using a constant  $\Delta t$  and a worst case value of  $c(x)$  (with a safety factor  $\beta$ ).

## 2.7.5 Neumann condition and a variable coefficient

Consider a Neumann condition  $\partial u / \partial x = 0$  at  $x = L = N_x \Delta x$ , discretized as

$$[D_{2x}u]_i^n = \frac{u_{i+1}^n - u_{i-1}^n}{2\Delta x} = 0 \Rightarrow u_{i+1}^n = u_{i-1}^n,$$

for  $i = N_x$ . Using the scheme (2.50) at the end point  $i = N_x$  with  $u_{i+1}^n = u_{i-1}^n$  results in

$$\begin{aligned} u_i^{n+1} &= -u_i^{n-1} + 2u_i^n + \\ &\left( \frac{\Delta t}{\Delta x} \right)^2 \left( q_{i+\frac{1}{2}}(u_{i-1}^n - u_i^n) - q_{i-\frac{1}{2}}(u_i^n - u_{i-1}^n) \right) + \Delta t^2 f_i^n \end{aligned} \quad (2.52)$$

$$\begin{aligned} &= -u_i^{n-1} + 2u_i^n + \left( \frac{\Delta t}{\Delta x} \right)^2 (q_{i+\frac{1}{2}} + q_{i-\frac{1}{2}})(u_{i-1}^n - u_i^n) + \Delta t^2 f_i^n \end{aligned} \quad (2.53)$$

$$\approx -u_i^{n-1} + 2u_i^n + \left( \frac{\Delta t}{\Delta x} \right)^2 2q_i(u_{i-1}^n - u_i^n) + \Delta t^2 f_i^n. \quad (2.54)$$

Here we used the approximation

$$\begin{aligned}
q_{i+\frac{1}{2}} + q_{i-\frac{1}{2}} &= q_i + \left( \frac{dq}{dx} \right)_i \Delta x + \left( \frac{d^2q}{dx^2} \right)_i \Delta x^2 + \dots + \\
&\quad q_i - \left( \frac{dq}{dx} \right)_i \Delta x + \left( \frac{d^2q}{dx^2} \right)_i \Delta x^2 + \dots \\
&= 2q_i + 2 \left( \frac{d^2q}{dx^2} \right)_i \Delta x^2 + \mathcal{O}(\Delta x^4) \\
&\approx 2q_i.
\end{aligned} \tag{2.55}$$

An alternative derivation may apply the arithmetic mean of  $q$  in (2.53), leading to the term

$$(q_i + \frac{1}{2}(q_{i+1} + q_{i-1}))(u_{i-1}^n - u_i^n).$$

Since  $\frac{1}{2}(q_{i+1} + q_{i-1}) = q_i + \mathcal{O}(\Delta x^2)$ , we can approximate with  $2q_i(u_{i-1}^n - u_i^n)$  for  $i = N_x$  and get the same term as we did above.

A common technique when implementing  $\partial u / \partial x = 0$  boundary conditions, is to assume  $dq/dx = 0$  as well. This implies  $q_{i+1} = q_{i-1}$  and  $q_{i+1/2} = q_{i-1/2}$  for  $i = N_x$ . The implications for the scheme are

$$\begin{aligned}
u_i^{n+1} &= -u_i^{n-1} + 2u_i^n + \\
&\quad \left( \frac{\Delta t}{\Delta x} \right)^2 \left( q_{i+\frac{1}{2}}(u_{i-1}^n - u_i^n) - q_{i-\frac{1}{2}}(u_i^n - u_{i-1}^n) \right) + \\
&\quad \Delta t^2 f_i^n
\end{aligned} \tag{2.56}$$

$$= -u_i^{n-1} + 2u_i^n + \left( \frac{\Delta t}{\Delta x} \right)^2 2q_{i-\frac{1}{2}}(u_{i-1}^n - u_i^n) + \Delta t^2 f_i^n. \tag{2.57}$$

## 2.7.6 Implementation of variable coefficients

The implementation of the scheme with a variable wave velocity  $q(x) = c^2(x)$  may assume that  $q$  is available as an array  $q[i]$  at the spatial mesh points. The following loop is a straightforward implementation of the scheme (2.50):

```

for i in range(1, Nx):
    u[i] = -u_2[i] + 2*u_1[i] + \
            C2*(0.5*(q[i] + q[i+1])*(u_1[i+1] - u_1[i]) - \
                  0.5*(q[i] + q[i-1])*(u_1[i] - u_1[i-1])) + \
            dt2*f(x[i], t[n])

```

The coefficient C2 is now defined as  $(dt/dx)^2$ , i.e., *not* as the squared Courant number, since the wave velocity is variable and appears inside the parenthesis.

With Neumann conditions  $u_x = 0$  at the boundary, we need to combine this scheme with the discrete version of the boundary condition, as shown in Section 2.7.5. Nevertheless, it would be convenient to reuse the formula for the interior points and just modify the indices  $ip1=i+1$  and  $im1=i-1$  as we did in Section 2.6.3. Assuming  $dq/dx = 0$  at the boundaries, we can implement the scheme at the boundary with the following code.

```
i = 0
ip1 = i+1
im1 = ip1
u[i] = - u_2[i] + 2*u_1[i] + \
    C2*(0.5*(q[i] + q[ip1])*(u_1[ip1] - u_1[i]) - \
        0.5*(q[i] + q[im1])*(u_1[i] - u_1[im1])) + \
    dt2*f(x[i], t[n])
```

With ghost cells we can just reuse the formula for the interior points also at the boundary, provided that the ghost values of both  $u$  and  $q$  are correctly updated to ensure  $u_x = 0$  and  $q_x = 0$ .

A vectorized version of the scheme with a variable coefficient at internal mesh points becomes

```
u[1:-1] = - u_2[1:-1] + 2*u_1[1:-1] + \
    C2*(0.5*(q[1:-1] + q[2:])*(u_1[2:] - u_1[1:-1]) - \
        0.5*(q[1:-1] + q[:-2])*(u_1[1:-1] - u_1[:-2])) + \
    dt2*f(x[1:-1], t[n])
```

## 2.7.7 A more general PDE model with variable coefficients

Sometimes a wave PDE has a variable coefficient in front of the time-derivative term:

$$\varrho(x) \frac{\partial^2 u}{\partial t^2} = \frac{\partial}{\partial x} \left( q(x) \frac{\partial u}{\partial x} \right) + f(x, t). \quad (2.58)$$

One example appears when modeling elastic waves in a rod with varying density, cf. (2.14.1) with  $\varrho(x)$ .

A natural scheme for (2.58) is

$$[\varrho D_t D_t u = D_x \bar{q}^x D_x u + f]_i^n. \quad (2.59)$$

We realize that the  $\varrho$  coefficient poses no particular difficulty, since  $\varrho$  enters the formula just as a simple factor in front of a derivative. There

is hence no need for any averaging of  $\varrho$ . Often,  $\varrho$  will be moved to the right-hand side, also without any difficulty:

$$[D_t D_t u = \varrho^{-1} D_x \bar{q}^x D_x u + f]_i^n. \quad (2.60)$$

## 2.7.8 Generalization: damping

Waves die out by two mechanisms. In 2D and 3D the energy of the wave spreads out in space, and energy conservation then requires the amplitude to decrease. This effect is not present in 1D. Damping is another cause of amplitude reduction. For example, the vibrations of a string die out because of damping due to air resistance and non-elastic effects in the string.

The simplest way of including damping is to add a first-order derivative to the equation (in the same way as friction forces enter a vibrating mechanical system):

$$\frac{\partial^2 u}{\partial t^2} + b \frac{\partial u}{\partial t} = c^2 \frac{\partial^2 u}{\partial x^2} + f(x, t), \quad (2.61)$$

where  $b \geq 0$  is a prescribed damping coefficient.

A typical discretization of (2.61) in terms of centered differences reads

$$[D_t D_t u + b D_{2t} u = c^2 D_x D_x u + f]_i^n. \quad (2.62)$$

Writing out the equation and solving for the unknown  $u_i^{n+1}$  gives the scheme

$$u_i^{n+1} = (1 + \frac{1}{2}b\Delta t)^{-1} ((\frac{1}{2}b\Delta t - 1)u_i^{n-1} + 2u_i^n + C^2(u_{i+1}^n - 2u_i^n + u_{i-1}^n) + \Delta t^2 f_i^n), \quad (2.63)$$

for  $i \in \mathcal{I}_x^i$  and  $n \geq 1$ . New equations must be derived for  $u_i^1$ , and for boundary points in case of Neumann conditions.

The damping is very small in many wave phenomena and thus only evident for very long time simulations. This makes the standard wave equation without damping relevant for a lot of applications.

## 2.8 Building a general 1D wave equation solver

The program `wave1D_dn_vc.py` is a fairly general code for 1D wave propagation problems that targets the following initial-boundary value problem

$$u_{tt} = (c^2(x)u_x)_x + f(x, t), \quad x \in (0, L), \quad t \in (0, T] \quad (2.64)$$

$$u(x, 0) = I(x), \quad x \in [0, L] \quad (2.65)$$

$$u_t(x, 0) = V(t), \quad x \in [0, L] \quad (2.66)$$

$$u(0, t) = U_0(t) \text{ or } u_x(0, t) = 0, \quad t \in (0, T] \quad (2.67)$$

$$u(L, t) = U_L(t) \text{ or } u_x(L, t) = 0, \quad t \in (0, T] \quad (2.68)$$

The only new feature here is the time-dependent Dirichlet conditions. These are trivial to implement:

```
i = Ix[0] # x=0
u[i] = U_0(t[n+1])

i = Ix[-1] # x=L
u[i] = U_L(t[n+1])
```

The `solver` function is a natural extension of the simplest `solver` function in the initial `wave1D_u0.py` program, extended with Neumann boundary conditions ( $u_x = 0$ ), time-varying Dirichlet conditions, as well as a variable wave velocity. The different code segments needed to make these extensions have been shown and commented upon in the preceding text. We refer to the `solver` function in the `wave1D_dn_vc.py` file for all the details.

The vectorization is only applied inside the time loop, not for the initial condition or the first time steps, since this initial work is negligible for long time simulations in 1D problems.

The following sections explain various more advanced programming techniques applied in the general 1D wave equation solver.

### 2.8.1 User action function as a class

A useful feature in the `wave1D_dn_vc.py` program is the specification of the `user_action` function as a class. This part of the program may need some motivation and explanation. Although the `plot_u_st` function (and the `PlotMatplotlib` class) in the `wave1D_u0.viz` function remembers

the local variables in the `viz` function, it is a cleaner solution to store the needed variables together with the function, which is exactly what a class offers.

**The code.** A class for flexible plotting, cleaning up files, making movie files, like the function `wave1D_u0.viz` did, can be coded as follows:

```
class PlotAndStoreSolution:
    """
    Class for the user_action function in solver.
    Visualizes the solution only.
    """
    def __init__(self,
                 casename='tmp',      # Prefix in filenames
                 umin=-1, umax=1,     # Fixed range of y axis
                 pause_between_frames=None, # Movie speed
                 backend='matplotlib',   # or 'gnuplot' or None
                 screen_movie=True,    # Show movie on screen?
                 title='',            # Extra message in title
                 skip_frame=1,        # Skip every skip_frame frame
                 filename=None):      # Name of file with solutions
        self.casename = casename
        self.yaxis = [umin, umax]
        self.pause = pause_between_frames
        self.backend = backend
        if backend is None:
            # Use native matplotlib
            import matplotlib.pyplot as plt
        elif backend in ('matplotlib', 'gnuplot'):
            module = 'scitools.easyviz.' + backend + '_'
            exec('import %s as plt' % module)
        self.plt = plt
        self.screen_movie = screen_movie
        self.title = title
        self.skip_frame = skip_frame
        self.filename = filename
        if filename is not None:
            # Store time points when u is written to file
            self.t = []
            filenames = glob.glob('..' + self.filename + '*.dat.npz')
            for filename in filenames:
                os.remove(filename)

            # Clean up old movie frames
            for filename in glob.glob('frame_*.png'):
                os.remove(filename)

    def __call__(self, u, x, t, n):
        """
        Callback function user_action, call by solver:
        Store solution, plot on screen and save to file.
        """

```

```

# Save solution u to a file using numpy.savetxt
if self.filename is not None:
    name = 'u%04d' % n # array name
    kwargs = {'name: u'}
    fname = '.' + self.filename + '_' + name + '.dat'
    np.savetxt(fname, **kwargs)
    self.t.append(t[n]) # store corresponding time value
    if n == 0:           # save x once
        np.savetxt('.' + self.filename + '_x.dat', x=x)

# Animate
if n % self.skip_frame != 0:
    return
title = 't=% .3f' % t[n]
if self.title:
    title = self.title + ' ' + title
if self.backend is None:
    # native matplotlib animation
    if n == 0:
        self.plt.ion()
        self.lines = self.plt.plot(x, u, 'r-')
        self.plt.axis([x[0], x[-1],
                      self.yaxis[0], self.yaxis[1]])
        self.plt.xlabel('x')
        self.plt.ylabel('u')
        self.plt.title(title)
        self.plt.legend(['t=% .3f' % t[n]])
    else:
        # Update new solution
        self.lines[0].set_ydata(u)
        self.plt.legend(['t=% .3f' % t[n]])
        self.plt.draw()
else:
    # scitools.easyviz animation
    self.plt.plot(x, u, 'r-',
                  xlabel='x', ylabel='u',
                  axis=[x[0], x[-1],
                        self.yaxis[0], self.yaxis[1]],
                  title=title,
                  show=self.screen_movie)

# pause
if t[n] == 0:
    time.sleep(2) # let initial condition stay 2 s
else:
    if self.pause is None:
        pause = 0.2 if u.size < 100 else 0
    time.sleep(pause)

self.plt.savefig('frame_%04d.png' % (n))

```

**Dissection.** Understanding this class requires quite some familiarity with Python in general and class programming in particular. The

class supports plotting with Matplotlib (`backend=None`) or SciTools (`backend=matplotlib` or `backend=gnuplot`) for maximum flexibility.

The constructor shows how we can flexibly import the plotting engine as (typically) `scitools.easyviz.gnuplot_` or `scitools.easyviz.matplotlib_` (note the trailing underscore - it is required). With the `screen_movie` parameter we can suppress displaying each movie frame on the screen. Alternatively, for slow movies associated with fine meshes, one can set `skip_frame=10`, causing every 10 frames to be shown.

The `__call__` method makes `PlotAndStoreSolution` instances behave like functions, so we can just pass an instance, say `p`, as the `user_action` argument in the `solver` function, and any call to `user_action` will be a call to `p.__call__`. The `__call__` method plots the solution on the screen, saves the plot to file, and stores the solution in a file for later retrieval.

More details on storing the solution in files appear in Section C.2.

### 2.8.2 Pulse propagation in two media

The function `pulse` in `wave1D_dn_vc.py` demonstrates wave motion in heterogeneous media where  $c$  varies. One can specify an interval where the wave velocity is decreased by a factor `slowness_factor` (or increased by making this factor less than one). Figure 2.5 shows a typical simulation scenario.

Four types of initial conditions are available:

1. a rectangular pulse (`plug`),
2. a Gaussian function (`gaussian`),
3. a “cosine hat” consisting of one period of the cosine function (`cosinehat`),
4. half a period of a “cosine hat” (`half-cosinehat`)

These peak-shaped initial conditions can be placed in the middle (`loc='center'`) or at the left end (`loc='left'`) of the domain. With the pulse in the middle, it splits in two parts, each with half the initial amplitude, traveling in opposite directions. With the pulse at the left end, centered at  $x = 0$ , and using the symmetry condition  $\partial u / \partial x = 0$ , only a right-going pulse is generated. There is also a left-going pulse, but it travels from  $x = 0$  in negative  $x$  direction and is not visible in the domain  $[0, L]$ .

The `pulse` function is a flexible tool for playing around with various wave shapes and location of a medium with a different wave velocity.

The code is shown to demonstrate how easy it is to reach this flexibility with the building blocks we have already developed:

```

def pulse(C=1,                  # Maximum Courant number
          Nx=200,             # spatial resolution
          animate=True,
          version='vectorized',
          T=2,                 # end time
          loc='left',           # location of initial condition
          pulse_tp='gaussian', # pulse/init.cond. type
          slowness_factor=2,   # inverse of wave vel. in right medium
          medium=[0.7, 0.9],    # interval for right medium
          skip_frame=1,         # skip frames in animations
          sigma=0.05            # width measure of the pulse
          ):
    """
    Various peaked-shaped initial conditions on [0,1].
    Wave velocity is decreased by the slowness_factor inside
    medium. The loc parameter can be 'center' or 'left',
    depending on where the initial pulse is to be located.
    The sigma parameter governs the width of the pulse.
    """
    # Use scaled parameters: L=1 for domain length, c_0=1
    # for wave velocity outside the domain.
    L = 1.0
    c_0 = 1.0
    if loc == 'center':
        xc = L/2
    elif loc == 'left':
        xc = 0

    if pulse_tp in ('gaussian','Gaussian'):
        def I(x):
            return np.exp(-0.5*((x-xc)/sigma)**2)
    elif pulse_tp == 'plug':
        def I(x):
            return 0 if abs(x-xc) > sigma else 1
    elif pulse_tp == 'cosinehat':
        def I(x):
            # One period of a cosine
            w = 2
            a = w*sigma
            return 0.5*(1 + np.cos(np.pi*(x-xc)/a)) \
                   if xc - a <= x <= xc + a else 0

    elif pulse_tp == 'half-cosinehat':
        def I(x):
            # Half a period of a cosine
            w = 4
            a = w*sigma
            return np.cos(np.pi*(x-xc)/a) \

```

```

        if xc - 0.5*a <= x <= xc + 0.5*a else 0
    else:
        raise ValueError('Wrong pulse_tp=%s' % pulse_tp)

def c(x):
    return c_0/slowness_factor \
        if medium[0] <= x <= medium[1] else c_0

umin=-0.5; umax=1.5*I(xc)
casename = '%s_Nx%s_sf%s' % \
    (pulse_tp, Nx, slowness_factor)
action = PlotMediumAndSolution(
    medium, casename=casename, umin=umin, umax=umax,
    skip_frame=skip_frame, screen_movie=animate,
    backend=None, filename='tmpdata')

# Choose the stability limit with given Nx, worst case c
# (lower C will then use this dt, but smaller Nx)
dt = (L/Nx)/c_0
cpu, hashed_input = solver(I=I, V=None, f=None, c=c,
                           U_0=None, U_L=None,
                           L=L, dt=dt, C=C, T=T,
                           user_action=action,
                           version=version,
                           stability_safety_factor=1)

if cpu > 0: # did we generate new data?
    action.close_file(hashed_input)
    action.make_movie_file()
print 'cpu (-1 means no new data generated):', cpu

```

The `PlotMediumAndSolution` class used here is a subclass of `PlotAndStoreSolution` where the medium with reduced  $c$  value, as specified by the `medium` interval, is visualized in the plots.

### Comment on the choices of discretization parameters

The argument  $N_x$  in the `pulse` function does not correspond to the actual spatial resolution of  $C < 1$ , since the `solver` function takes a fixed  $\Delta t$  and  $C$ , and adjusts  $\Delta x$  accordingly. As seen in the `pulse` function, the specified  $\Delta t$  is chosen according to the limit  $C = 1$ , so if  $C < 1$ ,  $\Delta t$  remains the same, but the `solver` function operates with a larger  $\Delta x$  and smaller  $N_x$  than was specified in the call to `pulse`. The practical reason is that we always want to keep  $\Delta t$  fixed such that plot frames and movies are synchronized in time regardless of the value of  $C$  (i.e.,  $\Delta x$  is varied when the Courant number varies).

The reader is encouraged to play around with the `pulse` function:

```
>>> import wave1D_dn_vc as w
>>> w.pulse(loc='left', pulse_tp='cosinehat', Nx=50, every_frame=10)
```

To easily kill the graphics by Ctrl-C and restart a new simulation it might be easier to run the above two statements from the command line with

---

Terminal

```
Terminal> python -c 'import wave1D_dn_vc as w; w.pulse(...)',
```

---

## 2.9 Exercises

### Exercise 2.6: Find the analytical solution to a damped wave equation

Consider the wave equation with damping (2.61). The goal is to find an exact solution to a wave problem with damping and zero source term. A starting point is the standing wave solution from Exercise 2.1. It becomes necessary to include a damping term  $e^{-\beta t}$  and also have both a sine and cosine component in time:

$$u_e(x, t) = e^{-\beta t} \sin kx (A \cos \omega t + B \sin \omega t).$$

Find  $k$  from the boundary conditions  $u(0, t) = u(L, t) = 0$ . Then use the PDE to find constraints on  $\beta$ ,  $\omega$ ,  $A$ , and  $B$ . Set up a complete initial-boundary value problem and its solution.

Filename: `damped_waves`.

### Problem 2.7: Explore symmetry boundary conditions

Consider the simple "plug" wave where  $\Omega = [-L, L]$  and

$$I(x) = \begin{cases} 1, & x \in [-\delta, \delta], \\ 0, & \text{otherwise} \end{cases}$$

for some number  $0 < \delta < L$ . The other initial condition is  $u_t(x, 0) = 0$  and there is no source term  $f$ . The boundary conditions can be set to  $u = 0$ . The solution to this problem is symmetric around  $x = 0$ . This means that we can simulate the wave process in only the half of the domain  $[0, L]$ .

**a)** Argue why the symmetry boundary condition is  $u_x = 0$  at  $x = 0$ .

**Hint.** Symmetry of a function about  $x = x_0$  means that  $f(x_0 + h) = f(x_0 - h)$ .

**b)** Perform simulations of the complete wave problem on  $[-L, L]$ . Thereafter, utilize the symmetry of the solution and run a simulation in half of the domain  $[0, L]$ , using a boundary condition at  $x = 0$ . Compare plots from the two solutions and confirm that they are the same.

**c)** Prove the symmetry property of the solution by setting up the complete initial-boundary value problem and showing that if  $u(x, t)$  is a solution, then also  $u(-x, t)$  is a solution.

**d)** If the code works correctly, the solution  $u(x, t) = x(L - x)(1 + \frac{t}{2})$  should be reproduced exactly. Write a test function `test_quadratic` that checks whether this is the case. Simulate for  $x$  in  $[0, \frac{L}{2}]$  with a symmetry condition at the end  $x = \frac{L}{2}$ .

Filename: `wave1D_symmetric`.

### Exercise 2.8: Send pulse waves through a layered medium

Use the `pulse` function in `wave1D_dn_vc.py` to investigate sending a pulse, located with its peak at  $x = 0$ , through two media with different wave velocities. The (scaled) velocity in the left medium is 1 while it is  $s_f$  in the right medium. Report what happens with a Gaussian pulse, a “cosine hat” pulse, half a “cosine hat” pulse, and a plug pulse for resolutions  $N_x = 40, 80, 160$ , and  $s_f = 2, 4$ . Simulate until  $T = 2$ .  
Filename: `pulse1D`.

### Exercise 2.9: Explain why numerical noise occurs

The experiments performed in Exercise 2.8 shows considerable numerical noise in the form of non-physical waves, especially for  $s_f = 4$  and the plug pulse or the half a “cosinehat” pulse. The noise is much less visible for a Gaussian pulse. Run the case with the plug and half a “cosinehat” pulses for  $s_f = 1, C = 0.9, 0.25$ , and  $N_x = 40, 80, 160$ . Use the numerical dispersion relation to explain the observations. Filename: `pulse1D_analysis`.

### Exercise 2.10: Investigate harmonic averaging in a 1D model

Harmonic means are often used if the wave velocity is non-smooth or discontinuous. Will harmonic averaging of the wave velocity give less numerical noise for the case  $s_f = 4$  in Exercise 2.8? Filename: `pulse1D_harmonic`.

### Problem 2.11: Implement open boundary conditions

To enable a wave to leave the computational domain and travel undisturbed through the boundary  $x = L$ , one can in a one-dimensional problem impose the following condition, called a *radiation condition* or *open boundary condition*:

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0. \quad (2.69)$$

The parameter  $c$  is the wave velocity.

Show that (2.69) accepts a solution  $u = g_R(x - ct)$  (right-going wave), but not  $u = g_L(x + ct)$  (left-going wave). This means that (2.69) will allow any right-going wave  $g_R(x - ct)$  to pass through the boundary undisturbed.

A corresponding open boundary condition for a left-going wave through  $x = 0$  is

$$\frac{\partial u}{\partial t} - c \frac{\partial u}{\partial x} = 0. \quad (2.70)$$

**a)** A natural idea for discretizing the condition (2.69) at the spatial end point  $i = N_x$  is to apply centered differences in time and space:

$$[D_{2t}u + cD_{2x}u = 0]_i^n, \quad i = N_x. \quad (2.71)$$

Eliminate the fictitious value  $u_{N_x+1}^n$  by using the discrete equation at the same point.

The equation for the first step,  $u_i^1$ , is in principle also affected, but we can then use the condition  $u_{N_x} = 0$  since the wave has not yet reached the right boundary.

**b)** A much more convenient implementation of the open boundary condition at  $x = L$  can be based on an explicit discretization

$$[D_t^+u + cD_x^-u = 0]_i^n, \quad i = N_x. \quad (2.72)$$

From this equation, one can solve for  $u_{N_x}^{n+1}$  and apply the formula as a Dirichlet condition at the boundary point. However, the finite difference approximations involved are of first order.

Implement this scheme for a wave equation  $u_{tt} = c^2 u_{xx}$  in a domain  $[0, L]$ , where you have  $u_x = 0$  at  $x = 0$ , the condition (2.69) at  $x = L$ , and an initial disturbance in the middle of the domain, e.g., a plug profile like

$$u(x, 0) = \begin{cases} 1, & L/2 - \ell \leq x \leq L/2 + \ell, \\ 0, & \text{otherwise} \end{cases}$$

Observe that the initial wave is split in two, the left-going wave is reflected at  $x = 0$ , and both waves travel out of  $x = L$ , leaving the solution as  $u = 0$  in  $[0, L]$ . Use a unit Courant number such that the numerical solution is exact. Make a movie to illustrate what happens.

Because this simplified implementation of the open boundary condition works, there is no need to pursue the more complicated discretization in a).

**Hint.** Modify the solver function in `wave1D_dn.py`.

**c)** Add the possibility to have either  $u_x = 0$  or an open boundary condition at the left boundary. The latter condition is discretized as

$$[D_t^+ u - c D_x^+ u = 0]_i^n, \quad i = 0, \quad (2.73)$$

leading to an explicit update of the boundary value  $u_0^{n+1}$ .

The implementation can be tested with a Gaussian function as initial condition:

$$g(x; m, s) = \frac{1}{\sqrt{2\pi}s} e^{-\frac{(x-m)^2}{2s^2}}.$$

Run two tests:

1. Disturbance in the middle of the domain,  $I(x) = g(x; L/2, s)$ , and open boundary condition at the left end.
2. Disturbance at the left end,  $I(x) = g(x; 0, s)$ , and  $u_x = 0$  as symmetry boundary condition at this end.

Make test functions for both cases, testing that the solution is zero after the waves have left the domain.

**d)** In 2D and 3D it is difficult to compute the correct wave velocity normal to the boundary, which is needed in generalizations of the open

boundary conditions in higher dimensions. Test the effect of having a slightly wrong wave velocity in (2.72). Make movies to illustrate what happens.

Filename: `wave1D_open_BC`.

**Remarks.** The condition (2.69) works perfectly in 1D when  $c$  is known. In 2D and 3D, however, the condition reads  $u_t + c_x u_x + c_y u_y = 0$ , where  $c_x$  and  $c_y$  are the wave speeds in the  $x$  and  $y$  directions. Estimating these components (i.e., the direction of the wave) is often challenging. Other methods are normally used in 2D and 3D to let waves move out of a computational domain.

### Exercise 2.12: Implement periodic boundary conditions

It is frequently of interest to follow wave motion over large distances and long times. A straightforward approach is to work with a very large domain, but that might lead to a lot of computations in areas of the domain where the waves cannot be noticed. A more efficient approach is to let a right-going wave out of the domain and at the same time let it enter the domain on the left. This is called a *periodic boundary condition*.

The boundary condition at the right end  $x = L$  is an open boundary condition (see Exercise 2.11) to let a right-going wave out of the domain. At the left end,  $x = 0$ , we apply, in the beginning of the simulation, either a symmetry boundary condition (see Exercise 2.7)  $u_x = 0$ , or an open boundary condition.

This initial wave will split in two and either be reflected or transported out of the domain at  $x = 0$ . The purpose of the exercise is to follow the right-going wave. We can do that with a *periodic boundary condition*. This means that when the right-going wave hits the boundary  $x = L$ , the open boundary condition lets the wave out of the domain, but at the same time we use a boundary condition on the left end  $x = 0$  that feeds the outgoing wave into the domain again. This periodic condition is simply  $u(0) = u(L)$ . The switch from  $u_x = 0$  or an open boundary condition at the left end to a periodic condition can happen when  $u(L, t) > \epsilon$ , where  $\epsilon = 10^{-4}$  might be an appropriate value for determining when the right-going wave hits the boundary  $x = L$ .

The open boundary conditions can conveniently be discretized as explained in Exercise 2.11. Implement the described type of boundary conditions and test them on two different initial shapes: a plug  $u(x, 0) = 1$  for  $x \leq 0.1$ ,  $u(x, 0) = 0$  for  $x > 0.1$ , and a Gaussian function in the

middle of the domain:  $u(x, 0) = \exp(-\frac{1}{2}(x - 0.5)^2/0.05)$ . The domain is the unit interval  $[0, 1]$ . Run these two shapes for Courant numbers 1 and 0.5. Assume constant wave velocity. Make movies of the four cases. Reason why the solutions are correct. Filename: `periodic`.

### Exercise 2.13: Compare discretizations of a Neumann condition

We have a 1D wave equation with variable wave velocity:  $u_{tt} = (qu_x)_x$ . A Neumann condition  $u_x$  at  $x = 0, L$  can be discretized as shown in (2.54) and (2.57).

The aim of this exercise is to examine the rate of the numerical error when using different ways of discretizing the Neumann condition.

- a)** As a test problem,  $q = 1 + (x - L/2)^4$  can be used, with  $f(x, t)$  adapted such that the solution has a simple form, say  $u(x, t) = \cos(\pi x/L) \cos(\omega t)$  for, e.g.,  $\omega = 1$ . Perform numerical experiments and find the convergence rate of the error using the approximation (2.54).
- b)** Switch to  $q(x) = 1 + \cos(\pi x/L)$ , which is symmetric at  $x = 0, L$ , and check the convergence rate of the scheme (2.57). Now,  $q_{i-1/2}$  is a 2nd-order approximation to  $q_i$ ,  $q_{i-1/2} = q_i + 0.25q''_i \Delta x^2 + \dots$ , because  $q'_i = 0$  for  $i = N_x$  (a similar argument can be applied to the case  $i = 0$ ).
- c)** A third discretization can be based on a simple and convenient, but less accurate, one-sided difference:  $u_i - u_{i-1} = 0$  at  $i = N_x$  and  $u_{i+1} - u_i = 0$  at  $i = 0$ . Derive the resulting scheme in detail and implement it. Run experiments with  $q$  from a) or b) to establish the rate of convergence of the scheme.
- d)** A fourth technique is to view the scheme as

$$[D_t D_t u]_i^n = \frac{1}{\Delta x} \left( [q D_x u]_{i+\frac{1}{2}}^n - [q D_x u]_{i-\frac{1}{2}}^n \right) + [f]_i^n,$$

and place the boundary at  $x_{i+\frac{1}{2}}$ ,  $i = N_x$ , instead of exactly at the physical boundary. With this idea of approximating (moving) the boundary, we can just set  $[q D_x u]_{i+\frac{1}{2}}^n = 0$ . Derive the complete scheme using this technique. The implementation of the boundary condition at  $L - \Delta x/2$  is  $\mathcal{O}(\Delta x^2)$  accurate, but the interesting question is what impact the movement of the boundary has on the convergence rate. Compute the errors as usual over the entire mesh and use  $q$  from a) or b).

Filename: `Neumann_discr`.

### Exercise 2.14: Verification by a cubic polynomial in space

The purpose of this exercise is to verify the implementation of the `solver` function in the program `wave1D_n0.py` by using an exact numerical solution for the wave equation  $u_{tt} = c^2 u_{xx} + f$  with Neumann boundary conditions  $u_x(0, t) = u_x(L, t) = 0$ .

A similar verification is used in the file `wave1D_u0.py`, which solves the same PDE, but with Dirichlet boundary conditions  $u(0, t) = u(L, t) = 0$ . The idea of the verification test in function `test_quadratic` in `wave1D_u0.py` is to produce a solution that is a lower-order polynomial such that both the PDE problem, the boundary conditions, and all the discrete equations are exactly fulfilled. Then the `solver` function should reproduce this exact solution to machine precision. More precisely, we seek  $u = X(x)T(t)$ , with  $T(t)$  as a linear function and  $X(x)$  as a parabola that fulfills the boundary conditions. Inserting this  $u$  in the PDE determines  $f$ . It turns out that  $u$  also fulfills the discrete equations, because the truncation error of the discretized PDE has derivatives in  $x$  and  $t$  of order four and higher. These derivatives all vanish for a quadratic  $X(x)$  and linear  $T(t)$ .

It would be attractive to use a similar approach in the case of Neumann conditions. We set  $u = X(x)T(t)$  and seek lower-order polynomials  $X$  and  $T$ . To force  $u_x$  to vanish at the boundary, we let  $X_x$  be a parabola. Then  $X$  is a cubic polynomial. The fourth-order derivative of a cubic polynomial vanishes, so  $u = X(x)T(t)$  will fulfill the discretized PDE also in this case, if  $f$  is adjusted such that  $u$  fulfills the PDE.

However, the discrete boundary condition is not exactly fulfilled by this choice of  $u$ . The reason is that

$$[D_{2x}u]_i^n = u_x(x_i, t_n) + \frac{1}{6}u_{xxx}(x_i, t_n)\Delta x^2 + \mathcal{O}(\Delta x^4). \quad (2.74)$$

At the boundary two boundary points, we must demand that the derivative  $X_x(x) = 0$  such that  $u_x = 0$ . However,  $u_{xxx}$  is a constant and not zero when  $X(x)$  is a cubic polynomial. Therefore, our  $u = X(x)T(t)$  fulfills

$$[D_{2x}u]_i^n = \frac{1}{6}u_{xxx}(x_i, t_n)\Delta x^2,$$

and not

$$[D_{2x}u]_i^n = 0, \quad i = 0, N_x,$$

as it should. (Note that all the higher-order terms  $\mathcal{O}(\Delta x^4)$  also have higher-order derivatives that vanish for a cubic polynomial.) So to summarize, the fundamental problem is that  $u$  as a product of a cubic polynomial and a linear or quadratic polynomial in time is not an exact solution of the discrete boundary conditions.

To make progress, we assume that  $u = X(x)T(t)$ , where  $T$  for simplicity is taken as a prescribed linear function  $1 + \frac{1}{2}t$ , and  $X(x)$  is taken as an *unknown* cubic polynomial  $\sum_{j=0}^3 a_j x^j$ . There are two different ways of determining the coefficients  $a_0, \dots, a_3$  such that both the discretized PDE and the discretized boundary conditions are fulfilled, under the constraint that we can specify a function  $f(x, t)$  for the PDE to feed to the `solver` function in `wave1D_n0.py`. Both approaches are explained in the subexercises.

**a)** One can insert  $u$  in the discretized PDE and find the corresponding  $f$ . Then one can insert  $u$  in the discretized boundary conditions. This yields two equations for the four coefficients  $a_0, \dots, a_3$ . To find the coefficients, one can set  $a_0 = 0$  and  $a_1 = 1$  for simplicity and then determine  $a_2$  and  $a_3$ . This approach will make  $a_2$  and  $a_3$  depend on  $\Delta x$  and  $f$  will depend on both  $\Delta x$  and  $\Delta t$ .

Use `sympy` to perform analytical computations. A starting point is to define  $u$  as follows:

```
def test_cubic1():
    import sympy as sm
    x, t, c, L, dx, dt = sm.symbols('x t c L dx dt')
    i, n = sm.symbols('i n', integer=True)

    # Assume discrete solution is a polynomial of degree 3 in x
    T = lambda t: 1 + sm.Rational(1,2)*t  # Temporal term
    a = sm.symbols('a_0 a_1 a_2 a_3')
    X = lambda x: sum(a[q]*x**q for q in range(4))  # Spatial term
    u = lambda x, t: X(x)*T(t)
```

The symbolic expression for  $u$  is reached by calling `u(x, t)` with `x` and `t` as `sympy` symbols.

Define `DxDx(u, i, n)`, `DtDt(u, i, n)`, and `D2x(u, i, n)` as Python functions for returning the difference approximations  $[D_x D_x u]_i^n$ ,  $[D_t D_t u]_i^n$ , and  $[D_{2x} u]_i^n$ . The next step is to set up the residuals for the equations  $[D_{2x} u]_0^n = 0$  and  $[D_{2x} u]_{N_x}^n = 0$ , where  $N_x = L/\Delta x$ . Call the residuals `R_0` and `R_L`. Substitute  $a_0$  and  $a_1$  by 0 and 1, respectively, in `R_0`, `R_L`, and `a`:

```
R_0 = R_0.subs(a[0], 0).subs(a[1], 1)
R_L = R_L.subs(a[0], 0).subs(a[1], 1)
```

```
a = list(a) # enable in-place assignment
a[0:2] = 0, 1
```

Determining  $a_2$  and  $a_3$  from the discretized boundary conditions is then about solving two equations with respect to  $a_2$  and  $a_3$ , i.e.,  $a[2:]$ :

```
s = sm.solve([R_0, R_L], a[2:])
# s is dictionary with the unknowns a[2] and a[3] as keys
a[2:] = s[a[2]], s[a[3]]
```

Now,  $a$  contains computed values and  $u$  will automatically use these new values since  $X$  accesses  $a$ .

Compute the source term  $f$  from the discretized PDE:  $f_i^n = [D_t D_t u - c^2 D_x D_x u]_i^n$ . Turn  $u$ , the time derivative  $u_t$  (needed for the initial condition  $V(x)$ ), and  $f$  into Python functions. Set numerical values for  $L$ ,  $N_x$ ,  $C$ , and  $c$ . Prescribe the time interval as  $\Delta t = CL/(N_x c)$ , which imply  $\Delta x = c\Delta t/C = L/N_x$ . Define new functions  $I(x)$ ,  $V(x)$ , and  $f(x, t)$  as wrappers of the ones made above, where fixed values of  $L$ ,  $c$ ,  $\Delta x$ , and  $\Delta t$  are inserted, such that  $I$ ,  $V$ , and  $f$  can be passed on to the `solver` function. Finally, call `solver` with a `user_action` function that compares the numerical solution to this exact solution  $u$  of the discrete PDE problem.

**Hint.** To turn a `sympy` expression  $e$ , depending on a series of symbols, say  $x$ ,  $t$ ,  $dx$ ,  $dt$ ,  $L$ , and  $c$ , into a plain Python function `e_exact(x,t,L,dx,dt,c)`, one can write

```
e_exact = sm.lambdify([x,t,L,dx,dt,c], e, 'numpy')
```

The '`numpy`' argument is a good habit as the `e_exact` function will then work with array arguments if it contains mathematical functions (but here we only do plain arithmetics, which automatically work with arrays).

**b)** An alternative way of determining  $a_0, \dots, a_3$  is to reason as follows. We first construct  $X(x)$  such that the boundary conditions are fulfilled:  $X = x(L - x)$ . However, to compensate for the fact that this choice of  $X$  does not fulfill the discrete boundary condition, we seek  $u$  such that

$$u_x = \frac{\partial}{\partial x} x(L - x)T(t) - \frac{1}{6}u_{xxx}\Delta x^2,$$

since this  $u$  will fit the discrete boundary condition. Assuming  $u = T(t) \sum_{j=0}^3 a_j x^j$ , we can use the above equation to determine the coefficients  $a_1, a_2, a_3$ . A value, e.g., 1 can be used for  $a_0$ . The following `sympy` code computes this  $u$ :

```

def test_cubic2():
    import sympy as sm
    x, t, c, L, dx = sm.symbols('x t c L dx')
    T = lambda t: 1 + sm.Rational(1,2)*t # Temporal term
    # Set u as a 3rd-degree polynomial in space
    X = lambda x: sum(a[i]*x**i for i in range(4))
    a = sm.symbols('a_0 a_1 a_2 a_3')
    u = lambda x, t: X(x)*T(t)
    # Force discrete boundary condition to be zero by adding
    # a correction term the analytical suggestion x*(L-x)*T
    # u_x = x*(L-x)*T(t) - 1/6*u_xxx*dx**2
    R = sm.diff(u(x,t), x) -
        x*(L-x) - sm.Rational(1,6)*sm.diff(u(x,t), x, x, x)*dx**2
    # R is a polynomial: force all coefficients to vanish.
    # Turn R to Poly to extract coefficients:
    R = sm.poly(R, x)
    coeff = R.all_coeffs()
    s = sm.solve(coeff, a[1:]) # a[0] is not present in R
    # s is dictionary with a[i] as keys
    # Fix a[0] as 1
    s[a[0]] = 1
    X = lambda x: sm.simplify(sum(s[a[i]]*x**i for i in range(4)))
    u = lambda x, t: X(x)*T(t)
    print 'u:', u(x,t)

```

The next step is to find the source term  $f_e$  by inserting  $u_e$  in the PDE. Thereafter, turn  $u$ ,  $f$ , and the time derivative of  $u$  into plain Python functions as in a), and then wrap these functions in new functions  $I$ ,  $V$ , and  $f$ , with the right signature as required by the `solver` function. Set parameters as in a) and check that the solution is exact to machine precision at each time level using an appropriate `user_action` function. Filename: `wave1D_n0_test_cubic`.

## 2.10 Analysis of the difference equations

### 2.10.1 Properties of the solution of the wave equation

The wave equation

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2}$$

has solutions of the form

$$u(x, t) = g_R(x - ct) + g_L(x + ct), \quad (2.75)$$

for any functions  $g_R$  and  $g_L$  sufficiently smooth to be differentiated twice. The result follows from inserting (2.75) in the wave equation. A function

of the form  $g_R(x - ct)$  represents a signal moving to the right in time with constant velocity  $c$ . This feature can be explained as follows. At time  $t = 0$  the signal looks like  $g_R(x)$ . Introducing a moving horizontal coordinate  $\xi = x - ct$ , we see the function  $g_R(\xi)$  is “at rest” in the  $\xi$  coordinate system, and the shape is always the same. Say the  $g_R(\xi)$  function has a peak at  $\xi = 0$ . This peak is located at  $x = ct$ , which means that it moves with the velocity  $dx/dt = c$  in the  $x$  coordinate system. Similarly,  $g_L(x + ct)$  is a function, initially with shape  $g_L(x)$ , that moves in the negative  $x$  direction with constant velocity  $c$  (introduce  $\xi = x + ct$ , look at the point  $\xi = 0$ ,  $x = -ct$ , which has velocity  $dx/dt = -c$ ).

With the particular initial conditions

$$u(x, 0) = I(x), \quad \frac{\partial}{\partial t} u(x, 0) = 0,$$

we get, with  $u$  as in (2.75),

$$g_R(x) + g_L(x) = I(x), \quad -cg'_R(x) + cg'_L(x) = 0.$$

The former suggests  $g_R = g_L$ , and the former then leads to  $g_R = g_L = I/2$ . Consequently,

$$u(x, t) = \frac{1}{2}I(x - ct) + \frac{1}{2}I(x + ct). \quad (2.76)$$

The interpretation of (2.76) is that the initial shape of  $u$  is split into two parts, each with the same shape as  $I$  but half of the initial amplitude. One part is traveling to the left and the other one to the right.

The solution has two important physical features: constant amplitude of the left and right wave, and constant velocity of these two waves. It turns out that the numerical solution will also preserve the constant amplitude, but the velocity depends on the mesh parameters  $\Delta t$  and  $\Delta x$ .

The solution (2.76) will be influenced by boundary conditions when the parts  $\frac{1}{2}I(x - ct)$  and  $\frac{1}{2}I(x + ct)$  hit the boundaries and get, e.g., reflected back into the domain. However, when  $I(x)$  is nonzero only in a small part in the middle of the spatial domain  $[0, L]$ , which means that the boundaries are placed far away from the initial disturbance of  $u$ , the solution (2.76) is very clearly observed in a simulation.

A useful representation of solutions of wave equations is a linear combination of sine and/or cosine waves. Such a sum of waves is a solution if the governing PDE is linear and each sine or cosine wave fulfills the equation. To ease analytical calculations by hand we shall work with complex exponential functions instead of real-valued sine or

cosine functions. The real part of complex expressions will typically be taken as the physical relevant quantity (whenever a physical relevant quantity is strictly needed). The idea now is to build  $I(x)$  of complex wave components  $e^{ikx}$ :

$$I(x) \approx \sum_{k \in K} b_k e^{ikx}. \quad (2.77)$$

Here,  $k$  is the frequency of a component,  $K$  is some set of all the discrete  $k$  values needed to approximate  $I(x)$  well, and  $b_k$  are constants that must be determined. We will very seldom need to compute the  $b_k$  coefficients: most of the insight we look for, and the understanding of the numerical methods we want to establish, come from investigating how the PDE and the scheme treat a single component  $e^{ikx}$  wave.

Letting the number of  $k$  values in  $K$  tend to infinity, makes the sum (2.77) converge to  $I(x)$ . This sum is known as a *Fourier series* representation of  $I(x)$ . Looking at (2.76), we see that the solution  $u(x, t)$ , when  $I(x)$  is represented as in (2.77), is also built of basic complex exponential wave components of the form  $e^{ik(x \pm ct)}$  according to

$$u(x, t) = \frac{1}{2} \sum_{k \in K} b_k e^{ik(x-ct)} + \frac{1}{2} \sum_{k \in K} b_k e^{ik(x+ct)}. \quad (2.78)$$

It is common to introduce the frequency in time  $\omega = kc$  and assume that  $u(x, t)$  is a sum of basic wave components written as  $e^{ikx - \omega t}$ . (Observe that inserting such a wave component in the governing PDE reveals that  $\omega^2 = k^2 c^2$ , or  $\omega = \pm kc$ , reflecting the two solutions: one  $(+kc)$  traveling to the right and the other  $(-kc)$  traveling to the left.)

## 2.10.2 More precise definition of Fourier representations

The above introduction to function representation by sine and cosine waves was quick and intuitive, but will suffice as background knowledge for the following material of single wave component analysis. However, to understand all details of how different wave components sum up to the analytical and numerical solutions, a more precise mathematical treatment is helpful and therefore summarized below.

It is well known that periodic functions can be represented by Fourier series. A generalization of the Fourier series idea to non-periodic functions defined on the real line is the *Fourier transform*:

$$I(x) = \int_{-\infty}^{\infty} A(k) e^{ikx} dk, \quad (2.79)$$

$$A(k) = \int_{-\infty}^{\infty} I(x) e^{-ikx} dx. \quad (2.80)$$

The function  $A(k)$  reflects the weight of each wave component  $e^{ikx}$  in an infinite sum of such wave components. That is,  $A(k)$  reflects the frequency content in the function  $I(x)$ . Fourier transforms are particularly fundamental for analyzing and understanding time-varying signals.

The solution of the linear 1D wave PDE can be expressed as

$$u(x, t) = \int_{-\infty}^{\infty} A(k) e^{i(kx - \omega(k)t)} dx.$$

In a finite difference method, we represent  $u$  by a mesh function  $u_q^n$ , where  $n$  counts temporal mesh points and  $q$  counts the spatial ones (the usual counter for spatial points,  $i$ , is here already used as imaginary unit). Similarly,  $I(x)$  is approximated by the mesh function  $I_q$ ,  $q = 0, \dots, N_x$ . On a mesh, it does not make sense to work with wave components  $e^{ikx}$  for very large  $k$ , because the shortest possible sine or cosine wave that can be represented uniquely on a mesh with spacing  $\Delta x$  is the wave with wavelength  $2\Delta x$ . This wave has its peaks and throughs at every two mesh points. That is, the wave “jumps up and down” between the mesh points.

The corresponding  $k$  value for the shortest possible wave in the mesh is  $k = 2\pi/(2\Delta x) = \pi/\Delta x$ . This maximum frequency is known as the *Nyquist frequency*. Within the range of relevant frequencies  $(0, \pi/\Delta x]$  one defines the [discrete Fourier transform](#), using  $N_x + 1$  discrete frequencies:

$$I_q = \frac{1}{N_x + 1} \sum_{k=0}^{N_x} A_k e^{i2\pi k j / (N_x + 1)}, \quad i = 0, \dots, N_x, \quad (2.81)$$

$$A_k = \sum_{q=0}^{N_x} I_q e^{-i2\pi k q / (N_x + 1)}, \quad k = 0, \dots, N_x + 1. \quad (2.82)$$

The  $A_k$  values represent the discrete Fourier transform of the  $I_q$  values, which themselves are the inverse discrete Fourier transform of the  $A_k$  values.

The discrete Fourier transform is efficiently computed by the *Fast Fourier transform* algorithm. For a real function  $I(x)$ , the relevant Python

code for computing and plotting the discrete Fourier transform appears in the example below.

```

import numpy as np
from numpy import sin, pi

def I(x):
    return sin(2*pi*x) + 0.5*sin(4*pi*x) + 0.1*sin(6*pi*x)

# Mesh
L = 10; Nx = 100
x = np.linspace(0, L, Nx+1)
dx = L/float(Nx)

# Discrete Fourier transform
A = np.fft.rfft(I(x))
A_amplitude = np.abs(A)

# Compute the corresponding frequencies
freqs = np.linspace(0, pi/dx, A_amplitude.size)

import matplotlib.pyplot as plt
plt.plot(freqs, A_amplitude)
plt.show()

```

### 2.10.3 Stability

The scheme

$$[D_t D_t u = c^2 D_x D_x u]_q^n \quad (2.83)$$

for the wave equation  $u_{tt} = c^2 u_{xx}$  allows basic wave components

$$u_q^n = e^{i(kx_q - \tilde{\omega}t_n)}$$

as solution, but it turns out that the frequency in time,  $\tilde{\omega}$ , is not equal to the exact frequency  $\omega = kc$ . The goal now is to find exactly what  $\tilde{\omega}$  is. We ask two key questions:

- How accurate is  $\tilde{\omega}$  compared to  $\omega$ ?
- Does the amplitude of such a wave component preserve its (unit) amplitude, as it should, or does it get amplified or damped in time (because of a complex  $\tilde{\omega}$ )?

The following analysis will answer these questions. We shall continue using  $q$  as an identifier for a certain mesh point in the  $x$  direction.

**Preliminary results.** A key result needed in the investigations is the finite difference approximation of a second-order derivative acting on a complex wave component:

$$[D_t D_t e^{i\omega t}]^n = -\frac{4}{\Delta t^2} \sin^2\left(\frac{\omega \Delta t}{2}\right) e^{i\omega n \Delta t}.$$

By just changing symbols ( $\omega \rightarrow k$ ,  $t \rightarrow x$ ,  $n \rightarrow q$ ) it follows that

$$[D_x D_x e^{ikx}]_q = -\frac{4}{\Delta x^2} \sin^2\left(\frac{k \Delta x}{2}\right) e^{ikq \Delta x}.$$

**Numerical wave propagation.** Inserting a basic wave component  $u_q^n = e^{i(kx_q - \tilde{\omega}t_n)}$  in (2.83) results in the need to evaluate two expressions:

$$\begin{aligned} [D_t D_t e^{ikx} e^{-i\tilde{\omega}t}]_q^n &= [D_t D_t e^{-i\tilde{\omega}t}]^n e^{ikq \Delta x} \\ &= -\frac{4}{\Delta t^2} \sin^2\left(\frac{\tilde{\omega} \Delta t}{2}\right) e^{-i\tilde{\omega}n \Delta t} e^{ikq \Delta x} \end{aligned} \quad (2.84)$$

$$\begin{aligned} [D_x D_x e^{ikx} e^{-i\tilde{\omega}t}]_q^n &= [D_x D_x e^{ikx}]_q e^{-i\tilde{\omega}n \Delta t} \\ &= -\frac{4}{\Delta x^2} \sin^2\left(\frac{k \Delta x}{2}\right) e^{ikq \Delta x} e^{-i\tilde{\omega}n \Delta t}. \end{aligned} \quad (2.85)$$

Then the complete scheme,

$$[D_t D_t e^{ikx} e^{-i\tilde{\omega}t}]_q^n = c^2 [D_x D_x e^{ikx} e^{-i\tilde{\omega}t}]_q^n$$

leads to the following equation for the unknown numerical frequency  $\tilde{\omega}$  (after dividing by  $-e^{ikx} e^{-i\tilde{\omega}t}$ ):

$$\frac{4}{\Delta t^2} \sin^2\left(\frac{\tilde{\omega} \Delta t}{2}\right) = c^2 \frac{4}{\Delta x^2} \sin^2\left(\frac{k \Delta x}{2}\right),$$

or

$$\sin^2\left(\frac{\tilde{\omega} \Delta t}{2}\right) = C^2 \sin^2\left(\frac{k \Delta x}{2}\right), \quad (2.86)$$

where

$$C = \frac{c \Delta t}{\Delta x} \quad (2.87)$$

is the Courant number. Taking the square root of (2.86) yields

$$\sin\left(\frac{\tilde{\omega}\Delta t}{2}\right) = C \sin\left(\frac{k\Delta x}{2}\right), \quad (2.88)$$

Since the exact  $\omega$  is real it is reasonable to look for a real solution  $\tilde{\omega}$  of (2.88). The right-hand side of (2.88) must then be in  $[-1, 1]$  because the sine function on the left-hand side has values in  $[-1, 1]$  for real  $\tilde{\omega}$ . The sine function on the right-hand side can attain the value 1 when

$$\frac{k\Delta x}{2} = m\frac{\pi}{2}, \quad m \in \mathbb{Z}.$$

With  $m = 1$  we have  $k\Delta x = \pi$ , which means that the wavelength  $\lambda = 2\pi/k$  becomes  $2\Delta x$ . This is the absolutely shortest wavelength that can be represented on the mesh: the wave jumps up and down between each mesh point. Larger values of  $|m|$  are irrelevant since these correspond to  $k$  values whose waves are too short to be represented on a mesh with spacing  $\Delta x$ . For the shortest possible wave in the mesh,  $\sin(k\Delta x/2) = 1$ , and we must require

$$C \leq 1. \quad (2.89)$$

Consider a right-hand side in (2.88) of magnitude larger than unity. The solution  $\tilde{\omega}$  of (2.88) must then be a complex number  $\tilde{\omega} = \tilde{\omega}_r + i\tilde{\omega}_i$  because the sine function is larger than unity for a complex argument. One can show that for any  $\omega_i$  there will also be a corresponding solution with  $-\omega_i$ . The component with  $\omega_i > 0$  gives an amplification factor  $e^{\omega_i t}$  that grows exponentially in time. We cannot allow this and must therefore require  $C \leq 1$  as a *stability criterion*.

### Remark on the stability requirement

For smoother wave components with longer wave lengths per length  $\Delta x$ , (2.89) can in theory be relaxed. However, small round-off errors are always present in a numerical solution and these vary arbitrarily from mesh point to mesh point and can be viewed as unavoidable noise with wavelength  $2\Delta x$ . As explained,  $C > 1$  will for this very small noise lead to exponential growth of the shortest possible wave component in the mesh. This noise will therefore grow with time and destroy the whole solution.

### 2.10.4 Numerical dispersion relation

Equation (2.88) can be solved with respect to  $\tilde{\omega}$ :

$$\tilde{\omega} = \frac{2}{\Delta t} \sin^{-1} \left( C \sin \left( \frac{k \Delta x}{2} \right) \right). \quad (2.90)$$

The relation between the numerical frequency  $\tilde{\omega}$  and the other parameters  $k$ ,  $c$ ,  $\Delta x$ , and  $\Delta t$  is called a *numerical dispersion relation*. Correspondingly,  $\omega = kc$  is the *analytical dispersion relation*. In general, dispersion refers to the phenomenon where the wave velocity depends on the spatial frequency ( $k$ , or the wave length  $\lambda = 2\pi/k$ ) of the wave. Since the wave velocity is  $\omega/k = c$ , we realize that the analytical dispersion relation reflects the fact that there is no dispersion. However, in a numerical scheme we have dispersive waves where the wave velocity depends on  $k$ .

The special case  $C = 1$  deserves attention since then the right-hand side of (2.90) reduces to

$$\frac{2}{\Delta t} \frac{k \Delta x}{2} = \frac{1}{\Delta t} \frac{\omega \Delta x}{c} = \frac{\omega}{C} = \omega.$$

That is,  $\tilde{\omega} = \omega$  and the numerical solution is exact at all mesh points regardless of  $\Delta x$  and  $\Delta t$ ! This implies that the numerical solution method is also an analytical solution method, at least for computing  $u$  at discrete points (the numerical method says nothing about the variation of  $u$  *between* the mesh points, and employing the common linear interpolation for extending the discrete solution gives a curve that in general deviates from the exact one).

For a closer examination of the error in the numerical dispersion relation when  $C < 1$ , we can study  $\tilde{\omega} - \omega$ ,  $\tilde{\omega}/\omega$ , or the similar error measures in wave velocity:  $\tilde{c} - c$  and  $\tilde{c}/c$ , where  $c = \omega/k$  and  $\tilde{c} = \tilde{\omega}/k$ . It appears that the most convenient expression to work with is  $\tilde{c}/c$ , since it can be written as a function of just two parameters:

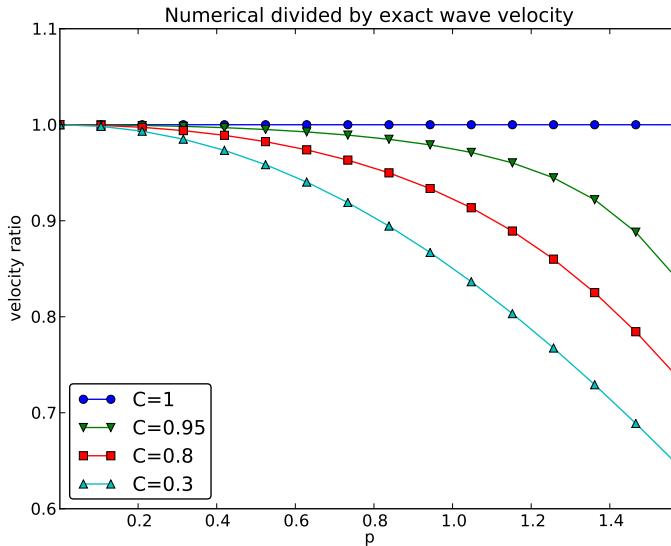
$$\frac{\tilde{c}}{c} = \frac{1}{Cp} \sin^{-1} (C \sin p),$$

with  $p = k \Delta x / 2$  as a non-dimensional measure of the spatial frequency. In essence,  $p$  tells how many spatial mesh points we have per wave length in space for the wave component with frequency  $k$  (recall that the wave length is  $2\pi/k$ ). That is,  $p$  reflects how well the spatial variation of the wave component is resolved in the mesh. Wave components with wave length less than  $2\Delta x$  ( $2\pi/k < 2\Delta x$ ) are not visible in the mesh, so it does not make sense to have  $p > \pi/2$ .

We may introduce the function  $r(C, p) = \tilde{c}/c$  for further investigation of numerical errors in the wave velocity:

$$r(C, p) = \frac{1}{Cp} \sin^{-1}(C \sin p), \quad C \in (0, 1], \quad p \in (0, \pi/2]. \quad (2.91)$$

This function is very well suited for plotting since it combines several parameters in the problem into a dependence on two dimensionless numbers,  $C$  and  $p$ .



**Fig. 2.6** The fractional error in the wave velocity for different Courant numbers.

Defining

```
def r(C, p):
    return 2/(C*p)*asin(C*sin(p))
```

we can plot  $r(C, p)$  as a function of  $p$  for various values of  $C$ , see Figure 2.6. Note that the shortest waves have the most erroneous velocity, and that short waves move more slowly than they should.

We can also easily make a Taylor series expansion in the discretization parameter  $p$ :

```
>>> import sympy as sym
>>> C, p = sym.symbols('C p')
```

```
>>> # Compute the 7 first terms around p=0 with no O() term
>>> rs = r(C, p).series(p, 0, 7).removeO()
>>> rs
p**6*(5*C**6/112 - C**4/16 + 13*C**2/720 - 1/5040) +
p**4*(3*C**4/40 - C**2/12 + 1/120) +
p**2*(C**2/6 - 1/6) + 1

>>> # Pick out the leading order term, but drop the constant 1
>>> rs_error_leading_order = (rs - 1).extract_leading_order(p)
>>> rs_error_leading_order
p**2*(C**2/6 - 1/6)

>>> # Turn the series expansion into a Python function
>>> rs_pyfunc = lambdify([C, p], rs, modules='numpy')

>>> # Check: rs_pyfunc is exact (=1) for C=1
>>> rs_pyfunc(1, 0.1)
1.0
```

Note that without the `.removeO()` call the series gets an  $O(x^{**7})$  term that makes it impossible to convert the series to a Python function (for, e.g., plotting).

From the `rs_error_leading_order` expression above, we see that the leading order term in the error of this series expansion is

$$\frac{1}{6} \left( \frac{k\Delta x}{2} \right)^2 (C^2 - 1) = \frac{k^2}{24} (c^2 \Delta t^2 - \Delta x^2), \quad (2.92)$$

pointing to an error  $\mathcal{O}(\Delta t^2, \Delta x^2)$ , which is compatible with the errors in the difference approximations ( $D_t D_t u$  and  $D_x D_x u$ ).

We can do more with a series expansion, e.g., factor it to see how the factor  $C - 1$  plays a significant role. To this end, we make a list of the terms, factor each term, and then sum the terms:

```
>>> rs = r(C, p).series(p, 0, 4).removeO().as_ordered_terms()
>>> rs
[1, C**2*p**2/6 - p**2/6,
 3*C**4*p**4/40 - C**2*p**4/12 + p**4/120,
 5*C**6*p**6/112 - C**4*p**6/16 + 13*C**2*p**6/720 - p**6/5040]
>>> rs = [factor(t) for t in rs]
>>> rs
[1, p**2*(C - 1)*(C + 1)/6,
 p**4*(C - 1)*(C + 1)*(3*C - 1)*(3*C + 1)/120,
 p**6*(C - 1)*(C + 1)*(225*C**4 - 90*C**2 + 1)/5040]
>>> rs = sum(rs) # Python's sum function sums the list
>>> rs
p**6*(C - 1)*(C + 1)*(225*C**4 - 90*C**2 + 1)/5040 +
p**4*(C - 1)*(C + 1)*(3*C - 1)*(3*C + 1)/120 +
p**2*(C - 1)*(C + 1)/6 + 1
```

We see from the last expression that  $C = 1$  makes all the terms in  $\mathbf{r}$ s vanish. Since we already know that the numerical solution is exact for  $C = 1$ , the remaining terms in the Taylor series expansion will also contain factors of  $C - 1$  and cancel for  $C = 1$ .

### 2.10.5 Extending the analysis to 2D and 3D

The typical analytical solution of a 2D wave equation

$$u_{tt} = c^2(u_{xx} + u_{yy}),$$

is a wave traveling in the direction of  $\mathbf{k} = k_x \mathbf{i} + k_y \mathbf{j}$ , where  $\mathbf{i}$  and  $\mathbf{j}$  are unit vectors in the  $x$  and  $y$  directions, respectively ( $\mathbf{i}$  must not be confused with  $i = \sqrt{-1}$ ). Such a wave can be expressed by

$$u(x, y, t) = g(k_x x + k_y y - kct)$$

for some twice differentiable function  $g$ , or with  $\omega = kc$ ,  $k = |\mathbf{k}|$ :

$$u(x, y, t) = g(k_x x + k_y y - \omega t).$$

We can, in particular, build a solution by adding complex Fourier components of the form

$$e^{(i(k_x x + k_y y - \omega t))}.$$

A discrete 2D wave equation can be written as

$$[D_t D_t u = c^2(D_x D_x u + D_y D_y u)]_{q,r}^n. \quad (2.93)$$

This equation admits a Fourier component

$$u_{q,r}^n = e^{(i(k_x q \Delta x + k_y r \Delta y - \tilde{\omega} n \Delta t))}, \quad (2.94)$$

as solution. Letting the operators  $D_t D_t$ ,  $D_x D_x$ , and  $D_y D_y$  act on  $u_{q,r}^n$  from (2.94) transforms (2.93) to

$$\frac{4}{\Delta t^2} \sin^2 \left( \frac{\tilde{\omega} \Delta t}{2} \right) = c^2 \frac{4}{\Delta x^2} \sin^2 \left( \frac{k_x \Delta x}{2} \right) + c^2 \frac{4}{\Delta y^2} \sin^2 \left( \frac{k_y \Delta y}{2} \right). \quad (2.95)$$

or

$$\sin^2 \left( \frac{\tilde{\omega} \Delta t}{2} \right) = C_x^2 \sin^2 p_x + C_y^2 \sin^2 p_y, \quad (2.96)$$

where we have eliminated the factor 4 and introduced the symbols

$$C_x = \frac{c\Delta t}{\Delta x}, \quad C_y = \frac{c\Delta t}{\Delta y}, \quad p_x = \frac{k_x \Delta x}{2}, \quad p_y = \frac{k_y \Delta y}{2}.$$

For a real-valued  $\tilde{\omega}$  the right-hand side must be less than or equal to unity in absolute value, requiring in general that

$$C_x^2 + C_y^2 \leq 1. \quad (2.97)$$

This gives the stability criterion, more commonly expressed directly in an inequality for the time step:

$$\Delta t \leq \frac{1}{c} \left( \frac{1}{\Delta x^2} + \frac{1}{\Delta y^2} \right)^{-1/2} \quad (2.98)$$

A similar, straightforward analysis for the 3D case leads to

$$\Delta t \leq \frac{1}{c} \left( \frac{1}{\Delta x^2} + \frac{1}{\Delta y^2} + \frac{1}{\Delta z^2} \right)^{-1/2} \quad (2.99)$$

In the case of a variable coefficient  $c^2 = c^2(\mathbf{x})$ , we must use the worst-case value

$$\bar{c} = \sqrt{\max_{\mathbf{x} \in \Omega} c^2(\mathbf{x})} \quad (2.100)$$

in the stability criteria. Often, especially in the variable wave velocity case, it is wise to introduce a safety factor  $\beta \in (0, 1]$  too:

$$\Delta t \leq \beta \frac{1}{\bar{c}} \left( \frac{1}{\Delta x^2} + \frac{1}{\Delta y^2} + \frac{1}{\Delta z^2} \right)^{-1/2} \quad (2.101)$$

The exact numerical dispersion relations in 2D and 3D becomes, for constant  $c$ ,

$$\tilde{\omega} = \frac{2}{\Delta t} \sin^{-1} \left( \left( C_x^2 \sin^2 p_x + C_y^2 \sin^2 p_y \right)^{\frac{1}{2}} \right), \quad (2.102)$$

$$\tilde{\omega} = \frac{2}{\Delta t} \sin^{-1} \left( \left( C_x^2 \sin^2 p_x + C_y^2 \sin^2 p_y + C_z^2 \sin^2 p_z \right)^{\frac{1}{2}} \right). \quad (2.103)$$

We can visualize the numerical dispersion error in 2D much like we did in 1D. To this end, we need to reduce the number of parameters in  $\tilde{\omega}$ . The direction of the wave is parameterized by the polar angle  $\theta$ , which means that

$$k_x = k \sin \theta, \quad k_y = k \cos \theta.$$

A simplification is to set  $\Delta x = \Delta y = h$ . Then  $C_x = C_y = c\Delta t/h$ , which we call  $C$ . Also,

$$p_x = \frac{1}{2}kh \cos \theta, \quad p_y = \frac{1}{2}kh \sin \theta.$$

The numerical frequency  $\tilde{\omega}$  is now a function of three parameters:

- $C$ , reflecting the number of cells a wave is displaced during a time step,
- $p = \frac{1}{2}kh$ , reflecting the number of cells per wave length in space,
- $\theta$ , expressing the direction of the wave.

We want to visualize the error in the numerical frequency. To avoid having  $\Delta t$  as a free parameter in  $\tilde{\omega}$ , we work with  $\tilde{c}/c = \tilde{\omega}/(kc)$ . The coefficient in front of the  $\sin^{-1}$  factor is then

$$\frac{2}{kc\Delta t} = \frac{2}{2kc\Delta th/h} = \frac{1}{Ckh} = \frac{2}{Cp},$$

and

$$\frac{\tilde{c}}{c} = \frac{2}{Cp} \sin^{-1} \left( C \left( \sin^2(p \cos \theta) + \sin^2(p \sin \theta) \right)^{\frac{1}{2}} \right).$$

We want to visualize this quantity as a function of  $p$  and  $\theta$  for some values of  $C \leq 1$ . It is instructive to make color contour plots of  $1 - \tilde{c}/c$  in *polar coordinates* with  $\theta$  as the angular coordinate and  $p$  as the radial coordinate.

The stability criterion (2.97) becomes  $C \leq C_{\max} = 1/\sqrt{2}$  in the present 2D case with the  $C$  defined above. Let us plot  $1 - \tilde{c}/c$  in polar coordinates for  $C_{\max}$ ,  $0.9C_{\max}$ ,  $0.5C_{\max}$ ,  $0.2C_{\max}$ . The program below does the somewhat tricky work in Matplotlib, and the result appears in Figure 2.7. From the figure we clearly see that the maximum  $C$  value gives the best results, and that waves whose propagation direction makes an angle of 45 degrees with an axis are the most accurate.

```
def dispersion_relation_2D(p, theta, C):
    arg = C*sqrt(sin(p*cos(theta))**2 +
                  sin(p*sin(theta))**2)
    c_frac = 2. / (C*p)*arcsin(arg)

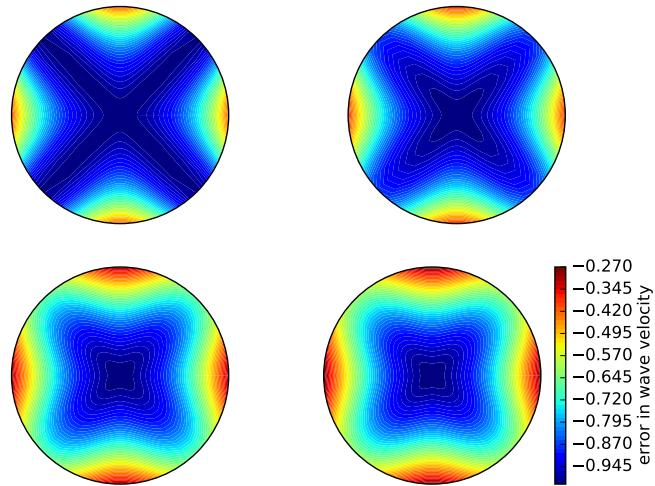
    return c_frac
```

```
import numpy as np
from numpy import \
    cos, sin, arcsin, sqrt, pi # for nicer math formulas

r = p = np.linspace(0.001, pi/2, 101)
theta = np.linspace(0, 2*pi, 51)
r, theta = np.meshgrid(r, theta)

# Make 2x2 filled contour plots for 4 values of C
import matplotlib.pyplot as plt
C_max = 1/sqrt(2)
C = [[C_max, 0.9*C_max], [0.5*C_max, 0.2*C_max]]
fix, axes = plt.subplots(2, 2, subplot_kw=dict(polar=True))
for row in range(2):
    for column in range(2):
        error = 1 - dispersion_relation_2D(
            p, theta, C[row][column])
        print error.min(), error.max()
        # use vmin=error.min(), vmax=error.max()
        cax = axes[row][column].contourf(
            theta, r, error, 50, vmin=-1, vmax=-0.28)
        axes[row][column].set_xticks([])
        axes[row][column].set_yticks([])

# Add colorbar to the last plot
cbar = plt.colorbar(cax)
cbar.ax.set_ylabel('error in wave velocity')
plt.savefig('disprel2D.png'); plt.savefig('disprel2D.pdf')
plt.show()
```



**Fig. 2.7** Error in numerical dispersion in 2D.

## 2.11 Finite difference methods for 2D and 3D wave equations

A natural next step is to consider extensions of the methods for various variants of the one-dimensional wave equation to two-dimensional (2D) and three-dimensional (3D) versions of the wave equation.

### 2.11.1 Multi-dimensional wave equations

The general wave equation in  $d$  space dimensions, with constant wave velocity  $c$ , can be written in the compact form

$$\frac{\partial^2 u}{\partial t^2} = c^2 \nabla^2 u \text{ for } \mathbf{x} \in \Omega \subset \mathbb{R}^d, \quad t \in (0, T], \quad (2.104)$$

where

$$\nabla^2 u = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2},$$

in a 2D problem ( $d = 2$ ) and

$$\nabla^2 u = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + \frac{\partial^2 u}{\partial z^2},$$

in three space dimensions ( $d = 3$ ).

Many applications involve variable coefficients, and the general wave equation in  $d$  dimensions is in this case written as

$$\varrho \frac{\partial^2 u}{\partial t^2} = \nabla \cdot (q \nabla u) + f \text{ for } \mathbf{x} \in \Omega \subset \mathbb{R}^d, \quad t \in (0, T], \quad (2.105)$$

which in, e.g., 2D becomes

$$\varrho(x, y) \frac{\partial^2 u}{\partial t^2} = \frac{\partial}{\partial x} \left( q(x, y) \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( q(x, y) \frac{\partial u}{\partial y} \right) + f(x, y, t). \quad (2.106)$$

To save some writing and space we may use the index notation, where subscript  $t$ ,  $x$ , or  $y$  means differentiation with respect to that coordinate. For example,

$$\begin{aligned} \frac{\partial^2 u}{\partial t^2} &= u_{tt}, \\ \frac{\partial}{\partial y} \left( q(x, y) \frac{\partial u}{\partial y} \right) &= (qu_y)_y. \end{aligned}$$

These comments extend straightforwardly to 3D, which means that the 3D versions of the two wave PDEs, with and without variable coefficients, can be stated as

$$u_{tt} = c^2(u_{xx} + u_{yy} + u_{zz}) + f, \quad (2.107)$$

$$\varrho u_{tt} = (qu_x)_x + (qu_z)_z + (qu_z)_z + f. \quad (2.108)$$

At *each point* of the boundary  $\partial\Omega$  (of  $\Omega$ ) we need *one* boundary condition involving the unknown  $u$ . The boundary conditions are of three principal types:

1.  $u$  is prescribed ( $u = 0$  or a known time variation of  $u$  at the boundary points, e.g., modeling an incoming wave),
2.  $\partial u / \partial n = \mathbf{n} \cdot \nabla u$  is prescribed (zero for reflecting boundaries),

3. an open boundary condition (also called radiation condition) is specified to let waves travel undisturbed out of the domain, see Exercise 2.11 for details.

All the listed wave equations with *second-order* derivatives in time need *two* initial conditions:

1.  $u = I$ ,
2.  $u_t = V$ .

## 2.11.2 Mesh

We introduce a mesh in time and in space. The mesh in time consists of time points

$$t_0 = 0 < t_1 < \dots < t_{N_t},$$

often with a constant spacing  $\Delta t = t_{n+1} - t_n$ ,  $n \in \mathcal{I}_t^-$ .

Finite difference methods are easy to implement on simple rectangle- or box-shaped domains. More complicated shapes of the domain require substantially more advanced techniques and implementational efforts. On a rectangle- or box-shaped domain, mesh points are introduced separately in the various space directions:

$x_0 < x_1 < \dots < x_{N_x}$  in the  $x$  direction,

$y_0 < y_1 < \dots < y_{N_y}$  in the  $y$  direction,

$z_0 < z_1 < \dots < z_{N_z}$  in the  $z$  direction .

We can write a general mesh point as  $(x_i, y_j, z_k, t_n)$ , with  $i \in \mathcal{I}_x$ ,  $j \in \mathcal{I}_y$ ,  $k \in \mathcal{I}_z$ , and  $n \in \mathcal{I}_t$ .

It is a very common choice to use constant mesh spacings:  $\Delta x = x_{i+1} - x_i$ ,  $i \in \mathcal{I}_x^-$ ,  $\Delta y = y_{j+1} - y_j$ ,  $j \in \mathcal{I}_y^-$ , and  $\Delta z = z_{k+1} - z_k$ ,  $k \in \mathcal{I}_z^-$ . With equal mesh spacings one often introduces  $h = \Delta x = \Delta y = \Delta z$ .

The unknown  $u$  at mesh point  $(x_i, y_j, z_k, t_n)$  is denoted by  $u_{i,j,k}^n$ . In 2D problems we just skip the  $z$  coordinate (by assuming no variation in that direction:  $\partial/\partial z = 0$ ) and write  $u_{i,j}^n$ .

### 2.11.3 Discretization

Two- and three-dimensional wave equations are easily discretized by assembling building blocks for discretization of 1D wave equations, because the multi-dimensional versions just contain terms of the same type as those in 1D.

**Discretizing the PDEs.** Equation (2.107) can be discretized as

$$[D_t D_t u = c^2(D_x D_x u + D_y D_y u + D_z D_z u) + f]_{i,j,k}^n. \quad (2.109)$$

A 2D version might be instructive to write out in detail:

$$[D_t D_t u = c^2(D_x D_x u + D_y D_y u) + f]_{i,j,k}^n,$$

which becomes

$$\frac{u_{i,j}^{n+1} - 2u_{i,j}^n + u_{i,j}^{n-1}}{\Delta t^2} = c^2 \frac{u_{i+1,j}^n - 2u_{i,j}^n + u_{i-1,j}^n}{\Delta x^2} + c^2 \frac{u_{i,j+1}^n - 2u_{i,j}^n + u_{i,j-1}^n}{\Delta y^2} + f_{i,j}^n,$$

Assuming, as usual, that all values at time levels  $n$  and  $n - 1$  are known, we can solve for the only unknown  $u_{i,j}^{n+1}$ . The result can be compactly written as

$$u_{i,j}^{n+1} = 2u_{i,j}^n + u_{i,j}^{n-1} + c^2 \Delta t^2 [D_x D_x u + D_y D_y u]_{i,j}^n. \quad (2.110)$$

As in the 1D case, we need to develop a special formula for  $u_{i,j}^1$  where we combine the general scheme for  $u_{i,j}^{n+1}$ , when  $n = 0$ , with the discretization of the initial condition:

$$[D_{2t} u = V]_{i,j}^0 \Rightarrow u_{i,j}^{-1} = u_{i,j}^1 - 2\Delta t V_{i,j}.$$

The result becomes, in compact form,

$$u_{i,j}^1 = u_{i,j}^0 - 2\Delta t V_{i,j} + \frac{1}{2} c^2 \Delta t^2 [D_x D_x u + D_y D_y u]_{i,j}^0. \quad (2.111)$$

The PDE (2.108) with variable coefficients is discretized term by term using the corresponding elements from the 1D case:

$$[\varrho D_t D_t u = (D_x \bar{q}^x D_x u + D_y \bar{q}^y D_y u + D_z \bar{q}^z D_z u) + f]_{i,j,k}^n. \quad (2.112)$$

When written out and solved for the unknown  $u_{i,j,k}^{n+1}$ , one gets the scheme

$$\begin{aligned}
 u_{i,j,k}^{n+1} = & -u_{i,j,k}^{n-1} + 2u_{i,j,k}^n + \\
 & \frac{1}{\varrho_{i,j,k}} \frac{1}{\Delta x^2} \left( \frac{1}{2}(q_{i,j,k} + q_{i+1,j,k})(u_{i+1,j,k}^n - u_{i,j,k}^n) - \right. \\
 & \quad \left. \frac{1}{2}(q_{i-1,j,k} + q_{i,j,k})(u_{i,j,k}^n - u_{i-1,j,k}^n) \right) + \\
 & \frac{1}{\varrho_{i,j,k}} \frac{1}{\Delta x^2} \left( \frac{1}{2}(q_{i,j,k} + q_{i,j+1,k})(u_{i,j+1,k}^n - u_{i,j,k}^n) - \right. \\
 & \quad \left. \frac{1}{2}(q_{i,j-1,k} + q_{i,j,k})(u_{i,j,k}^n - u_{i,j-1,k}^n) \right) + \\
 & \frac{1}{\varrho_{i,j,k}} \frac{1}{\Delta x^2} \left( \frac{1}{2}(q_{i,j,k} + q_{i,j,k+1})(u_{i,j,k+1}^n - u_{i,j,k}^n) - \right. \\
 & \quad \left. \frac{1}{2}(q_{i,j,k-1} + q_{i,j,k})(u_{i,j,k}^n - u_{i,j,k-1}^n) \right) + \\
 & \Delta t^2 f_{i,j,k}^n.
 \end{aligned}$$

Also here we need to develop a special formula for  $u_{i,j,k}^1$  by combining the scheme for  $n = 0$  with the discrete initial condition, which is just a matter of inserting  $u_{i,j,k}^{-1} = u_{i,j,k}^1 - 2\Delta t V_{i,j,k}$  in the scheme and solving for  $u_{i,j,k}^1$ .

**Handling boundary conditions where  $u$  is known.** The schemes listed above are valid for the internal points in the mesh. After updating these, we need to visit all the mesh points at the boundaries and set the prescribed  $u$  value.

**Discretizing the Neumann condition.** The condition  $\partial u / \partial n = 0$  was implemented in 1D by discretizing it with a  $D_{2x}u$  centered difference, followed by eliminating the fictitious  $u$  point outside the mesh by using the general scheme at the boundary point. Alternatively, one can introduce ghost cells and update a ghost value for use in the Neumann condition. Exactly the same ideas are reused in multiple dimensions.

Consider the condition  $\partial u / \partial n = 0$  at a boundary  $y = 0$  of a rectangular domain  $[0, L_x] \times [0, L_y]$  in 2D. The normal direction is then in  $-y$  direction, so

$$\frac{\partial u}{\partial n} = -\frac{\partial u}{\partial y},$$

and we set

$$[-D_{2y}u = 0]_{i,0}^n \Rightarrow \frac{u_{i,1}^n - u_{i,-1}^n}{2\Delta y} = 0.$$

From this it follows that  $u_{i,-1}^n = u_{i,1}^n$ . The discretized PDE at the boundary point  $(i, 0)$  reads

$$\frac{u_{i,0}^{n+1} - 2u_{i,0}^n + u_{i,0}^{n-1}}{\Delta t^2} = c^2 \frac{u_{i+1,0}^n - 2u_{i,0}^n + u_{i-1,0}^n}{\Delta x^2} + c^2 \frac{u_{i,1}^n - 2u_{i,0}^n + u_{i,-1}^n}{\Delta y^2} + f_{i,0}^n,$$

We can then just insert  $u_{i,1}^n$  for  $u_{i,-1}^n$  in this equation and solve for the boundary value  $u_{i,0}^{n+1}$ , just as was done in 1D.

From these calculations, we see a pattern: the general scheme applies at the boundary  $j = 0$  too if we just replace  $j - 1$  by  $j + 1$ . Such a pattern is particularly useful for implementations. The details follow from the explained 1D case in Section 2.6.3.

The alternative approach to eliminating fictitious values outside the mesh is to have  $u_{i,-1}^n$  available as a ghost value. The mesh is extended with one extra line (2D) or plane (3D) of ghost cells at a Neumann boundary. In the present example it means that we need a line with ghost cells below the  $y$  axis. The ghost values must be updated according to  $u_{i,-1}^{n+1} = u_{i,1}^{n+1}$ .

## 2.12 Implementation

We shall now describe in detail various Python implementations for solving a standard 2D, linear wave equation with constant wave velocity and  $u = 0$  on the boundary. The wave equation is to be solved in the space-time domain  $\Omega \times (0, T]$ , where  $\Omega = (0, L_x) \times (0, L_y)$  is a rectangular spatial domain. More precisely, the complete initial-boundary value problem is defined by

$$u_{tt} = c^2(u_{xx} + u_{yy}) + f(x, y, t), \quad (x, y) \in \Omega, \quad t \in (0, T], \quad (2.113)$$

$$u(x, y, 0) = I(x, y), \quad (x, y) \in \Omega, \quad (2.114)$$

$$u_t(x, y, 0) = V(x, y), \quad (x, y) \in \Omega, \quad (2.115)$$

$$u = 0, \quad (x, y) \in \partial\Omega, \quad t \in (0, T], \quad (2.116)$$

where  $\partial\Omega$  is the boundary of  $\Omega$ , in this case the four sides of the rectangle  $\Omega = [0, L_x] \times [0, L_y]$ :  $x = 0$ ,  $x = L_x$ ,  $y = 0$ , and  $y = L_y$ .

The PDE is discretized as

$$[D_t D_t u = c^2(D_x D_x u + D_y D_y u) + f]_{i,j}^n,$$

which leads to an explicit updating formula to be implemented in a program:

$$\begin{aligned} u_{i,j}^{n+1} = & -u_{i,j}^{n-1} + 2u_{i,j}^n + \\ & C_x^2(u_{i+1,j}^n - 2u_{i,j}^n + u_{i-1,j}^n) + C_y^2(u_{i,j+1}^n - 2u_{i,j}^n + u_{i,j-1}^n) + \Delta t^2 f_{i,j}^n, \end{aligned} \quad (2.117)$$

for all interior mesh points  $i \in \mathcal{I}_x^i$  and  $j \in \mathcal{I}_y^i$ , for  $n \in \mathcal{I}_t^+$ . The constants  $C_x$  and  $C_y$  are defined as

$$C_x = c \frac{\Delta t}{\Delta x}, \quad C_y = c \frac{\Delta t}{\Delta y}.$$

At the boundary, we simply set  $u_{i,j}^{n+1} = 0$  for  $i = 0, j = 0, \dots, N_y$ ;  $i = N_x, j = 0, \dots, N_y$ ;  $j = 0, i = 0, \dots, N_x$ ; and  $j = N_y, i = 0, \dots, N_x$ . For the first step,  $n = 0$ , (2.117) is combined with the discretization of the initial condition  $u_t = V$ ,  $[D_{2t}u = V]_{i,j}^0$  to obtain a special formula for  $u_{i,j}^1$  at the interior mesh points:

$$\begin{aligned} u_{i,j}^1 = & u_{i,j}^0 + \Delta t V_{i,j} + \\ & \frac{1}{2} C_x^2 (u_{i+1,j}^0 - 2u_{i,j}^0 + u_{i-1,j}^0) + \frac{1}{2} C_y^2 (u_{i,j+1}^0 - 2u_{i,j}^0 + u_{i,j-1}^0) + \\ & \frac{1}{2} \Delta t^2 f_{i,j}^n, \end{aligned} \quad (2.118)$$

The algorithm is very similar to the one in 1D:

1. Set initial condition  $u_{i,j}^0 = I(x_i, y_j)$
2. Compute  $u_{i,j}^1$  from (2.117)
3. Set  $u_{i,j}^1 = 0$  for the boundaries  $i = 0, N_x, j = 0, N_y$
4. For  $n = 1, 2, \dots, N_t$ :
  - a. Find  $u_{i,j}^{n+1}$  from (2.117) for all internal mesh points,  $i \in \mathcal{I}_x^i, j \in \mathcal{I}_y^i$
  - b. Set  $u_{i,j}^{n+1} = 0$  for the boundaries  $i = 0, N_x, j = 0, N_y$

### 2.12.1 Scalar computations

The `solver` function for a 2D case with constant wave velocity and boundary condition  $u = 0$  is analogous to the 1D case with similar parameter values (see `wave1D_u0.py`), apart from a few necessary extensions. The code is found in the program `wave2D_u0.py`.

**Domain and mesh.** The spatial domain is now  $[0, L_x] \times [0, L_y]$ , specified by the arguments `Lx` and `Ly`. Similarly, the number of mesh points in the  $x$  and  $y$  directions,  $N_x$  and  $N_y$ , become the arguments `Nx` and `. In multi-dimensional problems it makes less sense to specify a Courant number since the wave velocity is a vector and mesh spacings may differ in the various spatial directions. We therefore give  $\Delta t$  explicitly. The signature of the solver function is then`

```
def solver(I, V, f, c, Lx, Ly, Nx, Ny, dt, T,
          user_action=None, version='scalar'):
```

Key parameters used in the calculations are created as

```
x = linspace(0, Lx, Nx+1)                      # mesh points in x dir
y = linspace(0, Ly, Ny+1)                      # mesh points in y dir
dx = x[1] - x[0]
dy = y[1] - y[0]
Nt = int(round(T/float(dt)))
t = linspace(0, Nt*dt, Nt+1)                  # mesh points in time
Cx2 = (c*dt/dx)**2; Cy2 = (c*dt/dy)**2      # help variables
dt2 = dt**2
```

**Solution arrays.** We store  $u_{i,j}^{n+1}$ ,  $u_{i,j}^n$ , and  $u_{i,j}^{n-1}$  in three two-dimensional arrays,

```
u   = zeros((Nx+1,Ny+1))    # solution array
u_1 = zeros((Nx+1,Ny+1))    # solution at t-dt
u_2 = zeros((Nx+1,Ny+1))    # solution at t-2*dt
```

where  $u_{i,j}^{n+1}$  corresponds to `u[i,j]`,  $u_{i,j}^n$  to `u_1[i,j]`, and  $u_{i,j}^{n-1}$  to `u_2[i,j]`

**Index sets.** It is also convenient to introduce the index sets (cf. Section 2.6.4)

```
Ix = range(0, u.shape[0])
Iy = range(0, u.shape[1])
It = range(0, t.shape[0])
```

**Computing the solution.** Inserting the initial condition `I` in `u_1` and making a callback to the user in terms of the `user_action` function is a straightforward generalization of the 1D code from Section 2.1.6:

```

for i in Ix:
    for j in Iy:
        u_1[i,j] = I(x[i], y[j])

if user_action is not None:
    user_action(u_1, x, xv, y, yv, t, 0)

```

The `user_action` function has additional arguments compared to the 1D case. The arguments `xv` and `yv` will be commented upon in Section 2.12.2.

The key finite difference formula (2.110) for updating the solution at a time level is implemented in a separate function as

```

def advance_scalar(u, u_1, u_2, f, x, y, t, n, Cx2, Cy2, dt2,
                   V=None, step1=False):
    Ix = range(0, u.shape[0]); Iy = range(0, u.shape[1])
    if step1:
        dt = sqrt(dt2) # save
        Cx2 = 0.5*Cx2; Cy2 = 0.5*Cy2; dt2 = 0.5*dt2 # redefine
        D1 = 1; D2 = 0
    else:
        D1 = 2; D2 = 1
    for i in Ix[1:-1]:
        for j in Iy[1:-1]:
            u_xx = u_1[i-1,j] - 2*u_1[i,j] + u_1[i+1,j]
            u_yy = u_1[i,j-1] - 2*u_1[i,j] + u_1[i,j+1]
            u[i,j] = D1*u_1[i,j] - D2*u_2[i,j] + \
                      Cx2*u_xx + Cy2*u_yy + dt2*f(x[i], y[j], t[n])
            if step1:
                u[i,j] += dt*V(x[i], y[j])
    # Boundary condition u=0
    j = Iy[0]
    for i in Ix: u[i,j] = 0
    j = Iy[-1]
    for i in Ix: u[i,j] = 0
    i = Ix[0]
    for j in Iy: u[i,j] = 0
    i = Ix[-1]
    for j in Iy: u[i,j] = 0
    return u

```

The `step1` variable has been introduced to allow the formula to be reused for first step  $u_{i,j}^1$ :

```

u = advance_scalar(u, u_1, u_2, f, x, y, t,
                   n, Cx2, Cy2, dt, V, step1=True)

```

Below, we will make many alternative implementations of the `advance_scalar` function to speed up the code since most of the CPU time in simulations is spent in this function.

Finally, we remark that the `solver` function in the `wave2D_u0.py` code updates arrays for the next time step by switching references as described

in Section 2.4.5. If the solution  $\mathbf{u}$  is returned from `solver`, which it is not, it is important to set  $\mathbf{u} = \mathbf{u}_1$  after the time loop, otherwise  $\mathbf{u}$  actually contains  $\mathbf{u}_2$ .

## 2.12.2 Vectorized computations

The scalar code above turns out to be extremely slow for large 2D meshes, and probably useless in 3D beyond debugging of small test cases. Vectorization is therefore a must for multi-dimensional finite difference computations in Python. For example, with a mesh consisting of  $30 \times 30$  cells, vectorization brings down the CPU time by a factor of 70 (!). Equally important, vectorized code can also easily be parallelized to take (usually) optimal advantage of parallel computer platforms.

In the vectorized case, we must be able to evaluate user-given functions like  $I(x, y)$  and  $f(x, y, t)$  for the entire mesh in one operation (without loops). These user-given functions are provided as Python functions `I(x, y)` and `f(x, y, t)`, respectively. Having the one-dimensional coordinate arrays `x` and `y` is not sufficient when calling `I` and `f` in a vectorized way. We must extend `x` and `y` to their vectorized versions `xv` and `yv`:

```
from numpy import newaxis
xv = x[:,newaxis]
yv = y[newaxis,:]
# or
xv = x.reshape((x.size, 1))
yv = y.reshape((1, y.size))
```

This is a standard required technique when evaluating functions over a 2D mesh, say `sin(xv)*cos(xv)`, which then gives a result with shape  $(Nx+1, Ny+1)$ . Calling `I(xv, yv)` and `f(xv, yv, t[n])` will now return `I` and `f` values for the entire set of mesh points.

With the `xv` and `yv` arrays for vectorized computing, setting the initial condition is just a matter of

```
u_1[:, :] = I(xv, yv)
```

One could also have written `u_1 = I(xv, yv)` and let `u_1` point to a new object, but vectorized operations often make use of direct insertion in the original array through `u_1[:, :]`, because sometimes not all of the array is to be filled by such a function evaluation. This is the case with the computational scheme for  $u_{i,j}^{n+1}$ :

```
def advance_vectorized(u, u_1, u_2, f_a, Cx2, Cy2, dt2,
```

```

V=None, step1=False):
    if step1:
        dt = sqrt(dt2) # save
        Cx2 = 0.5*Cx2; Cy2 = 0.5*Cy2; dt2 = 0.5*dt2 # redefine
        D1 = 1; D2 = 0
    else:
        D1 = 2; D2 = 1
    u_xx = u_1[:-2,1:-1] - 2*u_1[1:-1,1:-1] + u_1[2:,1:-1]
    u_yy = u_1[1:-1,:-2] - 2*u_1[1:-1,1:-1] + u_1[1:-1,2:]
    u[1:-1,1:-1] = D1*u_1[1:-1,1:-1] - D2*u_2[1:-1,1:-1] + \
                    Cx2*u_xx + Cy2*u_yy + dt2*f_a[1:-1,1:-1]
    if step1:
        u[1:-1,1:-1] += dt*V[1:-1, 1:-1]
    # Boundary condition u=0
    j = 0
    u[:,j] = 0
    j = u.shape[1]-1
    u[:,j] = 0
    i = 0
    u[i,:] = 0
    i = u.shape[0]-1
    u[i,:] = 0
    return u

```

Array slices in 2D are more complicated to understand than those in 1D, but the logic from 1D applies to each dimension separately. For example, when doing  $u_{i,j}^n - u_{i-1,j}^n$  for  $i \in \mathcal{I}_x^+$ , we just keep  $j$  constant and make a slice in the first index:  $u_1[1:,:,j] - u_1[:-1,:,j]$ , exactly as in 1D. The  $1:$  slice specifies all the indices  $i = 1, 2, \dots, N_x$  (up to the last valid index), while  $:-1$  specifies the relevant indices for the second term:  $0, 1, \dots, N_x - 1$  (up to, but not including the last index).

In the above code segment, the situation is slightly more complicated, because each displaced slice in one direction is accompanied by a  $1:-1$  slice in the other direction. The reason is that we only work with the internal points for the index that is kept constant in a difference.

The boundary conditions along the four sides makes use of a slice consisting of all indices along a boundary:

```

u[:,0] = 0
u[:,Ny] = 0
u[0,:] = 0
u[Nx,:] = 0

```

In the vectorized update of  $\mathbf{u}$  (above), the function  $\mathbf{f}$  is first computed as an array over all mesh points:

```
f_a = f(xv, yv, t[n])
```

We could, alternatively, have used the call  $\mathbf{f}(\mathbf{xv}, \mathbf{yv}, \mathbf{t}[n])[1:-1,1:-1]$  in the last term of the update statement, but other implementations in

compiled languages benefit from having `f` available in an array rather than calling our Python function `f(x,y,t)` for every point.

Also in the `advance_vectorized` function we have introduced a boolean `step1` to reuse the formula for the first time step in the same way as we did with `advance_scalar`. We refer to the `solver` function in `wave2D_u0.py` for the details on how the overall algorithm is implemented.

The callback function now has the arguments `u`, `x`, `xv`, `y`, `yv`, `t`, `n`. The inclusion of `xv` and `yv` makes it easy to, e.g., compute an exact 2D solution in the callback function and compute errors, through an expression like `u - u_exact(xv, yv, t[n])`.

### 2.12.3 Verification

**Testing a quadratic solution.** The 1D solution from Section 2.2.4 can be generalized to multi-dimensions and provides a test case where the exact solution also fulfills the discrete equations, such that we know (to machine precision) what numbers the solver function should produce. In 2D we use the following generalization of (2.30):

$$u_e(x, y, t) = x(L_x - x)y(L_y - y)(1 + \frac{1}{2}t). \quad (2.119)$$

This solution fulfills the PDE problem if  $I(x, y) = u_e(x, y, 0)$ ,  $V = \frac{1}{2}u_e(x, y, 0)$ , and  $f = 2c^2(1 + \frac{1}{2}t)(y(L_y - y) + x(L_x - x))$ . To show that  $u_e$  also solves the discrete equations, we start with the general results  $[D_t D_t 1]^n = 0$ ,  $[D_t D_t t]^n = 0$ , and  $[D_t D_t t^2] = 2$ , and use these to compute

$$\begin{aligned} [D_x D_x u_e]_{i,j}^n &= [y(L_y - y)(1 + \frac{1}{2}t) D_x D_x x (L_x - x)]_{i,j}^n \\ &= y_j (L_y - y_j) (1 + \frac{1}{2}t_n) (-2). \end{aligned}$$

A similar calculation must be carried out for the  $[D_y D_y u_e]_{i,j}^n$  and  $[D_t D_t u_e]_{i,j}^n$  terms. One must also show that the quadratic solution fits the special formula for  $u_{i,j}^1$ . The details are left as Exercise 2.15. The `test_quadratic` function in the `wave2D_u0.py` program implements this verification as a proper test function for the pytest and nose frameworks.

## 2.13 Exercises

### Exercise 2.15: Check that a solution fulfills the discrete model

Carry out all mathematical details to show that (2.119) is indeed a solution of the discrete model for a 2D wave equation with  $u = 0$  on the boundary. One must check the boundary conditions, the initial conditions, the general discrete equation at a time level and the special version of this equation for the first time level. Filename: `check_quadratic_solution`.

### Project 2.16: Calculus with 2D mesh functions

The goal of this project is to redo Project 2.5 with 2D mesh functions ( $f_{i,j}$ ).

**Differentiation.** The differentiation results in a discrete gradient function, which in the 2D case can be represented by a three-dimensional array `df[d,i,j]` where `d` represents the direction of the derivative, and `i,j` is a mesh point in 2D. Use centered differences for the derivative at inner points and one-sided forward or backward differences at the boundary points. Construct unit tests and write a corresponding test function.

**Integration.** The integral of a 2D mesh function  $f_{i,j}$  is defined as

$$F_{i,j} = \int_{y_0}^{y_j} \int_{x_0}^{x_i} f(x, y) dx dy,$$

where  $f(x, y)$  is a function that takes on the values of the discrete mesh function  $f_{i,j}$  at the mesh points, but can also be evaluated in between the mesh points. The particular variation between mesh points can be taken as bilinear, but this is not important as we will use a product Trapezoidal rule to approximate the integral over a cell in the mesh and then we only need to evaluate  $f(x, y)$  at the mesh points.

Suppose  $F_{i,j}$  is computed. The calculation of  $F_{i+1,j}$  is then

$$\begin{aligned} F_{i+1,j} &= F_{i,j} + \int_{x_i}^{x_{i+1}} \int_{y_0}^{y_j} f(x, y) dy dx \\ &\approx \Delta x \frac{1}{2} \left( \int_{y_0}^{y_j} f(x_i, y) dy + \int_{y_0}^{y_j} f(x_{i+1}, y) dy \right) \end{aligned}$$

The integrals in the  $y$  direction can be approximated by a Trapezoidal rule. A similar idea can be used to compute  $F_{i,j+1}$ . Thereafter,  $F_{i+1,j+1}$  can be computed by adding the integral over the final corner cell to  $F_{i+1,j} + F_{i,j+1} - F_{i,j}$ . Carry out the details of these computations and implement a function that can return  $F_{i,j}$  for all mesh indices  $i$  and  $j$ . Use the fact that the Trapezoidal rule is exact for linear functions and write a test function. Filename: `mesh_calculus_2D`.

### Exercise 2.17: Implement Neumann conditions in 2D

Modify the `wave2D_u0.py` program, which solves the 2D wave equation  $u_{tt} = c^2(u_{xx} + u_{yy})$  with constant wave velocity  $c$  and  $u = 0$  on the boundary, to have Neumann boundary conditions:  $\partial u / \partial n = 0$ . Include both scalar code (for debugging and reference) and vectorized code (for speed).

To test the code, use  $u = 1.2$  as solution ( $I(x, y) = 1.2$ ,  $V = f = 0$ , and  $c$  arbitrary), which should be exactly reproduced with any mesh as long as the stability criterion is satisfied. Another test is to use the plug-shaped pulse in the `pulse` function from Section 2.8 and the `wave1D_dn_vc.py` program. This pulse is exactly propagated in 1D if  $c\Delta t / \Delta x = 1$ . Check that also the 2D program can propagate this pulse exactly in  $x$  direction ( $c\Delta t / \Delta x = 1$ ,  $\Delta y$  arbitrary) and  $y$  direction ( $c\Delta t / \Delta y = 1$ ,  $\Delta x$  arbitrary). Filename: `wave2D_dn`.

### Exercise 2.18: Test the efficiency of compiled loops in 3D

Extend the `wave2D_u0.py` code and the Cython, Fortran, and C versions to 3D. Set up an efficiency experiment to determine the relative efficiency of pure scalar Python code, vectorized code, Cython-compiled loops, Fortran-compiled loops, and C-compiled loops. Normalize the CPU time for each mesh by the fastest version. Filename: `wave3D_u0`.

## 2.14 Applications of wave equations

This section presents a range of wave equation models for different physical phenomena. Although many wave motion problems in physics can be modeled by the standard linear wave equation, or a similar formulation

with a system of first-order equations, there are some exceptions. Perhaps the most important is water waves: these are modeled by the Laplace equation with time-dependent boundary conditions at the water surface (long water waves, however, can be approximated by a standard wave equation, see Section 2.14.7). Quantum mechanical waves constitute another example where the waves are governed by the Schrödinger equation, i.e., not by a standard wave equation. Many wave phenomena also need to take nonlinear effects into account when the wave amplitude is significant. Shock waves in the air is a primary example.

The derivations in the following are very brief. Those with a firm background in continuum mechanics will probably have enough knowledge to fill in the details, while other readers will hopefully get some impression of the physics and approximations involved when establishing wave equation models.

### 2.14.1 Waves on a string

Figure 2.8 shows a model we may use to derive the equation for waves on a string. The string is modeled as a set of discrete point masses (at mesh points) with elastic strings in between. The string has a large constant tension  $T$ . We let the mass at mesh point  $x_i$  be  $m_i$ . The displacement of this mass point in the  $y$  direction is denoted by  $u_i(t)$ .

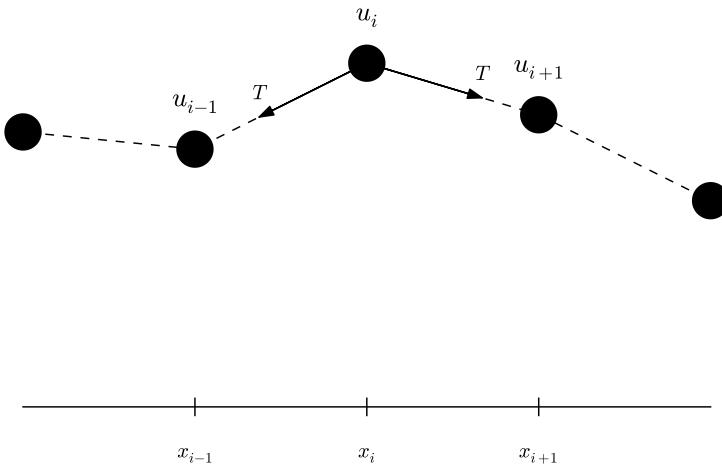
The motion of mass  $m_i$  is governed by Newton's second law of motion. The position of the mass at time  $t$  is  $x_i\mathbf{i} + u_i(t)\mathbf{j}$ , where  $\mathbf{i}$  and  $\mathbf{j}$  are unit vectors in the  $x$  and  $y$  direction, respectively. The acceleration is then  $u_i''(t)\mathbf{j}$ . Two forces are acting on the mass as indicated in Figure 2.8. The force  $\mathbf{T}^-$  acting toward the point  $x_{i-1}$  can be decomposed as

$$\mathbf{T}^- = -T \sin \phi \mathbf{i} - T \cos \phi \mathbf{j},$$

where  $\phi$  is the angle between the force and the line  $x = x_i$ . Let  $\Delta u_i = u_i - u_{i-1}$  and let  $\Delta s_i = \sqrt{\Delta u_i^2 + (x_i - x_{i-1})^2}$  be the distance from mass  $m_{i-1}$  to mass  $m_i$ . It is seen that  $\cos \phi = \Delta u_i / \Delta s_i$  and  $\sin \phi = (x_i - x_{i-1}) / \Delta s$  or  $\Delta x / \Delta s_i$  if we introduce a constant mesh spacing  $\Delta x = x_i - x_{i-1}$ . The force can then be written

$$\mathbf{T}^- = -T \frac{\Delta x}{\Delta s_i} \mathbf{i} - T \frac{\Delta u_i}{\Delta s_i} \mathbf{j}.$$

The force  $\mathbf{T}^+$  acting toward  $x_{i+1}$  can be calculated in a similar way:



**Fig. 2.8** Discrete string model with point masses connected by elastic strings.

$$\mathbf{T}^+ = T \frac{\Delta x}{\Delta s_{i+1}} \mathbf{i} + T \frac{\Delta u_{i+1}}{\Delta s_{i+1}} \mathbf{j}.$$

Newton's second law becomes

$$m_i u''_i(t) \mathbf{j} = \mathbf{T}^+ + \mathbf{T}^-,$$

which gives the component equations

$$T \frac{\Delta x}{\Delta s_i} = T \frac{\Delta x}{\Delta s_{i+1}}, \quad (2.120)$$

$$m_i u''_i(t) = T \frac{\Delta u_{i+1}}{\Delta s_{i+1}} - T \frac{\Delta u_i}{\Delta s_i}. \quad (2.121)$$

A basic reasonable assumption for a string is small displacements  $u_i$  and small displacement gradients  $\Delta u_i / \Delta x$ . For small  $g = \Delta u_i / \Delta x$  we have that

$$\Delta s_i = \sqrt{\Delta u_i^2 + \Delta x^2} = \Delta x \sqrt{1 + g^2} + \Delta x \left(1 + \frac{1}{2}g^2 + \mathcal{O}(g^4)\right) \approx \Delta x.$$

Equation (2.120) is then simply the identity  $T = T$ , while (2.121) can be written as

$$m_i u''_i(t) = T \frac{\Delta u_{i+1}}{\Delta x} - T \frac{\Delta u_i}{\Delta x},$$

which upon division by  $\Delta x$  and introducing the density  $\varrho_i = m_i / \Delta x$  becomes

$$\varrho_i u''_i(t) = T \frac{1}{\Delta x^2} (u_{i+1} - 2u_i + u_{i-1}). \quad (2.122)$$

We can now choose to approximate  $u''_i$  by a finite difference in time and get the discretized wave equation,

$$\varrho_i \frac{1}{\Delta t^2} (u_i^{n+1} - 2u_i^n + u_i^{n-1}) = T \frac{1}{\Delta x^2} (u_{i+1} - 2u_i + u_{i-1}). \quad (2.123)$$

On the other hand, we may go to the continuum limit  $\Delta x \rightarrow 0$  and replace  $u_i(t)$  by  $u(x, t)$ ,  $\varrho_i$  by  $\varrho(x)$ , and recognize that the right-hand side of (2.122) approaches  $\partial^2 u / \partial x^2$  as  $\Delta x \rightarrow 0$ . We end up with the continuous model for waves on a string:

$$\varrho \frac{\partial^2 u}{\partial t^2} = T \frac{\partial^2 u}{\partial x^2}. \quad (2.124)$$

Note that the density  $\varrho$  may change along the string, while the tension  $T$  is a constant. With variable wave velocity  $c(x) = \sqrt{T/\varrho(x)}$  we can write the wave equation in the more standard form

$$\frac{\partial^2 u}{\partial t^2} = c^2(x) \frac{\partial^2 u}{\partial x^2}. \quad (2.125)$$

Because of the way  $\varrho$  enters the equations, the variable wave velocity does *not* appear inside the derivatives as in many other versions of the wave equation. However, most strings of interest have constant  $\varrho$ .

The end points of a string are fixed so that the displacement  $u$  is zero. The boundary conditions are therefore  $u = 0$ .

**Damping.** Air resistance and non-elastic effects in the string will contribute to reduce the amplitudes of the waves so that the motion dies out after some time. This damping effect can be modeled by a term  $b u_t$  on the left-hand side of the equation

$$\varrho \frac{\partial^2 u}{\partial t^2} + b \frac{\partial u}{\partial t} = T \frac{\partial^2 u}{\partial x^2}. \quad (2.126)$$

The parameter  $b \geq 0$  is small for most wave phenomena, but the damping effect may become significant in long time simulations.

**External forcing.** It is easy to include an external force acting on the string. Say we have a vertical force  $\tilde{f}_i j$  acting on mass  $m_i$ . This force affects the vertical component of Newton's law and gives rise to an extra term  $\tilde{f}(x, t)$  on the right-hand side of (2.124). In the model (2.125) we would add a term  $f(x, t) = \tilde{f}(x, y)/\varrho(x)$ .

**Modeling the tension via springs.** We assumed, in the derivation above, that the tension in the string,  $T$ , was constant. It is easy to check this assumption by modeling the string segments between the masses as standard springs, where the force (tension  $T$ ) is proportional to the elongation of the spring segment. Let  $k$  be the spring constant, and set  $T_i = k\Delta\ell$  for the tension in the spring segment between  $x_{i-1}$  and  $x_i$ , where  $\Delta\ell$  is the elongation of this segment from the tension-free state. A basic feature of a string is that it has high tension in the equilibrium position  $u = 0$ . Let the string segment have an elongation  $\Delta\ell_0$  in the equilibrium position. After deformation of the string, the elongation is  $\Delta\ell = \Delta\ell_0 + \Delta s_i$ :  $T_i = k(\Delta\ell_0 + \Delta s_i) \approx k(\Delta\ell_0 + \Delta x)$ . This shows that  $T_i$  is independent of  $i$ . Moreover, the extra approximate elongation  $\Delta x$  is very small compared to  $\Delta\ell_0$ , so we may well set  $T_i = T = k\Delta\ell_0$ . This means that the tension is completely dominated by the initial tension determined by the tuning of the string. The additional deformations of the spring during the vibrations do not introduce significant changes in the tension.

## 2.14.2 Waves on a membrane

**hpl 9:** Must write the memberane model. Easiest to use Navier.

### 2.14.3 Elastic waves in a rod

Consider an elastic rod subject to a hammer impact at the end. This experiment will give rise to an elastic deformation pulse that travels through the rod. A mathematical model for longitudinal waves along an elastic rod starts with the general equation for deformations and stresses in an elastic medium,

$$\varrho \mathbf{u}_{tt} = \nabla \cdot \boldsymbol{\sigma} + \varrho \mathbf{f}, \quad (2.127)$$

where  $\varrho$  is the density,  $\mathbf{u}$  the displacement field,  $\boldsymbol{\sigma}$  the stress tensor, and  $\mathbf{f}$  body forces. The latter has normally no impact on elastic waves.

For stationary deformation of an elastic rod, one has that  $\sigma_{xx} = Eu_x$ , with all other stress components being zero. The parameter  $E$  is known as Young's modulus. Moreover, we set  $\mathbf{u} = u(x, t)\mathbf{i}$  and neglect the radial contraction and expansion (where Poisson's ratio is the important parameter). Assuming that this simple stress and deformation field is a good approximation, (2.127) simplifies to

$$\varrho \frac{\partial^2 u}{\partial t^2} = \frac{\partial}{\partial x} \left( E \frac{\partial u}{\partial x} \right). \quad (2.128)$$

The associated boundary conditions are  $u$  or  $\sigma_{xx} = Eu_x$  known, typically  $u = 0$  for a fixed end and  $\sigma_{xx} = 0$  for a free end.

### 2.14.4 The acoustic model for seismic waves

Seismic waves are used to infer properties of subsurface geological structures. The physical model is a heterogeneous elastic medium where sound is propagated by small elastic vibrations. The general mathematical model for deformations in an elastic medium is based on Newton's second law,

$$\varrho \mathbf{u}_{tt} = \nabla \cdot \boldsymbol{\sigma} + \varrho \mathbf{f}, \quad (2.129)$$

and a constitutive law relating  $\boldsymbol{\sigma}$  to  $\mathbf{u}$ , often Hooke's generalized law,

$$\boldsymbol{\sigma} = K \nabla \cdot \mathbf{u} \mathbf{I} + G (\nabla \mathbf{u} + (\nabla \mathbf{u})^T - \frac{2}{3} \nabla \cdot \mathbf{u} \mathbf{I}). \quad (2.130)$$

Here,  $\mathbf{u}$  is the displacement field,  $\boldsymbol{\sigma}$  is the stress tensor,  $\mathbf{I}$  is the identity tensor,  $\varrho$  is the medium's density,  $\mathbf{f}$  are body forces (such as gravity),  $K$  is the medium's bulk modulus and  $G$  is the shear modulus. All these quantities may vary in space, while  $\mathbf{u}$  and  $\boldsymbol{\sigma}$  will also show significant variation in time during wave motion.

The acoustic approximation to elastic waves arises from a basic assumption that the second term in Hooke's law, representing the deformations that give rise to shear stresses, can be neglected. This assumption can be interpreted as approximating the geological medium by a fluid. Neglecting also the body forces  $\mathbf{f}$ , (2.129) becomes

$$\varrho \mathbf{u}_{tt} = \nabla(K \nabla \cdot \mathbf{u}) \quad (2.131)$$

Introducing  $p$  as a pressure via

$$p = -K \nabla \cdot \mathbf{u}, \quad (2.132)$$

and dividing (2.131) by  $\varrho$ , we get

$$\mathbf{u}_{tt} = -\frac{1}{\varrho} \nabla p. \quad (2.133)$$

Taking the divergence of this equation, using  $\nabla \cdot \mathbf{u} = -p/K$  from (2.132), gives the *acoustic approximation to elastic waves*:

$$p_{tt} = K \nabla \cdot \left( \frac{1}{\varrho} \nabla p \right). \quad (2.134)$$

This is a standard, linear wave equation with variable coefficients. It is common to add a source term  $s(x, y, z, t)$  to model the generation of sound waves:

$$p_{tt} = K \nabla \cdot \left( \frac{1}{\varrho} \nabla p \right) + s. \quad (2.135)$$

A common additional approximation of (2.135) is based on using the chain rule on the right-hand side,

$$K \nabla \cdot \left( \frac{1}{\varrho} \nabla p \right) = \frac{K}{\varrho} \nabla^2 p + K \nabla \left( \frac{1}{\varrho} \right) \cdot \nabla p \approx \frac{K}{\varrho} \nabla^2 p,$$

under the assumption that the relative spatial gradient  $\nabla \varrho^{-1} = -\varrho^{-2} \nabla \varrho$  is small. This approximation results in the simplified equation

$$p_{tt} = \frac{K}{\varrho} \nabla^2 p + s. \quad (2.136)$$

The acoustic approximations to seismic waves are used for sound waves in the ground, and the Earth's surface is then a boundary where  $p$  equals the atmospheric pressure  $p_0$  such that the boundary condition becomes  $p = p_0$ .

**Anisotropy.** Quite often in geological materials, the effective wave velocity  $c = \sqrt{K/\varrho}$  is different in different spatial directions because geological layers are compacted, and often twisted, in such a way that the properties in the horizontal and vertical direction differ. With  $z$  as the vertical coordinate, we can introduce a vertical wave velocity  $c_z$  and a horizontal wave velocity  $c_h$ , and generalize (2.136) to

$$p_{tt} = c_z^2 p_{zz} + c_h^2 (p_{xx} + p_{yy}) + s. \quad (2.137)$$

### 2.14.5 Sound waves in liquids and gases

Sound waves arise from pressure and density variations in fluids. The starting point of modeling sound waves is the basic equations for a compressible fluid where we omit viscous (frictional) forces, body forces (gravity, for instance), and temperature effects:

$$\varrho_t + \nabla \cdot (\varrho \mathbf{u}) = 0, \quad (2.138)$$

$$\varrho \mathbf{u}_t + \varrho \mathbf{u} \cdot \nabla \mathbf{u} = -\nabla p, \quad (2.139)$$

$$\varrho = \varrho(p). \quad (2.140)$$

These equations are often referred to as the Euler equations for the motion of a fluid. The parameters involved are the density  $\varrho$ , the velocity  $\mathbf{u}$ , and the pressure  $p$ . Equation (2.139) reflects mass balance, (2.138) is Newton's second law for a fluid, with frictional and body forces omitted, and (2.140) is a constitutive law relating density to pressure by thermodynamic considerations. A typical model for (2.140) is the so-called [isentropic relation](#), valid for adiabatic processes where there is no heat transfer:

$$\varrho = \varrho_0 \left( \frac{p}{p_0} \right)^{1/\gamma}. \quad (2.141)$$

Here,  $p_0$  and  $\varrho_0$  are reference values for  $p$  and  $\varrho$  when the fluid is at rest, and  $\gamma$  is the ratio of specific heat at constant pressure and constant volume ( $\gamma = 5/3$  for air).

The key approximation in a mathematical model for sound waves is to assume that these waves are small perturbations to the density, pressure, and velocity. We therefore write

$$p = p_0 + \hat{p},$$

$$\varrho = \varrho_0 + \hat{\varrho},$$

$$\mathbf{u} = \hat{\mathbf{u}},$$

where we have decomposed the fields in a constant equilibrium value, corresponding to  $\mathbf{u} = 0$ , and a small perturbation marked with a hat symbol. By inserting these decompositions in (2.138) and (2.139), neglecting all product terms of small perturbations and/or their derivatives, and dropping the hat symbols, one gets the following linearized PDE system for the small perturbations in density, pressure, and velocity:

$$\varrho_t + \varrho_0 \nabla \cdot \mathbf{u} = 0, \quad (2.142)$$

$$\varrho_0 \mathbf{u}_t = -\nabla p. \quad (2.143)$$

Now we can eliminate  $\varrho_t$  by differentiating the relation  $\varrho(p)$ ,

$$\varrho_t = \varrho_0 \frac{1}{\gamma} \left( \frac{p}{p_0} \right)^{1/\gamma-1} \frac{1}{p_0} p_t = \frac{\varrho_0}{\gamma p_0} \left( \frac{p}{p_0} \right)^{1/\gamma-1} p_t.$$

The product term  $p^{1/\gamma-1} p_t$  can be linearized as  $p_0^{1/\gamma-1} p_t$ , resulting in

$$\varrho_t \approx \frac{\varrho_0}{\gamma p_0} p_t.$$

We then get

$$p_t + \gamma p_0 \nabla \cdot \mathbf{u} = 0, \quad (2.144)$$

$$\mathbf{u}_t = -\frac{1}{\varrho_0} \nabla p, . \quad (2.145)$$

Taking the divergence of (2.145) and differentiating (2.144) with respect to time gives the possibility to easily eliminate  $\nabla \cdot \mathbf{u}_t$  and arrive at a standard, linear wave equation for  $p$ :

$$p_{tt} = c^2 \nabla^2 p, \quad (2.146)$$

where  $c = \sqrt{\gamma p_0 / \rho_0}$  is the speed of sound in the fluid.

### 2.14.6 Spherical waves

Spherically symmetric three-dimensional waves propagate in the radial direction  $r$  only so that  $u = u(r, t)$ . The fully three-dimensional wave equation

$$\frac{\partial^2 u}{\partial t^2} = \nabla \cdot (c^2 \nabla u) + f$$

then reduces to the spherically symmetric wave equation

$$\frac{\partial^2 u}{\partial t^2} = \frac{1}{r^2} \frac{\partial}{\partial r} \left( c^2(r) r^2 \frac{\partial u}{\partial r} \right) + f(r, t), \quad r \in (0, R), \quad t > 0. \quad (2.147)$$

One can easily show that the function  $v(r, t) = ru(r, t)$  fulfills a standard wave equation in Cartesian coordinates if  $c$  is constant. To this end, insert  $u = v/r$  in

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( c^2(r) r^2 \frac{\partial u}{\partial r} \right)$$

to obtain

$$r \left( \frac{dc^2}{dr} \frac{\partial v}{\partial r} + c^2 \frac{\partial^2 v}{\partial r^2} \right) - \frac{dc^2}{dr} v.$$

The two terms in the parenthesis can be combined to

$$r \frac{\partial}{\partial r} \left( c^2 \frac{\partial v}{\partial r} \right),$$

which is recognized as the variable-coefficient Laplace operator in one Cartesian coordinate. The spherically symmetric wave equation in terms of  $v(r, t)$  now becomes

$$\frac{\partial^2 v}{\partial t^2} = \frac{\partial}{\partial r} \left( c^2(r) \frac{\partial v}{\partial r} \right) - \frac{1}{r} \frac{dc^2}{dr} v + rf(r, t), \quad r \in (0, R), \quad t > 0. \quad (2.148)$$

In the case of constant wave velocity  $c$ , this equation reduces to the wave equation in a single Cartesian coordinate called  $r$ :

$$\frac{\partial^2 v}{\partial t^2} = c^2 \frac{\partial^2 v}{\partial r^2} + rf(r, t), \quad r \in (0, R), \quad t > 0. \quad (2.149)$$

That is, any program for solving the one-dimensional wave equation in a Cartesian coordinate system can be used to solve (2.149), provided the source term is multiplied by the coordinate, and that we divide the Cartesian mesh solution by  $r$  to get the spherically symmetric solution. Moreover, if  $r = 0$  is included in the domain, spherical symmetry demands that  $\partial u / \partial r = 0$  at  $r = 0$ , which means that

$$\frac{\partial u}{\partial r} = \frac{1}{r^2} \left( r \frac{\partial v}{\partial r} - v \right) = 0, \quad r = 0,$$

implying  $v(0, t) = 0$  as a necessary condition. For practical applications, we exclude  $r = 0$  from the domain and assume that some boundary condition is assigned at  $r = \epsilon$ , for some  $\epsilon > 0$ .

### 2.14.7 The linear shallow water equations

The next example considers water waves whose wavelengths are much larger than the depth and whose wave amplitudes are small. This class of waves may be generated by catastrophic geophysical events, such as earthquakes at the sea bottom, landslides moving into water, or underwater slides (or a combination, as earthquakes frequently release avalanches of masses). For example, a subsea earthquake will normally have an extension of many kilometers but lift the water only a few meters. The wave length will have a size dictated by the earthquake area, which is much larger than the water depth, and compared to this wave length, an amplitude of a few meters is very small. The water is essentially a thin film, and mathematically we can average the problem in the vertical direction and approximate the 3D wave phenomenon by 2D PDEs. Instead of a moving water domain in three space dimensions, we get a horizontal 2D domain with an unknown function for the surface elevation and the water depth as a variable coefficient in the PDEs.

Let  $\eta(x, y, t)$  be the elevation of the water surface,  $H(x, y)$  the water depth corresponding to a flat surface ( $\eta = 0$ ),  $u(x, y, t)$  and  $v(x, y, t)$  the

depth-averaged horizontal velocities of the water. Mass and momentum balance of the water volume give rise to the PDEs involving these quantities:

$$\eta_t = -(Hu)_x - (Hv)_x \quad (2.150)$$

$$u_t = -g\eta_x, \quad (2.151)$$

$$v_t = -g\eta_y, \quad (2.152)$$

where  $g$  is the acceleration of gravity. Equation (2.150) corresponds to mass balance while the other two are derived from momentum balance (Newton's second law).

The initial conditions associated with (2.150)-(2.152) are  $\eta$ ,  $u$ , and  $v$  prescribed at  $t = 0$ . A common condition is to have some water elevation  $\eta = I(x, y)$  and assume that the surface is at rest:  $u = v = 0$ . A subsea earthquake usually means a sufficiently rapid motion of the bottom and the water volume to say that the bottom deformation is mirrored at the water surface as an initial lift  $I(x, y)$  and that  $u = v = 0$ .

Boundary conditions may be  $\eta$  prescribed for incoming, known waves, or zero normal velocity at reflecting boundaries (steep mountains, for instance):  $un_x + vn_y = 0$ , where  $(n_x, n_y)$  is the outward unit normal to the boundary. More sophisticated boundary conditions are needed when waves run up at the shore, and at open boundaries where we want the waves to leave the computational domain undisturbed.

Equations (2.150), (2.151), and (2.152) can be transformed to a standard, linear wave equation. First, multiply (2.151) and (2.152) by  $H$ , differentiate (2.151) with respect to  $x$  and (2.152) with respect to  $y$ . Second, differentiate (2.150) with respect to  $t$  and use that  $(Hu)_{xt} = (Hu_t)_x$  and  $(Hv)_{yt} = (Hv_t)_y$  when  $H$  is independent of  $t$ . Third, eliminate  $(Hu_t)_x$  and  $(Hv_t)_y$  with the aid of the other two differentiated equations. These manipulations result in a standard, linear wave equation for  $\eta$ :

$$\eta_{tt} = (gH\eta_x)_x + (gH\eta_y)_y = \nabla \cdot (gH\nabla\eta). \quad (2.153)$$

In the case we have an initial non-flat water surface at rest, the initial conditions become  $\eta = I(x, y)$  and  $\eta_t = 0$ . The latter follows from (2.150) if  $u = v = 0$ , or simply from the fact that the vertical velocity of the surface is  $\eta_t$ , which is zero for a surface at rest.

The system (2.150)-(2.152) can be extended to handle a time-varying bottom topography, which is relevant for modeling long waves generated

by underwater slides. In such cases the water depth function  $H$  is also a function of  $t$ , due to the moving slide, and one must add a time-derivative term  $H_t$  to the left-hand side of (2.150). A moving bottom is best described by introducing  $z = H_0$  as the still-water level,  $z = B(x, y, t)$  as the time- and space-varying bottom topography, so that  $H = H_0 - B(x, y, t)$ . In the elimination of  $u$  and  $v$  one may assume that the dependence of  $H$  on  $t$  can be neglected in the terms  $(Hu)_{xt}$  and  $(Hv)_{yt}$ . We then end up with a source term in (2.153), because of the moving (accelerating) bottom:

$$\eta_{tt} = \nabla \cdot (gH\nabla\eta) + B_{tt}. \quad (2.154)$$

The reduction of (2.154) to 1D, for long waves in a straight channel, or for approximately plane waves in the ocean, is trivial by assuming no change in  $y$  direction ( $\partial/\partial y = 0$ ):

$$\eta_{tt} = (gH\eta_x)_x + B_{tt}. \quad (2.155)$$

**Wind drag on the surface.** Surface waves are influenced by the drag of the wind, and if the wind velocity some meters above the surface is  $(U, V)$ , the wind drag gives contributions  $C_V\sqrt{U^2 + V^2}U$  and  $C_V\sqrt{U^2 + V^2}V$  to (2.151) and (2.152), respectively, on the right-hand sides.

**Bottom drag.** The waves will experience a drag from the bottom, often roughly modeled by a term similar to the wind drag:  $C_B\sqrt{u^2 + v^2}u$  on the right-hand side of (2.151) and  $C_B\sqrt{u^2 + v^2}v$  on the right-hand side of (2.152). Note that in this case the PDEs (2.151) and (2.152) become nonlinear and the elimination of  $u$  and  $v$  to arrive at a 2nd-order wave equation for  $\eta$  is not possible anymore.

**Effect of the Earth's rotation.** Long geophysical waves will often be affected by the rotation of the Earth because of the Coriolis force. This force gives rise to a term  $fv$  on the right-hand side of (2.151) and  $-fu$  on the right-hand side of (2.152). Also in this case one cannot eliminate  $u$  and  $v$  to work with a single equation for  $\eta$ . The Coriolis parameter is  $f = 2\Omega \sin \phi$ , where  $\Omega$  is the angular velocity of the earth and  $\phi$  is the latitude.

## 2.14.8 Waves in blood vessels

The flow of blood in our bodies is basically fluid flow in a network of pipes. Unlike rigid pipes, the walls in the blood vessels are elastic and

will increase their diameter when the pressure rises. The elastic forces will then push the wall back and accelerate the fluid. This interaction between the flow of blood and the deformation of the vessel wall results in waves traveling along our blood vessels.

A model for one-dimensional waves along blood vessels can be derived from averaging the fluid flow over the cross section of the blood vessels. Let  $x$  be a coordinate along the blood vessel and assume that all cross sections are circular, though with different radii  $R(x, t)$ . The main quantities to compute is the cross section area  $A(x, t)$ , the averaged pressure  $P(x, t)$ , and the total volume flux  $Q(x, t)$ . The area of this cross section is

$$A(x, t) = 2\pi \int_0^{R(x, t)} r dr, \quad (2.156)$$

Let  $v_x(x, t)$  be the velocity of blood averaged over the cross section at point  $x$ . The volume flux, being the total volume of blood passing a cross section per time unit, becomes

$$Q(x, t) = A(x, t)v_x(x, t) \quad (2.157)$$

Mass balance and Newton's second law lead to the PDEs

$$\frac{\partial A}{\partial t} + \frac{\partial Q}{\partial x} = 0, \quad (2.158)$$

$$\frac{\partial Q}{\partial t} + \frac{\gamma+2}{\gamma+1} \frac{\partial}{\partial x} \left( \frac{Q^2}{A} \right) + \frac{A}{\varrho} \frac{\partial P}{\partial x} = -2\pi(\gamma+2) \frac{\mu}{\varrho} \frac{Q}{A}, \quad (2.159)$$

where  $\gamma$  is a parameter related to the velocity profile,  $\varrho$  is the density of blood, and  $\mu$  is the dynamic viscosity of blood.

We have three unknowns  $A$ ,  $Q$ , and  $P$ , and two equations (2.158) and (2.159). A third equation is needed to relate the flow to the deformations of the wall. A common form for this equation is

$$\frac{\partial P}{\partial t} + \frac{1}{C} \frac{\partial Q}{\partial x} = 0, \quad (2.160)$$

where  $C$  is the compliance of the wall, given by the constitutive relation

$$C = \frac{\partial A}{\partial P} + \frac{\partial A}{\partial t}, \quad (2.161)$$

which requires a relationship between  $A$  and  $P$ . One common model is to view the vessel wall, locally, as a thin elastic tube subject to an internal pressure. This gives the relation

$$P = P_0 + \frac{\pi h E}{(1 - \nu^2) A_0} (\sqrt{A} - \sqrt{A_0}),$$

where  $P_0$  and  $A_0$  are corresponding reference values when the wall is not deformed,  $h$  is the thickness of the wall, and  $E$  and  $\nu$  are Young's modulus and Poisson's ratio of the elastic material in the wall. The derivative becomes

$$C = \frac{\partial A}{\partial P} = \frac{2(1 - \nu^2) A_0}{\pi h E} \sqrt{A_0} + 2 \left( \frac{(1 - \nu^2) A_0}{\pi h E} \right)^2 (P - P_0). \quad (2.162)$$

Another (nonlinear) deformation model of the wall, which has a better fit with experiments, is

$$P = P_0 \exp(\beta(A/A_0 - 1)),$$

where  $\beta$  is some parameter to be estimated. This law leads to

$$C = \frac{\partial A}{\partial P} = \frac{A_0}{\beta P}. \quad (2.163)$$

**Reduction to the standard wave equation.** It is not uncommon to neglect the viscous term on the right-hand side of (2.159) and also the quadratic term with  $Q^2$  on the left-hand side. The reduced equations (2.159) and (2.160) form a first-order linear wave equation system:

$$C \frac{\partial P}{\partial t} = -\frac{\partial Q}{\partial x}, \quad (2.164)$$

$$\frac{\partial Q}{\partial t} = -\frac{A}{\varrho} \frac{\partial P}{\partial x}. \quad (2.165)$$

These can be combined into standard 1D wave PDE by differentiating the first equation with respect to  $t$  and the second with respect to  $x$ ,

$$\frac{\partial}{\partial t} \left( C \frac{\partial P}{\partial t} \right) = \frac{\partial}{\partial x} \left( \frac{A}{\varrho} \frac{\partial P}{\partial x} \right),$$

which can be approximated by

$$\frac{\partial^2 Q}{\partial t^2} = c^2 \frac{\partial^2 Q}{\partial x^2}, \quad c = \sqrt{\frac{A}{\varrho C}}, \quad (2.166)$$

where the  $A$  and  $C$  in the expression for  $c$  are taken as constant reference values.

### 2.14.9 Electromagnetic waves

Light and radio waves are governed by standard wave equations arising from Maxwell's general equations. When there are no charges and no currents, as in a vacuum, Maxwell's equations take the form

$$\begin{aligned}\nabla \cdot \mathbf{E} &= 0, \\ \nabla \cdot \mathbf{B} &= 0, \\ \nabla \times \mathbf{E} &= -\frac{\partial \mathbf{B}}{\partial t}, \\ \nabla \times \mathbf{B} &= \mu_0 \epsilon_0 \frac{\partial \mathbf{E}}{\partial t},\end{aligned}$$

where  $\epsilon_0 = 8.854187817620 \cdot 10^{-12}$  (F/m) is the permittivity of free space, also known as the electric constant, and  $\mu_0 = 1.2566370614 \cdot 10^{-6}$  (H/m) is the permeability of free space, also known as the magnetic constant. Taking the curl of the two last equations and using the mathematical identity

$$\nabla \times (\nabla \times \mathbf{E}) = \nabla(\nabla \cdot \mathbf{E}) - \nabla^2 \mathbf{E} = -\nabla^2 \mathbf{E} \text{ when } \nabla \cdot \mathbf{E} = 0,$$

gives the wave equation governing the electric and magnetic field:

$$\frac{\partial^2 \mathbf{E}}{\partial t^2} = c^2 \nabla^2 \mathbf{E}, \quad (2.167)$$

$$\frac{\partial^2 \mathbf{B}}{\partial t^2} = c^2 \nabla^2 \mathbf{B}, \quad (2.168)$$

with  $c = 1/\sqrt{\mu_0 \epsilon_0}$  as the velocity of light. Each component of  $\mathbf{E}$  and  $\mathbf{B}$  fulfills a wave equation and can hence be solved independently.

## 2.15 Exercises

### Exercise 2.19: Simulate waves on a non-homogeneous string

Simulate waves on a string that consists of two materials with different density. The tension in the string is constant, but the density has a jump at the middle of the string. Experiment with different sizes of the jump and produce animations that visualize the effect of the jump on the wave motion.

**Hint.** According to Section 2.14.1, the density enters the mathematical model as  $\varrho$  in  $\varrho u_{tt} = Tu_{xx}$ , where  $T$  is the string tension. Modify, e.g., the `wave1D_u0v.py` code to incorporate the tension and two density values. Make a mesh function `rho` with density values at each spatial mesh point. A value for the tension may be 150 N. Corresponding density values can be computed from the wave velocity estimations in the `guitar` function in the `wave1D_u0v.py` file.

Filename: `wave1D_u0_sv_discont`.

### Exercise 2.20: Simulate damped waves on a string

Formulate a mathematical model for damped waves on a string. Use data from Section 2.3.5, and tune the damping parameter so that the string is very close to the rest state after 15 s. Make a movie of the wave motion. Filename: `wave1D_u0_sv_damping`.

### Exercise 2.21: Simulate elastic waves in a rod

A hammer hits the end of an elastic rod. The exercise is to simulate the resulting wave motion using the model (2.128) from Section 2.14.3. Let the rod have length  $L$  and let the boundary  $x = L$  be stress free so that  $\sigma_{xx} = 0$ , implying that  $\partial u / \partial x = 0$ . The left end  $x = 0$  is subject to a strong stress pulse (the hammer), modeled as

$$\sigma_{xx}(t) = \begin{cases} S, & 0 < t \leq t_s, \\ 0, & t > t_s \end{cases}$$

The corresponding condition on  $u$  becomes  $u_x = S/E$  for  $t \leq t_s$  and zero afterwards (recall that  $\sigma_{xx} = Eu_x$ ). This is a non-homogeneous Neumann condition, and you will need to approximate this condition and

combine it with the scheme (the ideas and manipulations follow closely the handling of a non-zero initial condition  $u_t = V$  in wave PDEs or the corresponding second-order ODEs for vibrations). Filename: `wave_rod`.

### Exercise 2.22: Simulate spherical waves

Implement a model for spherically symmetric waves using the method described in Section 2.14.6. The boundary condition at  $r = 0$  must be  $\partial u / \partial r = 0$ , while the condition at  $r = R$  can either be  $u = 0$  or a radiation condition as described in Problem 2.11. The  $u = 0$  condition is sufficient if  $R$  is so large that the amplitude of the spherical wave has become insignificant. Make movie(s) of the case where the source term is located around  $r = 0$  and sends out pulses

$$f(r, t) = \begin{cases} Q \exp\left(-\frac{r^2}{2\Delta r^2}\right) \sin \omega t, & \sin \omega t \geq 0 \\ 0, & \sin \omega t < 0 \end{cases}$$

Here,  $Q$  and  $\omega$  are constants to be chosen.

**Hint.** Use the program `wave1D_u0v.py` as a starting point. Let `solver` compute the  $v$  function and then set  $u = v/r$ . However,  $u = v/r$  for  $r = 0$  requires special treatment. One possibility is to compute  $u[1:] = v[1:]/r[1:]$  and then set  $u[0]=u[1]$ . The latter makes it evident that  $\partial u / \partial r = 0$  in a plot.

Filename: `wave1D_spherical`.

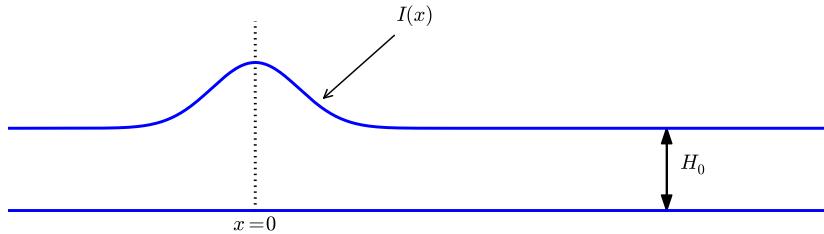
### Problem 2.23: Earthquake-generated tsunami over a subsea hill

A subsea earthquake leads to an immediate lift of the water surface, see Figure 2.9. The lifted water surface splits into two tsunamis, one traveling to the right and one to the left, as depicted in Figure 2.10. Since tsunamis are normally very long waves, compared to the depth, with a small amplitude, compared to the wave length, the wave equation model described in Section 2.14.7 is relevant:

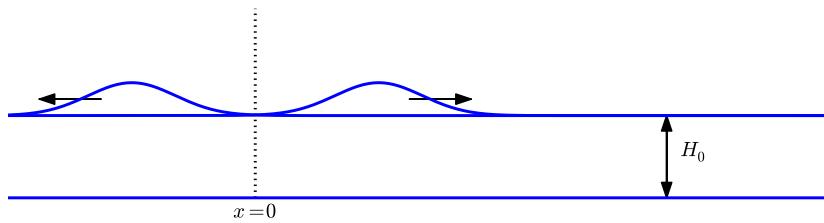
$$\eta_{tt} = (gH(x)\eta_x)_x,$$

where  $g$  is the acceleration of gravity, and  $H(x)$  is the still water depth.

To simulate the right-going tsunami, we can impose a symmetry boundary at  $x = 0$ :  $\partial \eta / \partial x = 0$ . We then simulate the wave motion



**Fig. 2.9** Sketch of initial water surface due to a subsea earthquake.



**Fig. 2.10** An initial surface elevation is split into two waves.

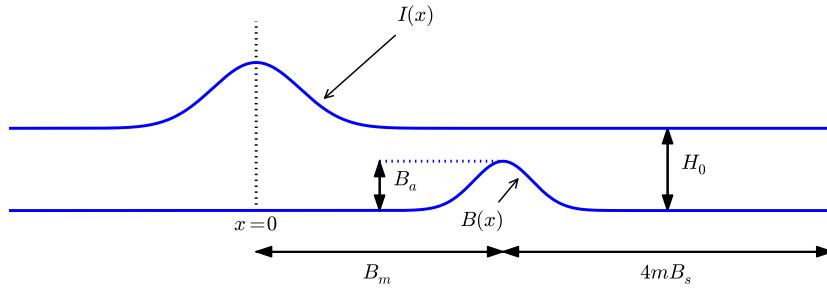
in  $[0, L]$ . Unless the ocean ends at  $x = L$ , the waves should travel undisturbed through the boundary  $x = L$ . A radiation condition as explained in Problem 2.11 can be used for this purpose. Alternatively, one can just stop the simulations before the wave hits the boundary at  $x = L$ . In that case it does not matter what kind of boundary condition we use at  $x = L$ . Imposing  $\eta = 0$  and stopping the simulations when  $|\eta_i^n| > \epsilon$ ,  $i = N_x - 1$ , is a possibility ( $\epsilon$  is a small parameter).

The shape of the initial surface can be taken as a Gaussian function,

$$I(x; I_0, I_a, I_m, I_s) = I_0 + I_a \exp\left(-\left(\frac{x - I_m}{I_s}\right)^2\right), \quad (2.169)$$

with  $I_m = 0$  reflecting the location of the peak of  $I(x)$  and  $I_s$  being a measure of the width of the function  $I(x)$  ( $I_s$  is  $\sqrt{2}$  times the standard deviation of the familiar normal distribution curve).

Now we extend the problem with a hill at the sea bottom, see Figure 2.11. The wave speed  $c = \sqrt{gH(x)} = \sqrt{g(H_0 - B(x))}$  will then be reduced in the shallow water above the hill.



**Fig. 2.11** Sketch of an earthquake-generated tsunami passing over a subsea hill.

One possible form of the hill is a Gaussian function,

$$B(x; B_0, B_a, B_m, B_s) = B_0 + B_a \exp\left(-\left(\frac{x - B_m}{B_s}\right)^2\right), \quad (2.170)$$

but many other shapes are also possible, e.g., a "cosine hat" where

$$B(x; B_0, B_a, B_m, B_s) = B_0 + B_a \cos\left(\pi \frac{x - B_m}{2B_s}\right), \quad (2.171)$$

when  $x \in [B_m - B_s, B_m + B_s]$  while  $B = B_0$  outside this interval.

Also an abrupt construction may be tried:

$$B(x; B_0, B_a, B_m, B_s) = B_0 + B_a, \quad (2.172)$$

for  $x \in [B_m - B_s, B_m + B_s]$  while  $B = B_0$  outside this interval.

The `wave1D_dn_vc.py` program can be used as starting point for the implementation. Visualize both the bottom topography and the water

surface elevation in the same plot. Allow for a flexible choice of bottom shape: (2.170), (2.171), (2.172), or  $B(x) = B_0$  (flat).

The purpose of this problem is to explore the quality of the numerical solution  $\eta_i^n$  for different shapes of the bottom obstruction. The "cosine hat" and the box-shaped hills have abrupt changes in the derivative of  $H(x)$  and are more likely to generate numerical noise than the smooth Gaussian shape of the hill. Investigate if this is true. Filename: `tsunami1D_hill`.

### Problem 2.24: Earthquake-generated tsunami over a 3D hill

This problem extends Problem 2.23 to a three-dimensional wave phenomenon, governed by the 2D PDE (2.153). We assume that the earthquake arises from a fault along the line  $x = 0$  in the  $xy$ -plane so that the initial lift of the surface can be taken as  $I(x)$  in Problem 2.23. That is, a plane wave is propagating to the right, but will experience bending because of the bottom.

The bottom shape is now a function of  $x$  and  $y$ . An "elliptic" Gaussian function in two dimensions, with its peak at  $(B_{mx}, B_{my})$ , generalizes (2.170):

$$B(x; B_0, B_a, B_{mx}, B_{my}, B_s, b) = B_0 + B_a \exp \left( - \left( \frac{x - B_{mx}}{B_s} \right)^2 - \left( \frac{y - B_{my}}{bB_s} \right)^2 \right), \quad (2.173)$$

where  $b$  is a scaling parameter:  $b = 1$  gives a circular Gaussian function with circular contour lines, while  $b \neq 1$  gives an elliptic shape with elliptic contour lines.

The "cosine hat" (2.171) can also be generalized to

$$B(x; B_0, B_a, B_{mx}, B_{my}, B_s) = B_0 + B_a \cos \left( \pi \frac{x - B_{mx}}{2B_s} \right) \cos \left( \pi \frac{y - B_{my}}{2B_s} \right), \quad (2.174)$$

when  $0 \leq \sqrt{x^2 + y^2} \leq B_s$  and  $B = B_0$  outside this circle.

A box-shaped obstacle means that

$$B(x; B_0, B_a, B_m, B_s, b) = B_0 + B_a \quad (2.175)$$

for  $x$  and  $y$  inside a rectangle

$$B_{mx} - B_s \leq x \leq B_{mx} + B_s, \quad B_{my} - bB_s \leq y \leq B_{my} + bB_s,$$

and  $B = B_0$  outside this rectangle. The  $b$  parameter controls the rectangular shape of the cross section of the box.

Note that the initial condition and the listed bottom shapes are symmetric around the line  $y = B_{my}$ . We therefore expect the surface elevation also to be symmetric with respect to this line. This means that we can halve the computational domain by working with  $[0, L_x] \times [0, B_{my}]$ . Along the upper boundary,  $y = B_{my}$ , we must impose the symmetry condition  $\partial\eta/\partial n = 0$ . Such a symmetry condition ( $-\eta_x = 0$ ) is also needed at the  $x = 0$  boundary because the initial condition has a symmetry here. At the lower boundary  $y = 0$  we also set a Neumann condition (which becomes  $-\eta_y = 0$ ). The wave motion is to be simulated until the wave hits the reflecting boundaries where  $\partial\eta/\partial n = \eta_x = 0$  (one can also set  $\eta = 0$  - the particular condition does not matter as long as the simulation is stopped before the wave is influenced by the boundary condition).

Visualize the surface elevation. Investigate how different hill shapes, different sizes of the water gap above the hill, and different resolutions  $\Delta x = \Delta y = h$  and  $\Delta t$  influence the numerical quality of the solution. Filename: `tsunami2D_hill`.

### **Problem 2.25: Investigate Matplotlib for visualization**

Play with native Matplotlib code for visualizing 2D solutions of the wave equation with variable wave velocity. See if there are effective ways to visualize both the solution and the wave velocity. Filename: `tsunami2D_hill_mpl`.

### **Problem 2.26: Investigate visualization packages**

Create some fancy 3D visualization of the water waves *and* the sub-sea hill in Problem 2.24. Try to make the hill transparent. Possible visualization tools are [Mayavi](#), [Paraview](#), and [OpenDX](#). Filename: `tsunami2D_hill_viz`.

### **Problem 2.27: Implement loops in compiled languages**

Extend the program from Problem 2.24 such that the loops over mesh points, inside the time loop, are implemented in compiled languages.

Consider implementations in Cython, Fortran via `f2py`, C via Cython, C via `f2py`, C/C++ via Instant, and C/C++ via `scipy.weave`. Perform efficiency experiments to investigate the relative performance of the various implementations. It is often advantageous to normalize CPU times by the fastest method on a given mesh. Filename: `tsunami2D_hill_compiled`.

### Exercise 2.28: Simulate seismic waves in 2D

The goal of this exercise is to simulate seismic waves using the PDE model (2.137) in a 2D  $xz$  domain with geological layers. Introduce  $m$  horizontal layers of thickness  $h_i$ ,  $i = 0, \dots, m - 1$ . Inside layer number  $i$  we have a vertical wave velocity  $c_{z,i}$  and a horizontal wave velocity  $c_{h,i}$ . Make a program for simulating such 2D waves. Test it on a case with 3 layers where

$$c_{z,0} = c_{z,1} = c_{z,2}, \quad c_{h,0} = c_{h,2}, \quad c_{h,1} \ll c_{h,0}.$$

Let  $s$  be a localized point source at the middle of the Earth's surface (the upper boundary) and investigate how the resulting wave travels through the medium. The source can be a localized Gaussian peak that oscillates in time for some time interval. Place the boundaries far enough from the expanding wave so that the boundary conditions do not disturb the wave. Then the type of boundary condition does not matter, except that we physically need to have  $p = p_0$ , where  $p_0$  is the atmospheric pressure, at the upper boundary. Filename: `seismic2D`.

### Project 2.29: Model 3D acoustic waves in a room

The equation for sound waves in air is derived in Section 2.14.5 and reads

$$p_{tt} = c^2 \nabla^2 p,$$

where  $p(x, y, z, t)$  is the pressure and  $c$  is the speed of sound, taken as 340 m/s. However, sound is absorbed in the air due to relaxation of molecules in the gas. A model for simple relaxation, valid for gases consisting only of one type of molecules, is a term  $c^2 \tau_s \nabla^2 p_t$  in the PDE, where  $\tau_s$  is the relaxation time. If we generate sound from, e.g., a loudspeaker in the room, this sound source must also be added to the governing equation.

The PDE with the mentioned type of damping and source then becomes

$$p_t t = c^2 \nabla^p + c^2 \tau_s \nabla^2 p_t + f, \quad (2.176)$$

where  $f(x, y, z, t)$  is the source term.

The walls can absorb some sound. A possible model is to have a "wall layer" (thicker than the physical wall) outside the room where  $c$  is changed such that some of the wave energy is reflected and some is absorbed in the wall. The absorption of energy can be taken care of by adding a damping term  $b p_t$  in the equation:

$$p_t t + b p_t = c^2 \nabla^p + c^2 \tau_s \nabla^2 p_t + f. \quad (2.177)$$

Typically,  $b = 0$  in the room and  $b > 0$  in the wall. A discontinuity in  $b$  or  $c$  will give rise to reflections. It can be wise to use a constant  $c$  in the wall to control reflections because of the discontinuity between  $c$  in the air and in the wall, while  $b$  is gradually increased as we go into the wall to avoid reflections because of rapid changes in  $b$ . At the outer boundary of the wall the condition  $p = 0$  or  $\partial p / \partial n = 0$  can be imposed. The waves should anyway be approximately damped to  $p = 0$  this far out in the wall layer.

There are two strategies for discretizing the  $\nabla^2 p_t$  term: using a center difference between times  $n + 1$  and  $n - 1$  (if the equation is sampled at level  $n$ ), or use a one-sided difference based on levels  $n$  and  $n - 1$ . The latter has the advantage of not leading to any equation system, while the former is second-order accurate as the scheme for the simple wave equation  $p_t t = c^2 \nabla^2 p$ . To avoid an equation system, go for the one-sided difference such that the overall scheme becomes explicit and only of first order in time.

Develop a 3D solver for the specified PDE and introduce a wall layer. Test the solver with the method of manufactured solutions. Make some demonstrations where the wall reflects and absorbs the waves (reflection because of discontinuity in  $b$  and absorption because of growing  $b$ ). Experiment with the impact of the  $\tau_s$  parameter. Filename: `acoustics`.

### Project 2.30: Solve a 1D transport equation

We shall study the wave equation

$$u_t + c u_x = 0, \quad x \in (0, L], \quad t \in (0, T], \quad (2.178)$$

with initial condition

$$u(x, 0) = I(x), \quad x \in [0, L], \quad (2.179)$$

and *one* periodic boundary condition

$$u(0, t) = u(L, t). \quad (2.180)$$

This boundary condition means that what goes out of the domain at  $x = L$  comes in at  $x = 0$ . Roughly speaking, we need only one boundary condition because of the spatial derivative is of first order only.

**Physical interpretation.** The parameter  $c$  can be constant or variable,  $c = c(x)$ . The equation (2.178) arises in *transport* problems where a quantity  $u$ , which could be temperature or concentration of some contaminant, is transported with the velocity  $c$  of a fluid. In addition to the transport imposed by "travelling with the fluid",  $u$  may also be transported by diffusion (such as heat conduction or Fickian diffusion), but we have in the model  $u_t + cu_x$  assumed that diffusion effects are negligible, which they often are.

**a)** Show that under the assumption of  $a = \text{const}$ ,

$$u(x, t) = I(x - ct) \quad (2.181)$$

fulfills the PDE as well as the initial and boundary condition (provided  $I(0) = I(L)$ ).

A widely used numerical scheme for (2.178) applies a forward difference in time and a backward difference in space when  $c > 0$ :

$$[D_t^+ u + c D_x^- u = 0]_i^n. \quad (2.182)$$

For  $c < 0$  we use a forward difference in space:  $[c D_x^+ u]_i^n$ .

**b)** Set up a computational algorithm and implement it in a function. Assume  $a$  is constant and positive.

**c)** Test the implementation by using the remarkable property that the numerical solution is exact at the mesh points if  $\Delta t = c^{-1} \Delta x$ .

**d)** Make a movie comparing the numerical and exact solution for the following two choices of initial conditions:

$$I(x) = \left[ \sin\left(\pi \frac{x}{L}\right) \right]^{2n} \quad (2.183)$$

where  $n$  is an integer, typically  $n = 5$ , and

$$I(x) = \exp\left(-\frac{(x - L/2)^2}{2\sigma^2}\right). \quad (2.184)$$

Choose  $\Delta t = c^{-1} \Delta x, 0.9c^{-1} \Delta x, 0.5c^{-1} \Delta x$ .

- e)** The performance of the suggested numerical scheme can be investigated by analyzing the numerical dispersion relation. Analytically, we have that the *Fourier component*

$$u(x, t) = e^{i(kx - \omega t)},$$

is a solution of the PDE if  $\omega = kc$ . This is the *analytical dispersion relation*. A complete solution of the PDE can be built by adding up such Fourier components with different amplitudes, where the initial condition  $I$  determines the amplitudes. The solution  $u$  is then represented by a Fourier series.

A similar discrete Fourier component at  $(x_p, t_n)$  is

$$u_p^q = e^{i(kp\Delta x - \tilde{\omega}n\Delta t)},$$

where in general  $\tilde{\omega}$  is a function of  $k$ ,  $\Delta t$ , and  $\Delta x$ , and differs from the exact  $\omega = kc$ .

Insert the discrete Fourier component in the numerical scheme and derive an expression for  $\tilde{\omega}$ , i.e., the discrete dispersion relation. Show in particular that if  $\Delta t/(c\Delta x) = 1$ , the discrete solution coincides with the exact solution at the mesh points, regardless of the mesh resolution (!). Show that if the stability condition

$$\frac{\Delta t}{c\Delta x} \leq 1,$$

the discrete Fourier component cannot grow (i.e.,  $\tilde{\omega}$  is real).

- f)** Write a test for your implementation where you try to use information from the numerical dispersion relation.

We shall hereafter assume that  $c(x) > 0$ .

- g)** Set up a computational algorithm for the variable coefficient case and implement it in a function. Make a test that the function works for constant  $a$ .

- h)** It can be shown that for an observer moving with velocity  $c(x)$ ,  $u$  is constant. This can be used to derive an exact solution when  $a$  varies with  $x$ . Show first that

$$u(x, t) = f(C(x) - t), \quad (2.185)$$

where

$$C'(x) = \frac{1}{c(x)},$$

is a solution of (2.178) for any differentiable function  $f$ .

- i) Use the initial condition to show that an exact solution is

$$u(x, t) = I(C^{-1}(C(x) - t)),$$

with  $C^{-1}$  being the inverse function of  $C = \int c^1 dx$ . Since  $C(x)$  is an integral  $\int_0^x (1/c)dx$ ,  $C(x)$  is monotonically increasing and there exists hence an inverse function  $C^{-1}$  with values in  $[0, L]$ .

To compute (2.185) we need to integrate  $1/c$  to obtain  $C$  and then compute the inverse of  $C$ .

The inverse function computation can be easily done if we first think discretely. Say we have some function  $y = g(x)$  and seek its inverse. Plotting  $(x_i, y_i)$ , where  $y_i = g(x_i)$  for some mesh points  $x_i$ , displays  $g$  as a function of  $x$ . The inverse function is simply  $x$  as a function of  $g$ , i.e., the curve with points  $(y_i, x_i)$ . We can therefore quickly compute points at the curve of the inverse function. One way of extending these points to a continuous function is to assume a linear variation (known as linear interpolation) between the points (which actually means to draw straight lines between the points, exactly as done by a plotting program).

The function `wrap2callable` in `scitools.std` can take a set of points and return a continuous function that corresponds to linear variation between the points. The computation of the inverse of a function  $g$  on  $[0, L]$  can then be done by

```
def inverse(g, domain, resolution=101):
    x = linspace(domain[0], domain[L], resolution)
    y = g(x)
    from scitools.std import wrap2callable
    g_inverse = wrap2callable((y, x))
    return g_inverse
```

To compute  $C(x)$  we need to integrate  $1/c$ , which can be done by a Trapezoidal rule. Suppose we have computed  $C(x_i)$  and need to compute  $C(x_{i+1})$ . Using the Trapezoidal rule with  $m$  subintervals over the integration domain  $[x_i, x_{i+1}]$  gives

$$C(x_{i+1}) = C(x_i) + \int_{x_i}^{x_{i+1}} \frac{dx}{c} \approx h \left( \frac{1}{2} \frac{1}{c(x_i)} + \frac{1}{2} \frac{1}{c(x_{i+1})} + \sum_{j=1}^{m-1} \frac{1}{c(x_i + jh)} \right), \quad (2.186)$$

where  $h = (x_{i+1} - x_i)/m$  is the length of the subintervals used for the integral over  $[x_i, x_{i+1}]$ . We observe that (2.186) is a *difference equation* which we can solve by repeatedly applying (2.186) for  $i = 0, 1, \dots, N_x - 1$  if a mesh  $x_0, x_1, \dots, x_{N_x}$  is prescribed. Note that  $C(0) = 0$ .

- j)** Implement a function for computing  $C(x_i)$  and one for computing  $C^{-1}(x)$  for any  $x$ . Use these two functions for computing the exact solution  $I(C^{-1}(C(x) - t))$ . End up with a function `u_exact_variable_c(x, n, c, I)` that returns the value of  $I(C^{-1}(C(x) - t_n))$ .

- k)** Make movies showing a comparison of the numerical and exact solutions for the two initial conditions (2.183) and (2.30). Choose  $\Delta t = \Delta x / \max_{0,L} c(x)$  and the velocity of the medium as

1.  $c(x) = 1 + \epsilon \sin(k\pi x/L)$ ,  $\epsilon < 1$ ,
2.  $c(x) = 1 + I(x)$ , where  $I$  is given by (2.183) or (2.30).

The PDE  $u_t + cu_x = 0$  expresses that the initial condition  $I(x)$  is transported with velocity  $c(x)$ .

Filename: `advec1D`.

### Problem 2.31: General analytical solution of a 1D damped wave equation

We consider an initial-boundary value problem for the damped wave equation:

$$\begin{aligned} u_{tt} + bu_t - c^2 u_{xx} &= 0, & x \in (0, L), t \in (0, T] \\ u(0, t) &= 0, \\ u(L, t) &= 0, \\ u(x, 0) &= I(x), \\ u_t(x, 0) &= V(x). \end{aligned}$$

Here,  $b \geq 0$  and  $c$  are given constants. The aim is to derive a general analytical solution of this problem. Familiarity with the method of separation of variables for solving PDEs will be assumed.

- a)** Seek a solution on the form  $u(x, t) = X(x)T(t)$ . Insert this solution in the PDE and show that it leads to two differential equations for  $X$  and  $T$ :

$$T'' + bT' + \lambda T = 0, \quad c^2 X'' + \lambda X = 0,$$

with  $X(0) = X(L) = 0$  as boundary conditions, and  $\lambda$  as a constant to be determined.

- b)** Show that  $X(x)$  is on the form

$$X_n(x) = C_n \sin kx, \quad k = \frac{n\pi}{L}, \quad n = 1, 2, \dots$$

where  $C_n$  is an arbitrary constant.

- c)** Under the assumption that  $(b/2)^2 < k^2$ , show that  $T(t)$  is on the form

$$T_n(t) = e^{-\frac{1}{2}bt}(a_n \cos \omega t + b_n \sin \omega t), \quad \omega = \sqrt{k^2 - \frac{1}{4}b^2}, \quad n = 1, 2, \dots$$

The complete solution is then

$$u(x, t) = \sum_{n=1}^{\infty} \sin kx e^{-\frac{1}{2}bt} (A_n \cos \omega t + B_n \sin \omega t),$$

where the constants  $A_n$  and  $B_n$  must be computed from the initial conditions.

- d)** Derive a formula for  $A_n$  from  $u(x, 0) = I(x)$  and developing  $I(x)$  as a sine Fourier series on  $[0, L]$ .

- e)** Derive a formula for  $B_n$  from  $u_t(x, 0) = V(x)$  and developing  $V(x)$  as a sine Fourier series on  $[0, L]$ .

- f)** Calculate  $A_n$  and  $B_n$  from vibrations of a string where  $V(x) = 0$  and

$$I(x) = \begin{cases} ax/x_0, & x < x_0, \\ a(L-x)/(L-x_0), & \text{otherwise} \end{cases} \quad (2.187)$$

- g)** Implement the series for  $u(x, t)$  in a function `u_series(x, t, tol=1E-10)`, where `tol` is a tolerance for truncating the series. Simply sum the terms until  $|a_n|$  and  $|b_n|$  both are less than `tol`.

- h)** What will change in the derivation of the analytical solution if we have  $u_x(0, t) = u_x(L, t) = 0$  as boundary conditions? And how will you solve the problem with  $u(0, t) = 0$  and  $u_x(L, t) = 0$ ?

Filename: `damped_wave1D`.

**Problem 2.32: General analytical solution of a 2D damped wave equation**

Carry out Problem 2.31 in the 2D case:  $u_{tt} + bu_t = c^2(u_{xx} + u_{yy})$ , where  $(x, y) \in (0, L_x) \times (0, L_y)$ . Assume a solution on the form  $u(x, y, t) = X(x)Y(y)T(t)$ . Filename: `damped_wave2D`.

# Diffusion equations

3

**hpl 10:** Here is a todo list.

- There are few index entries, see the “wave equation” index entries and how many sub-entries we have there...
- Check if the material on diffusion equation with radial symmetry is okay.
- Implement a diffusion equation with axi-symmetry such we can solve pulsating flow in a circular blood vessel, see Exercise 3.10.

The famous *diffusion equation*, also known as the *heat equation*, reads

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2},$$

where  $u(x, t)$  is the unknown function to be solved for,  $x$  is a coordinate in space, and  $t$  is time. The coefficient  $\alpha$  is the *diffusion coefficient* and determines how fast  $u$  changes in time. A quick short form for the diffusion equation is  $u_t = \alpha u_{xx}$ .

Compared to the wave equation,  $u_{tt} = c^2 u_{xx}$ , which looks very similar, the diffusion equation features solutions that are very different from those of the wave equation. Also, the diffusion equation makes quite different demands to the numerical methods.

Typical diffusion problems may experience rapid change in the very beginning, but then the evolution of  $u$  becomes slower and slower. The solution is usually very smooth, and after some time, one cannot recognize

the initial shape of  $u$ . This is in sharp contrast to solutions of the wave equation where the initial shape is preserved - the solution is basically a moving initial condition. The standard wave equation  $u_{tt} = c^2 u_{xx}$  has solutions that propagates with speed  $c$  forever, without changing shape, while the diffusion equation converges to a *stationary solution*  $\bar{u}(x)$  as  $t \rightarrow \infty$ . In this limit,  $u_t = 0$ , and  $\bar{u}$  is governed by  $\bar{u}''(x) = 0$ . This stationary limit of the diffusion equation is called the *Laplace* equation and arises in a very wide range of applications throughout the sciences.

It is possible to solve for  $u(x, t)$  using an explicit scheme, as we do in Section 3.1, but the time step restrictions soon become much less favorable than for an explicit scheme for the wave equation. And of more importance, since the solution  $u$  of the diffusion equation is very smooth and changes slowly, small time steps are not convenient and not required by accuracy as the diffusion process converges to a stationary state. Therefore, implicit schemes as described in Section 3.2 are popular, but these require solutions of systems of algebraic equations. We shall use ready-made software for this purpose, but also program some simple iterative methods.

The exposition below assumes that the reader is familiar with the basic ideas of discretization and implementation of wave equations from Chapter 2. Readers not familiar with the Forward Euler, Backward Euler, and Crank-Nicolson (or centered or midpoint) discretization methods in time should consult, e.g., Section 1.1 in [4].

## 3.1 An explicit method for the 1D diffusion equation

### 3.1.1 The initial-boundary value problem for 1D diffusion

To obtain a unique solution of the diffusion equation, or equivalently, to apply numerical methods, we need initial and boundary conditions. The diffusion equation goes with one initial condition  $u(x, 0) = I(x)$ , where  $I$  is a prescribed function. One boundary condition is required at each point on the boundary, which in 1D means that  $u$  must be known,  $u_x$  must be known, or some combination of them.

We shall start with the simplest boundary condition:  $u = 0$ . The complete initial-boundary value diffusion problem in one space dimension can then be specified as

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2} + f, \quad x \in (0, L), \quad t \in (0, T] \quad (3.1)$$

$$u(x, 0) = I(x), \quad x \in [0, L] \quad (3.2)$$

$$u(0, t) = 0, \quad t > 0, \quad (3.3)$$

$$u(L, t) = 0, \quad t > 0. \quad (3.4)$$

With only a first-order derivative in time, only one *initial condition* is needed, while the second-order derivative in space leads to a demand for two *boundary conditions*. We have added a source term  $f = f(x, t)$  for convenience when testing implementations.

Diffusion equations like (3.1) have a wide range of applications throughout physical, biological, and financial sciences. One of the most common applications is propagation of heat, where  $u(x, t)$  represents the temperature of some substance at point  $x$  and time  $t$ .

### 3.1.2 Forward Euler scheme

The first step in the discretization procedure is to replace the domain  $[0, L] \times [0, T]$  by a set of mesh points. Here we apply equally spaced mesh points

$$x_i = i \Delta x, \quad i = 0, \dots, N_x,$$

and

$$t_n = n \Delta t, \quad n = 0, \dots, N_t.$$

Moreover,  $u_i^n$  denotes the mesh function that approximates  $u(x_i, t_n)$  for  $i = 0, \dots, N_x$  and  $n = 0, \dots, N_t$ . Requiring the PDE (3.1) to be fulfilled at a mesh point  $(x_i, t_n)$  leads to the equation

$$\frac{\partial}{\partial t} u(x_i, t_n) = \alpha \frac{\partial^2}{\partial x^2} u(x_i, t_n) + f(x_i, t_n), \quad (3.5)$$

The next step is to replace the derivatives by finite difference approximations. The computationally simplest method arises from using a forward difference in time and a central difference in space:

$$[D_t^+ u = \alpha D_x D_x u + f]_i^n. \quad (3.6)$$

Written out,

$$\frac{u_i^{n+1} - u_i^n}{\Delta t} = \alpha \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2} + f_i^n. \quad (3.7)$$

We have turned the PDE into algebraic equations, also often called discrete equations. The key property of the equations is that they are algebraic, which makes them easy to solve. As usual, we anticipate that  $u_i^n$  is already computed such that  $u_i^{n+1}$  is the only unknown in (3.7). Solving with respect to this unknown is easy:

$$u_i^{n+1} = u_i^n + F (u_{i+1}^n - 2u_i^n + u_{i-1}^n) + \Delta t f_i^n, \quad (3.8)$$

where we have introduced the *mesh Fourier number*:

$$F = \alpha \frac{\Delta t}{\Delta x^2}. \quad (3.9)$$

### F is the key parameter in the discrete diffusion equation

Note that  $F$  is a *dimensionless* number that lumps the key physical parameter in the problem,  $\alpha$ , and the discretization parameters  $\Delta x$  and  $\Delta t$  into a single parameter. All the properties of the numerical method are critically dependent upon the value of  $F$  (see Section 3.3 for details).

The computational algorithm then becomes

1. compute  $u_i^0 = I(x_i)$  for  $i = 0, \dots, N_x$
2. for  $n = 0, 1, \dots, N_t$ :
  - a. apply (3.8) for all the internal spatial points  $i = 1, \dots, N_x - 1$
  - b. set the boundary values  $u_i^{n+1} = 0$  for  $i = 0$  and  $i = N_x$

The algorithm is compactly and fully specified in Python:

```
import numpy as np
x = np.linspace(0, L, Nx+1)      # mesh points in space
dx = x[1] - x[0]
t = np.linspace(0, T, Nt+1)      # mesh points in time
dt = t[1] - t[0]
F = alpha*dt/dx**2
u   = np.zeros(Nx+1)              # unknown u at new time level
u_1 = np.zeros(Nx+1)              # u at the previous time level

# Set initial condition u(x,0) = I(x)
for i in range(0, Nx+1):
    u_1[i] = I(x[i])
```

```

for n in range(0, Nt):
    # Compute u at inner mesh points
    for i in range(1, Nx):
        u[i] = u_1[i] + F*(u_1[i-1] - 2*u_1[i] + u_1[i+1]) + \
            f(x[i], t[n])

    # Insert boundary conditions
    u[0] = 0; u[Nx] = 0

    # Update u_1 before next step
    u_1[:] = u

```

Note that we use `a` for  $\alpha$  in the code, motivated by easy visual mapping between the variable name and the mathematical symbol in formulas.

We need to state already now that the shown algorithm does not produce meaningful results unless  $F \leq 1/2$ . Why is explained in Section 3.3.

### 3.1.3 Implementation

The file `diffu1D_u0.py` contains a complete function `solver_FE_simple` for solving the 1D diffusion equation with  $u = 0$  on the boundary as specified in the algorithm above:

```

import numpy as np

def solver_FE_simple(I, a, f, L, dt, F, T):
    """
    Simplest expression of the computational algorithm
    using the Forward Euler method and explicit Python loops.
    For this method F <= 0.5 for stability.
    """
    import time; t0 = time.clock() # For measuring the CPU time

    Nt = int(round(T/float(dt)))
    t = np.linspace(0, Nt*dt, Nt+1) # Mesh points in time
    dx = np.sqrt(a*dt/F)
    Nx = int(round(L/dx))
    x = np.linspace(0, L, Nx+1)      # Mesh points in space
    # Make sure dx and dt are compatible with x and t
    dx = x[1] - x[0]
    dt = t[1] - t[0]

    u = np.zeros(Nx+1)
    u_1 = np.zeros(Nx+1)

    # Set initial condition u(x,0) = I(x)
    for i in range(0, Nx+1):
        u_1[i] = I(x[i])

```

```

for n in range(0, Nt):
    # Compute u at inner mesh points
    for i in range(1, Nx):
        u[i] = u_1[i] + F*(u_1[i-1] - 2*u_1[i] + u_1[i+1]) + \
                dt*f(x[i], t[n])

    # Insert boundary conditions
    u[0] = 0; u[Nx] = 0

    # Switch variables before next step
    #u_1[:] = u # safe, but slow
    u_1, u = u, u_1

t1 = time.clock()
return u_1, x, t, t1-t0 # u_1 holds latest u

```

A faster version, based on vectorization of the finite difference scheme, is available in the function `solver_FE`. The vectorized version replaces the explicit loop

```

for i in range(1, Nx):
    u[i] = u_1[i] + F*(u_1[i-1] - 2*u_1[i] + u_1[i+1]) \
            + dt*f(x[i], t[n])

```

by arithmetics on displaced slices of the `u` array:

```

u[1:Nx] = u_1[1:Nx] + F*(u_1[0:Nx-1] - 2*u_1[1:Nx] + u_1[2:Nx+1]) \
            + dt*f(x[1:Nx], t[n])
# or
u[1:-1] = u_1[1:-1] + F*(u_1[0:-2] - 2*u_1[1:-1] + u_1[2:]) \
            + dt*f(x[1:-1], t[n])

```

For example, the vectorized version runs 70 times faster than the scalar version in a case with 100 time steps and a spatial mesh of  $10^5$  cells.

The `solver_FE` function also features a callback function such that the user can process the solution at each time level. The callback function looks like `user_action(u, x, t, n)`, where `u` is the array containing the solution at time level `n`, `x` holds all the spatial mesh points, while `t` holds all the temporal mesh points. Apart from the vectorized loop over the spatial mesh points, the callback function, and a bit more complicated setting of the source `f` it is not specified (`None`), the `solver_FE` is identical to `solver_FE_simple` above:

```

def solver_FE(I, a, f, L, dt, F, T,
              user_action=None, version='scalar'):
    """
    Vectorized implementation of solver_FE_simple.
    """
    import time; t0 = time.clock() # for measuring the CPU time

    Nt = int(round(T/float(dt)))

```

```

t = np.linspace(0, Nt*dt, Nt+1)    # Mesh points in time
dx = np.sqrt(a*dt/F)
Nx = int(round(L/dx))
x = np.linspace(0, L, Nx+1)        # Mesh points in space
# Make sure dx and dt are compatible with x and t
dx = x[1] - x[0]
dt = t[1] - t[0]

u   = np.zeros(Nx+1)    # solution array
u_1 = np.zeros(Nx+1)    # solution at t-dt

# Set initial condition
for i in range(0,Nx+1):
    u_1[i] = I(x[i])

if user_action is not None:
    user_action(u_1, x, t, 0)

for n in range(0, Nt):
    # Update all inner points
    if version == 'scalar':
        for i in range(1, Nx):
            u[i] = u_1[i] + F*(u_1[i-1] - 2*u_1[i] + u_1[i+1]) +\
                dt*f(x[i], t[n])

    elif version == 'vectorized':
        u[1:Nx] = u_1[1:Nx] + \
            F*(u_1[0:Nx-1] - 2*u_1[1:Nx] + u_1[2:Nx+1]) +\
            dt*f(x[1:Nx], t[n])
    else:
        raise ValueError('version=%s' % version)

    # Insert boundary conditions
    u[0] = 0; u[Nx] = 0
    if user_action is not None:
        user_action(u, x, t, n+1)

    # Switch variables before next step
    u_1, u = u, u_1

t1 = time.clock()
return t1-t0

```

### 3.1.4 Verification

Before thinking about running the functions in the previous section, we need to construct a suitable test example for verification. It appears that a manufactured solution that is linear in time and at most quadratic in space fulfills the Forward Euler scheme exactly. With the restriction that  $u = 0$  for  $x = 0, L$ , we can try the solution

$$u(x, t) = 5tx(L - x).$$

Inserted in the PDE, it requires a source term

$$f(x, t) = 10\alpha t + 5x(L - x).$$

With the formulas from Appendix A.4 we can easily check that the manufactured  $u$  fulfills the scheme:

$$\begin{aligned} [D_t^+ u = \alpha D_x D_x u + f]_i^n &= [5x(L - x)D_t^+ t = 5t\alpha D_x D_x(xL - x^2) + \\ &\quad 10\alpha t + 5x(L - x)]_i^n \\ &= [5x(L - x) = 5t\alpha(-2) + 10\alpha t + 5x(L - x)]_i^n. \end{aligned}$$

The computation of the source term, given any  $u$ , is easily automated with `sympy`:

```
import sympy as sym
x, t, a, L = sym.symbols('x t a L')
u = x*(L-x)*5*t

def pde(u):
    return sym.diff(u, t) - a*sym.diff(u, x, x)

f = sym.simplify(pde(u))
```

Now we can choose any expression for  $u$  and automatically get the suitable source term  $f$ . However, the manufactured solution  $u$  will in general not be exactly reproduced by the scheme: only constant and linear functions are differentiated correctly by a forward difference, while only constant, linear, and quadratic functions are differentiated exactly by a  $[D_x D_x u]_i^n$  difference.

The numerical code will need to access the  $u$  and  $f$  above as Python functions. The exact solution is wanted as a Python function  $u_{\text{exact}}(x, t)$ , while the source term is wanted as  $f(x, t)$ . The parameters  $a$  and  $L$  in  $u$  and  $f$  above are symbols and must be replaced by `float` objects in a Python function. This can be done by redefining  $a$  and  $L$  as `float` objects and performing substitutions of symbols by numbers in  $u$  and  $f$ . The appropriate code looks like this:

```
a = 0.5
L = 1.5
u_exact = sym.lambdify(
    [x, t], u.subs('L', L).subs('a', a), modules='numpy')
f = sym.lambdify(
```

```
[x, t], f.subs('L', L).subs('a', a), modules='numpy')
I = lambda x: u_exact(x, 0)
```

Here we also make a function `I` for the initial condition.

The idea now is that our manufactured solution should be exactly reproduced by the code (to machine precision). For this purpose we make a test function for comparing the exact and numerical solutions at the end of the time interval:

```
def test_solver_FE():
    # Define u_exact, f, I as explained above

    dx = L/3  # 3 cells
    F = 0.5
    dt = F*dx**2

    u, x, t, cpu = solver_FE_simple(
        I=I, a=a, f=f, L=L, dt=dt, F=F, T=2)
    u_e = u_exact(x, t[-1])
    diff = abs(u_e - u).max()
    tol = 1E-14
    assert diff < tol, 'max diff solver_FE_simple: %g' % diff

    u, x, t, cpu = solver_FE(
        I=I, a=a, f=f, L=L, dt=dt, F=F, T=2,
        user_action=None, version='scalar')
    u_e = u_exact(x, t[-1])
    diff = abs(u_e - u).max()
    tol = 1E-14
    assert diff < tol, 'max diff solver_FE, scalar: %g' % diff

    u, x, t, cpu = solver_FE(
        I=I, a=a, f=f, L=L, dt=dt, F=F, T=2,
        user_action=None, version='vectorized')
    u_e = u_exact(x, t[-1])
    diff = abs(u_e - u).max()
    tol = 1E-14
    assert diff < tol, 'max diff solver_FE, vectorized: %g' % diff
```

We emphasize that the value  $F=0.5$  is critical: the tests above will fail if  $F$  has a larger value (this is because the Forward Euler scheme is unstable for  $F > 1/2$ ).

### 3.1.5 Numerical experiments

When a test function like the one above runs silently without errors, we have some evidence for a correct implementation of the numerical method. The next step is to do some experiments with more interesting solutions.

We target a scaled diffusion problem where  $x/L$  is a new spatial coordinate and  $\alpha t/L^2$  is a new time coordinate. The source term  $f$  is omitted, and  $u$  is scaled by  $\max_{x \in [0,L]} |I(x)|$  (see Section 3.2 in [5] for details). The governing PDE is then

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2},$$

in the spatial domain  $[0, L]$ , with boundary conditions  $u(0) = u(1) = 0$ . Two initial conditions will be tested: a discontinuous plug,

$$I(x) = \begin{cases} 0, & |x - L/2| > 0.1 \\ 1, & \text{otherwise} \end{cases}$$

and a smooth Gaussian function,

$$I(x) = e^{-\frac{1}{2\sigma^2}(x-L/2)^2}.$$

The functions `plug` and `gaussian` in `diffu1D_u0.py` run the two cases, respectively:

```
def plug(scheme='FE', F=0.5, Nx=50):
    L = 1.
    a = 1.
    T = 0.1
    # Compute dt from Nx and F
    dx = L/Nx; dt = F/a*dx**2

    def I(x):
        """Plug profile as initial condition."""
        if abs(x-L/2.0) > 0.1:
            return 0
        else:
            return 1

    cpu = viz(I, a, L, dt, F, T,
              umin=-0.1, umax=1.1,
              scheme=scheme, animate=True, framefiles=True)
    print 'CPU time:', cpu

def gaussian(scheme='FE', F=0.5, Nx=50, sigma=0.05):
    L = 1.
    a = 1.
    T = 0.1
    # Compute dt from Nx and F
    dx = L/Nx; dt = F/a*dx**2

    def I(x):
        """Gaussian profile as initial condition."""
        return exp(-0.5*((x-L/2.0)**2)/sigma**2)
```

```

u, cpu = viz(I, a, L, dt, F, T,
             umin=-0.1, umax=1.1,
             scheme=scheme, animate=True, framefiles=True)
print 'CPU time:', cpu

```

These functions make use of the function `viz` for running the solver and visualizing the solution using a callback function with plotting:

```

def viz(I, a, L, dt, F, T, umin, umax,
       scheme='FE', animate=True, framefiles=True):

    def plot_u(u, x, t, n):
        plt.plot(x, u, 'r-', axis=[0, L, umin, umax],
                  title='t=%f' % t[n])
        if framefiles:
            plt.savefig('tmp_frame%04d.png' % n)
        if t[n] == 0:
            time.sleep(2)
        elif not framefiles:
            # It takes time to write files so pause is needed
            # for screen only animation
            time.sleep(0.2)

    user_action = plot_u if animate else lambda u,x,t,n: None

    cpu = eval('solver_'+scheme)(I, a, L, dt, F, T,
                                 user_action=user_action)
    return cpu

```

Notice that this `viz` function stores all the solutions in a list `solutions` in the callback function. Modern computers have hardly any problem with storing a lot of such solutions for moderate values of  $N_x$  in 1D problems, but for 2D and 3D problems, this technique cannot be used and solutions must be stored in files.

### **hpl 11: Better to show the scalable file solution here?**

Our experiments employ a time step  $\Delta t = 0.0002$  and simulate for  $t \in [0, 0.1]$ . First we try the highest value of  $F$ :  $F = 0.5$ . This resolution corresponds to  $N_x = 50$ . A possible terminal command is

---

Terminal

```
Terminal> python -c 'from diffu1D_u0 import gaussian
> gaussian("solver_FE", F=0.5, dt=0.0002)'
```

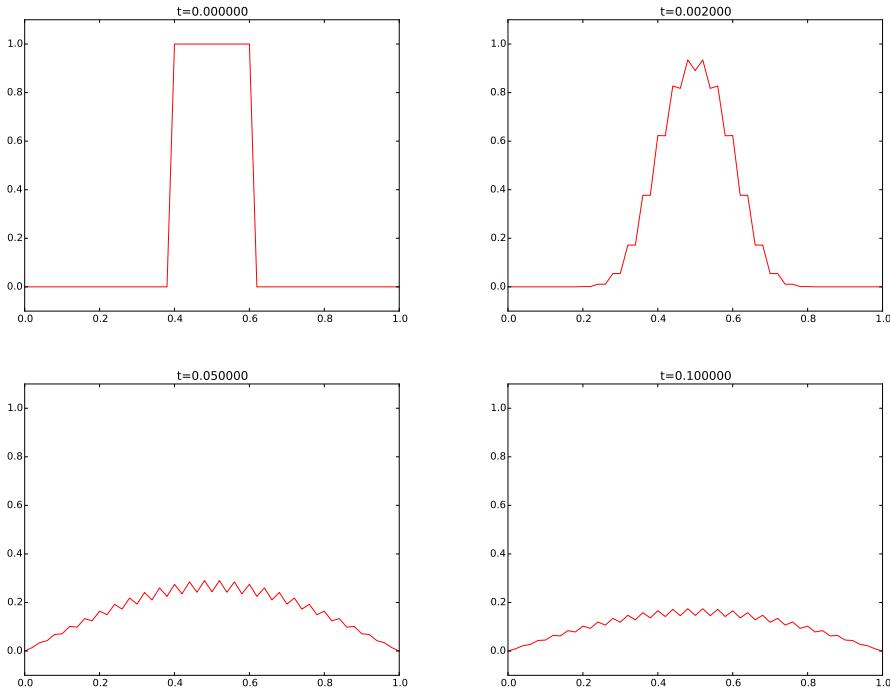
---

The  $u(x, t)$  curve as a function of  $x$  is shown in Figure 3.1 at four time levels (see also a [movie](#)).

We see that the curves have saw-tooth waves in the beginning of the simulation. This non-physical noise is smoothed out with time, but solutions of the diffusion equations are known to be smooth, and this

numerical solution is definitely not smooth. Lowering  $F$  helps:  $F \leq 0.25$  gives a smooth solution, see Figure 3.2 (and a [movie](#)).

Increasing  $F$  slightly beyond the limit 0.5, to  $F = 0.51$ , gives growing, non-physical instabilities, as seen in Figure 3.3.

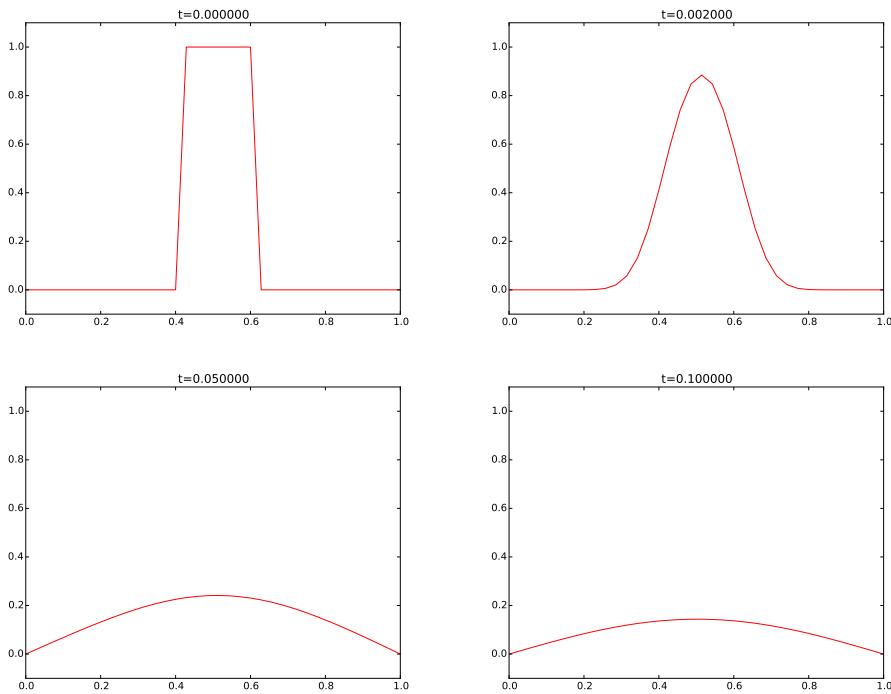


**Fig. 3.1** Forward Euler scheme for  $F = 0.5$ .

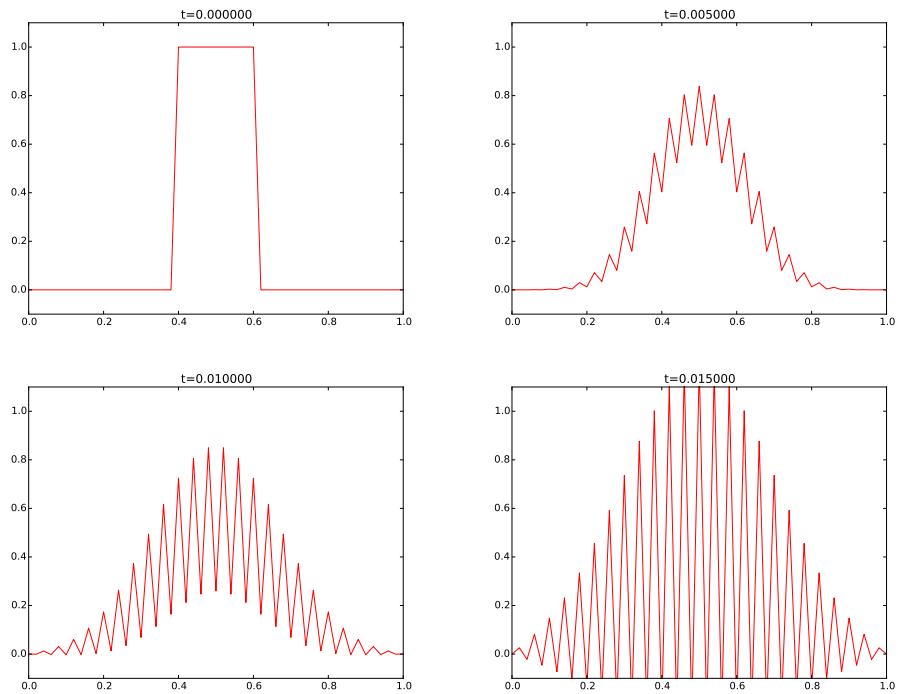
Instead of a discontinuous initial condition we now try the smooth Gaussian function for  $I(x)$ . A simulation for  $F = 0.5$  is shown in Figure 3.4. Now the numerical solution is smooth for all times, and this is true for any  $F \leq 0.5$ .

Experiments with these two choices of  $I(x)$  reveal some important observations:

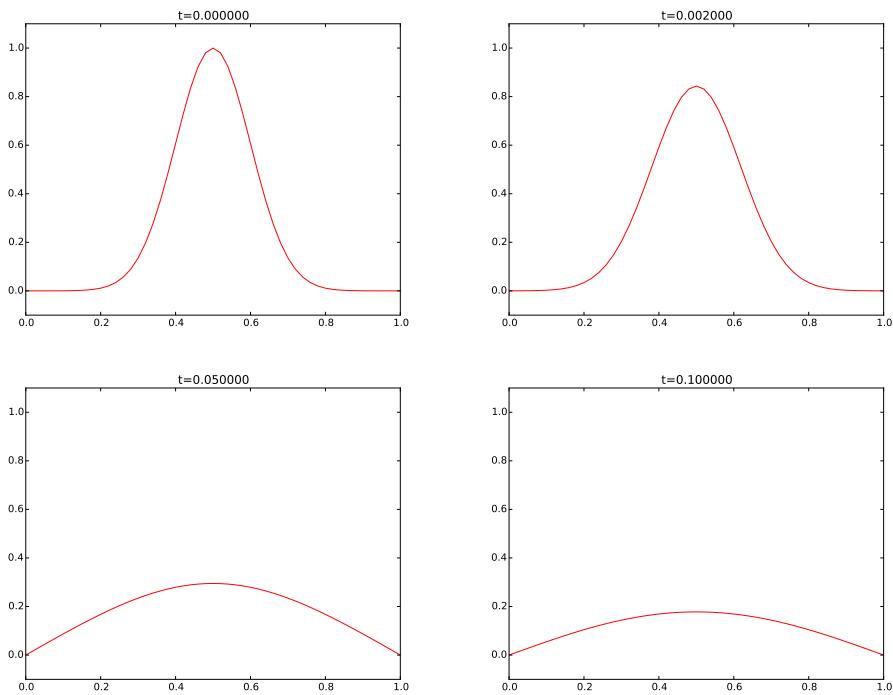
- The Forward Euler scheme leads to growing solutions if  $F > \frac{1}{2}$ .
- $I(x)$  as a discontinuous plug leads to a saw tooth-like noise for  $F = \frac{1}{2}$ , which is absent for  $F \leq \frac{1}{4}$ .
- The smooth Gaussian initial function leads to a smooth solution for all relevant  $F$  values ( $F \geq \frac{1}{2}$ ).



**Fig. 3.2** Forward Euler scheme for  $F = 0.25$ .



**Fig. 3.3** Forward Euler scheme for  $F = 0.51$ .



**Fig. 3.4** Forward Euler scheme for  $F = 0.5$ .

## 3.2 Implicit methods for the 1D diffusion equation

Simulations with the Forward Euler scheme shows that the time step restriction,  $F \leq \frac{1}{2}$ , which means  $\Delta t \leq \Delta x^2/(2\alpha)$ , may be relevant in the beginning of the diffusion process, when the solution changes quite fast, but as time increases, the process slows down, and a small  $\Delta t$  may be inconvenient. By using *implicit schemes*, which lead to a coupled system of linear equations to be solved at each time level, any size of  $\Delta t$  is possible (but the accuracy decreases with increasing  $\Delta t$ ). The Backward Euler scheme, derived and implemented below, is the simplest implicit scheme for the diffusion equation.

### 3.2.1 Backward Euler scheme

We now apply a backward difference in time in (3.5), but the same central difference in space:

$$[D_t^- u = D_x D_x u + f]_i^n, \quad (3.10)$$

which written out reads

$$\frac{u_i^n - u_i^{n-1}}{\Delta t} = \alpha \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2} + f_i^n. \quad (3.11)$$

Now we assume  $u_i^{n-1}$  is already computed, but all quantities at the “new” time level  $n$  are unknown. This time it is not possible to solve with respect to  $u_i^n$  because this value couples to its neighbors in space,  $u_{i-1}^n$  and  $u_{i+1}^n$ , which are also unknown. Let us examine this fact for the case when  $N_x = 3$ . Equation (3.11) written for  $i = 1, \dots, Nx - 1 = 1, 2$  becomes

$$\frac{u_1^n - u_1^{n-1}}{\Delta t} = \alpha \frac{u_2^n - 2u_1^n + u_0^n}{\Delta x^2} + f_1^n \quad (3.12)$$

$$\frac{u_2^n - u_2^{n-1}}{\Delta t} = \alpha \frac{u_3^n - 2u_2^n + u_1^n}{\Delta x^2} + f_2^n \quad (3.13)$$

The boundary values  $u_0^n$  and  $u_3^n$  are known as zero. Collecting the unknown new values  $u_1^n$  and  $u_2^n$  on the left-hand side and multiplying by  $\Delta t$  gives

$$(1 + 2F) u_1^n - Fu_2^n = u_1^{n-1} + \Delta t f_1^n, \quad (3.14)$$

$$-Fu_1^n + (1 + 2F) u_2^n = u_2^{n-1} + \Delta t f_2^n. \quad (3.15)$$

This is a coupled  $2 \times 2$  system of algebraic equations for the unknowns  $u_1^n$  and  $u_2^n$ . The equivalent matrix form is

$$\begin{pmatrix} 1 + 2F & -F \\ -F & 1 + 2F \end{pmatrix} \begin{pmatrix} u_1^n \\ u_2^n \end{pmatrix} = \begin{pmatrix} u_1^{n-1} + \Delta t f_1^n \\ u_2^{n-1} + \Delta t f_2^n \end{pmatrix}$$

### Implicit vs. explicit methods

Discretization methods that lead to a coupled system of equations for the unknown function at a new time level are said to be *implicit methods*. The counterpart, *explicit methods*, refers to discretization methods where there is a simple explicit formula for the values of the unknown function at each of the spatial mesh points at the new time level. From an implementational point of view, implicit methods are more comprehensive to code since they require the solution of coupled equations, i.e., a matrix system, at each time level.

In the general case, (3.11) gives rise to a coupled  $(Nx - 1) \times (Nx - 1)$  system of algebraic equations for all the unknown  $u_i^n$  at the interior spatial points  $i = 1, \dots, Nx - 1$ . Collecting the unknowns on the left-hand side, (3.11) can be written

$$-Fu_{i-1}^n + (1 + 2F) u_i^n - Fu_{i+1}^n = u_{i-1}^{n-1}, \quad (3.16)$$

for  $i = 1, \dots, Nx - 1$ . One can either view these equations as a system for where the  $u_i^n$  values at the internal mesh points,  $i = 1, \dots, N_x - 1$ , are unknown, or we may append the boundary values  $u_0^n$  and  $u_{N_x}^n$  to the system. In the latter case, all  $u_i^n$  for  $i = 0, \dots, N_x$  are unknown and we must add the boundary equations to the  $N_x - 1$  equations in (3.16):

$$u_0^n = 0, \quad (3.17)$$

$$u_{N_x}^n = 0. \quad (3.18)$$

A coupled system of algebraic equations can be written on matrix form, and this is important if we want to call up ready-made software for

solving the system. The equations (3.16) and (3.17)–(3.18) correspond to the matrix equation

$$AU = b$$

where  $U = (u_0^n, \dots, u_{N_x}^n)$ , and the matrix  $A$  has the following structure:

$$A = \begin{pmatrix} A_{0,0} & A_{0,1} & 0 & \cdots & \cdots & \cdots & \cdots & \cdots & 0 \\ A_{1,0} & A_{1,1} & A_{1,2} & \ddots & & & & & \vdots \\ 0 & A_{2,1} & A_{2,2} & A_{2,3} & \ddots & & & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 & & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \ddots & \ddots & & \vdots \\ \vdots & & & 0 & A_{i,i-1} & A_{i,i} & A_{i,i+1} & \ddots & \vdots \\ \vdots & & & & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & & & & & \ddots & \ddots & \ddots & A_{N_x-1,N_x} \\ 0 & \cdots & \cdots & \cdots & \cdots & 0 & A_{N_x,N_x-1} & A_{N_x,N_x} & \end{pmatrix} \quad (3.19)$$

The nonzero elements are given by

$$A_{i,i-1} = -F \quad (3.20)$$

$$A_{i,i} = 1 + 2F \quad (3.21)$$

$$A_{i,i+1} = -F \quad (3.22)$$

for the equations for internal points,  $i = 1, \dots, N_x - 1$ . The first and last equation correspond to the boundary condition, where we know the solution, and therefore we must have

$$A_{0,0} = 1, \quad (3.23)$$

$$A_{0,1} = 0, \quad (3.24)$$

$$A_{N_x,N_x-1} = 0, \quad (3.25)$$

$$A_{N_x,N_x} = 1. \quad (3.26)$$

The right-hand side  $b$  is written as

$$b = \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_i \\ \vdots \\ b_{N_x} \end{pmatrix} \quad (3.27)$$

with

$$b_0 = 0, \quad (3.28)$$

$$b_i = u_i^{n-1}, \quad i = 1, \dots, N_x - 1, \quad (3.29)$$

$$b_{N_x} = 0. \quad (3.30)$$

We observe that the matrix  $A$  contains quantities that do not change in time. Therefore,  $A$  can be formed once and for all before we enter the recursive formulas for the time evolution. The right-hand side  $b$ , however, must be updated at each time step. This leads to the following computational algorithm, here sketched with Python code:

```
x = np.linspace(0, L, Nx+1)    # mesh points in space
dx = x[1] - x[0]
t = np.linspace(0, T, Nt+1)    # mesh points in time
u  = np.zeros(Nx+1)            # unknown u at new time level
u_1 = np.zeros(Nx+1)           # u at the previous time level

# Data structures for the linear system
A = np.zeros((Nx+1, Nx+1))
b = np.zeros(Nx+1)

for i in range(1, Nx):
    A[i,i-1] = -F
    A[i,i+1] = -F
    A[i,i] = 1 + 2*F
A[0,0] = A[Nx,Nx] = 1

# Set initial condition u(x,0) = I(x)
for i in range(0, Nx+1):
    u_1[i] = I(x[i])

import scipy.linalg

for n in range(0, Nt):
    # Compute b and solve linear system
    for i in range(1, Nx):
        b[i] = -u_1[i]
    b[0] = b[Nx] = 0
```

```

u[:] = scipy.linalg.solve(A, b)

# Update u_1 before next step
u_1[:] = u

```

### 3.2.2 Sparse matrix implementation

We have seen from (3.19) that the matrix  $A$  is tridiagonal. The code segment above used a full, dense matrix representation of  $A$ , which stores a lot of values we know are zero beforehand, and worse, the solution algorithm computes with all these zeros. With  $N_x + 1$  unknowns, the work by the solution algorithm is  $\frac{1}{3}(N_x + 1)^3$  and the storage requirements  $(N_x + 1)^2$ . By utilizing the fact that  $A$  is tridiagonal and employing corresponding software tools that work with the three diagonals, the work and storage demands can be proportional to  $N_x$  only. This leads to a dramatic improvement: with  $N_x = 200$ , which is an often needed resolution, the code runs about 40,000 times faster and reduces the storage to just 1.5%! It is no doubt that we need to take advantage of the fact that  $A$  is tridiagonal.

The key idea is to apply a data structure for a tridiagonal or sparse matrix. The `scipy.sparse` package has relevant utilities. For example, we can store the nonzero diagonals of a matrix. The package also has linear system solvers that operate on sparse matrix data structures. The code below illustrates how we can store only the main diagonal and the upper and lower diagonals.

```

# Representation of sparse matrix and right-hand side
main = np.zeros(Nx+1)
lower = np.zeros(Nx)
upper = np.zeros(Nx)
b     = np.zeros(Nx+1)

# Precompute sparse matrix
main[:] = 1 + 2*F
lower[:] = -F
upper[:] = -F
# Insert boundary conditions
main[0] = 1
main[Nx] = 1

A = scipy.sparse.diags(
    diagonals=[main, lower, upper],
    offsets=[0, -1, 1], shape=(Nx+1, Nx+1),
    format='csr')
print A.todense() # Check that A is correct

```

```
# Set initial condition
for i in range(0,Nx+1):
    u_1[i] = I(x[i])

for n in range(0, Nt):
    b = u_1
    b[0] = b[-1] = 0.0 # boundary conditions
    u[:] = scipy.sparse.linalg.spsolve(A, b)
    u_1[:] = u
```

The `scipy.sparse.linalg.spsolve` function utilizes the sparse storage structure of  $A$  and performs in this case a very efficient Gaussian elimination solve.

The program `diffu1D_u0.py` contains a function `solver_BE`, which implements the Backward Euler scheme sketched above. As mentioned in Section 3.1.2, the functions `plug` and `gaussian` runs the case with  $I(x)$  as a discontinuous plug or a smooth Gaussian function. All experiments point to two characteristic features of the Backward Euler scheme: 1) it is always stable, and 2) it always gives a smooth, decaying solution.

### 3.2.3 Crank-Nicolson scheme

The idea in the Crank-Nicolson scheme is to apply centered differences in space and time, combined with an average in time. We demand the PDE to be fulfilled at the spatial mesh points, but midway between the points in the time mesh:

$$\frac{\partial}{\partial t}u(x_i, t_{n+\frac{1}{2}}) = \alpha \frac{\partial^2}{\partial x^2}u(x_i, t_{n+\frac{1}{2}}) + f(x_i, t_{n+\frac{1}{2}}),$$

for  $i = 1, \dots, N_x - 1$  and  $n = 0, \dots, N_t - 1$ .

With centered differences in space and time, we get

$$[D_t u = \alpha D_x D_x u + f]_i^{n+\frac{1}{2}}.$$

On the right-hand side we get an expression

$$\frac{1}{\Delta x^2} \left( u_{i-1}^{n+\frac{1}{2}} - 2u_i^{n+\frac{1}{2}} + u_{i+1}^{n+\frac{1}{2}} \right) + f_i^{n+\frac{1}{2}}.$$

This expression is problematic since  $u_i^{n+\frac{1}{2}}$  is not one of the unknowns we compute. A possibility is to replace  $u_i^{n+\frac{1}{2}}$  by an arithmetic average:

$$u_i^{n+\frac{1}{2}} \approx \frac{1}{2} (u_i^n + u_i^{n+1}) .$$

In the compact notation, we can use the arithmetic average notation  $\bar{u}^t$ :

$$[D_t u = \alpha D_x D_x \bar{u}^t + f]_i^{n+\frac{1}{2}} .$$

We can also use an average for  $f_i^{n+\frac{1}{2}}$ :

$$[D_t u = \alpha D_x D_x \bar{u}^t + \bar{f}^t]_i^{n+\frac{1}{2}} .$$

After writing out the differences and average, multiplying by  $\Delta t$ , and collecting all unknown terms on the left-hand side, we get

$$\begin{aligned} u_i^{n+1} - \frac{1}{2} F(u_{i-1}^{n+1} - 2u_i^{n+1} + u_{i+1}^{n+1}) &= u_i^n + \frac{1}{2} F(u_{i-1}^n - 2u_i^n + u_{i+1}^n) \\ &\quad \frac{1}{2} f_i^{n+1} + \frac{1}{2} f_i^n . \end{aligned} \quad (3.31)$$

Also here, as in the Backward Euler scheme, the new unknowns  $u_{i-1}^{n+1}$ ,  $u_i^{n+1}$ , and  $u_{i+1}^{n+1}$  are coupled in a linear system  $AU = b$ , where  $A$  has the same structure as in (3.19), but with slightly different entries:

$$A_{i,i-1} = -\frac{1}{2} F \quad (3.32)$$

$$A_{i,i} = \frac{1}{2} + F \quad (3.33)$$

$$A_{i,i+1} = -\frac{1}{2} F \quad (3.34)$$

for the equations for internal points,  $i = 1, \dots, N_x - 1$ . The equations for the boundary points correspond to

$$A_{0,0} = 1, \quad (3.35)$$

$$A_{0,1} = 0, \quad (3.36)$$

$$A_{N_x,N_x-1} = 0, \quad (3.37)$$

$$A_{N_x,N_x} = 1. \quad (3.38)$$

The right-hand side  $b$  has entries

$$b_0 = 0, \quad (3.39)$$

$$b_i = u_i^{n-1} + \frac{1}{2}(f_i^n + f_i^{n+1}), \quad i = 1, \dots, N_x - 1, \quad (3.40)$$

$$b_{N_x} = 0. \quad (3.41)$$

### 3.2.4 The unifying $\theta$ rule

For the equation

$$\frac{\partial u}{\partial t} = G(u),$$

where  $G(u)$  is some a spatial differential operator, the  $\theta$ -rule looks like

$$\frac{u_i^{n+1} - u_i^n}{\Delta t} = \theta G(u_i^{n+1}) + (1 - \theta)G(u_i^n).$$

The important feature of this time discretization scheme is that we can implement one formula and then generate a family of well-known and widely used schemes:

- $\theta = 0$  gives the Forward Euler scheme in time
- $\theta = 1$  gives the Backward Euler scheme in time
- $\theta = \frac{1}{2}$  gives the Crank-Nicolson scheme in time

Applied to the 1D diffusion problem, the  $\theta$ -rule gives

$$\begin{aligned} \frac{u_i^{n+1} - u_i^n}{\Delta t} &= \alpha \left( \theta \frac{u_{i+1}^{n+1} - 2u_i^{n+1} + u_{i-1}^{n+1}}{\Delta x^2} + (1 - \theta) \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2} \right) \\ &\quad + \theta f_i^{n+1} + (1 - \theta) f_i^n. \end{aligned}$$

This scheme also leads to a matrix system with entries

$$A_{i,i-1} = -F\theta, \quad A_{i,i} = 1 + 2F\theta, \quad A_{i,i+1} = -F\theta,$$

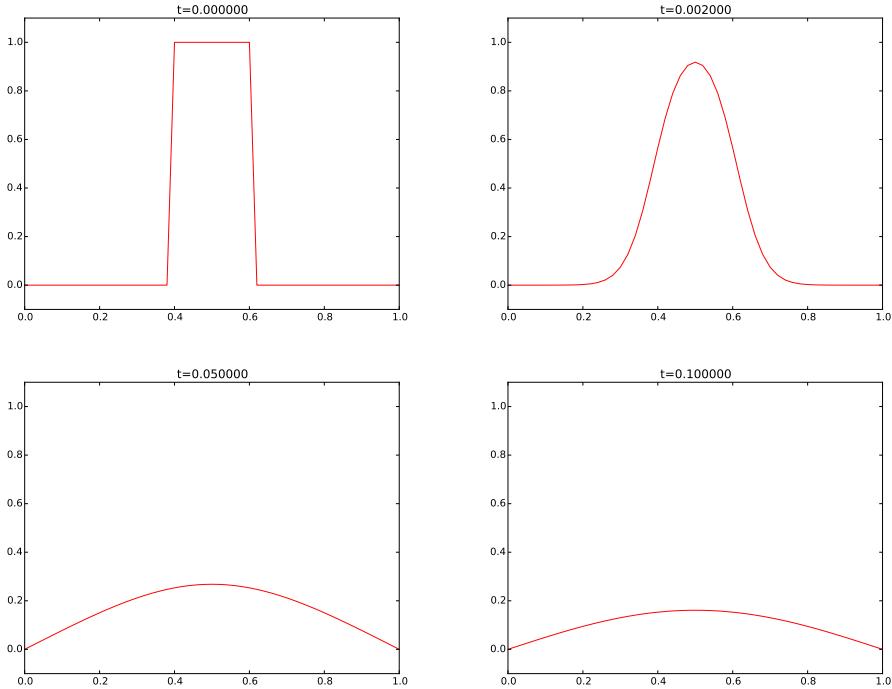
while right-hand side entry  $b_i$  is

$$b_i = u_i^n + F(1 - \theta) \frac{u_{i+1}^n - 2u_i^n + u_{i-1}^n}{\Delta x^2} + \Delta t \theta f_i^{n+1} + \Delta t (1 - \theta) f_i^n.$$

The corresponding entries for the boundary points are as in the Backward Euler and Crank-Nicolson schemes listed earlier.

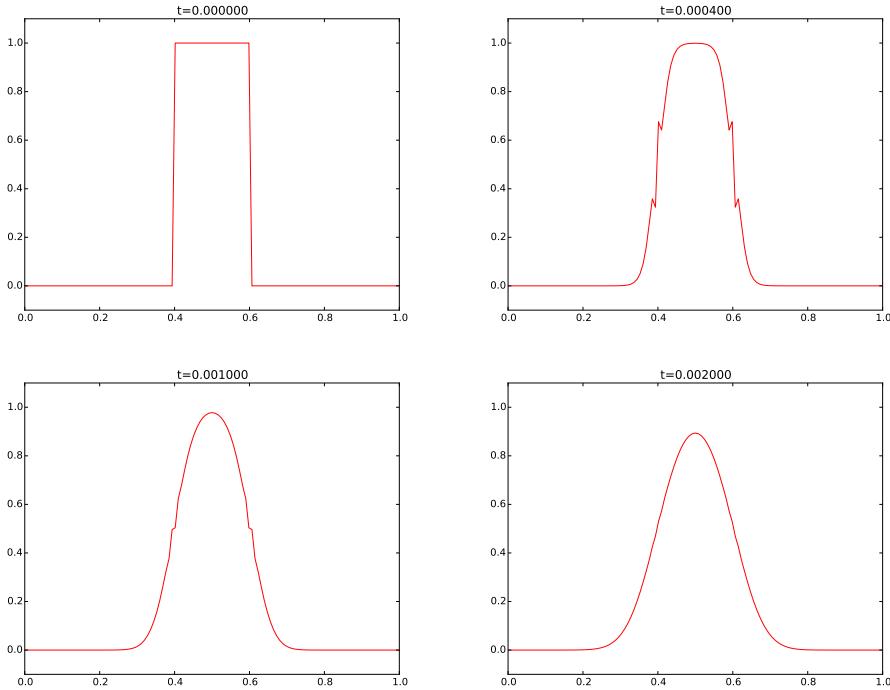
### 3.2.5 Experiments

We can repeat the experiments from Section 3.1.5 to see if the Backward Euler or Crank-Nicolson schemes have problems with sawtooth-like noise when starting with a discontinuous initial condition. We can also verify that we can have  $F > \frac{1}{2}$ , which allows larger time steps than in the Forward Euler method.



**Fig. 3.5** Backward Euler scheme for  $F = 0.5$ .

The Backward Euler scheme always produces smooth solutions for any  $F$ . Figure 3.5 shows one example. The Crank-Nicolson method produces smooth solutions for small  $F$ ,  $F \leq \frac{1}{2}$ , but small noise is more and more evident as  $F$  increases. Figures 3.6 and 3.7 demonstrate the effect for  $F = 3$  and  $F = 10$ , respectively. Section 3.3 explains why such noise occur.



**Fig. 3.6** Crank-Nicolson scheme for  $F = 3$ .

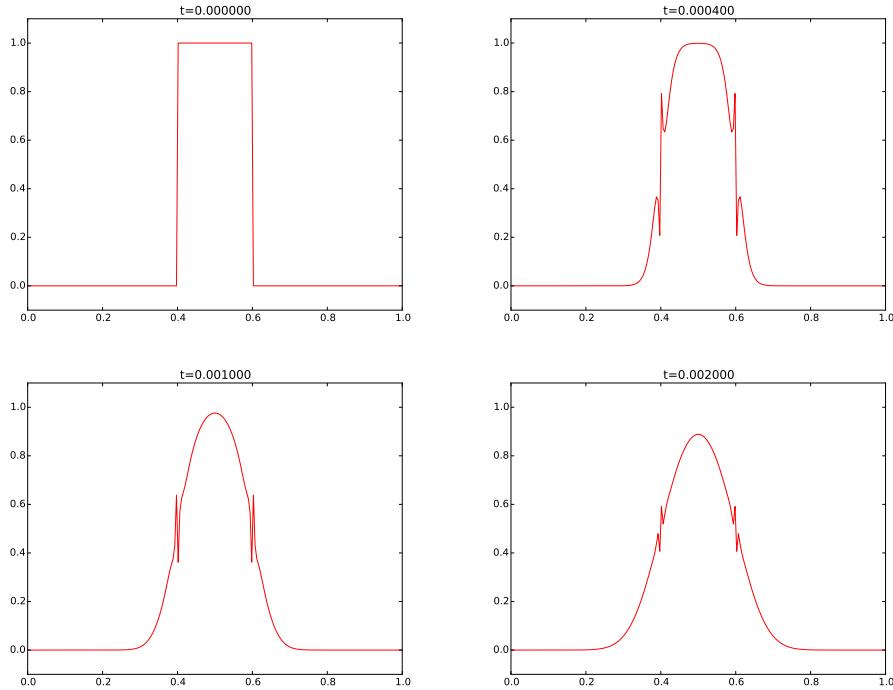
### 3.2.6 The Laplace and Poisson equation

The Laplace equation,  $\nabla^2 u = 0$ , and the Poisson equation,  $-\nabla^2 u = f$ , occur in numerous applications throughout science and engineering. In 1D these equations read  $u''(x) = 0$  and  $-u''(x) = f(x)$ , respectively. We can solve 1D variants of the Laplace equations with the listed software, because we can interpret  $u_{xx} = 0$  as the limiting solution of  $u_t = \alpha u_{xx}$  when  $u$  reaches a steady state limit where  $u_t \rightarrow 0$ . Similarly, Poisson's equation  $-u_{xx} = f$  arises from solving  $u_t = u_{xx} + f$  and letting  $t \rightarrow \infty$  so  $u_t \rightarrow 0$ .

Technically in a program, we can simulate  $t \rightarrow \infty$  by just taking one large time step:  $\Delta t \rightarrow \infty$ . In the limit, the Backward Euler scheme gives

$$-\frac{u_{i+1}^{n+1} - 2u_i^{n+1} + u_{i-1}^{n+1}}{\Delta x^2} = f_i^{n+1},$$

which is nothing but the discretization  $[-D_x D_x u = f]_i^{n+1} = 0$  of  $-u_{xx} = f$ .



**Fig. 3.7** Crank-Nicolson scheme for  $F = 10$ .

The result above means that the Backward Euler scheme can solve the limit equation directly and hence produce a solution of the 1D Laplace equation. With the Forward Euler scheme we must do the time stepping since  $\Delta t > \Delta x^2/\alpha$  is illegal and leads to instability. We may interpret this time stepping as solving the equation system from  $-u_{xx} = f$  by iterating on a pseudo time variable.

**hpl 12:** Better to say the last sentence when we treat iterative methods.

### 3.3 Analysis of schemes for the diffusion equation

The numerical experiments in Sections 3.1.5 and 3.2.5 reveal that there are some numerical problems with the Forward Euler and Crank-Nicolson schemes: sawtooth-like noise is sometimes present in solutions that are, from a mathematical point of view, expected to be smooth. This section presents a mathematical analysis that explains the observed behavior and arrives at criteria for obtaining numerical solutions that reproduce

the qualitative properties of the exact solutions. In short, we shall explain what is observed in Figures 3.1, 3.2, 3.3, 3.4, 3.5, 3.6, and 3.7.

### 3.3.1 Properties of the solution

A particular characteristic of diffusive processes, governed by an equation like

$$u_t = \alpha u_{xx}, \quad (3.42)$$

is that the initial shape  $u(x, 0) = I(x)$  spreads out in space with time, along with a decaying amplitude. Three different examples will illustrate the spreading of  $u$  in space and the decay in time.

**Similarity solution.** The diffusion equation (3.42) admits solutions that depend on  $\eta = (x - c)/\sqrt{4\alpha t}$  for a given value of  $c$ . One particular solution is

$$u(x, t) = a \operatorname{erf}(\eta) + b, \quad (3.43)$$

where

$$\operatorname{erf}(\eta) = \frac{2}{\sqrt{\pi}} \int_0^\eta e^{-\zeta^2} d\zeta, \quad (3.44)$$

is the *error function*, and  $a$  and  $b$  are arbitrary constants. The error function lies in  $(-1, 1)$ , is odd around  $\eta = 0$ , and goes relatively quickly to  $\pm 1$ :

$$\begin{aligned} \lim_{\eta \rightarrow -\infty} \operatorname{erf}(\eta) &= -1, \\ \lim_{\eta \rightarrow \infty} \operatorname{erf}(\eta) &= 1, \\ \operatorname{erf}(\eta) &= -\operatorname{erf}(-\eta), \\ \operatorname{erf}(0) &= 0, \\ \operatorname{erf}(2) &= 0.99532227, \\ \operatorname{erf}(3) &= 0.99997791. \end{aligned}$$

As  $t \rightarrow 0$ , the error function approaches a step function centered at  $x = c$ . For a diffusion problem posed on the unit interval  $[0, 1]$ , we may choose the step at  $x = 1/2$  (meaning  $c = 1/2$ ),  $a = -1/2$ ,  $b = 1/2$ . Then

$$u(x, t) = \frac{1}{2} \left( 1 - \operatorname{erf} \left( \frac{x - \frac{1}{2}}{\sqrt{4\alpha t}} \right) \right) = \frac{1}{2} \operatorname{erfc} \left( \frac{x - \frac{1}{2}}{\sqrt{4\alpha t}} \right), \quad (3.45)$$

where we have introduced the *complementary error function*  $\operatorname{erfc}(\eta) = 1 - \operatorname{erf}(\eta)$ . The solution (3.45) implies the boundary conditions

$$u(0, t) = \frac{1}{2} \left( 1 - \operatorname{erf} \left( \frac{-1/2}{\sqrt{4\alpha t}} \right) \right), \quad (3.46)$$

$$u(1, t) = \frac{1}{2} \left( 1 - \operatorname{erf} \left( \frac{1/2}{\sqrt{4\alpha t}} \right) \right). \quad (3.47)$$

For small enough  $t$ ,  $u(0, t) \approx 1$  and  $u(1, t) \approx 1$ , but as  $t \rightarrow \infty$ ,  $u(x, t) \rightarrow 1/2$  on  $[0, 1]$ .

**Solution for a Gaussian pulse.** The standard diffusion equation  $u_t = \alpha u_{xx}$  admits a Gaussian function as solution:

$$u(x, t) = \frac{1}{\sqrt{4\pi\alpha t}} \exp \left( -\frac{(x - c)^2}{4\alpha t} \right). \quad (3.48)$$

At  $t = 0$  this is a Dirac delta function, so for computational purposes one must start to view the solution at some time  $t = t_\epsilon > 0$ . Replacing  $t$  by  $t_\epsilon + t$  in (3.48) makes it easy to operate with a (new)  $t$  that starts at  $t = 0$  with an initial condition with a finite width. The important feature of (3.48) is that the standard deviation  $\sigma$  of a sharp initial Gaussian pulse increases in time according to  $\sigma = \sqrt{2\alpha t}$ , making the pulse diffuse and flatten out.

**Solution for a sine component.** For example, (3.42) admits a solution of the form

$$u(x, t) = Q e^{-at} \sin(kx). \quad (3.49)$$

The parameters  $Q$  and  $k$  can be freely chosen, while inserting (3.49) in (3.42) gives the constraint

$$a = -\alpha k^2.$$

A very important feature is that the initial shape  $I(x) = Q \sin kx$  undergoes a damping  $\exp(-\alpha k^2 t)$ , meaning that rapid oscillations in space, corresponding to large  $k$ , are very much faster damped than slow oscillations in space, corresponding to small  $k$ . This feature leads to a smoothing of the initial condition with time.

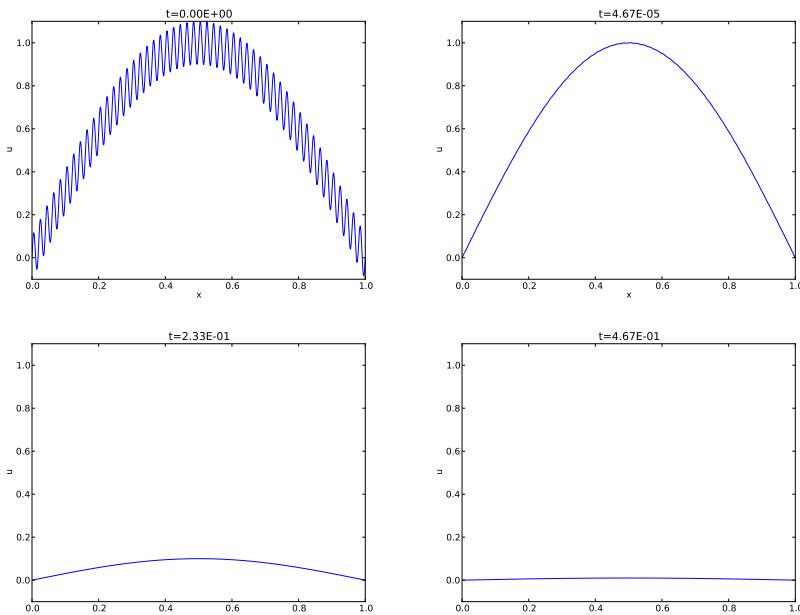
The following examples illustrates the damping properties of (3.49). We consider the specific problem

$$\begin{aligned} u_t &= u_{xx}, \quad x \in (0, 1), \quad t \in (0, T], \\ u(0, t) &= u(1, t) = 0, \quad t \in (0, T], \\ u(x, 0) &= \sin(\pi x) + 0.1 \sin(100\pi x). \end{aligned}$$

The initial condition has been chosen such that adding two solutions like (3.49) constructs an analytical solution to the problem:

$$u(x, t) = e^{-\pi^2 t} \sin(\pi x) + 0.1e^{-\pi^2 10^4 t} \sin(100\pi x). \quad (3.50)$$

Figure 3.8 illustrates the rapid damping of rapid oscillations  $\sin(100\pi x)$  and the very much slower damping of the slowly varying  $\sin(\pi x)$  term. After about  $t = 0.5 \cdot 10^{-4}$  the rapid oscillations do not have a visible amplitude, while we have to wait until  $t \sim 0.5$  before the amplitude of the long wave  $\sin(\pi x)$  becomes very small.



**Fig. 3.8** Evolution of the solution of a diffusion problem: initial condition (upper left),  $1/100$  reduction of the small waves (upper right),  $1/10$  reduction of the long wave (lower left), and  $1/100$  reduction of the long wave (lower right).

### 3.3.2 Analysis of discrete equations

A counterpart to (3.49) is the complex representation of the same function:

$$u(x, t) = Q e^{-at} e^{ikx},$$

where  $i = \sqrt{-1}$  is the imaginary unit. We can add such functions, often referred to as wave components, to make a Fourier representation of a general solution of the diffusion equation:

$$u(x, t) \approx \sum_{k \in K} b_k e^{-\alpha k^2 t} e^{ikx}, \quad (3.51)$$

where  $K$  is a set of an infinite number of  $k$  values needed to construct the solution. In practice, however, the series is truncated and  $K$  is a finite set of  $k$  values needed to build a good approximate solution. Note that (3.50) is a special case of (3.51) where  $K = \{\pi, 100\pi\}$ ,  $b_\pi = 1$ , and  $b_{100\pi} = 0.1$ .

The amplitudes  $b_k$  of the individual Fourier waves must be determined from the initial condition. At  $t = 0$  we have  $u \approx \sum_k b_k \exp(ikx)$  and find  $K$  and  $b_k$  such that

$$I(x) \approx \sum_{k \in K} b_k e^{ikx}. \quad (3.52)$$

(The relevant formulas for  $b_k$  come from Fourier analysis, or equivalently, a least-squares method for approximating  $I(x)$  in a function space with basis  $\exp(ikx)$ .)

Much insight about the behavior of numerical methods can be obtained by investigating how a wave component  $\exp(-\alpha k^2 t) \exp(ikx)$  is treated by the numerical scheme. It appears that such wave components are also solutions of the schemes, but the damping factor  $\exp(-\alpha k^2 t)$  varies among the schemes. To ease the forthcoming algebra, we write the damping factor as  $A^n$ . The exact amplification factor corresponding to  $A$  is  $A_e = \exp(-\alpha k^2 \Delta t)$ .

### 3.3.3 Analysis of the finite difference schemes

We have seen that a general solution of the diffusion equation can be built as a linear combination of basic components

$$e^{-\alpha k^2 t} e^{ikx}.$$

A fundamental question is whether such components are also solutions of the finite difference schemes. This is indeed the case, but the amplitude  $\exp(-\alpha k^2 t)$  might be modified (which also happens when solving the ODE counterpart  $u' = -\alpha u$ ). We therefore look for numerical solutions of the form

$$u_q^n = A^n e^{ikq \Delta x} = A^n e^{ikx}, \quad (3.53)$$

where the amplification factor  $A$  must be determined by inserting the component into an actual scheme.

**Stability.** The exact amplification factor is  $A_e = \exp(-\alpha^2 k^2 \Delta t)$ . We should therefore require  $|A| < 1$  to have a decaying numerical solution as well. If  $-1 \leq A < 0$ ,  $A^n$  will change sign from time level to time level, and we get stable, non-physical oscillations in the numerical solutions that are not present in the exact solution.

**Accuracy.** To determine how accurately a finite difference scheme treats one wave component (3.53), we see that the basic deviation from the exact solution is reflected in how well  $A^n$  approximates  $A_e^n$ , or how well  $A$  approximates  $A_e$ . We can plot  $A_e$  and the various expressions for  $A$ , and we can make Taylor expansions of  $A/A_e$  to see the error more analytically.

### 3.3.4 Analysis of the Forward Euler scheme

The Forward Euler finite difference scheme for  $u_t = \alpha u_{xx}$  can be written as

$$[D_t^+ u = \alpha D_x D_x u]_q^n.$$

Inserting a wave component (3.53) in the scheme demands calculating the terms

$$e^{ikq \Delta x} [D_t^+ A]^n = e^{ikq \Delta x} A^n \frac{A - 1}{\Delta t},$$

and

$$A^n D_x D_x [e^{ikx}]_q = A^n \left( -e^{ikq \Delta x} \frac{4}{\Delta x^2} \sin^2 \left( \frac{k \Delta x}{2} \right) \right).$$

Inserting these terms in the discrete equation and dividing by  $A^n e^{ikq \Delta x}$  leads to

$$\frac{A - 1}{\Delta t} = -\alpha \frac{4}{\Delta x^2} \sin^2 \left( \frac{k \Delta x}{2} \right),$$

and consequently

$$A = 1 - 4F \sin^2 p \quad (3.54)$$

where

$$F = \frac{\alpha \Delta t}{\Delta x^2} \quad (3.55)$$

is the *numerical Fourier number*, and  $p = k \Delta x / 2$ . The complete numerical solution is then

$$u_q^n = \left( 1 - 4F \sin^2 p \right)^n e^{ikq \Delta x}. \quad (3.56)$$

**Stability.** We easily see that  $A \leq 1$ . However, the  $A$  can be less than  $-1$ , which will lead to growth of a numerical wave component. The criterion  $A \geq -1$  implies

$$4F \sin^2(p/2) \leq 2.$$

The worst case is when  $\sin^2(p/2) = 1$ , so a sufficient criterion for stability is

$$F \leq \frac{1}{2}, \quad (3.57)$$

or expressed as a condition on  $\Delta t$ :

$$\Delta t \leq \frac{\Delta x^2}{2\alpha}. \quad (3.58)$$

Note that halving the spatial mesh size,  $\Delta x \rightarrow \frac{1}{2}\Delta x$ , requires  $\Delta t$  to be reduced by a factor of  $1/4$ . The method hence becomes very expensive for fine spatial meshes.

**Accuracy.** Since  $A$  is expressed in terms of  $F$  and the parameter we now call  $p = k \Delta x / 2$ , we should also express  $A_e$  by  $F$  and  $p$ . The exponent in  $A_e$  is  $-\alpha k^2 \Delta t$ , which equals  $-F k^2 \Delta x^2 = -F 4p^2$ . Consequently,

$$A_e = \exp(-\alpha k^2 \Delta t) = \exp(-4F p^2).$$

All our  $A$  expressions as well as  $A_e$  are now functions of the two dimensionless parameters  $F$  and  $p$ .

Computing the Taylor series expansion of  $A/A_e$  in terms of  $F$  can easily be done with aid of `sympy`:

```
def A_exact(F, p):
    return exp(-4*F*p**2)

def A_FE(F, p):
    return 1 - 4*F*sin(p)**2

from sympy import *
F, p = symbols('F p')
A_err_FE = A_FE(F, p)/A_exact(F, p)
print A_err_FE.series(F, 0, 6)
```

The result is

$$\frac{A}{A_e} = 1 - 4F \sin^2 p + 2Fp^2 - 16F^2 p^2 \sin^2 p + 8F^2 p^4 + \dots$$

Recalling that  $F = \alpha \Delta t / \Delta x$ ,  $p = k \Delta x / 2$ , and that  $\sin^2 p \leq 1$ , we realize that the dominating terms in  $A/A_e$  are at most

$$1 - 4\alpha \frac{\Delta t}{\Delta x^2} + \alpha \Delta t - 4\alpha^2 \Delta t^2 + \alpha^2 \Delta t^2 \Delta x^2 + \dots.$$

### 3.3.5 Analysis of the Backward Euler scheme

Discretizing  $u_t = \alpha u_{xx}$  by a Backward Euler scheme,

$$[D_t^- u = \alpha D_x D_x u]_q^n,$$

and inserting a wave component (3.53), leads to calculations similar to those arising from the Forward Euler scheme, but since

$$e^{ikq \Delta x} [D_t^- A]^n = A^n e^{ikq \Delta x} \frac{1 - A^{-1}}{\Delta t},$$

we get

$$\frac{1 - A^{-1}}{\Delta t} = -\alpha \frac{4}{\Delta x^2} \sin^2 \left( \frac{k \Delta x}{2} \right),$$

and then

$$A = \left( 1 + 4F \sin^2 p \right)^{-1}. \quad (3.59)$$

The complete numerical solution can be written

$$u_q^n = \left(1 + 4F \sin^2 p\right)^{-n} e^{ikq\Delta x}. \quad (3.60)$$

**Stability.** We see from (3.59) that  $0 < A < 1$ , which means that all numerical wave components are stable and non-oscillatory for any  $\Delta t > 0$ .

### 3.3.6 Analysis of the Crank-Nicolson scheme

The Crank-Nicolson scheme can be written as

$$[D_t u]_q = \alpha [D_x D_x \bar{u}^x]_q^{n+\frac{1}{2}},$$

or

$$[D_t u]_q^{n+\frac{1}{2}} = \frac{1}{2} \alpha \left( [D_x D_x u]_q^n + [D_x D_x u]_q^{n+1} \right).$$

Inserting (3.53) in the time derivative approximation leads to

$$[D_t A^n e^{ikq\Delta x}]^{n+\frac{1}{2}} = A^{n+\frac{1}{2}} e^{ikq\Delta x} \frac{A^{\frac{1}{2}} - A^{-\frac{1}{2}}}{\Delta t} = A^n e^{ikq\Delta x} \frac{A - 1}{\Delta t}.$$

Inserting (3.53) in the other terms and dividing by  $A^n e^{ikq\Delta x}$  gives the relation

$$\frac{A - 1}{\Delta t} = -\frac{1}{2} \alpha \frac{4}{\Delta x^2} \sin^2 \left( \frac{k\Delta x}{2} \right) (1 + A),$$

and after some more algebra,

$$A = \frac{1 - 2F \sin^2 p}{1 + 2F \sin^2 p}. \quad (3.61)$$

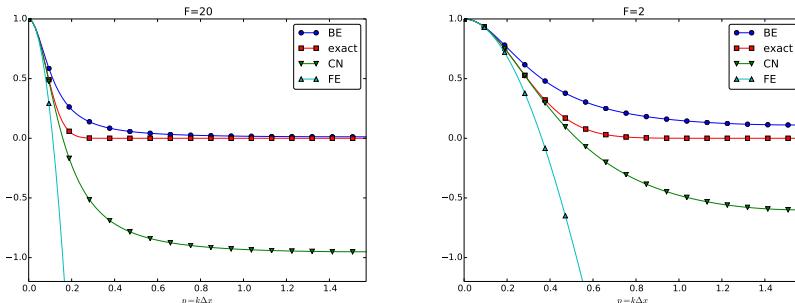
The exact numerical solution is hence

$$u_q^n = \left( \frac{1 - 2F \sin^2 p}{1 + 2F \sin^2 p} \right)^n e^{ikp\Delta x}. \quad (3.62)$$

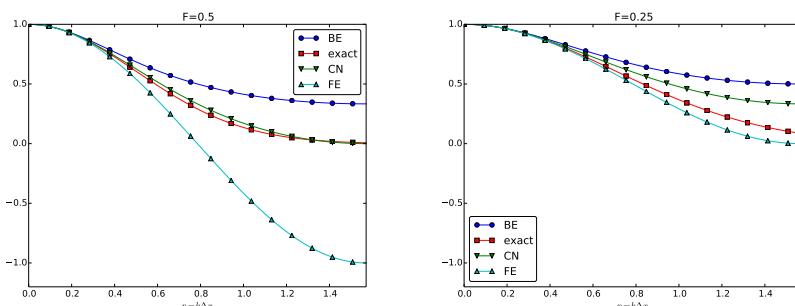
**Stability.** The criteria  $A > -1$  and  $A < 1$  are fulfilled for any  $\Delta t > 0$ . Therefore, the solution cannot grow, but it will oscillate if  $1 - 2F \sin^2 p < 0$ . To avoid such non-physical oscillations, we must demand  $F \leq \frac{1}{2}$ .

### 3.3.7 Summary of accuracy of amplification factors

We can plot the various amplification factors against  $p = k\Delta x/2$  for different choices of the  $F$  parameter. Figures 3.9, 3.10, and 3.11 show how long and small waves are damped by the various schemes compared to the exact damping. As long as all schemes are stable, the amplification factor is positive, except for Crank-Nicolson when  $F > 0.5$ .

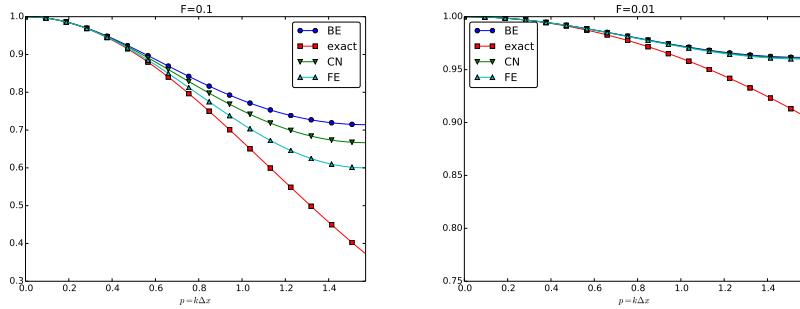


**Fig. 3.9** Amplification factors for large time steps.



**Fig. 3.10** Amplification factors for time steps around the Forward Euler stability limit.

The effect of negative amplification factors is that  $A^n$  changes sign from one time level to the next, thereby giving rise to oscillations in time in an animation of the solution. We see from Figure 3.9 that for  $F = 20$ , waves with  $p \geq \pi/2$  undergo a damping close to  $-1$ , which means that the amplitude does not decay and that the wave component jumps up and down (flips amplitude) in time. For  $F = 2$  we have a damping of a factor of 0.5 from one time level to the next, which is very much smaller than the exact damping. Short waves will therefore fail to be effectively damped. These waves will manifest themselves as high frequency oscillatory noise in the solution.



**Fig. 3.11** Amplification factors for small time steps.

A value  $p = \pi/4$  corresponds to four mesh points per wave length of  $e^{ikx}$ , while  $p = \pi/2$  implies only two points per wave length, which is the smallest number of points we can have to represent the wave on the mesh.

To demonstrate the oscillatory behavior of the Crank-Nicolson scheme, we choose an initial condition that leads to short waves with significant amplitude. A discontinuous  $I(x)$  will in particular serve this purpose.

**hpl 13:** Run example!!

### 3.3.8 Analysis of the 2D diffusion equation

We first consider the 2D diffusion equation

$$u_t = \alpha(u_{xx} + u_{yy}),$$

which has Fourier component solutions of the form

$$u(x, y, t) = Ae^{-\alpha k^2 t} e^{i(k_x x + k_y y)},$$

and the schemes have discrete versions of this Fourier component:

$$u_{q,r}^n = A\xi^n e^{i(k_x q \Delta x + k_y r \Delta y)}.$$

**The Forward Euler scheme.** For the Forward Euler discretization,

$$[D_t^+ u = \alpha(D_x D_x u + D_y D_y u)]_{i,j}^n,$$

we get

$$\frac{\xi - 1}{\Delta t} = -\alpha \frac{4}{\Delta x^2} \sin^2 \left( \frac{k_x \Delta x}{2} \right) - \alpha \frac{4}{\Delta y^2} \sin^2 \left( \frac{k_y \Delta y}{2} \right).$$

Introducing

$$p_x = \frac{k_x \Delta x}{2}, \quad p_y = \frac{k_y \Delta y}{2},$$

we can write the equation for  $\xi$  more compactly as

$$\frac{\xi - 1}{\Delta t} = -\alpha \frac{4}{\Delta x^2} \sin^2 p_x - \alpha \frac{4}{\Delta y^2} \sin^2 p_y,$$

and solve for  $\xi$ :

$$\xi = 1 - 4F_x \sin^2 p_x - 4F_y \sin^2 p_y. \quad (3.63)$$

The complete numerical solution for a wave component is

$$u_{q,r}^n = A(1 - 4F_x \sin^2 p_x - 4F_y \sin^2 p_y)^n e^{i(k_x p \Delta x + k_y q \Delta y)}. \quad (3.64)$$

For stability we demand  $-1 \leq \xi \leq 1$ , and  $-1 \leq \xi$  is the critical limit, since clearly  $\xi \leq 1$ , and the worst case happens when the sines are at their maximum. The stability criterion becomes

$$F_x + F_y \leq \frac{1}{2}. \quad (3.65)$$

For the special, yet common, case  $\Delta x = \Delta y = h$ , the stability criterion can be written as

$$\Delta t \leq \frac{h^2}{2d\alpha},$$

where  $d$  is the number of space dimensions:  $d = 1, 2, 3$ .

**The Backward Euler scheme.** The Backward Euler method,

$$[D_t^- u = \alpha(D_x D_x u + D_y D_y u)]_{i,j}^n,$$

results in

$$1 - \xi^{-1} = -4F_x \sin^2 p_x - 4F_y \sin^2 p_y,$$

and

$$\xi = (1 + 4F_x \sin^2 p_x + 4F_y \sin^2 p_y)^{-1},$$

which is always in  $(0, 1]$ . The solution for a wave component becomes

$$u_{q,r}^n = A(1 + 4F_x \sin^2 p_x + 4F_y \sin^2 p_y)^{-n} e^{i(k_x q \Delta x + k_y r \Delta y)}. \quad (3.66)$$

**The Crank-Nicolson scheme.** With a Crank-Nicolson discretization,

$$[D_t u]_{i,j}^{n+\frac{1}{2}} = \frac{1}{2}[\alpha(D_x D_x u + D_y D_y u)]_{i,j}^{n+1} + \frac{1}{2}[\alpha(D_x D_x u + D_y D_y u)]_{i,j}^n,$$

we have, after some algebra,

$$\xi = \frac{1 - 2(F_x \sin^2 p_x + F_x \sin^2 p_y)}{1 + 2(F_x \sin^2 p_x + F_x \sin^2 p_y)}.$$

The fraction on the right-hand side is always less than 1, so stability in the sense of non-growing wave components is guaranteed for all physical and numerical parameters. However, the fraction can become negative and result in non-physical oscillations. This phenomenon happens when

$$F_x \sin^2 p_x + F_x \sin^2 p_y > \frac{1}{2}.$$

A criterion against non-physical oscillations is therefore

$$F_x + F_y \leq \frac{1}{2},$$

which is the same limit as the stability criterion for the Forward Euler scheme.

The exact discrete solution is

$$u_{q,r}^n = A \left( \frac{1 - 2(F_x \sin^2 p_x + F_x \sin^2 p_y)}{1 + 2(F_x \sin^2 p_x + F_x \sin^2 p_y)} \right)^n e^{i(k_x q \Delta x + k_y r \Delta y)}. \quad (3.67)$$

### 3.3.9 Explanation of numerical artifacts

The behavior of the Forward Euler discretization in time (and centered differences in space) is summarized at the end of Section 3.1.5. Can we from the analysis above explain the behavior?

We may start by looking at Figure 3.3 where  $F = 0.51$ . The figure shows that the solution is unstable and grows in time. The stability limit for such growth is  $F = 0.5$  and since the  $F$  in this simulation is slightly larger, growth is unavoidable.

Figure 3.1 has unexpected features: we would expect the solution of the diffusion equation to be smooth, but the graphs in Figure 3.1 contain non-smooth noise. Turning to Figure 3.4, which has a quite similar initial condition, we see that the curves are indeed smooth. The problem with the results in Figure 3.1 is that the initial condition is discontinuous. To represent it, we need a significant amplitude on the shortest waves in the mesh. However, for  $F = 0.5$ , the shortest wave ( $p = \pi/2$ ) gives the amplitude in the numerical solution as  $(1 - 4F)^n$ , which oscillates between negative and positive values at subsequent time levels for  $F > \frac{1}{4}$ . Since the shortest waves have visible amplitudes in the solution profile, the oscillations becomes visible. The smooth initial condition in Figure 3.4, on the other hand, lead to very small amplitudes of the shortest waves. That these waves then oscillate in a non-physical way for  $F = 0.5$  is not a visible effect. The oscillations in time in the amplitude  $(1 - 4F)^n$  disappear for  $F \leq \frac{1}{4}$ , and that is why also the discontinuous initial condition leads to always smooth solutions in Figure refrefdiffu:pde1:FE:fig:F=0.25, where  $F = \frac{1}{4}$ .

Turning the attention to the Backward Euler scheme and the experiments in Figure 3.5, we see that even the discontinuous initial condition gives smooth solutions for  $F = 0.5$  (and in fact all other  $F$  values). From the exact expression of the numerical amplitude,  $(1 + 4F \sin^2 p)^{-1}$ , we realize that this factor can never flip between positive and negative values, and no instabilities can occur. The conclusion is that the Backward Euler scheme always produces smooth solutions. Also, the Backward Euler scheme guarantees that the solution cannot grow in time (unless we add a source term to the PDE, but that is meant to represent a physically relevant growth).

Finally, we have some small, strange artifacts when simulating the development of the initial plug profile with the Crank-Nicolson scheme, see Figure 3.7, where  $F = 3$ . The Crank-Nicolson scheme cannot give growing amplitudes, but it may give oscillating amplitudes in time. The critical factor is  $1 - 2F \sin^2 p$ , which for the shortest waves ( $p = \pi/2$ ) indicates a stability limit  $F = 0.5$ . With the discontinuous initial condition, we have enough amplitude on the shortest waves so their wrong behavior is visible, and this is what we see as small instabilities in Figure 3.7. The only remedy is to lower the  $F$  value.

## 3.4 Exercises

### Exercise 3.1: Explore symmetry in a 1D problem

This exercise simulates the exact solution (3.48). Suppose for simplicity that  $c = 0$ .

- a) Formulate an initial-boundary value problem that has (3.48) as solution in the domain  $[-L, L]$ . Use the exact solution (3.48) as Dirichlet condition at the boundaries. Simulate the diffusion of the Gaussian peak. Observe that the solution is symmetric around  $x = 0$ .
- b) Show from (3.48) that  $u_x(c, t) = 0$ . Since the solution is symmetric around  $x = c = 0$ , we can solve the numerical problem in half of the domain, using a *symmetry boundary condition*  $u_x = 0$  at  $x = 0$ . Set up the initial-boundary value problem in this case. Simulate the diffusion problem in  $[0, L]$  and compare with the solution in a).

Filename: `diffu_symmetric_gaussian`.

### Exercise 3.2: Investigate approximation errors from a $u_x = 0$ boundary condition

We consider the problem solved in Exercise 3.1 part b). The boundary condition  $u_x(0, t) = 0$  can be implemented in two ways: 1) by a standard symmetric finite difference  $[D_{2x}u]_i^n = 0$ , or 2) by a one-sided difference  $[D^+u = 0]_i^n = 0$ . Investigate the effect of these two conditions on the convergence rate in space.

**Hint.** If you use a Forward Euler scheme, choose a discretization parameter  $h = \Delta t = \Delta x^2$  and assume the error goes like  $E \sim h^r$ . The error in the scheme is  $\mathcal{O}(\Delta t, \Delta x^2)$  so one should expect that the estimated  $r$  approaches 1. The question is if a one-sided difference approximation to  $u_x(0, t) = 0$  destroys this convergence rate.

Filename: `diffu_onesided_fd`.

### Exercise 3.3: Experiment with open boundary conditions in 1D

We address diffusion of a Gaussian function as in Exercise 3.1, in the domain  $[0, L]$ , but now we shall explore different types of boundary conditions on  $x = L$ . In real-life problems we do not know the exact solution on  $x = L$  and must use something simpler.

**a)** Imagine that we want to solve the problem numerically on  $[0, L]$ , with a symmetry boundary condition  $u_x = 0$  at  $x = 0$ , but we do not know the exact solution and cannot of that reason assign a correct Dirichlet condition at  $x = L$ . One idea is to simply set  $u(L, t) = 0$  since this will be an accurate approximation before the diffused pulse reaches  $x = L$  and even thereafter it might be a satisfactory condition if the exact  $u$  has a small value. Let  $u_e$  be the exact solution and let  $u$  be the solution of  $u_t = \alpha u_{xx}$  with an initial Gaussian pulse and the boundary conditions  $u_x(0, t) = u(L, t) = 0$ . Derive a diffusion problem for the error  $e = u_e - u$ . Solve this problem numerically using an exact Dirichlet condition at  $x = L$ . Animate the evolution of the error and make a curve plot of the error measure

$$E(t) = \sqrt{\frac{\int_0^L e^2 dx}{\int_0^L u dx}}.$$

Is this a suitable error measure for the present problem?

**b)** Instead of using  $u(L, t) = 0$  as approximate boundary condition for letting the diffused Gaussian pulse move out of our finite domain, one may try  $u_x(L, t) = 0$  since the solution for large  $t$  is quite flat. Argue that this condition gives a completely wrong asymptotic solution as  $t \rightarrow 0$ . To do this, integrate the diffusion equation from 0 to  $L$ , integrate  $u_{xx}$  by parts (or use Gauss' divergence theorem in 1D) to arrive at the important property

$$\frac{d}{dt} \int_0^L u(x, t) dx = 0,$$

implying that  $\int_0^L u dx$  must be constant in time, and therefore

$$\int_0^L u(x, t) dx = \int_0^L I(x) dx.$$

The integral of the initial pulse is 1.

**c)** Another idea for an artificial boundary condition at  $x = L$  is to use a cooling law

$$-\alpha u_x = q(u - u_S), \quad (3.68)$$

where  $q$  is an unknown heat transfer coefficient and  $u_S$  is the surrounding temperature in the medium outside of  $[0, L]$ . (Note that arguing that  $u_S$  is approximately  $u(L, t)$  gives the  $u_x = 0$  condition from the previous

subexercise that is qualitatively wrong for large  $t$ .) Develop a diffusion problem for the error in the solution using (3.68) as boundary condition. Assume one can take  $u_S = 0$  “outside the domain” since  $u_e \rightarrow 0$  as  $x \rightarrow \infty$ . Find a function  $q = q(t)$  such that the exact solution obeys the condition (3.68). Test some constant values of  $q$  and animate how the corresponding error function behaves. Also compute  $E(t)$  curves as defined above.

Filename: `diffu_open_BC`.

### Exercise 3.4: Simulate a diffused Gaussian peak in 2D/3D

- a)** Generalize (3.48) to multi dimensions by assuming that one-dimensional solutions can be multiplied to solve  $u_t = \alpha \nabla^2 u$ . Set  $c = 0$  such that the peak of the Gaussian is at the origin.
- b)** One can from the exact solution show that  $u_x = 0$  on  $x = 0$ ,  $u_y = 0$  on  $y = 0$ , and  $u_z = 0$  on  $z = 0$ . The approximately correct condition  $u = 0$  can be set on the remaining boundaries (say  $x = L$ ,  $y = L$ ,  $z = L$ ), cf. Exercise 3.3. Simulate a 2D case and make an animation of the diffused Gaussian peak.
- c)** The formulation in b) makes use of symmetry of the solution such that we can solve the problem in the first quadrant (2D) or octant (3D) only. To check that the symmetry assumption is correct, formulate the problem without symmetry in a domain  $[-L, L] \times [L, L]$  in 2D. Use  $u = 0$  as approximately correct boundary condition. Simulate the same case as in b), but in a four times as large domain. Make an animation and compare it with the one in b).

Filename: `diffu_symmetric_gaussian_2D`.

### Exercise 3.5: Examine stability of a diffusion model with a source term

Consider a diffusion equation with a linear  $u$  term:

$$u_t = \alpha u_{xx} + \beta u .$$

- a)** Derive in detail a Forward Euler scheme, a Backward Euler scheme, and a Crank-Nicolson for this type of diffusion model. Thereafter, formulate a  $\theta$ -rule to summarize the three schemes.

**b)** Assume a solution like (3.49) and find the relation between  $a$ ,  $k$ ,  $\alpha$ , and  $\beta$ .

**Hint.** Insert (3.49) in the PDE problem.

**c)** Calculate the stability of the Forward Euler scheme. Design numerical experiments to confirm the results.

**Hint.** Insert the discrete counterpart to (3.49) in the numerical scheme. Run experiments at the stability limit and slightly above.

**d)** Repeat c) for the Backward Euler scheme.

**e)** Repeat c) for the Crank-Nicolson scheme.

**f)** How does the extra term  $bu$  impact the accuracy of the three schemes?

**Hint.** For analysis of the accuracy, compare the numerical and exact amplification factors, in graphs and/or by Taylor series expansion.

Filename: `diffu_stability_uterm`.

## 3.5 Diffusion in heterogeneous media

Diffusion in heterogeneous media normally implies a non-constant diffusion coefficient  $\alpha = \alpha(x)$ . A 1D diffusion model with such a variable diffusion coefficient reads

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left( \alpha(x) \frac{\partial u}{\partial x^2} \right) + f(x, t), \quad x \in (0, L), \quad t \in (0, T], \quad (3.69)$$

$$u(x, 0) = I(x), \quad x \in [0, L], \quad (3.70)$$

$$u(0, t) = U_0, \quad t > 0, \quad (3.71)$$

$$u(L, t) = U_L, \quad t > 0. \quad (3.72)$$

A short form of the diffusion equation with variable coefficients is  $u_t = (\alpha u_x)_x$ .

### 3.5.1 Discretization

We can discretize (3.69) by a  $\theta$ -rule in time and centered differences in space:

$$[D_t u]^{n+\frac{1}{2}} = \theta [D_x \bar{\alpha}^x D_x u + f]^{n+1} + (1 - \theta) [D_x \bar{\alpha}^x D_x u + f]^n.$$

Written out, this becomes

$$\begin{aligned} \frac{u^{n+1} - u^n}{\Delta t} &= \theta \frac{1}{\Delta x^2} (\alpha_{i+\frac{1}{2}} (u_{i+1}^{n+1} - u_i^{n+1}) - \alpha_{i-\frac{1}{2}} (u_i^{n+1} - u_{i+1}^{n+1})) + \\ &\quad (1 - \theta) \frac{1}{\Delta x^2} (\alpha_{i+\frac{1}{2}} (u_{i+1}^n - u_i^n) - \alpha_{i-\frac{1}{2}} (u_i^n - u_{i+1}^n)) + \\ &\quad \theta f_i^{n+1} + (1 - \theta) f_i^n, \end{aligned}$$

where, e.g., an arithmetic mean can be used for  $\alpha_{i+\frac{1}{2}}$ :

$$\alpha_{i+\frac{1}{2}} = \frac{1}{2} (\alpha_i + \alpha_{i+1}).$$

### 3.5.2 Implementation

Suitable code for solving the discrete equations is very similar to what we created for a constant  $\alpha$ . Since the Fourier number has no meaning for varying  $\alpha$ , we introduce a related parameter  $D = \Delta t / \Delta x^2$ .

```
def solver_theta(I, a, L, Nx, D, T, theta=0.5, u_L=1, u_R=0,
                 user_action=None):
    x = linspace(0, L, Nx+1)    # mesh points in space
    dx = x[1] - x[0]
    dt = D*dx**2
    Nt = int(round(T/float(dt)))
    t = linspace(0, T, Nt+1)    # mesh points in time

    u    = zeros(Nx+1)    # solution array at t[n+1]
    u_1 = zeros(Nx+1)    # solution at t[n]

    Dl = 0.5*D*theta
    Dr = 0.5*D*(1-theta)

    # Representation of sparse matrix and right-hand side
    diagonal = zeros(Nx+1)
    lower    = zeros(Nx)
    upper    = zeros(Nx)
    b        = zeros(Nx+1)

    # Precompute sparse matrix (scipy format)
    diagonal[1:-1] = 1 + Dl*(a[2:] + 2*a[1:-1] + a[:-2])
    lower[:-1] = -Dl*(a[1:-1] + a[:-2])
    upper[1:]  = -Dl*(a[2:] + a[1:-1])
    # Insert boundary conditions
    diagonal[0] = 1
```

```

upper[0] = 0
diagonal[Nx] = 1
lower[-1] = 0

A = scipy.sparse.diags(
    diagonals=[diagonal, lower, upper],
    offsets=[0, -1, 1],
    shape=(Nx+1, Nx+1),
    format='csr')

# Set initial condition
for i in range(0,Nx+1):
    u_1[i] = I(x[i])

if user_action is not None:
    user_action(u_1, x, t, 0)

# Time loop
for n in range(0, Nt):
    b[1:-1] = u_1[1:-1] + Dr*(
        (a[2:] + a[1:-1])*(u_1[2:] - u_1[1:-1]) -
        (a[1:-1] + a[0:-2])*(u_1[1:-1] - u_1[:-2]))
    # Boundary conditions
    b[0] = u_L(t[n+1])
    b[-1] = u_R(t[n+1])
    # Solve
    u[:] = scipy.sparse.linalg.spsolve(A, b)

    if user_action is not None:
        user_action(u, x, t, n+1)

    # Switch variables before next step
    u_1, u = u, u_1

```

The code is found in the file [diffu1D\\_vc.py](#).

### 3.5.3 Stationary solution

As  $t \rightarrow \infty$ , the solution of the problem (3.69)-(3.72) will approach a stationary limit where  $\partial u / \partial t = 0$ . The governing equation is then

$$\frac{d}{dx} \left( \alpha \frac{du}{dx} \right) = 0, \quad (3.73)$$

with boundary conditions  $u(0) = U_0$  and  $u(L) = u_L$ . It is possible to obtain an exact solution of (3.73) for any  $\alpha$ . Integrating twice and applying the boundary conditions to determine the integration constants gives

$$u(x) = U_0 + (U_L - U_0) \frac{\int_0^x (\alpha(\xi))^{-1} d\xi}{\int_0^L (\alpha(\xi))^{-1} d\xi}. \quad (3.74)$$

### 3.5.4 Piecewise constant medium

Consider a medium built of  $M$  layers. The boundaries between the layers are denoted by  $b_0, \dots, b_M$ , where  $b_0 = 0$  and  $b_M = L$ . If the material in each layer potentially differs from the others, but is otherwise constant, we can express  $\alpha$  as a *piecewise constant function* according to

$$\alpha(x) = \begin{cases} \alpha_0, & b_0 \leq x < b_1, \\ \vdots \\ \alpha_i, & b_i \leq x < b_{i+1}, \\ \vdots \\ \alpha_0, & b_{M-1} \leq x \leq b_M. \end{cases} \quad (3.75)$$

The exact solution (3.74) in case of such a piecewise constant  $\alpha$  function is easy to derive. Assume that  $x$  is in the  $m$ -th layer:  $x \in [b_m, b_{m+1}]$ . In the integral  $\int_0^x (\alpha(\xi))^{-1} d\xi$  we must integrate through the first  $m - 1$  layers and then add the contribution from the remaining part  $x - b_m$  into the  $m$ -th layer:

$$u(x) = U_0 + (U_L - U_0) \frac{\sum_{j=0}^{m-1} (b_{j+1} - b_j)/\alpha(b_j) + (x - b_m)/\alpha(b_m)}{\sum_{j=0}^{M-1} (b_{j+1} - b_j)/\alpha(b_j)} \quad (3.76)$$

**Remark.** It may sound strange to have a discontinuous  $\alpha$  in a differential equation where one is to differentiate, but a discontinuous  $\alpha$  is compensated by a discontinuous  $u_x$  such that  $\alpha u_x$  is continuous and therefore can be differentiated as  $(\alpha u_x)_x$ .

### 3.5.5 Implementation of diffusion in a piecewise constant medium

Programming with piecewise function definition quickly becomes cumbersome as the most naive approach is to test for which interval  $x$  lies, and then start evaluating a formula like (3.76). In Python, vectorized expressions may help to speed up the computations. The convenience classes `PiecewiseConstant` and `IntegratedPiecewiseConstant` in the

`Heaviside` module were made to simplify programming with functions like (3.5.4) and expressions like (3.76). These utilities not only represent piecewise constant functions, but also *smoothed* versions of them where the discontinuities can be smoothed out in a controlled fashion. This is advantageous in many computational contexts (although seldom for pure finite difference computations of the solution  $u$ ).

The `PiecewiseConstant` class is created by sending in the domain as a 2-tuple or 2-list and a `data` object describing the boundaries  $b_0, \dots, b_M$  and the corresponding function values  $\alpha_0, \dots, \alpha_{M-1}$ . More precisely, `data` is a nested list, where `data[i][0]` holds  $b_i$  and `data[i][1]` holds the corresponding value  $\alpha_i$ , for  $i = 0, \dots, M - 1$ . Given  $b_i$  and  $\alpha_i$  in arrays `b` and `a`, it is easy to fill out the nested list `data`.

In our application, we want to represent  $\alpha$  and  $1/\alpha$  as piecewise constant function, in addition to the  $u(x)$  function which involves the integrals of  $1/\alpha$ . A class creating the functions we need and a method for evaluating  $u$ , can take the form

```
class SerialLayers:
    """
        b: coordinates of boundaries of layers, b[0] is left boundary
        and b[-1] is right boundary of the domain [0,L].
        a: values of the functions in each layer (len(a) = len(b)-1).
        U_0: u(x) value at left boundary x=0=b[0].
        U_L: u(x) value at right boundary x=L=b[0].
    """

    def __init__(self, a, b, U_0, U_L, eps=0):
        self.a, self.b = np.asarray(a), np.asarray(b)
        self.eps = eps # smoothing parameter for smoothed a
        self.U_0, self.U_L = U_0, U_L

        a_data = [[bi, ai] for bi, ai in zip(self.b, self.a)]
        domain = [b[0], b[-1]]
        self.a_func = PiecewiseConstant(domain, a_data, eps)

        # inv_a = 1/a is needed in formulas
        inv_a_data = [[bi, 1./ai] for bi, ai in zip(self.b, self.a)]
        self.inv_a_func = \
            PiecewiseConstant(domain, inv_a_data, eps)
        self.integral_of_inv_a_func = \
            IntegratedPiecewiseConstant(domain, inv_a_data, eps)
        # Denominator in the exact formula is constant
        self.inv_a_OL = self.integral_of_inv_a_func(b[-1])

    def __call__(self, x):
        solution = self.U_0 + (self.U_L-self.U_0)*\
                   self.integral_of_inv_a_func(x)/self.inv_a_OL
        return solution
```

A visualization method is also convenient to have. Below we plot  $u(x)$  along with  $\alpha(x)$  (which works well as long as  $\max \alpha(x)$  is of the same size as  $\max u = \max(U_0, U_L)$ ).

```
class SerialLayers:
    ...

    def plot(self):
        x, y_a = self.a_func.plot()
        x = np.asarray(x); y_a = np.asarray(y_a)
        y_u = self.u_exact(x)
        import matplotlib.pyplot as plt
        plt.figure()
        plt.plot(x, y_u, 'b')
        plt.hold('on') # Matlab style
        plt.plot(x, y_a, 'r')
        ymin = -0.1
        ymax = 1.2 * max(y_u.max(), y_a.max())
        plt.axis([x[0], x[-1], ymin, ymax])
        plt.legend(['solution $u$', 'coefficient $a$'], loc='upper left')
        if self.eps > 0:
            plt.title('Smoothing eps: %s' % self.eps)
        plt.savefig('tmp.pdf')
        plt.savefig('tmp.png')
        plt.show()
```

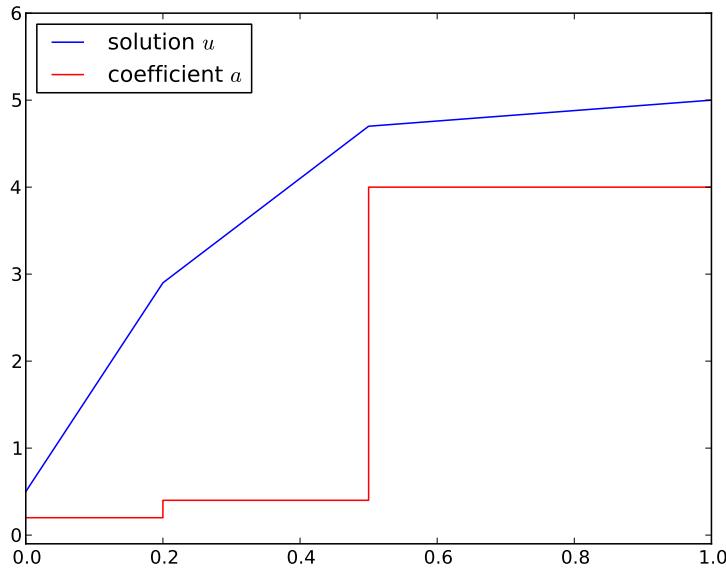
Figure 3.12 shows the case where

```
b = [0, 0.25, 0.5, 1] # material boundaries
a = [0.2, 0.4, 4]      # material values
U_0 = 0.5; U_L = 5     # boundary conditions
```

By adding the `eps` parameter to the constructor of the `SerialLayers` class, we can experiment with smoothed versions of  $\alpha$  and see the (small) impact on  $u$ . Figure 3.13 shows the result.

### 3.5.6 Axi-symmetric diffusion

Suppose we have a diffusion process taking place in a straight tube with radius  $R$ . We assume axi-symmetry such that  $u$  is just a function of  $r$  and  $t$ ,  $r$  being the radial distance from a point to the center axis of the tube. With such axi-symmetry it is advantageous to introduce *cylindrical coordinates*  $r$ ,  $\theta$ , and  $z$ , where  $z$  is in the direction of the tube and  $(r, \theta)$  are polar coordinates in a cross section. Axi-symmetry means that all quantities are independent of  $\theta$ . From the relations  $x = \cos \theta$ ,  $y = \sin \theta$ , and  $z = z$ , between Cartesian and cylindrical coordinates, one can (after some work) derive the diffusion equation in cylindrical coordinates, which with axi-symmetry takes the form



**Fig. 3.12** Solution of the stationary diffusion equation corresponding to a piecewise constant diffusion coefficient.

$$\frac{\partial u}{\partial t} = \frac{1}{r} \frac{\partial}{\partial r} \left( r \alpha(r, z) \frac{\partial u}{\partial r} \right) + \frac{\partial}{\partial z} \left( \alpha(r, z) \frac{\partial u}{\partial z} \right) + f(r, z, t).$$

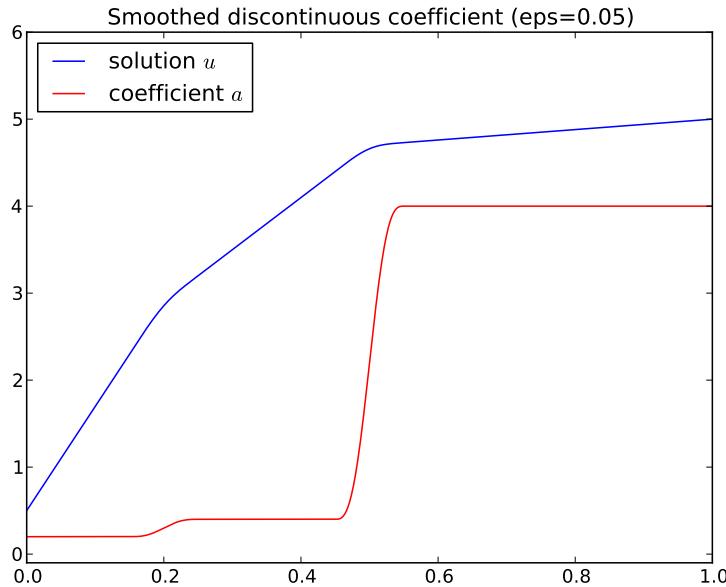
Let us assume that  $u$  does not change along the tube axis so it suffices to compute variations in a cross section. Then  $\partial u / \partial z = 0$  and we have a 1D diffusion equation in the radial coordinate  $r$  and time  $t$ . In particular, we shall address the initial-boundary value problem

$$\frac{\partial u}{\partial t} = \frac{1}{r} \frac{\partial}{\partial r} \left( r \alpha(r) \frac{\partial u}{\partial r} \right) + f(t), \quad r \in (0, R), \quad t \in (0, T], \quad (3.77)$$

$$\frac{\partial u}{\partial r}(0, t) = 0, \quad t \in (0, T], \quad (3.78)$$

$$u(R, t) = 0, \quad t \in (0, T], \quad (3.79)$$

$$u(r, 0) = I(r), \quad r \in [0, R]. \quad (3.80)$$



**Fig. 3.13** Solution of the stationary diffusion equation corresponding to a *smoothed* piecewise constant diffusion coefficient.

The condition (3.78) is a necessary symmetry condition at  $r = 0$ , while (3.79) could be any Dirichlet or Neumann condition (or Robin condition in case of cooling or heating).

The finite difference approximation will need the discretized version of the PDE for  $r = 0$  (just as we use the PDE at the boundary when implementing Neumann conditions). However, discretizing the PDE at  $r = 0$  poses a problem because of the  $1/r$  factor. We therefore need to work out the PDE for discretization at  $r = 0$  with care. Let us in case of constant  $\alpha$  expand the spatial derivative to

$$\text{alpha} \frac{\partial^2 u}{\partial r^2} + \alpha \frac{1}{r} \frac{\partial u}{\partial r}.$$

The last term faces a difficulty at  $r = 0$  since it becomes a  $0/0$  expression because of the symmetry condition at  $r = 0$ . L'Hospital's rule can be used:

$$\lim_{r \rightarrow 0} \frac{1}{r} \frac{\partial u}{\partial r} = \frac{\partial^2 u}{\partial r^2}.$$

The PDE at  $r = 0$  therefore becomes

$$\frac{\partial u}{\partial t} = 2\alpha \frac{\partial^2 u}{\partial r^2} + f(t). \quad (3.81)$$

For a variable coefficient  $\alpha(r)$  the expanded derivative reads

$$\alpha(r) \frac{\partial^2 u}{\partial r^2} + \frac{1}{r}(\alpha(r) + r\alpha'(r)) \frac{\partial u}{\partial r}.$$

We are interested in this expression for  $r = 0$ . A necessary condition for  $u$  to be axi-symmetric is that all input data, including  $\alpha$ , must also be axi-symmetric, implying that  $\alpha'(0) = 0$  (the second term vanishes anyway because of  $r = 0$ ). The limit of interest is

$$\lim_{r \rightarrow 0} \frac{1}{r} \alpha(r) \frac{\partial u}{\partial r} = \alpha(0) \frac{\partial^2 u}{\partial r^2}.$$

The PDE at  $r = 0$  now looks like

$$\frac{\partial u}{\partial t} = 2\alpha(0) \frac{\partial^2 u}{\partial r^2} + f(t), \quad (3.82)$$

so there is no essential difference between the constant coefficient and variable coefficient cases.

The second-order derivative in (3.81) and (3.82) are discretized in the usual way.

$$2\alpha \frac{\partial^2}{\partial r^2} u(r_0, t_n) \approx [2\alpha 2D_r D_r u]_0^n = 2\alpha \frac{u_1^n - 2u_0^n + u_{-1}^n}{\Delta r^2}.$$

The fictitious value  $u_{-1}^n$  can be eliminated using the discrete symmetry condition

$$[D_{2r} u = 0]_0^n \Rightarrow u_{-1}^n = u_1^n,$$

which then gives the modified approximation to the second-order derivative of  $u$  in  $r$  at  $r = 0$ :

$$4\alpha \frac{u_1^n - u_0^n}{\Delta r^2}. \quad (3.83)$$

The discretization of the second-order derivative in  $r$  at another internal mesh point is straightforward:

$$\begin{aligned} \left[ \frac{1}{r} \frac{\partial}{\partial r} \left( r \alpha \frac{\partial u}{\partial r} \right) \right]_i^n &\approx [r^{-1} D_r(r \alpha D_r u)]_i^n \\ &= \frac{1}{\Delta r^2} \left( r_{i+\frac{1}{2}} \alpha_{i+\frac{1}{2}} (u_{i+1}^n - u_i^n) - r_{i-\frac{1}{2}} \alpha_{i-\frac{1}{2}} (u_i^n - u_{i-1}^n) \right). \end{aligned}$$

To complete the discretization, we need a scheme in time, but that can be done as before and does not interfere with the discretization in space.

### 3.5.7 Spherically-symmetric diffusion

**Discretization in spherical coordinates.** Let us now pose the problem from Section 3.5.6 in spherical coordinates, where  $u$  only depends on the radial coordinate  $r$  and time  $t$ . That is, we have spherical symmetry. For simplicity we restrict the diffusion coefficient  $\alpha$  to be a constant. The PDE reads

$$\frac{\partial u}{\partial t} = \frac{\alpha}{r^\gamma} \frac{\partial}{\partial r} \left( r^\gamma \frac{\partial u}{\partial r} \right) + f(t), \quad (3.84)$$

for  $r \in (0, R)$  and  $t \in (0, T]$ . The parameter  $\gamma$  is 2 for spherically-symmetric problems and 1 for axi-symmetric problems. The boundary and initial conditions have the same mathematical form as in (3.77)-(3.80).

Since the PDE in spherical coordinates has the same form as the PDE in Section 3.5.6, just with the  $\gamma$  parameter being different, we can use the same discretization approach. At the origin  $r = 0$  we get problems with the term

$$\frac{\gamma}{r} \frac{\partial u}{\partial t},$$

but L'Hospital's rule shows that this term equals  $\gamma \partial^2 u / \partial r^2$ , and the PDE at  $r = 0$  becomes

$$\frac{\partial u}{\partial t} = (\gamma + 1) \alpha \frac{\partial^2 u}{\partial r^2} + f(t). \quad (3.85)$$

The associated discrete form is then

$$[D_t u]_i = \frac{1}{2} (\gamma + 1) \alpha ([D_r D_r \bar{u}^t + \bar{f}^t]_i^n, \quad (3.86)$$

for a Crank-Nicolson scheme.

**Discretization in Cartesian coordinates.** The spherically-symmetric spatial derivative can be transformed to the Cartesian counterpart by introducing

$$v(r, t) = ru(r, t).$$

Inserting  $u = v/r$  in the PDE yields

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left( \alpha(r) r^2 \frac{\partial u}{\partial r} \right),$$

and then

$$r \left( \frac{dc^2}{dr} \frac{\partial v}{\partial r} + \alpha \frac{\partial^2 v}{\partial r^2} \right) - \frac{dc^2}{dr} v.$$

The two terms in the parenthesis can be combined to

$$r \frac{\partial}{\partial r} \left( \alpha \frac{\partial v}{\partial r} \right),$$

which is recognized as the variable-coefficient Laplace operator in one Cartesian coordinate.

## 3.6 Diffusion in 2D

We now address a diffusion in two space dimensions:

$$\frac{\partial u}{\partial t} = \alpha \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) + f(x, y), \quad (3.87)$$

in a domain

$$(x, y) \in (0, L_x) \times (0, L_y), \quad t \in (0, T],$$

with  $u = 0$  on the boundary and  $u(x, y, 0) = I(x, y)$  as initial condition.

### 3.6.1 Discretization

For generality, it is natural to use a  $\theta$ -rule for the time discretization. Standard, second-order accurate finite differences are used for the spatial

derivatives. We sample the PDE at a space-time point  $(i, j, n + \frac{1}{2})$  and apply the difference approximations:

$$[D_t u]^{n+\frac{1}{2}} = \theta[\alpha(D_x D_x u + D_y D_y u) + f]^{n+1} + \\ (1 - \theta)[\alpha(D_x D_x u + D_y D_y u) + f]^n. \quad (3.88)$$

Written out,

$$\frac{u_{i,j}^{n+1} - u_{i,j}^n}{\Delta t} = \theta(\alpha(\frac{u_{i-1,j}^{n+1} - 2u_{i,j}^{n+1} + u_{i+1,j}^{n+1}}{\Delta x^2}) + (\frac{u_{i,j-1}^{n+1} - 2u_{i,j}^{n+1} + u_{i,j+1}^{n+1}}{\Delta y^2})) + f_{i,j}^{n+1}) + \\ (1 - \theta)(\alpha(\frac{u_{i-1,j}^n - 2u_{i,j}^n + u_{i+1,j}^n}{\Delta x^2}) + (\frac{u_{i,j-1}^n - 2u_{i,j}^n + u_{i,j+1}^n}{\Delta y^2})) + f_{i,j}^n) \quad (3.89)$$

We collect the unknowns on the left-hand side

$$u_{i,j}^{n+1} - \theta(F_x(u_{i-1,j}^{n+1} - 2u_{i,j}^{n+1} + u_{i+1,j}^{n+1}) + F_y(u_{i,j-1}^{n+1} - 2u_{i,j}^{n+1} + u_{i,j+1}^{n+1})) = u_{i,j}^n + \\ (1 - \theta)(F_x(u_{i-1,j}^n - 2u_{i,j}^n + u_{i+1,j}^n) + F_y(u_{i,j-1}^n - 2u_{i,j}^n + u_{i,j+1}^n)) + \\ \theta \Delta t f_{i,j}^{n+1} + (1 - \theta) \Delta t f_{i,j}^n, \quad (3.90)$$

where

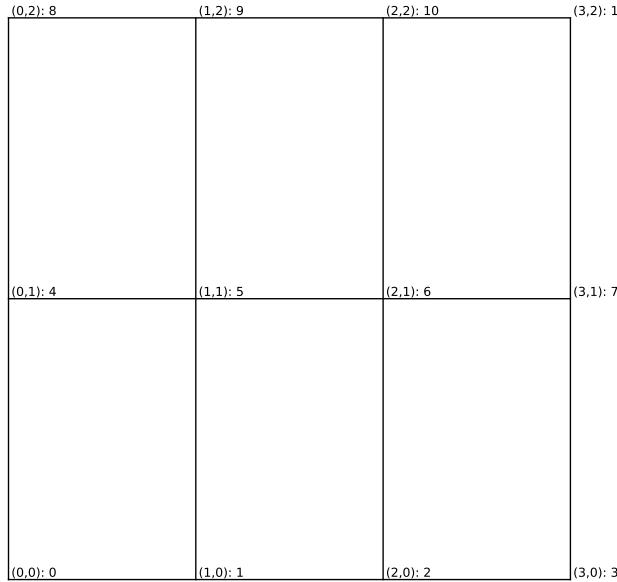
$$F_x = \frac{\alpha \Delta t}{\Delta x^2}, \quad F_y = \frac{\alpha \Delta t}{\Delta y^2},$$

are the Fourier numbers in  $x$  and  $y$  direction, respectively.

### 3.6.2 Numbering of mesh points versus equations and unknowns

The equations (3.90) are coupled at the new time level  $n + 1$ . That is, we must solve a system of (linear) algebraic equations, which we will write as  $Ac = b$ , where  $A$  is the coefficient matrix,  $c$  is the vector of unknowns, and  $b$  is the right-hand side.

Let us examine the equations in  $Ac = b$  on a mesh with  $N_x = 3$  and  $N_y = 2$  cells in each direction. The spatial mesh is depicted in Figure 3.14. The equations at the boundary just implement the boundary condition  $u = 0$ :



**Fig. 3.14** 3x2 2D mesh.

$$u_{0,0}^{n+1} = u_{1,0}^{n+1} = u_{2,0}^{n+1} = u_{3,0}^{n+1} = u_{0,1}^{n+1} = u_{3,1}^{n+1} = u_{0,2}^{n+1} = u_{1,2}^{n+1} = u_{2,2}^{n+1} = u_{3,2}^{n+1} = 0.$$

We are left with two interior points, with  $i = 1, j = 1$  and  $i = 2, j = 1$ . The corresponding equations are

$$\begin{aligned} u_{i,j}^{n+1} - \theta \left( F_x(u_{i-1,j}^{n+1} - 2u_{i,j}^{n+1} + u_{i+1,j}^{n+1}) + F_y(u_{i,j-1}^{n+1} - 2u_{i,j}^{n+1} + u_{i,j+1}^{n+1}) \right) = u_{i,j}^n + \\ (1 - \theta) \left( F_x(u_{i-1,j}^n - 2u_{i,j}^n + u_{i+1,j}^n) + F_y(u_{i,j-1}^n - 2u_{i,j}^n + u_{i,j+1}^n) \right) + \\ \theta \Delta t f_{i,j}^{n+1} + (1 - \theta) \Delta t f_{i,j}^n, \end{aligned}$$

There are in total 12 unknowns  $u_{i,j}^{n+1}$  for  $i = 0, 1, 2, 3$  and  $j = 0, 1, 2$ . To solve the equations, we need to form a matrix system  $Ac = b$ . In that system, the solution vector  $c$  can only one index. Thus, we need a numbering of the unknowns with one index, not two as used in the mesh. We introduce a mapping  $m(i, j)$  from a mesh point with indices  $(i, j)$  to the corresponding unknown  $p$  in the equation system:

$$p = m(i, j) = j(N_x + 1) + i.$$

When  $i$  and  $j$  runs through their values we see the following mapping to  $p$ :

$$\begin{aligned}(0,0) &\rightarrow 0, (0,1) \rightarrow 1, (0,2) \rightarrow 2, (0,3) \rightarrow 3, \\ (1,0) &\rightarrow 4, (1,1) \rightarrow 5, (1,2) \rightarrow 6, (1,3) \rightarrow 7, \\ (2,0) &\rightarrow 8, (2,1) \rightarrow 9, (2,2) \rightarrow 10, (2,3) \rightarrow 11.\end{aligned}$$

That is, we number the points along the  $x$  axis, starting with  $y = 0$ , and the progress one horizontal mesh line at a time. In Figure 3.14 you can see that the  $(i,j)$  and the corresponding single index  $(p)$  are listed for each mesh point.

We could equally well numbered the equations in other ways, e.g., let the  $j$  index be the fastest varying index:  $p = m(i, j) = i(N_y + 1) + j$ .

Let us form the coefficient matrix  $A$ , or more precisely, insert matrix element (according Python's convention with zero as base index) for each of the nonzero elements in  $A$  (the indices run through the values of  $p$ , i.e.,  $p = 0, \dots, 11$ ):

$$\left( \begin{array}{cccccccccccc} (0,0) & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & (1,1) & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & (2,2) & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & (3,3) & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & (4,4) & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & (5,5) & (5,6) & 0 & 0 & (5,9) & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & (6,6) & (6,7) & 0 & 0 & (6,10) & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & (7,7) & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & (8,8) & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & (9,9) & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & (10,10) & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & (11,11) \end{array} \right)$$

Here is a more compact visualization of the coefficient matrix where we insert dots for zeros and bullets for non-zero elements:

$$\left( \begin{array}{cccccccccccc} \bullet & \cdot \\ \cdot & \bullet & \cdot \\ \cdot & \cdot & \bullet & \cdot \\ \cdot & \cdot & \cdot & \bullet & \cdot \\ \cdot & \cdot & \bullet & \cdot & \bullet & \bullet & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \bullet & \cdot & \bullet & \bullet & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \bullet & \cdot & \bullet & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \bullet & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \bullet & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \bullet & \cdot & \cdot & \cdot & \cdot \end{array} \right)$$

It is clearly seen that most of the elements are zero. This is a general feature of coefficient matrices arising from discretizing PDEs by finite difference methods. We say that the matrix is *sparse*.

Let  $A_{p,q}$  be the value of element  $(p,q)$  in the coefficient matrix  $A$ , where  $p$  and  $q$  now correspond to the numbering of the unknowns in the equation system. We have  $A_{p,q} = 1$  for  $p = q = 0, 1, 2, 3, 4, 7, 8, 9, 10, 11$ , corresponding to all the known boundary values. Let  $p$  be  $m(i,j)$ , i.e., the single index corresponding to mesh point  $(i,j)$ . Then we have

$$A_{m(i,j),m(i,j)} = A_{p,p} = 1 + \theta(F_x + F_y), \quad (3.91)$$

$$A_{p,m(i-1,j)} = A_{p,p-1} = -\theta F_x, \quad (3.92)$$

$$A_{p,m(i+1,j)} = A_{p,p+1} = -\theta F_x, \quad (3.93)$$

$$A_{p,m(i,j-1)} = A_{p,p-(N_x+1)} = -\theta F_y, \quad (3.94)$$

$$A_{p,m(i,j+1)} = A_{p,p+(N_x+1)} = -\theta F_y, \quad (3.95)$$

$$(3.96)$$

for the equations associated with the two interior mesh points. At these interior points, the single index  $p$  takes on the specific values  $p = 5, 6$ , corresponding to the values  $(1, 1)$  and  $(1, 2)$  of the pair  $(i, j)$ .

The above values for  $A_{p,q}$  can be inserted in the matrix:

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -\theta F_y & 0 & 0 & -\theta F_x & 1 + 2\theta F_x & -\theta F_x & 0 & 0 & -\theta F_y & 0 & 0 \\ 0 & 0 & -\theta F_y & 0 & 0 & -\theta F_x & 1 + 2\theta F_x & -\theta F_x & 0 & 0 & -\theta F_y & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

The corresponding right-hand side vector in the equation system has the entries  $b_p$ , where  $p$  numbers the equations. We have

$$b_0 = b_1 = b_2 = b_3 = b_4 = b_7 = b_8 = b_9 = b_{10} = b_{11} = 0,$$

for the boundary values. For the equations associated with the interior points, we get for  $p = 5, 6$ , corresponding to  $i = 1, 2$  and  $j = 1$ :

$$b_p = u_i + (1 - \theta) \left( F_x(u_{i-1,j}^n - 2u_{i,j}^n + u_{i,j+1}^n) + F_y(u_{i,j-1}^n - 2u_{i,j}^n + u_{i,j+1}^n) \right) + \theta \Delta t f_{i,j}^{n+1} + (1 - \theta) \Delta t f_{i,j}^n.$$

Recall that  $p = m(i,j) = j(N_x + 1) + j$  in this expression.

We can, as an alternative, leave the boundary mesh points out of the matrix system. For a mesh with  $N_x = 3$  and  $N_y = 2$  there are only two

internal mesh points whose unknowns will enter the matrix system. We must now number the unknowns at the interior points:

$$p = (j - 1)(N_x - 1) + i,$$

for  $i = 1, \dots, N_x - 1$ ,  $j = 1, \dots, N_y - 1$ .

**hpl 14:** Fill in details.

(0,3): 15	(1,3): 16	(2,3): 17	(3,3): 18	(4,3): 19
(0,2): 10	(1,2): 11	(2,2): 12	(3,2): 13	(4,2): 14
(0,1): 5	(1,1): 6	(2,1): 7	(3,1): 8	(4,1): 9
(0,0): 0	(1,0): 1	(2,0): 2	(3,0): 3	(4,0): 4

**Fig. 3.15** 4x3 2D mesh.

We can continue with illustrating a bit larger mesh,  $N_x = 4$  and  $N_y = 3$ , see Figure 3.15. The corresponding coefficient matrix with dots for zeros and bullets for non-zeroes look as follows (values at boundary points are included in the equation system):

$$\begin{pmatrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{pmatrix}$$

### The coefficient matrix is banded

Besides being sparse, we observe that the coefficient matrix is *banded*: it has five distinct bands. We have the diagonal  $A_{i,i}$ , the subdiagonal  $A_{i-1,j}$ , the superdiagonal  $A_{i,i+1}$ , a lower diagonal  $A_{i,i-(Nx+1)}$ , and an upper diagonal  $A_{i,i+(Nx+1)}$ . The other matrix entries are known to be zero. With  $N_x + 1 = N_y + 1 = N$ , only a fraction  $5N^{-2}$  of the matrix entries are nonzero, so the matrix is clearly very sparse for relevant  $N$  values. The more we can compute with the nonzeros only, the faster the solution methods will be.

### 3.6.3 Algorithm for setting up the coefficient matrix

We looked at a specific mesh in the previous section, formulated the equations, and saw what the corresponding coefficient matrix and right-hand side are. Now our aim is to set up a general algorithm, for any choice of  $N_x$  and  $N_y$ , that produces the coefficient matrix and the right-hand side vector. We start with a zero matrix and vector, run through each mesh point, and fill in the values depending on whether the mesh point is an interior point or on the boundary.

- for  $i = 0, \dots, N_x$ 
  - for  $j = 0, \dots, N_y$ 
    - $p = j(N_x + 1) + i$
    - if point  $(i, j)$  is on the boundary:
      - $A_{p,p} = 1, b_p = 0$
      - else:

- fill  $A_{p,m(i-1,j)}$ ,  $A_{p,m(i+1,j)}$ ,  $A_{p,m(i,j)}$ ,  $A_{p,m(i,j-1)}$ ,  $A_{p,m(i,j+1)}$ , and  $b_p$

To ease the test on whether  $(i, j)$  is on the boundary or not, we can split the loops a bit, starting with the boundary line  $j = 0$ , then treat the interior lines  $1 \leq j < N_y$ , and finally treat the boundary line  $j = N_y$ :

- for  $i = 0, \dots, N_x$ 
  - boundary  $j = 0$ :  $p = j(N_x + 1) + i$ ,  $A_{p,p} = 1$
- for  $j = 0, \dots, N_y$ 
  - boundary  $i = 0$ :  $p = j(N_x + 1) + i$ ,  $A_{p,p} = 1$
  - for  $i = 1, \dots, N_x - 1$ 
    - interior point  $p = j(N_x + 1) + i$
    - fill  $A_{p,m(i-1,j)}$ ,  $A_{p,m(i+1,j)}$ ,  $A_{p,m(i,j)}$ ,  $A_{p,m(i,j-1)}$ ,  $A_{p,m(i,j+1)}$ , and  $b_p$
  - boundary  $i = N_x$ :  $p = j(N_x + 1) + i$ ,  $A_{p,p} = 1$
- for  $i = 0, \dots, N_x$ 
  - boundary  $j = N_y$ :  $p = j(N_x + 1) + i$ ,  $A_{p,p} = 1$

The right-hand side is set up as follows.

- for  $i = 0, \dots, N_x$ 
  - boundary  $j = 0$ :  $p = j(N_x + 1) + i$ ,  $b_p = 0$
- for  $j = 0, \dots, N_y$ 
  - boundary  $i = 0$ :  $p = j(N_x + 1) + i$ ,  $b_p = 0$
  - for  $i = 1, \dots, N_x - 1$ 
    - interior point  $p = j(N_x + 1) + i$
    - fill  $b_p$
  - boundary  $i = N_x$ :  $p = j(N_x + 1) + i$ ,  $b_p = 0$
- for  $i = 0, \dots, N_x$ 
  - boundary  $j = N_y$ :  $p = j(N_x + 1) + i$ ,  $b_p = 0$

### 3.6.4 Implementation with a dense coefficient matrix

The goal now is to map the algorithms in the previous section to Python code. One should for computational efficiency reasons take advantage of the fact that the coefficient matrix is sparse and/or banded, i.e., take advantage of all the zeros; however, we first demonstrate how to fill an  $N \times N$  dense square matrix, where  $N$  is the number of unknowns, here  $N = (N_x + 1)(N_y + 1)$ . The dense matrix is much easier to understand than the sparse matrix case.

```
import numpy as np

def solver_dense(
    I, a, f, Lx, Ly, Nx, Ny, dt, T, theta=0.5, user_action=None):
    """
    Solve u_t = a*(u_xx + u_yy) + f, u(x,y,0)=I(x,y), with u=0
    on the boundary, on [0,Lx]x[0,Ly]x[0,T], with time step dt,
    using the theta-scheme.
    """
    x = np.linspace(0, Lx, Nx+1)      # mesh points in x dir
    y = np.linspace(0, Ly, Ny+1)      # mesh points in y dir
    dx = x[1] - x[0]
    dy = y[1] - y[0]

    dt = float(dt)                  # avoid integer division
    Nt = int(round(T/float(dt)))
    t = np.linspace(0, Nt*dt, Nt+1)  # mesh points in time

    # Mesh Fourier numbers in each direction
    Fx = a*dt/dx**2
    Fy = a*dt/dy**2
```

The  $u_{i,j}^{n+1}$  and  $u_{i,j}^n$  mesh functions are represented by their spatial values at the mesh points:

```
u   = np.zeros((Nx+1, Ny+1))      # unknown u at new time level
u_1 = np.zeros((Nx+1, Ny+1))      # u at the previous time level
```

It is a good habit (for extensions) to introduce index sets for all mesh points:

```
Ix = range(0, Nx+1)
Iy = range(0, Ny+1)
It = range(0, Nt+1)
```

The initial condition is easy to fill in:

```
# Load initial condition into u_1
for i in Ix:
    for j in Iy:
        u_1[i,j] = I(x[i], y[j])
```

The memory for the coefficient matrix and right-hand side vector is allocated by

```
N = (Nx+1)*(Ny+1) # no of unknowns
A = np.zeros((N, N))
b = np.zeros(N)
```

The filling of A goes like this:

```
m = lambda i, j: j*(Nx+1) + i

# Equations corresponding to j=0, i=0,1,... (u known)
j = 0
for i in Ix:
    p = m(i,j); A[p, p] = 1

# Loop over all internal mesh points in y direction
# and all mesh points in x direction
for j in Iy[1:-1]:
    i = 0; p = m(i,j); A[p, p] = 1 # Boundary
    for i in Ix[1:-1]: # Interior points
        p = m(i,j)
        A[p, m(i,j-1)] = - theta*Fy
        A[p, m(i-1,j)] = - theta*Fx
        A[p, p] = 1 + 2*theta*(Fx+Fy)
        A[p, m(i+1,j)] = - theta*Fx
        A[p, m(i,j+1)] = - theta*Fy
    i = Nx; p = m(i,j); A[p, p] = 1 # Boundary
# Equations corresponding to j=Ny, i=0,1,... (u known)
j = Ny
for i in Ix:
    p = m(i,j); A[p, p] = 1
```

Since A is independent of time, it can be filled once and for all before the time loop. The right-hand side vector must be filled at each time level inside the time loop:

```
import scipy.linalg

for n in It[0:-1]:
    # Compute b
    j = 0
    for i in Ix:
        p = m(i,j); b[p] = 0 # Boundary
    for j in Iy[1:-1]:
        i = 0; p = m(i,j); b[p] = 0 # Boundary
        for i in Ix[1:-1]: # Interior points
            p = m(i,j)
            b[p] = u_1[i,j] + \
                (1-theta)*(
                    Fx*(u_1[i+1,j] - 2*u_1[i,j] + u_1[i-1,j]) + \
                    Fy*(u_1[i,j+1] - 2*u_1[i,j] + u_1[i,j-1]))\ 
                + theta*dt*f(i*dx, j*dy, (n+1)*dt) + \
```

```

        (1-theta)*dt*f(i*dx,j*dy,n*dt)
        i = Nx;  p = m(i,j);  b[p] = 0 # Boundary
j = Ny
for i in Ix:
    p = m(i,j);  b[p] = 0           # Boundary

# Solve matrix system A*c = b
c = scipy.linalg.solve(A, b)

# Fill u with vector c
for i in Ix:
    for j in Iy:
        u[i,j] = c[m(i,j)]

# Update u_1 before next step
u_1, u = u, u_1

```

We use `solve` from `scipy.linalg` and not from `numpy.linalg`. The difference is stated below.

### scipy.linalg versus numpy.linalg

Quote from the SciPy documentation:

`scipy.linalg` contains all the functions in `numpy.linalg` plus some other more advanced ones not contained in `numpy.linalg`.

Another advantage of using `scipy.linalg` over `numpy.linalg` is that it is always compiled with BLAS/LAPACK support, while for NumPy this is optional. Therefore, the SciPy version might be faster depending on how NumPy was installed.

Therefore, unless you don't want to add SciPy as a dependency to your NumPy program, use `scipy.linalg` instead of `numpy.linalg`.

The code shown above is available in the `solver_dense` function in the file `diffu2D_u0.py`, differing only in the boundary conditions, which in the code can be an arbitrary function along each side of the domain.

We do not bother to look at vectorized versions of filling `A` since a dense matrix is just used of pedagogical reasons for the very first implementation. Vectorization will be treated when `A` has a sparse matrix representation, as in Section 3.6.7.

### How to debug the computation of $A$ and $b$

A good starting point for debugging the filling of  $A$  and  $b$  is to choose a very coarse mesh, say  $N_x = N_y = 2$ , where there is just one internal mesh point, compute the equations by hand, and print out  $\mathbf{A}$  and  $\mathbf{b}$  for comparison in the code. If wrong elements in  $\mathbf{A}$  or  $\mathbf{b}$  occur, print out each assignment to elements in  $\mathbf{A}$  and  $\mathbf{b}$  inside the loops and compare with what you expect.

To let the user store, analyze, or visualize the solution at each time level, we include a callback function, named `user_action`, to be called before the time loop and in each pass in that loop. The function has the signature

```
user_action(u, x, xv, y, yv, t, n)
```

where  $u$  is a two-dimensional array holding the solution at time level  $n$  and time  $t[n]$ . The  $x$  and  $y$  coordinates of the mesh points are given by the arrays  $x$  and  $y$ , respectively. The arrays  $xv$  and  $yv$  are vectorized representations of the mesh points such that vectorized function evaluations can be invoked. The  $xv$  and  $yv$  arrays are defined by

```
xv = x[:, np.newaxis]
yv = y[np.newaxis, :]
```

One can then evaluate, e.g.,  $f(x, y, t)$  at all internal mesh points at time level  $n$  by first evaluating  $f$  at all points,

```
f_a = f(xv, yv, t[n])
```

and then use slices to extract a view of the values at the internal mesh points:  $f_a[1:-1, 1:-1]$ . The next section features an example on writing a `user_action` callback function.

### 3.6.5 Verification: exact numerical solution

A good test example to start with is one that preserves the solution  $u = 0$ , i.e.,  $f = 0$  and  $I(x, y) = 0$ . This trivial solution can uncover some bugs.

The first real test example is based on having an exact solution of the discrete equations. This solution is linear in time and quadratic in space:

$$u(x, y, t) = 5tx(L_x - x)y(y - L_y).$$

Inserting this manufactured solution in the PDE shows that the source term  $f$  must be

$$f(x, y, t) = 5x(L_x - x)y(y - L_y) + 10\alpha t(x(L_x - x) + y(y - L_y)).$$

We can use the `user_action` function to compare the numerical solution with the exact solution at each time level. A suitable helper function for checking the solution goes like this:

```
def quadratic(theta, Nx, Ny):

    def u_exact(x, y, t):
        return 5*t*x*(Lx-x)*y*(Ly-y)
    def I(x, y):
        return u_exact(x, y, 0)
    def f(x, y, t):
        return 5*x*(Lx-x)*y*(Ly-y) + 10*a*t*(y*(Ly-y)+x*(Lx-x))

    # Use rectangle to detect errors in switching i and j in scheme
    Lx = 0.75
    Ly = 1.5
    a = 3.5
    dt = 0.5
    T = 2

    def assert_no_error(u, x, xv, y, yv, t, n):
        """Assert zero error at all mesh points."""
        u_e = u_exact(xv, yv, t[n])
        diff = abs(u - u_e).max()
        tol = 1E-12
        msg = 'diff=%g, step %d, time=%g' % (diff, n, t[n])
        print msg
        assert diff < tol, msg

    solver_dense(
        I, a, f, Lx, Ly, Nx, Ny,
        dt, T, theta, user_action=assert_no_error)
```

A true test function for checking the quadratic solution for several different meshes and  $\theta$  values can take the form

```
def test_quadratic():
    # For each of the three schemes (theta = 1, 0.5, 0), a series of
    # meshes are tested (Nx > Ny and Nx < Ny)
    for theta in [1, 0.5, 0]:
        for Nx in range(2, 6, 2):
            for Ny in range(2, 6, 2):
                print 'testing for %dx%d mesh' % (Nx, Ny)
                quadratic(theta, Nx, Ny)
```

### 3.6.6 Verification: convergence rates

For any manufactured solution of the PDE problem we can compute the numerical error and check that this error has the expected dependence on the discretization parameters. Truncation error analysis and other forms of error analysis point to a formula like

$$E = C_t \Delta t^p + C_x \Delta x^2 + C_y \Delta y^2$$

for the error in a 2D problem, where  $p$ ,  $C_t$ ,  $C_x$ , and  $C_y$  are unknown constants. A Crank-Nicolson method has  $p = 2$ , while the Forward and Backward Euler schemes have  $p = 1$ .

When checking the error formula empirically, we need to reduce it to a form  $E = Ch^r$  with a single discretization parameter  $h$  and some rate  $r$  to be estimated. For the Backward Euler method, where  $p = 1$ , we can introduce a single discretization parameter according to

$$h = \Delta x^2 = \Delta y^2, \quad h = K^{-1} \Delta t,$$

where  $K$  is a constant. The error formula then becomes

$$E = C_t K h + C_x h + C_y = \tilde{C} h, \quad \tilde{C} = C_t K + C_x + C_y.$$

The simplest choice is  $K = 1$ , but if we consider the Forward Euler method instead, stability requires  $\Delta t = hK \leq h/(4\alpha)$ , so  $K \leq 1/(4\alpha)$ .

For the Crank-Nicolson method,  $p = 2$ , and we can simply choose

$$h = \Delta x = \Delta y = \Delta t,$$

since there is no restriction on  $\Delta t$  in terms of  $\Delta x$  and  $\Delta y$ .

A frequently used error measure is the  $\ell^2$  norm of the error mesh point values. Section 2.2.3 and the formula (2.26) shows the error measure for a 1D time-dependent problem. The extension to the current 2D problem reads

$$E = \left( \Delta t \Delta x \Delta y \sum_{n=0}^{N_t} \sum_{i=0}^{N_x} \sum_{j=0}^{N_y} (u_e(x_i, y_j, t_n) - u_{i,j}^n)^2 \right)^{\frac{1}{2}}.$$

One attractive manufactured solution is

$$u_e = e^{-pt} \sin(k_x x) \sin(k_y y), \quad k_x = \frac{\pi}{L_x}, k_y = \frac{\pi}{L_y},$$

where  $p$  can be arbitrary. The required source term is

$$f = (\alpha(k_x^2 + k_y^2) - p)u_e$$

The function `convergence_rates` in `diffu2D_u0.py` implements a convergence rate test. Two potential difficulties are important to be aware of:

- The error formula is assumed to be correct when  $h \rightarrow 0$ , so for coarse meshes the estimated rate  $r$  may be somewhat away from the expected value. Fine meshes may lead to prohibitively long execution times.
  - Choosing  $p = \alpha(k_x^2 + k_y^2)$  in the manufactured solution above seems attractive ( $f = 0$ ), but leads to a slower approach to the asymptotic range where the error formula is valid (i.e.,  $r$  fluctuates and needs finer meshes to stabilize).

### 3.6.7 Implementation with a sparse coefficient matrix

We used a sparse matrix implementation in Section 3.2.2 for a 1D problem with a tridiagonal matrix. The present matrix, arising from a 2D problem, has five diagonals, but we can use the same sparse matrix data structure `scipy.sparse.diags`.

**Understanding the diagonals.** Let us look closer at the diagonals in the example with a  $4 \times 3$  mesh as depicted in Figure 3.15 and its associated matrix visualized by dots for zeros and bullets for nonzeros. From the example mesh, we may generalize to an  $N_x \times N_y$  mesh.

The main diagonal has  $N = (N_x + 1)(N_y + 1)$  elements, while the sub- and super-diagonals have  $N - 1$  elements. By looking at the matrix above, we realize that the lower diagonal starts in row  $N_x + 1$  and goes to row  $N$ , so its length is  $N - (N_x + 1)$ . Similarly, the upper diagonal

starts at row 0 and lasts to row  $N - (N_x + 1)$ , so it has the same length. Based on this information, we declare the diagonals by

```
main = np.zeros(N)          # diagonal
lower = np.zeros(N-1)        # subdiagonal
upper = np.zeros(N-1)        # superdiagonal
lower2 = np.zeros(N-(Nx+1))  # lower diagonal
upper2 = np.zeros(N-(Nx+1))  # upper diagonal
b     = np.zeros(N)          # right-hand side
```

**Filling the diagonals.** We run through all mesh points and fill in elements on the various diagonals. The line of mesh points corresponding to  $j = 0$  are all on the boundary, and only the main diagonal gets a contribution:

```
m = lambda i, j: j*(Nx+1) + i
j = 0; main[m(0,j):m(Nx+1,j)] = 1 # j=0 boundary line
```

Then we run through all interior  $j = \text{const}$  lines of mesh points. The first and the last point on each line,  $i = 0$  and  $i = N_x$ , correspond to boundary points:

```
for j in Iy[1:-1]:          # Interior mesh lines j=1,...,Ny-1
    i = 0; main[m(i,j)] = 1
    i = Nx; main[m(i,j)] = 1 # Boundary
```

For the interior mesh points  $i = 1, \dots, N_x - 1$  on a mesh line  $y = \text{const}$  we can start with the main diagonal. The entries to be filled go from  $i = 1$  to  $i = N_x - 1$  so the relevant slice in the `main` vector is `m(1,j):m(Nx,j)`:

```
main[m(1,j):m(Nx,j)] = 1 + 2*theta*(Fx+Fy)
```

The `upper` array for the superdiagonal has its index 0 corresponding to row 0 in the matrix, and the array entries to be set go from  $m(1,j)$  to  $m(N_x - 1, j)$ :

```
upper[m(1,j):m(Nx,j)] = - theta*Fx
```

The subdiagonal (`lower` array), however, has its index 0 corresponding to row 1, so there is an offset of 1 in indices compared to the matrix. The first nonzero occurs (interior point) at a mesh line  $j = \text{const}$  corresponds to matrix row  $m(1, j)$ , and the corresponding array index in `lower` is then  $m(1, j)$ . To fill the entries from  $m(1, j)$  to  $m(N_x - 1, j)$  we set the following slice in `lower`:

```
lower_offset = 1
lower[m(1,j)-lower_offset:m(Nx,j)-lower_offset] = - theta*Fx
```

For the upper diagonal, its index 0 corresponds to matrix row 0, so there is no offset and we can set the entries correspondingly to `upper`:

```
upper2[m(1,j):m(Nx,j)] = - theta*Fy
```

The lower2 diagonal, however, has its first index 0 corresponding to row  $N_x + 1$ , so here we need to subtract the offset  $N_x + 1$ :

```
lower2_offset = Nx+1
lower2[m(1,j)-lower2_offset:m(Nx,j)-lower2_offset] = - theta*Fy
```

We can now summarize the above code lines for setting the entries in the sparse matrix representation of the coefficient matrix:

```
lower_offset = 1
lower2_offset = Nx+1
m = lambda i, j: j*(Nx+1) + i

j = 0; main[m(0,j):m(Nx+1,j)] = 1 # j=0 boundary line
for j in Iy[1:-1]: # Interior mesh lines j=1,...,Ny-1
    i = 0; main[m(i,j)] = 1 # Boundary
    i = Nx; main[m(i,j)] = 1 # Boundary
    # Interior i points: i=1,...,N_x-1
    lower2[m(1,j)-lower2_offset:m(Nx,j)-lower2_offset] = - theta*Fy
    lower[m(1,j)-lower_offset:m(Nx,j)-lower_offset] = - theta*Fx
    main[m(1,j):m(Nx,j)] = 1 + 2*theta*(Fx+Fy)
    upper[m(1,j):m(Nx,j)] = - theta*Fx
    upper2[m(1,j):m(Nx,j)] = - theta*Fy
j = Ny; main[m(0,j):m(Nx+1,j)] = 1 # Boundary line
```

The next task is to create the sparse matrix from these diagonals:

```
import scipy.sparse

A = scipy.sparse.diags(
    diagonals=[main, lower, upper, lower2, upper2],
    offsets=[0, -lower_offset, lower_offset,
             -lower2_offset, lower2_offset],
    shape=(N, N), format='csr')
```

**Filling the right-hand side; scalar version.** Setting the entries in the right-hand side is easier since there are no offsets in the array to take into account. This is in fact similar to the one previously shown when we used a dense matrix representation (the right-hand side vector is, of course, independent of what type of representation we use for the coefficient matrix). The complete time loop goes as follows.

```
import scipy.sparse.linalg

for n in It[0:-1]:
    # Compute b
    j = 0
    for i in Ix:
        p = m(i,j); b[p] = 0 # Boundary
```

```

for j in Iy[1:-1]:
    i = 0; p = m(i,j); b[p] = 0 # Boundary
    for i in Ix[1:-1]:
        p = m(i,j) # Interior
        b[p] = u_1[i,j] + \
            (1-theta)*(
                Fx*(u_1[i+1,j] - 2*u_1[i,j] + u_1[i-1,j]) + \
                Fy*(u_1[i,j+1] - 2*u_1[i,j] + u_1[i,j-1]))\
                + theta*dt*f(i*dx, j*dy, (n+1)*dt) + \
                (1-theta)*dt*f(i*dx, j*dy, n*dt)
    i = Nx; p = m(i,j); b[p] = 0 # Boundary
j = Ny
for i in Ix:
    p = m(i,j); b[p] = 0 # Boundary

# Solve matrix system A*c = b
c = scipy.sparse.linalg.spsolve(A, b)

# Fill u with vector c
for i in Ix:
    for j in Iy:
        u[i,j] = c[m(i,j)]

# Update u_1 before next step
u_1, u = u, u_1

```

**Filling the right-hand side; vectorized version.** Since we use a sparse matrix and try to speed up the computations, we should examine the loops and see if some can be easily removed by vectorization. In the filling of  $A$  we have already used vectorized expressions at each  $j = \text{const}$  line of mesh points. We can very easily do the same in the code above and remove the need for loops over the  $i$  index:

```

for n in It[0:-1]:
    # Compute b, vectorized version

    # Precompute f in array so we can make slices
    f_a_np1 = f(xv, yv, t[n+1])
    f_a_n   = f(xv, yv, t[n])

    j = 0; b[m(0,j):m(Nx+1,j)] = 0 # Boundary
    for j in Iy[1:-1]:
        i = 0; p = m(i,j); b[p] = 0 # Boundary
        i = Nx; p = m(i,j); b[p] = 0 # Boundary
        imin = Ix[1]
        imax = Ix[-1] # for slice, max i index is Ix[-1]-1
        b[m(imin,j):m(imax,j)] = u_1[imin:imax,j] + \
            (1-theta)*(Fx*(
                u_1[imin+1:imax+1,j] -
                2*u_1[imin:imax,j] + \
                u_1[imin-1:imax-1,j]) + \
                Fy*(

```

```

        u_1[imin:imax,j+1] -
    2*u_1[imin:imax,j] +
    u_1[imin:imax,j-1])) + \
    theta*dt*f_a_np1[imin:imax,j] + \
    (1-theta)*dt*f_a_n[imin:imax,j]
j = Ny; b[m(0,j):m(Nx+1,j)] = 0 # Boundary

# Solve matrix system A*c = b
c = scipy.sparse.linalg.spsolve(A, b)

# Fill u with vector c
u[:, :] = c.reshape(Ny+1, Nx+1).T

# Update u_1 before next step
u_1, u = u, u_1

```

The most tricky part of this code snippet is the loading of values in the one-dimensional array `c` into the two-dimensional array `u`. With our numbering of unknowns from left to right along “horizontal” mesh lines, the correct reordering of the one-dimensional array `c` as a two-dimensional array requires first a reshaping as an  $(Ny+1, Nx+1)$  two-dimensional array and then taking the transpose. The result is an  $(Nx+1, Ny+1)$  array compatible with `u` both in size and appearance of the function values.

The `spsolve` function in `scipy.sparse.linalg` is an efficient version of Gaussian elimination suited for matrices described by diagonals. Actually, only the matrix elements within the bands (from `lower2` to `upper2`) are computed with, and these elements constitute only a fraction of all  $N^2$  matrix elements, a crucial property to exploit. The Gaussian elimination algorithm for banded matrices is therefore much faster and requires much less storage than standard Gaussian elimination for a dense matrix. More precisely, with  $b = N_x + 1$  as the *bandwidth* of the matrix [][

**hpl 15:** Problem: if  $N_x \gg N_y$  one should number the unknowns in *y* direction to get a smaller bandwidth.

The complete code is found in the `solver_sparse` function in the file `diffu2D_u0.py`.

**Verification.** We can easily extend the function `quadratic` from Section 3.6.5 to include a test of the `solver_sparse` function as well.

```

def quadratic(theta, Nx, Ny):
    ...
    t, cpu = solver_sparse(
        I, a, f, Lx, Ly, Nx, Ny,
        dt, T, theta, user_action=assert_no_error)

```

### 3.6.8 The Jacobi iterative method

So far we have created a matrix and right-hand side of a linear system  $Ac = b$  and solved the system for  $c$  by calling an exact algorithm based on Gaussian elimination. A much simpler implementation, which requires no memory for the coefficient matrix  $A$ , arises if we solve the system by *iterative* methods. These methods are only approximate, and the core algorithm is repeated many times until the solution is considered to be converged.

**Numerical scheme and linear system.** To illustrate the idea of the Jacobi method, we simplify the numerical scheme to the Backward Euler case,  $\theta = 1$ , so there are fewer terms to write:

$$\begin{aligned} u_{i,j}^{n+1} - \left( F_x(u_{i-1,j}^{n+1} - 2u_{i,j}^{n+1} + u_{i,j}^{n+1}) + F_y(u_{i,j-1}^{n+1} - 2u_{i,j}^{n+1} + u_{i,j+1}^{n+1}) \right) = \\ u_{i,j}^n + \Delta t f_{i,j}^{n+1} \end{aligned} \quad (3.97)$$

The idea of the *Jacobi* iterative method is to introduce an iteration, here with index  $r$ , where we in each iteration treat  $u_{i,j}^{n+1}$  as unknown, but use values from the previous iteration for the other unknowns  $u_{i\pm 1,j\pm 1}^{n+1}$ .

**Iterations.** Let  $u_{i,j}^{n+1,r}$  be the approximation to  $u_{i,j}^{n+1}$  in iteration  $r$ , for all relevant  $i$  and  $j$  indices. We first solve with respect to  $u_{i,j}^{n+1}$  to get the equation to solve:

$$\begin{aligned} u_{i,j}^{n+1} = (1 + 2F_x + 2F_y)^{-1} \left( F_x(u_{i-1,j}^{n+1} + u_{i,j}^{n+1}) + F_y(u_{i,j-1}^{n+1} + u_{i,j+1}^{n+1}) \right) + \\ u_{i,j}^n + \Delta t f_{i,j}^{n+1} \end{aligned} \quad (3.98)$$

The iteration is introduced by using iteration index  $r$ , for computed values, on the right-hand side and  $r + 1$  (unknown in this iteration) on the left-hand side:

$$\begin{aligned} u_{i,j}^{n+1,r+1} = (1 + 2F_x + 2F_y)^{-1} \left( (F_x(u_{i-1,j}^{n+1,r} + u_{i,j}^{n+1,r}) + F_y(u_{i,j-1}^{n+1,r} + u_{i,j+1}^{n+1,r})) + \right. \\ \left. u_{i,j}^n + \Delta t f_{i,j}^{n+1} \right) \end{aligned} \quad (3.99)$$

**Initial guess.** We start the iteration with the computed values at the previous time level:

$$u_{i,j}^{n+1,0} = u_{i,j}^n, \quad i = 0, \dots, N_x, \quad j = 0, \dots, N_y. \quad (3.100)$$

**Relaxation.** A common technique in iterative methods is to introduce a *relaxation*, which means that the new approximation is a weighted mean of the approximation as suggested by the algorithm and the previous approximation. Naming the quantity on the left-hand side of (3.99) as  $u_{i,j}^{n+1,*}$ , a new approximation based on relaxation reads

$$u^{n+1,r+1} = \omega u_{i,j}^{n+1,*} + (1 - \omega) u_{i,j}^{n+1,r}. \quad (3.101)$$

Under-relaxation means  $\omega < 1$ , while over-relaxation has  $\omega > 1$ .

**Stopping criteria.** The iteration can be stopped when the change from one iteration to the next is sufficiently small ( $\epsilon$ ), using either an infinity norm,

$$\max_{i,j} |u_{i,j}^{n+1,r+1} - u_{i,j}^{n+1,r}| \leq \epsilon, \quad (3.102)$$

or an  $L^2$  norm,

$$\left( \Delta x \Delta y \sum_{i,j} (u_{i,j}^{n+1,r+1} - u_{i,j}^{n+1,r})^2 \right)^{\frac{1}{2}} \leq \epsilon. \quad (3.103)$$

Another widely used criterion measures how well the equations are solved by looking at the residual (essentially  $b - Ac^{r+1}$  if  $c^{r+1}$  is the approximation to the solution in iteration  $r + 1$ ). The residual, defined in terms of the finite difference stencil, is

$$\begin{aligned} R_{i,j} = & u_{i,j}^{n+1,r+1} - (F_x(u_{i-1,j}^{n+1,r+1} - 2u_{i,j}^{n+1,r+1} + u_{i+1,j}^{n+1,r+1}) + \\ & F_y(u_{i,j-1}^{n+1,r+1} - 2u_{i,j}^{n+1,r+1} + u_{i,j+1}^{n+1,r+1})) - \\ & u_{i,j}^n - \Delta t f_{i,j}^{n+1} \end{aligned} \quad (3.104)$$

One can then iterate until the norm of the mesh function  $R_{i,j}$  is less than some tolerance:

$$\left( \Delta x \Delta y \sum_{i,j} R_{i,j}^2 \right)^{\frac{1}{2}} \leq \epsilon. \quad (3.105)$$

**Code-friendly notation.** To make the mathematics as close as possible to what we will write in a computer program, we may introduce some new notation:  $u_{i,j}$  is a short notation for  $u_{i,j}^{n+1,r+1}$ ,  $u_{i,j}^-$  is a short notation for  $u_{i,j}^{n+1,r}$ , and  $u_{i,j}^{(s)}$  denotes  $u_{i,j}^{n+1-s}$ . That is,  $u_{i,j}$  is the unknown,  $u_{i,j}^-$

is its most recently computed approximation, and  $s$  counts time levels backwards in time. The Jacobi method (ref(3.99)) takes the following form with the new notation:

$$u_{i,j}^* = (1 + 2F_x + 2F_y)^{-1}((F_x(u_{i-1,j}^- + u_{i,j}^-) + F_y(u_{i,j-1}^{n+1,r} + u_{i,j+1}^{n+1,r})) + \\ u_{i,j}^{(1)} + \Delta t f_{i,j}^{n+1}) \quad (3.106)$$

**Generalization of the scheme.** We can also quite easily introduce the  $\theta$  rule for discretization in time and write up the Jacobi iteration in that case as well:

$$u_{i,j}^* = (1 + 2\theta(F_x + F_y))^{-1}(\theta(F_x(u_{i-1,j}^- + u_{i,j}^-) + F_y(u_{i,j-1}^- + u_{i,j+1}^-)) + \\ u_{i,j}^{(1)} + \theta \Delta t f_{i,j}^{n+1} + (1 - \theta) \Delta t f_{i,j}^n + \\ (1 - \theta)(F_x(u_{i-1,j}^{(1)} - 2u_{i,j}^{(1)} + u_{i+1,j}^{(1)}) + F_y(u_{i,j-1}^{(1)} - 2u_{i,j}^{(1)} + u_{i,j+1}^{(1)}))) . \quad (3.107)$$

The final update of  $u$  applies relaxation:

$$u_{i,j} = \omega u_{i,j}^* + (1 - \omega)u_{i,j}^- .$$

### 3.6.9 Implementation of the Jacobi method

The Jacobi method needs no coefficient matrix and right-hand side vector, but it needs an array for  $u$  in the previous iteration. We call this array  $u_-$ , using the notation at the end of the previous section (at the same time level). The unknown itself is called  $u$ , while  $u_-$  is the computed solution one time level back in time. With a  $\theta$  rule in time, the time loop can be coded like this:

```
for n in It[0:-1]:
    # Solve linear system by Jacobi iteration at time level n+1
    u_[:, :] = u_- # Start value
    converged = False
    r = 0
    while not converged:
        if version == 'scalar':
            j = 0
            for i in Ix:
                u[i, j] = U_Oy(t[n+1]) # Boundary
                for j in Iy[1:-1]:
```

```

    i = 0;   u[i,j] = U_Ox(t[n+1])  # Boundary
    i = Nx;  u[i,j] = U_Lx(t[n+1])  # Boundary
# Interior points
    for i in Ix[1:-1]:
        u_new = 1.0/(1.0 + 2*theta*(Fx + Fy))*(theta*(
            Fx*(u_[i+1,j] + u_[i-1,j]) +
            Fy*(u_[i,j+1] + u_[i,j-1])) + \
            u_1[i,j] + \
            (1-theta)*(Fx*(
                u_1[i+1,j] - 2*u_1[i,j] + u_1[i-1,j]) + \
                Fy*(
                    u_1[i,j+1] - 2*u_1[i,j] + u_1[i,j-1])))\ \
            + theta*dt*f(i*dx,j*dy,(n+1)*dt) + \
            (1-theta)*dt*f(i*dx,j*dy,n*dt))
        u[i,j] = omega*u_new + (1-omega)*u_[i,j]

    j = Ny
    for i in Ix:
        u[i,j] = U_Ly(t[n+1])      # Boundary

elif version == 'vectorized':
    j = 0;  u[:,j] = U_Oy(t[n+1])  # Boundary
    i = 0;  u[i,:] = U_Ox(t[n+1])  # Boundary
    i = Nx; u[i,:] = U_Lx(t[n+1])  # Boundary
    j = Ny; u[:,j] = U_Ly(t[n+1])  # Boundary
# Internal points
    f_a_np1 = f(xv, yv, t[n+1])
    f_a_n   = f(xv, yv, t[n])
    u_new = 1.0/(1.0 + 2*theta*(Fx + Fy))*(theta*(Fx*(
        u_[2:,1:-1] + u_[:-2,1:-1]) +
        Fy*(
            u_[1:-1,2:] + u_[1:-1,:-2])) +\
        u_1[1:-1,1:-1] + \
        (1-theta)*(Fx*(
            u_1[2:,1:-1] - 2*u_1[1:-1,1:-1] + u_1[:-2,1:-1]) +\
            Fy*(
                u_1[1:-1,2:] - 2*u_1[1:-1,1:-1] + u_1[1:-1,:-2])))\ \
        + theta*dt*f_a_np1[1:-1,1:-1] + \
        (1-theta)*dt*f_a_n[1:-1,1:-1])
    u[1:-1,1:-1] = omega*u_new + (1-omega)*u_[1:-1,1:-1]
    r += 1
    converged = np.abs(u-u_).max() < tol or r >= max_iter
    u_[:, :] = u

# Update u_1 before next step
    u_1, u = u, u_1

```

The vectorized version should be quite straightforward to understand once one has an understanding of how a standard 2D finite stencil is vectorized.

### **hpl 16:** Make references to 1D and 2D wave equation vectorization.

The first natural verification is to use the test problem from in the function `quadratic` from Section 3.6.5. This problem is known to have no approximation error, but any iterative method will produce an ap-

proximate solution with unknown error. For a tolerance  $10^{-k}$  in the iterative method, we can, e.g., use a slightly larger tolerance  $10^{-(k-1)}$  for the difference between the exact and the computed solution.

```
def quadratic(theta, Nx, Ny):
    ...
    def assert_small_error(u, x, xv, y, yv, t, n):
        """Assert small error for iterative methods."""
        u_e = u_exact(xv, yv, t[n])
        diff = abs(u - u_e).max()
        tol = 1E-4
        msg = 'diff=%g, step %d, time=%g' % (diff, n, t[n])
        assert diff < tol, msg

    for version in 'scalar', 'vectorized':
        for theta in 1, 0.5:
            print 'testing Jacobi, %s version, theta=%g' % \
                (version, theta)
            t, cpu = solver_Jacobi(
                I=I, a=a, f=f, Lx=Lx, Ly=Ly, Nx=Nx, Ny=Ny,
                dt=dt, T=T, theta=theta,
                U_0x=0, U_0y=0, U_Lx=0, U_Ly=0,
                user_action=assert_small_error,
                version=version, iteration='Jacobi',
                omega=1.0, max_iter=100, tol=1E-5)
```

Even for a very coarse  $4 \times 4$  mesh, the Jacobi method requires 26 iterations to reach a tolerance of  $10^{-5}$ , which is quite many iterations, given that there are only 25 unknowns.

### 3.6.10 Test problem: diffusion of a sine hill

It can be shown that

$$u_e = A e^{-\alpha\pi^2(L_x^{-2} + L_y^{-2})} \sin\left(\frac{\pi}{L_x}x\right) \sin\left(\frac{\pi}{L_y}y\right), \quad (3.108)$$

is a solution of the 2D homogeneous diffusion equation  $u_t = \alpha(u_{xx} + u_{yy})$  in a rectangle  $[0, L_x] \times [0, L_y]$ , for any value of the amplitude  $A$ . This solution vanishes at the boundaries, and the initial condition is the product of two sines. We may choose  $A = 1$  for simplicity.

It is difficult to know if our Jacobi method works properly since we are faced with two sources of errors: one from the discretization,  $E_\Delta$ , and one from the iterative Jacobi method,  $E_i$ . The total error in the computed  $u$  can be represented as

$$E_u = E_\Delta + E_i.$$

One error measure is to look at the maximum value, which is obtained for the midpoint  $x = L_x/2$  and  $y = L_y/2$ . This midpoint is represented in the discrete  $\mathbf{u}$  if  $N_x$  and  $N_y$  are even numbers. We can then compute  $E_u$  as  $E_u = |\max u_e - \max u|$ , when we know an exact solution  $u_e$  of the problem.

What about  $E_\Delta$ ? If we use the maximum value as a measure of the error, we have in fact analytical insight into the approximation error in this particular problem. According to Section 3.3.8, the exact solution (3.108) of the PDE problem is also an exact solution of the discrete equations, except that the damping factor in time is different. More precisely, (3.66) and (3.67) are solutions of the discrete problem for  $\theta = 1$  (Backward Euler) and  $\theta = \frac{1}{2}$  (Crank-Nicolson), respectively. The factors raised to the power  $n$  is the numerical amplitude, and the errors in these factors become

$$E_\Delta = e^{-\alpha k^2 t} - \left( \frac{1 - 2(F_x \sin^2 p_x + F_y \sin^2 p_y)}{1 + 2(F_x \sin^2 p_x + F_y \sin^2 p_y)} \right)^n, \quad \theta = \frac{1}{2},$$

$$E_\Delta = e^{-\alpha k^2 t} - (1 + 4F_x \sin^2 p_x + 4F_y \sin^2 p_y)^{-n}, \quad \theta = 1.$$

We are now in a position to compute  $E_i$  numerically. That is, we can compute the error due to iterative solution of the linear system and see if it corresponds to the convergence tolerance used in the method. Note that the convergence is based on measuring the difference in two consecutive approximations, which is not exactly the error due to the iteration, but it is a kind of measure, and it should have about the same size as  $E_i$ .

The function `demo_classic_iterative` in `diffu2D_u0.py` implements the idea above (also for the methods in Section 3.6.12). The value of  $E_i$  is in particular printed at each time level. By changing the tolerance in the convergence criterion in the Jacobi method, we can see that  $E_i$  is of the same order of magnitude as the prescribed tolerance in the Jacobi method. For example:  $E_\Delta \sim 10^{-2}$  with  $N_x = N_y = 10$  and  $\theta = \frac{1}{2}$ , as long as  $\max u$  has some significant size ( $\max u > 0.02$ ). An appropriate value of the tolerance is then  $10^{-3}$ , such that the error in the Jacobi method does not become bigger than the discretization error. In that case,  $E_i$  is around  $5 \cdot 10^{-3}$ . The corresponding number of Jacobi iterations (with  $\omega = 1$ ) varies from 31 to 12 during the time

simulation (for  $\max u > 0.02$ ). Changing the tolerance to  $10^{-5}$  causes many more iterations (61 to 42) without giving any contribution to the overall accuracy, because the total error is dominated by  $E_\Delta$ .

Also, with a  $N_x = N_y = 20$ , the spatial accuracy increases and many more iterations are needed (143 to 45), but the dominating error is from the time discretization. However, with such a finer spatial mesh, a higher tolerance in the convergence criterion  $10^{-4}$  is needed to keep  $E_i \sim 10^{-3}$ . More experiments show the disadvantage of the very simple Jacobi iteration method: the number of iterations increases with the number of unknowns, keeping the tolerance fixed, but the tolerance should also be lowered to avoid the iteration error to dominate the total error. A small adjustment of the Jacobi method, as described in Section 3.6.12, provides a better method.

### 3.6.11 The relaxed Jacobi method and its relation to the Forward Euler method

We shall now show that solving the Poisson equation  $-\alpha \nabla^2 u = f$  by the Jacobi iterative method is in fact equivalent to using a Forward Euler scheme on  $u_t = \alpha \nabla^2 u + f$  and letting  $t \rightarrow \infty$ .

A Forward Euler discretization of the 2D diffusion equation,

$$[D_t^+ u = \alpha(D_x D_x u + D_y D_y u) + f]_{i,j}^n,$$

can be written out as

$$u^{n+1} = u^n + \frac{\Delta t}{\alpha h^2} (u_{i-1,j}^n + u_{i+1,j}^n + u_{i,j-1}^n + u_{i,j+1}^n - 4u_{i,j}^n + h^2 f_{i,j}),$$

where  $h = \Delta x = \Delta y$  has been introduced for simplicity. The scheme can be reordered as

$$u^{n+1} = (1 - \omega) u_{i,j}^n + \frac{1}{4} \omega (u_{i-1,j}^n + u_{i+1,j}^n + u_{i,j-1}^n + u_{i,j+1}^n - 4u_{i,j}^n + h^2 f_{i,j}),$$

with

$$\omega = 4 \frac{\Delta t}{\alpha h^2},$$

but this latter form is nothing but the relaxed Jacobi method applied to

$$[D_x D_x u + D_y D_y u = -f]_{i,j}^n.$$

From the equivalence above we know a couple of things about the Jacobi method for solving  $-\nabla^2 u = f$ :

1. The method is unstable if  $\omega > 1$  (since the Forward Euler method is then unstable).
2. The convergence is really slow as the iteration index increases (coming from the fact that the Forward Euler scheme requires many small time steps to reach the stationary solution).

These observations are quite disappointing: if we already have a time-dependent diffusion problem and want to take larger time steps by an implicit time discretization method, we will with the Jacobi method end up with something close to a slow Forward Euler simulation of the original problem at each time level. Nevertheless, there are two reasons for why the Jacobi method remains a fundamental building block for solving linear systems arising from PDEs: 1) a couple of iterations remove large parts of the error and this is effectively used in the very efficient class of multigrid methods; and 2) the idea of the Jacobi method can be developed into more efficient methods, especially the SOR method, which is treated next.

### 3.6.12 The Gauss-Seidel and SOR methods

If we update the mesh points according to the Jacobi method (3.98) for a Backward Euler discretization with a loop over  $i = 1, \dots, N_x - 1$  and  $j = 1, \dots, N_y - 1$ , we realize that when  $u_{i,j}^{n+1,r+1}$  is computed,  $u_{i-1,j}^{n+1,r+1}$  and  $u_{i,j-1}^{n+1,r+1}$  are already computed, so these new values can be used rather than  $u_{i-1,j}^{n+1,r}$  and  $u_{i,j-1}^{n+1,r}$  (respectively) in the formula for  $u_{i,j}^{n+1,r+1}$ . This idea gives rise to the *Gauss-Seidel* iteration method, which mathematically is just a small adjustment of (3.98):

$$\begin{aligned} u_{i,j}^{n+1,r+1} &= (1 + 2F_x + 2F_y)^{-1}(( \\ &\quad F_x(u_{i-1,j}^{n+1,r+1} + u_{i,j}^{n+1,r}) + F_y(u_{i,j-1}^{n+1,r+1} + u_{i,j+1}^{n+1,r})) + u_{i,j}^n + \Delta t f_{i,j}^{n+1}). \end{aligned} \tag{3.109}$$

Observe that the way we access the mesh points in the formula (3.109) is important: points with  $i - 1$  must be computed before points with  $i$ , and points with  $j - 1$  must be computed before points with  $j$ . Any sequence of mesh points can be used in the Gauss-Seidel method, but the particular math formula must distinguish between already visited points in the current iteration and the points not yet visited.

The idea of relaxation (3.101) can equally well be applied to the Gauss-Seidel method. Actually, the Gauss-Seidel method with an arbitrary  $0 < \omega \leq 2$  has its own name: the *Successive Over-Relaxation* method, abbreviated as SOR.

The SOR method for a  $\theta$  rule discretization, with the shortened  $u$  and  $u^-$  notation, can be written

$$\begin{aligned} u_{i,j}^* = & (1 + 2\theta(F_x + F_y))^{-1}(\theta(F_x(u_{i-1,j} + u_{i,j}^-) + F_y(u_{i,j-1} + u_{i,j+1}^-)) + \\ & u_{i,j}^{(1)} + \theta \Delta t f_{i,j}^{n+1} + (1 - \theta) \Delta t f_{i,j}^n + \\ & (1 - \theta)(F_x(u_{i-1,j}^{(1)} - 2u_{i,j}^{(1)} + u_{i+1,j}^{(1)}) + F_y(u_{i,j-1}^{(1)} - 2u_{i,j}^{(1)} + u_{i,j+1}^{(1)}))), \end{aligned} \quad (3.110)$$

$$u_{i,j} = \omega u_{i,j}^* + (1 - \omega)u_{i,j}^- \quad (3.111)$$

The sequence of mesh points in (3.110) is  $i = 1, \dots, N_x - 1$ ,  $j = 1, \dots, N_y - 1$  (but whether  $i$  runs faster or slower than  $j$  does not matter).

### 3.6.13 Scalar implementation of the SOR method

Since the Jacobi and Gauss-Seidel methods with relaxation are so similar, we can easily make a common code for the two:

```
for n in It[0:-1]:
    # Solve linear system by Jacobi/SOR iteration at time level n+1
    u_[:, :] = u_1  # Start value
    converged = False
    r = 0
    while not converged:
        if version == 'scalar':
            if iteration == 'Jacobi':
                u__ = u_
            elif iteration == 'SOR':
                u__ = u
            j = 0
            for i in Ix:
```

```

        u[i,j] = U_Oy(t[n+1]) # Boundary
    for j in Iy[1:-1]:
        i = 0; u[i,j] = U_Ox(t[n+1]) # Boundary
        i = Nx; u[i,j] = U_Lx(t[n+1]) # Boundary
    for i in Ix[1:-1]:
        u_new = 1.0/(1.0 + 2*theta*(Fx + Fy))*(theta*(  

            Fx*(u_[i+1,j] + u_1[i-1,j]) +  

            Fy*(u_[i,j+1] + u_1[i,j-1])) + \  

            u_1[i,j] + (1-theta)*(  

            Fx*(  

            u_1[i+1,j] - 2*u_1[i,j] + u_1[i-1,j]) +  

            Fy*(  

            u_1[i,j+1] - 2*u_1[i,j] + u_1[i,j-1]))\ \  

            + theta*dt*f(i*dx, j*dy, (n+1)*dt) + \  

            (1-theta)*dt*f(i*dx, j*dy, n*dt))  

        u[i,j] = omega*u_new + (1-omega)*u_[i,j]
    j = Ny
    for i in Ix:
        u[i,j] = U_Ly(t[n+1]) # boundary
    r += 1
    converged = np.abs(u-u_).max() < tol or r >= max_iter
    u_[:, :] = u

    u_1, u = u, u_1 # Get ready for next iteration

```

The idea here is to introduce  $u_{-}$  to be used for already computed values ( $u$ ) in the Gauss-Seidel/SOR version of the implementation, or just values from the previous iteration ( $u_-$ ) in case of the Jacobi method.

### 3.6.14 Vectorized implementation of the SOR method

Vectorizing the Gauss-Seidel iteration step turns out to be non-trivial. The problem is that vectorized operations typically imply operations on arrays where the sequence we visit the elements in does not matter. In particular, this principle makes vectorized code trivial to parallelize. However, in the Gauss-Seidel algorithm the sequence we visit the elements in the arrays does matter, and it is well known that the basic method as explained above cannot be parallelized. Therefore, also vectorization will require new thinking.

The strategy for vectorizing (and parallelizing) the Gauss-Seidel method is to use a special numbering of the mesh points called red-black numbering: every other point is red or black as in a checkerboard pattern. This numbering requires  $N_x$  and  $N_y$  to be even numbers. Here is an example of a  $6 \times 6$  mesh:

```
r b r b r b r  
b r b r b r b
```

```
r b r b r b r
b r b r b r b
r b r b r b r
b r b r b r b
r b r b r b r
```

The idea now is to first update all the red points. Each formula for updating a red point involves only the black neighbors. Thereafter, we update all the black points, and at each black point, only the recently computed red points are involved.

The scalar implementation of the red-black numbered Gauss-Seidel method is really compact, since we can update values directly in  $\mathbf{u}$  (that guarantees that we use the most recently computed values). Here is the relevant code for the Backward Euler scheme in time and without a source term:

```
# Update internal points
for sweep in 'red', 'black':
    for j in range(1, Ny, 1):
        if sweep == 'red':
            start = 1 if j % 2 == 1 else 2
        elif sweep == 'black':
            start = 2 if j % 2 == 1 else 1
    for i in range(start, Nx, 2):
        u[i,j] = 1.0/(1.0 + 2*(Fx + Fy))*(  

                    Fx*(u[i+1,j] + u[i-1,j]) +  

                    Fy*(u[i,j+1] + u[i,j-1]) + u_1[i,j])
```

The vectorized version must be based on slices. Looking at a typical red-black pattern, e.g.,

```
r b r b r b r
b r b r b r b
r b r b r b r
b r b r b r b
r b r b r b r
b r b r b r b
r b r b r b r
```

we want to update the internal points (marking boundary points with x):

```
x x x x x x x
x r b r b r x
x b r b r b x
x r b r b r x
x b r b r b x
x r b r b r x
x x x x x x x
```

It is impossible to make one slice that picks out all the internal red points. Instead, we need two slices. The first involves points marked with R:

```
x x x x x x x x
x R b R b R x
x b r b r b x
x R b R b R x
x b r b r b x
x R b R b R x
x x x x x x x
```

This slice is specified as `1::2` for `i` and `1::2` for `j`, or with `slice` objects:

```
i = slice(1, None, 2); j = slice(1, None, 2)
```

The second slice involves the red points with R:

```
x x x x x x x x
x r b r b r x
x b R b R b x
x r b r b r x
x b R b R b x
x r b r b r x
x x x x x x x
```

The slices are

```
i = slice(2, None, 2); j = slice(2, None, 2)
```

For the black points, the first slice involves the B points:

```
x x x x x x x x
x r B r B r x
x b r b r b x
x r B r B r x
x b r b r b x
x r B r B r x
x x x x x x x
```

with slice objects

```
i = slice(2, None, 2); j = slice(1, None, 2)
```

The second set of black points is shown here:

```
x x x x x x x x
x r b r b r x
x B r B r B x
x r b r b r x
x B r B r B x
x r b r b r x
x x x x x x x
```

with slice objects

```
i = slice(1, None, 2); j = slice(2, None, 2)
```

That is, we need four sets of slices. The simplest way of implementing the algorithm is to make a function with variables for the slices representing  $i$ ,  $i - 1$ ,  $i + 1$ ,  $j$ ,  $j - 1$ , and  $j + 1$ , here called `ic` (“ $i$  center”), `im1` (“ $i$  minus 1”), `ip1` (“ $i$  plus 1”), `jc`, `jm1`, and `jp1`, respectively.

```
def update(u_, u_1, ic, im1, ip1, jc, jm1, jp1):
    return \
        1.0/(1.0 + 2*theta*(Fx + Fy))*(theta*(
            Fx*(u_[ip1,jc] + u_[im1,jc]) +
            Fy*(u_[ic,jp1] + u_[ic,jm1])) +\
        u_1[ic,jc] + (1-theta)*(
            Fx*(u_1[ip1,jc] - 2*u_1[ic,jc] + u_1[im1,jc]) +\
            Fy*(u_1[ic,jp1] - 2*u_1[ic,jc] + u_1[ic,jm1]))+\
            theta*dt*f_a_np1[ic,jc] + \
            (1-theta)*dt*f_a_n[ic,jc])
```

The formula returned from `update` is to be compared with (3.110).

The relaxed Jacobi iteration can be implemented by

```
ic = jc = slice(1,-1)
im1 = jm1 = slice(0,-2)
ip1 = jp1 = slice(2,None)
u_new[ic,jc] = update(
    u_, u_1, ic, im1, ip1, jc, jm1, jp1)
u[ic,jc] = omega*u_new[ic,jc] + (1-omega)*u_[ic,jc]
```

The Gauss-Seidel (or SOR) updates need four different steps. The `ic` and `jc` slices are specified above. For each of these, we must specify the corresponding `im1`, `ip1`, `jm1`, and `jp1` slices. The code below contain the details.

```
# Red points
ic = slice(1,-1,2)
im1 = slice(0,-2,2)
ip1 = slice(2,None,2)
jc = slice(1,-1,2)
jm1 = slice(0,-2,2)
jp1 = slice(2,None,2)
u_new[ic,jc] = update(
    u_new, u_1, ic, im1, ip1, jc, jm1, jp1)

ic = slice(2,-1,2)
im1 = slice(1,-2,2)
ip1 = slice(3,None,2)
jc = slice(2,-1,2)
jm1 = slice(1,-2,2)
jp1 = slice(3,None,2)
u_new[ic,jc] = update(
    u_new, u_1, ic, im1, ip1, jc, jm1, jp1)
```

```

# Black points
ic  = slice(2,-1,2)
im1 = slice(1,-2,2)
ip1 = slice(3,None,2)
jc  = slice(1,-1,2)
jm1 = slice(0,-2,2)
jp1 = slice(2,None,2)
u_new[ic,jc] = update(
    u_new, u_1, ic, im1, ip1, jc, jm1, jp1)

ic  = slice(1,-1,2)
im1 = slice(0,-2,2)
ip1 = slice(2,None,2)
jc  = slice(2,-1,2)
jm1 = slice(1,-2,2)
jp1 = slice(3,None,2)
u_new[ic,jc] = update(
    u_new, u_1, ic, im1, ip1, jc, jm1, jp1)

# Relax
c = slice(1,-1)
u[c,c] = omega*u_new[c,c] + (1-omega)*u_[c,c]

```

The function `solver_classic_iterative` in `diffu2D_u0.py` contains a unified implementation of the relaxed Jacobi and SOR methods in scalar and vectorized versions using the techniques explained above.

**hpl 17:** Experiments with SOR, benefits. Rerun the Jacobi experiments. Just mention the most important conclusions. We can have more comprehensive experiments with comparison of all methods in a later section, after CG with preconditioning.

### 3.6.15 Direct versus iterative methods

**Direct methods.** There are two classes of methods for solving linear systems: direct methods and iterative methods. Direct methods are based on variants of the Gaussian elimination procedure and will produce an exact solution (in exact arithmetics) in an a priori known number of steps. Iterative methods, on the other hand, produce an approximate solution, and the amount of work for reaching a given accuracy is usually not known.

The most common direct method today is to use the *LU factorization* procedure to factor the coefficient matrix  $A$  as the product of a lower-triangular matrix  $L$  (with unit diagonal terms) and an upper-triangular matrix  $U$ :  $A = LU$ . As soon as we have  $L$  and  $U$ , a system of equations  $LUc = b$  is easy to solve because of the triangular nature of  $L$  and  $U$ . We

first solve  $Ly = b$  for  $y$  (forward substitution), and thereafter we find  $c$  from solving  $Uc = y$  (backward substitution). When  $A$  is a dense  $N \times N$  matrix, the LU factorization costs  $\frac{1}{3}N^3$  arithmetic operations, while the forward and backward substitution steps each require of the order  $N^2$  arithmetic operations. That is, factorization dominates the costs, while the substitution steps are cheap.

Symmetric, positive definite coefficient matrices often arise when discretizing PDEs. In this case, the LU factorization becomes  $A = LL^T$ , and the associated algorithm is known as *Cholesky factorization*. Most linear algebra software offers highly optimized implementations of LU and Cholesky factorization as well as forward and backward substitution (`scipy.linalg` is the relevant Python package).

Finite difference discretizations lead to sparse coefficient matrices. An extreme case arose in Section 3.2.1 where  $A$  is tridiagonal. For a tridiagonal matrix, the amount of arithmetic operations in the LU and Cholesky factorization algorithms is just of the order  $N$ , not  $N^3$ . Tridiagonal matrices are special cases of *banded matrices*, where the matrices contain just a set of diagonal bands. Finite difference methods on regularly numbered rectangular and box-shaped meshes give rise to such banded matrices, with 5 bands in 2D and 7 in 3D for diffusion problems. Gaussian elimination only needs to work within the bands, leading to much more efficient algorithms. For example, factorization of more general, sparse matrices can often take advantage of the sparsity through modified Gaussian elimination algorithms. The relevant Python package is `scipy.sparse.linalg`.

Although a direct method is an exact algorithm, rounding errors may in practice accumulate and pollute the solution. The effect grows with the size of the linear system, so both for accuracy and efficiency, iterative methods are better suited than direct methods for solving really large linear systems.

**Iterative methods.** The Jacobi and SOR iterative methods belong to a class of iterative methods where the idea is to solve  $Au = b$  by splitting  $A$  into two parts,  $A = M - N$ , such that solving systems  $Mu = c$  is easy and efficient. With the splitting, we get a system

$$Mu = Nu + b,$$

which suggests an iterative method

$$Mu^{r+1} = Nu^r + b, \quad r = 0, 1, 2, \dots,$$

where  $u^{r+1}$  is a new approximation to  $u$  in the  $r + 1$ -th iteration. To initiate the iteration, we need a start vector  $u^0$ .

The Jacobi and SOR methods are based on splitting  $A$  into a lower tridiagonal part  $L$ , the diagonal  $D$ , and an upper tridiagonal part  $U$ , such that  $A = L + D + U$ . The Jacobi method corresponds to  $M = D$  and  $N = -L - U$ . The Gauss-Seidel method employs  $M = L + D$  and  $N = -U$ , while the SOR method corresponds to

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

The relaxed Jacobi method has similar expressions:

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

With the matrix forms of the Jacobi and SOR methods as written above, we could in an implementation alternatively fill the matrix  $A$  with entries and call general implementations of the Jacobi or SOR methods that work on a system  $Au = b$ . However, this is almost never done since forming the matrix  $A$  requires quite some code and storing  $A$  in the computer's memory is unnecessary. It is much easier to just apply the Jacobi and SOR ideas to the finite difference stencils directly in an implementation, as we have shown in detail.

Nevertheless, the matrix formulation of the Jacobi and SOR methods have been important for analyzing their convergence behavior. One can show that the error  $u^r - u$  fulfills  $u^r - u = G^r(u^0 - u)$ , where  $G = M^{-1}N$  and  $G^k$  is a matrix exponential. For the method to converge,  $\lim_{r \rightarrow \infty} \|G^r\| = 0$  is a necessary and sufficient condition. This implies that the *spectral radius* of  $G$  must be less than one. Since  $G$  is directly related to the finite difference scheme for the underlying PDE problem, one can in principle compute the spectral radius. For a given PDE problem, however, this is not a practical strategy, since it is very difficult to develop useful formulas. Analysis of model problems, usually related to the Poisson equation, reveals some trends of interest: the convergence rate of the Jacobi method goes like  $h^2$ , while that of SOR with an optimal  $\omega$  goes like  $h$ , where  $h$  is the spatial spacing:  $h = \Delta x = \Delta y$ . That is, the efficiency of the Jacobi method quickly deteriorates with the increasing mesh resolution, and SOR is much to be preferred (even if the optimal  $\omega$  remains an open question). We refer to Chapter 4 of [8] for more information on the convergence theory. One important result is that if

$A$  is symmetric and positive definite, then SOR will converge for any  $0 < \omega < 2$ .

The optimal  $\omega$  parameter can be theoretically established for a Poisson problem as

$$\omega_o = \frac{2}{1 + \sqrt{1 - \varrho^2}}, \quad \varrho = \frac{\cos(\pi/N_x) + (\Delta x/\Delta y)^2 \cos(\pi/N_y)}{1 + (\Delta x/\Delta y)^2}. \quad (3.112)$$

This formula can be used as a guide also in other problems.

The Jacobi and the SOR methods have their great advantage of being trivial to implement, so they are obviously popular of this reason. However, the slow convergence of these methods limits the popularity to fairly small linear systems (i.e., coarse meshes). As soon as the matrix size grows, one is better off with more sophisticated iterative methods like the preconditioned Conjugate gradient method, which we now turn to.

Finally, we mention that there is a variant of the SOR method, called Symmetric Successive Overrelaxation method, known as SSOR, where one runs a standard SOR sweep through the mesh points and then a new sweep but visiting the points in reverse order.

**hpl 18:** Line Jacobi/SOR?

### 3.6.16 The Conjugate gradient method

There is no simple intuitive derivation of the Conjugate gradient method, so we refer to the many excellent expositions in the literature for the idea of the method and how the algorithm is derived. In particular, we recommend the books [1, 8, 2]. A brief overview is provided in the [Wikipedia article](#). Here, we just state the pros and cons of the method from a user's perspective and how we utilize it in code.

The original Conjugate gradient method is limited to linear systems  $Au = b$ , where  $A$  is a symmetric and positive definite matrix. There are, however, extensions of the method to non-symmetric matrices.

A major advantage of all conjugate gradient methods is that the matrix  $A$  is only used in matrix-vector products, so we do not need form and store  $A$  if we can provide code for computing a matrix-vector product  $Au$ . Another important feature is that the algorithm is very easy to vectorize and parallelize. The primary downside of the method is that it converges slowly unless one has an effective *preconditioner* for the

system. That is, instead of solving  $Au = b$ , we try to solve  $M^{-1}Au = M^{-1}b$  in the hope that the method works better for this *preconditioned* system. The matrix  $M$  is the *preconditioner* or preconditioning matrix. Now we need to perform matrix-vector products  $y = M^{-1}Au$ , which is done in two steps: first the matrix-vector product  $v = Au$  is carried out and then the system  $My = v$  must be solved. Therefore,  $M$  must be cheap to compute and systems  $My = v$  must be cheap to solve.

A perfect preconditioner is  $M = A$ , but in each iteration in the Conjugate gradient method one then has to solve a system with  $A$  as coefficient matrix! A key idea is to let  $M$  be some kind of *cheap approximation* to  $A$ . The simplest preconditioner is to set  $M = D$ , where  $D$  is the diagonal of  $A$ . This choice means running one Jacobi iteration as preconditioner. Exercise 3.8 shows that the Jacobi and SOR methods can also be viewed as preconditioners.

Constructing good preconditioners is a scientific field on its own. Here we shall treat the topic just very briefly. For a user having access to the `scipy.sparse.linalg` library, there are Conjugate gradient methods and preconditioners readily available:

- For positive definite, symmetric systems: `cg` (the Conjugate gradient method)
- For symmetric systems: `minres` (Minimum residual method)
- For non-symmetric systems:
  - `gmres` (GMRES: Generalized minimum residual method)
  - `bicg` (BiConjugate gradient method)
  - `bicgstab` (Stabilized BiConjugate gradient method)
  - `cgs` (Conjugate gradient squared method)
  - `qmr` (Quasi-minimal residual iteration)
- Preconditioner: `spilu` (Sparse, incomplete LU factorization)

The ILU preconditioner is an attractive all-round type of preconditioner that is suitable for most problems on serial computers. A more efficient preconditioner is the multigrid method, and algebraic multigrid is also an all-round choice as preconditioner. The Python package `PyAMG` offers efficient implementations of the algebraic multigrid method, to be used both as a preconditioner and as a stand-alone iterative method.

The matrix arising from implicit time discretization methods of the diffusion equation is symmetric and positive definite so we can use the Conjugate gradient method (`cg`), typically in combination with an ILU preconditioner. The code is very similar to the one we created when

solving the linear system by sparse Gaussian elimination, the main difference is that we now allow for calling up the Conjugate gradient function as an alternative solver.

```

def solver_sparse(
    I, a, f, Lx, Ly, Nx, Ny, dt, T, theta=0.5,
    U_Ox=0, U_Oy=0, U_Lx=0, U_Ly=0, user_action=None,
    method='direct', CG_prec='ILU', CG_tol=1E-5):
    """
        Full solver for the model problem using the theta-rule
        difference approximation in time. Sparse matrix with
        dedicated Gaussian elimination algorithm (method='direct')
        or ILU preconditioned Conjugate Gradients (method='CG' with
        tolerance CG_tol and preconditioner CG_prec ('ILU' or None)).
    """
    # Set up data structures as shown before

    # Precompute sparse matrix
    ...

    A = scipy.sparse.diags(
        diagonals=[main, lower, upper, lower2, upper2],
        offsets=[0, -lower_offset, lower_offset,
                  -lower2_offset, lower2_offset],
        shape=(N, N), format='csc')

    if method == 'CG':
        if CG_prec == 'ILU':
            # Find ILU preconditioner (constant in time)
            A_ilu = scipy.sparse.linalg.spilu(A) # SuperLU defaults
            M = scipy.sparse.linalg.LinearOperator(
                shape=(N, N), matvec=A_ilu.solve)
        else:
            M = None
        CG_iter = [] # No of CG iterations at time level n

    # Time loop
    for n in It[0:-1]:
        # Compute b, vectorized version

        # Solve matrix system A*c = b
        if method == 'direct':
            c = scipy.sparse.linalg.spsolve(A, b)
        elif method == 'CG':
            x0 = u_1.T.reshape(N) # Start vector is u_1
            CG_iter.append(0)

            def CG_callback(c_k):
                """Trick to count the no of iterations in CG."""
                CG_iter[-1] += 1

            c, info = scipy.sparse.linalg.cg(
                A, b, x0=x0, tol=CG_tol, maxiter=N, M=M,

```

```
callback=CG_callback)

# Fill u with vector c
# Update u_1 before next step
u_1, u = u, u_1
```

The number of iterations in the Conjugate gradient method is of interest, but unfortunately not available from the `cg` function, so we perform a trick: in each iteration a user function `CG_callback` is called where we accumulate the number of iteration in a list `CG_iter`.

## 3.7 Random walk

Models leading to diffusion equations, see Section 3.8, are usually based on reasoning with *averaged* physical quantities such as concentration, temperature, and velocity. The underlying physical processes involve complicated microscopic movement of atoms and molecules, but an average of a large number of molecules is performed in a small volume before the modeling starts, and the averaged quantity inside this volume is assigned as a point value at the centroid of the volume. This means that concentration, temperature, and velocity at a space-time point represent averages around the point in a small time interval and small spatial volume.

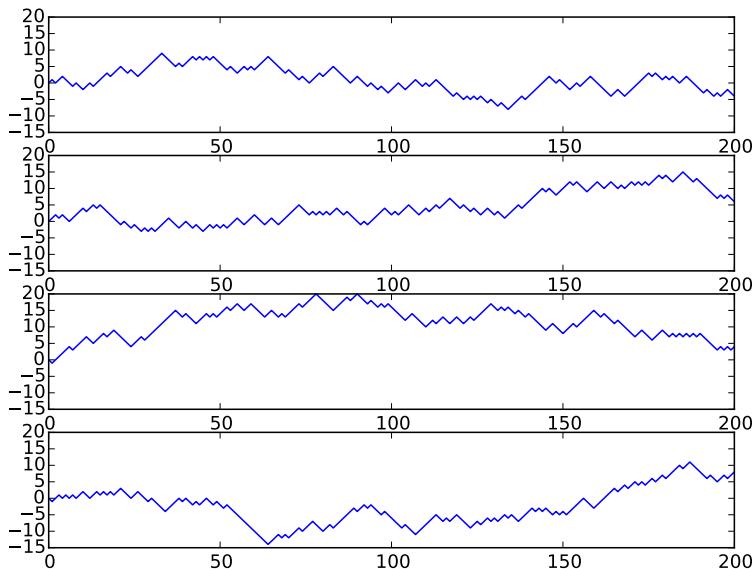
Random walk is a principally totally different kind of modeling procedure. The idea is to have a large number of “particles” that undergo random movements. Averaging can then be used afterwards to compute macroscopic quantities like concentration. The “particles” and their random movement represent a very simplified microscopic behavior of molecules, much simpler and computationally much more efficient than direct `molecular simulation`, yet the random walk model has been very powerful to describe a wide range of phenomena, including heat conduction, quantum mechanics, polymer chains, population genetics, neuroscience, hazard games, and pricing of financial instruments.

It can be shown that random walk, when averaged, produces models that are mathematically equivalent to diffusion equations. This is the primary reason why we treat random walk in this chapter: two very different algorithms (finite difference stencils and random walk) solve the same type of problems. The simplicity of the random walk algorithm makes it particularly attractive for solving diffusion equations on massively parallel computers.

### 3.7.1 Random walk in 1D

Imagine that we have some particles that perform random moves, either to the right or to the left. We may flip a coin to decide the movement of each particle, say head implies movement to the right and tail means movement to the left. Each move is one unit length. Physicists use the term *random walk* for this type of movement. The movement is also known as [drunkard's walk](#). You may try this yourself: flip the coin and make one step to the left or right, and repeat the process.

We introduce the symbol  $N$  for the number of steps in a random walk. Figure 3.16 shows four different random walks with  $N = 200$ .



**Fig. 3.16** Ensemble of 4 random walks, each with 200 steps.

### 3.7.2 Statistical considerations

**hpl 19:** Need to distinguish between scaled and unscaled position in the notation here!

Let  $S_k$  be the stochastic variable representing a step to the left or to the right in step number  $k$ . We have that  $S_k = -1$  with probability  $p$

and  $S_k = 1$  with probability  $q = 1 - p$ . The variable  $S_k$  is known as a **Bernoulli variable**. The expectation of  $S_k$  is

$$\mathbb{E}[S_k] = p \cdot (-1) + q \cdot 1 = 1 - 2p,$$

and the variance is

$$\text{Var}[S_k] = \mathbb{E}[S_k^2] - \mathbb{E}[S_k]^2 = 1 - (1 - 2p)^2 = 4p(1 - p).$$

The position after  $k$  steps is another stochastic variable

$$\bar{X}_k = \sum_{i=0}^{k-1} S_i.$$

The expected position is

$$\mathbb{E}[\bar{X}_k] = \mathbb{E}\left[\sum_{i=0}^{k-1} S_i\right] = \sum_{i=0}^{k-1} \mathbb{E}[S_i] = k(1 - 2p).$$

All the  $S_k$  variables are independent. The variance therefore becomes

$$\text{Var}[\bar{X}_k] = \text{Var}\left[\sum_{i=0}^{k-1} S_i\right] = \sum_{i=0}^{k-1} \text{Var}[S_i] = k4p(1 - p).$$

We see that  $\text{Var}[\bar{X}_k]$  is proportional with the number of steps  $k$ . For the very important case  $p = q = \frac{1}{2}$ ,  $\mathbb{E}[\bar{X}_k] = 0$  and  $\text{Var}[\bar{X}_k] = k$ .

How can we estimate  $\mathbb{E}[\bar{X}_k] = 0$  and  $\text{Var}[\bar{X}_k] = N$ ? We must have many random walks of the type in Figure 3.16. For a given  $k$ , say  $k = 100$ , we find all the values of  $\bar{X}_k$ , name them  $\bar{x}_{0,k}, \bar{x}_{1,k}, \bar{x}_{2,k}$ , and so on. The empirical estimate of  $\mathbb{E}[\bar{X}_k]$  is the average,

$$\mathbb{E}[\bar{X}_k] \approx \frac{1}{W} \sum_{j=0}^{W-1} \bar{x}_{j,k},$$

while an empirical estimate of  $\text{Var}[\bar{X}_k]$  is

$$\text{Var}[\bar{X}_k] \approx \frac{1}{W} \sum_{j=0}^{W-1} (\bar{x}_{j,k})^2 - \left( \frac{1}{W} \sum_{j=0}^{W-1} \bar{x}_{j,k} \right)^2.$$

That is, we take the statistics for a given  $K$  across the ensemble of random walks (“vertically” in Figure 3.16). The key quantities to record are  $\sum_i \bar{x}_{i,k}$  and  $\sum_i \bar{x}_{i,k}^2$ .

### 3.7.3 Playing around with some code

**Scalar code.** Python has a `random` module for drawing random numbers, and a function `uniform(a, b)` for drawing a uniformly distributed random number in the interval  $[a, b]$ . If an event happens with probability  $p$ , we can simulate this on the computer by drawing a random number  $r$  in  $[0, 1]$ , because then  $r \leq p$  with probability  $p$  and  $r > p$  with probability  $1 - p$ :

```
import random
r = random.uniform(0, 1)
if r <= p:
    # Event happens
else:
    # Event does not happen
```

A random walk with  $N$  steps, starting at  $x_0$ , where we move to the left with probability  $p$  and to the right with probability  $1 - p$  can now be implemented by

```
import random, numpy as np

def random_walk1D(x0, N, p):
    """1D random walk with 1 particle."""
    # Store position in step k in position[k]
    position = np.zeros(N)
    position[0] = x0
    current_pos = x0
    for k in range(N-1):
        r = random.uniform(0, 1)
        if r <= p:
            current_pos -= 1
        else:
            current_pos += 1
        position[k+1] = current_pos
    return position
```

**Vectorized code.** Since  $N$  is supposed to be large and we want to repeat the process for many particles, we should speed up the code as much as possible. Vectorization is the obvious technique here: we draw all the random numbers at once with aid of `numpy`, and then we formulate vector operations to get rid of the loop over the steps ( $k$ ). The `numpy.random` module has vectorized versions of the functions in Python's built-in `random` module. For example, `numpy.random.uniform(a, b, N)` returns  $N$  random numbers uniformly distributed between  $a$  (included) and  $b$  (not included).

We can then make an array of all the steps in a random walk: if the random number is less than or equal to  $p$ , the step is  $-1$ , otherwise the step is  $1$ :

```
r = np.random.uniform(0, 1, size=N)
steps = np.where(r <= p, -1, 1)
```

The value of `position[k]` is the sum of all steps up to step  $k$ . Such sums are often needed in vectorized algorithms and therefore available by the `numpy.cumsum` function:

```
>>> import numpy as np
>>> np.cumsum(np.array([1,3,4,6]))
array([ 1,  4,  8, 14])
```

The resulting array in this demo has elements  $1$ ,  $1 + 3 = 4$ ,  $1 + 3 + 4 = 8$ , and  $1 + 3 + 4 + 6 = 14$ .

We can now vectorize the `random_walk1D` function:

```
def random_walk1D_vec(x0, N, p):
    """Vectorized version of random_walk1D."""
    # Store position in step k in position[k]
    position = np.zeros(N+1)
    position[0] = x0
    r = np.random.uniform(0, 1, size=N)
    steps = np.where(r <= p, -1, 1)
    position[1:] = x0 + np.cumsum(steps)
    return position
```

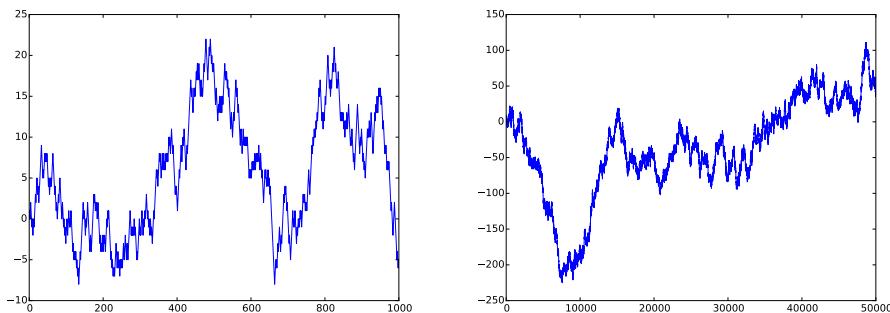
This code runs about 10 times faster than the scalar version. With a parallel `numpy` library, the code can also automatically take advantage of hardware for parallel computing because each of the four array operations can be trivially parallelized.

**Fixing the random sequence.** During software development with random numbers it is advantageous to always generate the same sequence of random numbers as this may help debugging processes. To fix the sequence, we set a *seed* of the random number generator to some chosen integer, e.g.,

```
np.random.seed(10)
```

Calls to `random_walk1D_vec` give positions of the particle as depicted in Figure 3.17. The particle starts at the origin and moves with  $p = \frac{1}{2}$ . Since the seed is the same, the plot to the left is just a magnification of the first 1/50 steps in the plot to the right.

**Verification.** When we have a scalar and a vectorized code, it is always a good idea to develop a unit test for checking that they produce the same



**Fig. 3.17** 1,000 (left) and 50,000 (right) steps of a random walk.

result. A problem in the present context is that the two versions apply to different random number generators. For a test to be meaningful, we need to fix the seed and use the same generator. This means that the scalar version must either use `np.random` or have this as an option. An option is the most flexible choice:

```
import random

def random_walk1D(x0, N, p, random=random):
    ...
    r = random.uniform(0, 1)
```

Using `random=np.random`, the `r` variable gets computed by `np.random.uniform`, and the sequence of random numbers will be the same as in the vectorized version that employs the same generator (given that the seed is also the same). A proper test function may be to check that the positions in the walk are the same in the scalar and vectorized implementations:

```
def test_random_walk1D():
    # For fixed seed, check that scalar and vectorized versions
    # produce the same result
    x0 = 2; N = 4; p = 0.6
    np.random.seed(10)
    scalar_computed = random_walk1D(x0, N, p, random=np.random)
    np.random.seed(10)
    vectorized_computed = random_walk1D_vec(x0, N, p)
    assert (scalar_computed == vectorized_computed).all()
```

Note that we employ `==` for arrays with real numbers, which is normally an inadequate test due to rounding errors, but in the present case, all arithmetics consists of adding or subtracting one, so these operations are expected to have no rounding errors. Comparing two `numpy` arrays with `==` results in a boolean array, so we need to call the `all()` method

to ensure that all elements are `True`, i.e., that all elements in the two arrays match each other pairwise.

### 3.7.4 Equivalence with diffusion

The original random walk algorithm can be said to work with dimensionless coordinates  $\bar{x}_i = -N + i$ ,  $i = 0, 1, \dots, 2N + 1$  ( $i \in [-N, N]$ ), and  $\bar{t}_n = n$ ,  $n = 0, 1, \dots, N$ . A mesh with spacings  $\Delta x$  and  $\Delta t$  with dimensions can be introduced by

$$x_i = X_0 + \bar{x}_i \Delta x, \quad t_n = \bar{t}_n \Delta t.$$

If we implement the algorithm with dimensionless coordinates, we can just use this rescaling to obtain the movement in a coordinate system without unit spacings.

Let  $P_i^{n+1}$  be the probability of finding the particle at mesh point  $\bar{x}_i$  at time  $\bar{t}_{n+1}$ . We can reach mesh point  $(i, n+1)$  in two ways: either coming in from the left from  $(i-1, n)$  or from the right  $(i+1, n)$ . Both has probability  $\frac{1}{2}$  (if we assume  $p = q = \frac{1}{2}$ ). The fundamental equation for  $P_i^{n+1}$  is

$$P_i^{n+1} = \frac{1}{2} P_{i-1}^n + \frac{1}{2} P_{i+1}^n. \quad (3.113)$$

(This equation is easiest to understand if one looks at the random walk as a Markov process and applies the transition probabilities, but this is beyond scope of the present text.)

Subtracting  $P_i^n$  from (3.7.1) results in

$$P_i^{n+1} - P_i^n = \frac{1}{2}(P_{i-1}^n - 2P_i^n + \frac{1}{2}P_{i+1}^n).$$

Readers who have seen the Forward Euler discretization of a 1D diffusion equation recognize this scheme as very close to such a discretization. We have

$$\frac{\partial}{\partial t} P(x_i, t_n) = \frac{P_i^{n+1} - P_i^n}{\Delta t} + \mathcal{O}(\Delta t),$$

or in dimensionless coordinates

$$\frac{\partial}{\partial \bar{t}} P(\bar{x}_i, \bar{t}_n) \approx P_i^{n+1} - P_i^n.$$

Similarly, we have

$$\frac{\partial^2}{\partial x^2} P(x_i, t_n) = \frac{P_{i-1}^n - 2P_i^n + \frac{1}{2}P_{i+1}^n}{\Delta x^2} + \mathcal{O}(\Delta x^2),$$

$$\frac{\partial^2}{\partial x^2} P(\bar{x}_i, \bar{t}_n) \approx P_{i-1}^n - 2P_i^n + \frac{1}{2}P_{i+1}^n.$$

Equation (3.7.1) is therefore equivalent with the dimensionless diffusion equation

$$\frac{\partial P}{\partial \bar{t}} = \frac{1}{2} \frac{\partial^2 P}{\partial \bar{x}^2}, \quad (3.114)$$

or the diffusion equation

$$\frac{\partial P}{\partial t} = D \frac{\partial^2 P}{\partial x^2}, \quad (3.115)$$

with diffusion coefficient

$$D = \frac{\Delta x^2}{2\Delta t}.$$

This derivation shows the tight link between random walk and diffusion. If we keep track of where the particle is, and repeat the process many times, or run the algorithms for lots of particles, the histogram of the positions will approximate the solution of the diffusion equation for the local probability  $P_i^n$ .

Suppose all the random walks start at the origin. Then the initial condition for the probability distribution is the Dirac delta function  $\delta(x)$ . The solution of (3.114) can be shown to be

$$\bar{P}(\bar{x}, \bar{t}) = \frac{1}{\sqrt{4\pi\alpha\bar{t}}} e^{-\frac{\bar{x}^2}{4\alpha\bar{t}}}, \quad (3.116)$$

where  $\alpha = \frac{1}{2}$ .

### 3.7.5 Implementation of multiple walks

Our next task is to implement an ensemble of walks (for statistics, see Section 3.7.2) and also provide data from the walks such that we can compute the probabilities of the positions as introduced in the previous section. An appropriate representation of probabilities  $P_i^n$  are histograms (with  $i$  along the  $x$  axis) for a few selected values of  $n$ .

To estimate the expectation and variance of the random walks, Section 3.7.2 points to recording  $\sum_j x_{j,k}$  and  $\sum_j x_{j,k}^2$ , where  $x_{j,k}$  is the position at time/step level  $k$  in random walk number  $j$ . The histogram of positions needs the individual values  $x_{i,k}$  for all  $i$  values and some selected  $k$  values.

We introduce `position[k]` to hold  $\sum_j x_{j,k}$ , `position2[k]` to hold  $\sum_j (x_{j,k})^2$ , and `pos_hist[i,k]` to hold  $x_{i,k}$ . A selection of  $k$  values can be specified by saying how many, `num_times`, and let them be equally spaced through time:

```
pos_hist_times = [(N//num_times)*i for i in range(num_times)]
```

This is one of the few situations we want integer division (`//`) or real division rounded to an integer.

**Scalar version.** Our scalar implementation of running `num_walks` random walks may go like this:

```
def random_walks1D(x0, N, p, num_walks=1, num_times=1,
                    random=random):
    """Simulate num_walks random walks from x0 with N steps."""
    position = np.zeros(N+1)      # Accumulated positions
    position[0] = x0*num_walks
    position2 = np.zeros(N+1)     # Accumulated positions**2
    position2[0] = x0**2*num_walks
    # Histogram at num_times selected time points
    pos_hist = np.zeros((num_walks, num_times))
    pos_hist_times = [(N//num_times)*i for i in range(num_times)]
    #print 'save hist:', pos_hist_times

    for n in range(num_walks):
        num_times_counter = 0
        current_pos = x0
        for k in range(N):
            if k in pos_hist_times:
                #print 'save, k:', k, num_times_counter, n
                pos_hist[n,num_times_counter] = current_pos
                num_times_counter += 1
            # current_pos corresponds to step k+1
            r = random.uniform(0, 1)
            if r <= p:
                current_pos -= 1
            else:
                current_pos += 1
            position [k+1] += current_pos
            position2[k+1] += current_pos**2
    return position, position2, pos_hist, np.array(pos_hist_times)
```

**Vectorized version.** We have already vectorized a single random walk. The additional challenge here is to vectorize the computation of the

data for the histogram, `pos_hist`, but given the selected steps in `pos_hist_times`, we can find the corresponding positions by indexing with the list `pos_hist_times`: `position[pos_hist_times]`, which are to be inserted in `pos_hist[n,:]`.

```
def random_walks1D_vec1(x0, N, p, num_walks=1, num_times=1):
    """Vectorized version of random_walks1D."""
    position = np.zeros(N+1)      # Accumulated positions
    position2 = np.zeros(N+1)     # Accumulated positions**2
    walk = np.zeros(N+1)          # Positions of current walk
    walk[0] = x0
    # Histogram at num_times selected time points
    pos_hist = np.zeros((num_walks, num_times))
    pos_hist_times = [(N//num_times)*i for i in range(num_times)]

    for n in range(num_walks):
        r = np.random.uniform(0, 1, size=N)
        steps = np.where(r <= p, -1, 1)
        walk[1:] = x0 + np.cumsum(steps) # Positions of this walk
        position += walk
        position2 += walk**2
        pos_hist[n,:] = walk[pos_hist_times]
    return position, position2, pos_hist, np.array(pos_hist_times)
```

**Improved vectorized version.** Looking at the vectorized version above, we still have one potentially long Python loop over `n`. Normally, `num_walks` will be much larger than `N`. The vectorization of the loop over `N` certainly speeds up the program, but if we think of vectorization as also a way to parallelize the code, all the independent walks (the `n` loop) can be executed in parallel. Therefore, we should include this loop as well in the vectorized expressions, at the expense of using more memory.

We introduce the array `walks` to hold the  $N + 1$  steps of all the walks: each row represents the steps in one walk.

```
walks = np.zeros((num_walks, N+1)) # Positions of each walk
walks[:,0] = x0
```

Since all the steps are independent, we can just make one long vector of enough random numbers ( $N * \text{num\_walks}$ ), translate these numbers to  $\pm 1$ , then we reshape the array such that the steps of each walk are stored in the rows.

```
r = np.random.uniform(0, 1, size=N*num_walks)
steps = np.where(r <= p, -1, 1).reshape(num_walks, N)
```

The next step is to sum up the steps in each walk. We need the `np.cumsum` function for this, with the argument `axis=1` for indicating a sum across the columns:

```
>>> a = np.arange(6).reshape(2,3)
>>> a
array([[0, 1, 2],
       [3, 4, 5]])
>>> np.cumsum(a, axis=1)
array([[ 0,  1,  3],
       [ 3,  7, 12]])
```

Now `walks` can be computed by

```
walks[:,1:] = x0 + np.cumsum(steps, axis=1)
```

The `position` vector is the sum of all the walks. That is, we want to sum all the rows, obtained by

```
position = np.sum(walks, axis=0)
```

A corresponding expression computes the squares of the positions. Finally, we need to compute `pos_hist`, but that is a matter of grabbing some of the walks (according to `pos_hist_times`):

```
pos_hist[:, :] = walks[:, pos_hist_times]
```

The complete vectorized algorithm without any loop can now be summarized:

```
def random_walks1D_vec2(x0, N, p, num_walks=1, num_times=1):
    """Vectorized version of random_walks1D; no loops."""
    position = np.zeros(N+1)      # Accumulated positions
    position2 = np.zeros(N+1)     # Accumulated positions**2
    walks = np.zeros((num_walks, N+1))  # Positions of each walk
    walks[:,0] = x0
    # Histogram at num_times selected time points
    pos_hist = np.zeros((num_walks, num_times))
    pos_hist_times = [(N//num_times)*i for i in range(num_times)]

    r = np.random.uniform(0, 1, size=N*num_walks)
    steps = np.where(r <= p, -1, 1).reshape(num_walks, N)
    walks[:,1:] = x0 + np.cumsum(steps, axis=1)
    position = np.sum(walks, axis=0)
    position2 = np.sum(walks**2, axis=0)
    pos_hist[:, :] = walks[:, pos_hist_times]
    return position, position2, pos_hist, np.array(pos_hist_times)
```

What is the gain of the vectorized implementations? One important gain is that each vectorized operation can be automatically parallelized if one applies a parallel `numpy` library like `Numba`. On a single CPU, however, the speed up of the vectorized operations is also significant. With  $N = 1,000$  and 50,000 repeated walks, the two vectorized versions run about 25 and 18 times faster than the scalar version, with `random_walks1D_vec1` being fastest.

**Remark on vectorized code and parallelization.** Our first attempt on vectorization removed the loop over the  $N$  steps in a single walk. However, the number of walks is usually much larger than  $N$ , because of the need for accurate statistics. Therefore, we should rather remove the loop over all walks. It turns out, from our efficiency experiments, that the function `random_walks1D_vec2` (with no loops) is slower than `random_walks1D_vec1`. This is a bit surprising and may be explained by less efficiency in the statements involving very large arrays, containing all steps for all walks at once.

From a parallelization and improved vectorization point of view, it would be more natural to switch the sequence of the loops in the serial code such that the shortest loop is the outer loop:

```
def random_walks1D2(x0, N, p, num_walks=1, num_times=1, ...):
    ...
    current_pos = x0 + np.zeros(num_walks)
    num_times_counter = -1

    for k in range(N):
        if k in pos_hist_times:
            num_times_counter += 1
            store_hist = True
        else:
            store_hist = False

        for n in range(num_walks):
            # current_pos corresponds to step k+1
            r = random.uniform(0, 1)
            if r <= p:
                current_pos[n] -= 1
            else:
                current_pos[n] += 1
            position[k+1] += current_pos[n]
            position2[k+1] += current_pos[n]**2
            if store_hist:
                pos_hist[n,num_times_counter] = current_pos[n]
    return position, position2, pos_hist, np.array(pos_hist_times)
```

The vectorized version of this code, where we just vectorize the loop over  $n$ , becomes

```
def random_walks1D2_vec1(x0, N, p, num_walks=1, num_times=1):
    """Vectorized version of random_walks1D2."""
    position = np.zeros(N+1)      # Accumulated positions
    position2 = np.zeros(N+1)     # Accumulated positions**2
    # Histogram at num_times selected time points
    pos_hist = np.zeros((num_walks, num_times))
    pos_hist_times = [(N//num_times)*i for i in range(num_times)]

    current_pos = np.zeros(num_walks)
```

```

current_pos[0] = x0
num_times_counter = -1

for k in range(N):
    if k in pos_hist_times:
        num_times_counter += 1
        store_hist = True # Store histogram data for this k
    else:
        store_hist = False

    # Move all walks one step
    r = np.random.uniform(0, 1, size=num_walks)
    steps = np.where(r <= p, -1, 1)
    current_pos += steps
    position[k+1] = np.sum(current_pos)
    position2[k+1] = np.sum(current_pos**2)
    if store_hist:
        pos_hist[:,num_times_counter] = current_pos
return position, position2, pos_hist, np.array(pos_hist_times)

```

This function runs significantly faster than the `random_walks1D_vec1` function above, typically 1.7 times faster. The code is also more appropriate in a parallel computing context since each vectorized statement can work with data of size `num_walks` over the compute units, repeated `N` times (compared with data of size `N`, repeated `num_walks` times, in `random_walks1D_vec1`).

The scalar code with switched loops, `random_walks1D2` runs a bit slower than the original code in `random_walks1D`, so with the longest loop as the inner loop, the vectorized function `random_walks1D2_vec1` is almost 60 times faster than the scalar counterpart, while the code `random_walks1D_vec2` without loops is only around 18 times faster. Taking into account the very large arrays required by the latter function, we end up with `random_walks1D2_vec1` as the preferred implementation.

**Test function.** During program development, it is highly recommended to carry out computations by hand for, e.g., `N=4` and `num_walks=3`. Normally, this is done by executing the program with these parameters and checking with pen and paper that the computations make sense. The next step is to use this test for correctness in a formal test function.

First, we need to check that the simulation of multiple random walks reproduces the results of `random_walk1D`, `random_walk1D_vec1`, and `random_walk1D_vec2` for the first walk, if the seed is the same. Second, we run three random walks (`N=4`) with the scalar and the two vectorized versions and check that the returned arrays are identical.

For this type of test to be successful, we must be sure that exactly the same set of random numbers are used in the three versions, a fact that

requires the same random number generator and the same seed, of course, but also the same sequence of computations. This is not obviously the case with the three `random_walk1D*` functions we have presented. The critical issue in `random_walk1D_vec1` is that the first random numbers are used for the first walk, the second set of random numbers is used for the second walk and so, to be compatible with how the random numbers are used in the function `random_walk1D`. For the function `random_walk1D_vec2` the situation is a bit more complicated since we generate all the random numbers at once. However, the critical step now is the reshaping of the array returned from `np.where`: we must reshape as `(num_walks, N)` to ensure that the first `N` random numbers are used for the first walk, the next `N` numbers are used for the second walk, and so on.

We arrive at the test function below.

```
def test_random_walks1D():
    # For fixed seed, check that scalar and vectorized versions
    # produce the same result
    x0 = 0;  N = 4;  p = 0.5

    # First, check that random_walks1D for 1 walk reproduces
    # the walk in random_walk1D
    num_walks = 1
    np.random.seed(10)
    computed = random_walks1D(
        x0, N, p, num_walks, random=np.random)
    np.random.seed(10)
    expected = random_walk1D(
        x0, N, p, random=np.random)
    assert (computed[0] == expected).all()

    # Same for vectorized versions
    np.random.seed(10)
    computed = random_walks1D_vec1(x0, N, p, num_walks)
    np.random.seed(10)
    expected = random_walk1D_vec(x0, N, p)
    assert (computed[0] == expected).all()
    np.random.seed(10)
    computed = random_walks1D_vec2(x0, N, p, num_walks)
    np.random.seed(10)
    expected = random_walk1D_vec(x0, N, p)
    assert (computed[0] == expected).all()

    # Second, check multiple walks: scalar == vectorized
    num_walks = 3
    num_times = N
    np.random.seed(10)
    serial_computed = random_walks1D(
        x0, N, p, num_walks, num_times, random=np.random)
    np.random.seed(10)
```

```

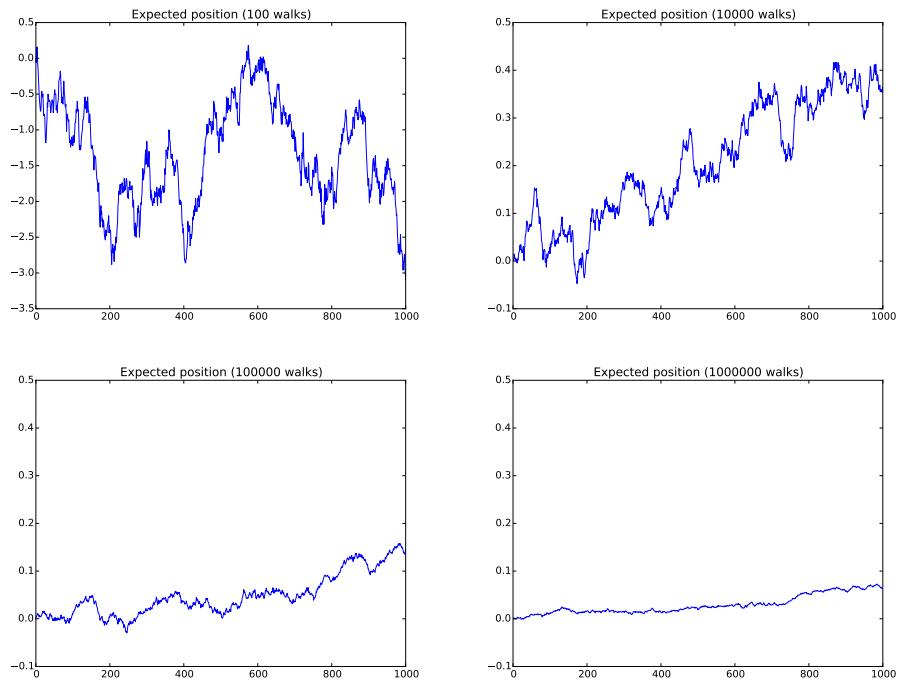
vectorized1_computed = random_walks1D_vec1(
    x0, N, p, num_walks, num_times)
np.random.seed(10)
vectorized2_computed = random_walks1D_vec2(
    x0, N, p, num_walks, num_times)
# positions: [0', 1, 0, 1, 2]
# Can test without tolerance since everything is +/- 1
return_values = ['pos', 'pos2', 'pos_hist', 'pos_hist_times']
for s, v, r in zip(serial_computed,
                    vectorized1_computed,
                    return_values):
    msg = '%s: %s (serial) vs %s (vectorized)' % (r, s, v)
    assert (s == v).all(), msg
for s, v, r in zip(serial_computed,
                    vectorized2_computed,
                    return_values):
    msg = '%s: %s (serial) vs %s (vectorized)' % (r, s, v)
    assert (s == v).all(), msg

```

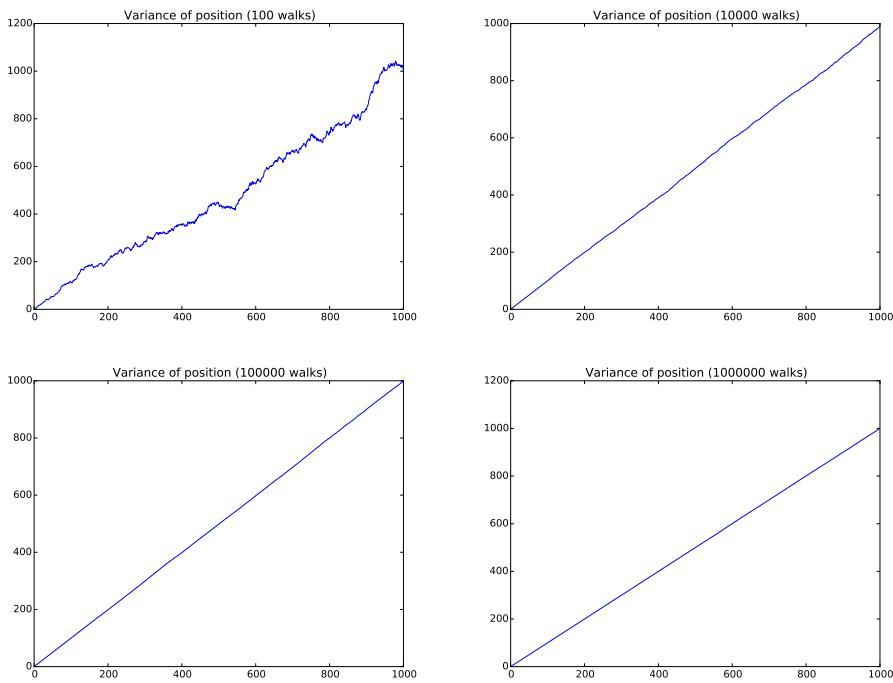
Such test functions are indispensable for further development of the code as we can at any time test whether the basic computations remain correct or not. This is particularly important in stochastic simulations since without test functions and fixed seeds, we always experience variations from run to run, and it can be very difficult to spot bugs through averaged statistical quantities.

### 3.7.6 Demonstration of multiple walks

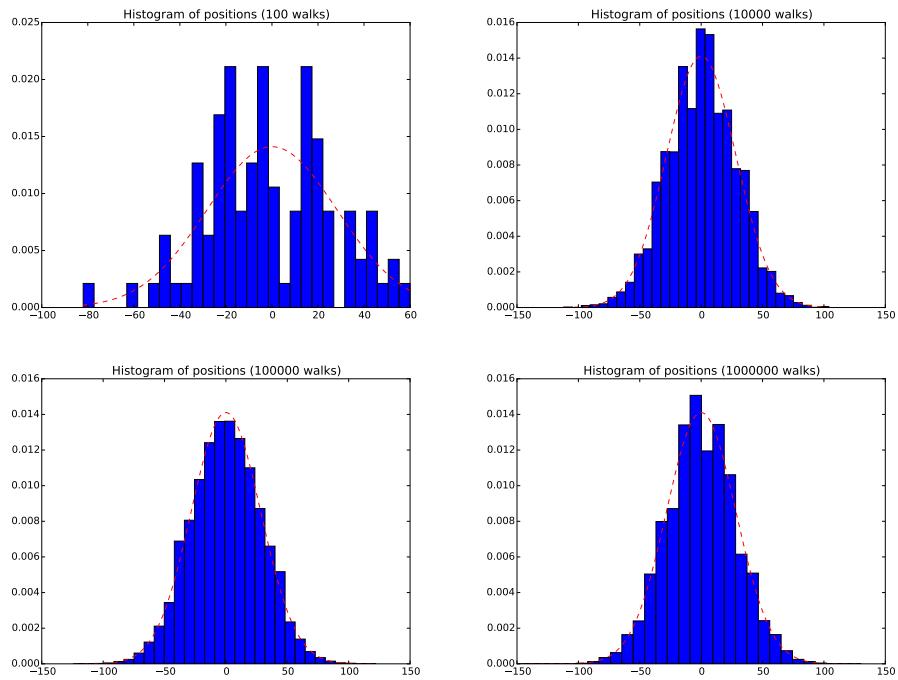
Assuming now that the code works, we can just scale up the number of steps in each walk and the number of walks. The latter influences the accuracy of the statistical estimates. Figure 3.18 shows the impact of the number of walks on the expectation, which should approach zero. Figure 3.19 displays the corresponding estimate of the variance of the position, which should grow linearly with the number of steps. It does, seemingly very accurately, but notice that the scale on the  $y$  axis is so much larger than for the expectation, so irregularities due to the stochastic nature of the process become so much less visible in the variance plots. The probability of finding a particle at a certain position at time (or step) 800 is shown in Figure 3.20. The dashed red line is the theoretical distribution (3.116) arising from solving the diffusion equation (3.114) instead. As always, we realize that one needs significantly more statistical samples to estimate a histogram accurately than the expectation or variance.



**Fig. 3.18** Estimated expected value for 1000 steps, using 100 walks (upper left), 10,000 (upper right), 100,000 (lower left), and 1,000,000 (lower right).



**Fig. 3.19** Estimated variance over 1000 steps, using 100 walks (upper left), 10,000 (upper right), 100,000 (lower left), and 1,000,000 (lower right).



**Fig. 3.20** Estimated probability distribution at step 800, using 100 walks (upper left), 10,000 (upper right), 100,000 (lower left), and 1,000,000 (lower right).

### 3.7.7 Ascii visualization of 1D random walk

If we want to study (very) long time series of random walks, it can be convenient to plot the position in a terminal window with the time axis pointing downwards. The module `avplotter` in SciTools has a class `Plotter` for plotting functions in the terminal window with the aid of ascii symbols only. Below is the code required to visualize a simple random walk, starting at the origin, and considered over when the point  $x = -1$  is reached. We use a spacing  $\Delta x = 0.05$  (so  $x = -1$  corresponds to  $i = -20$ ).

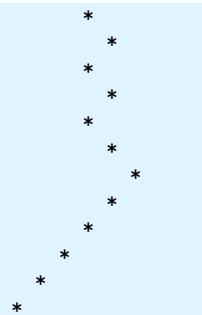
```
def run_random_walk():
    from scitools.avplotter import Plotter
    import time, numpy as np
    p = Plotter(-1, 1, width=75)    # Horizontal axis: 75 chars wide
    dx = 0.05
    np.random.seed(10)

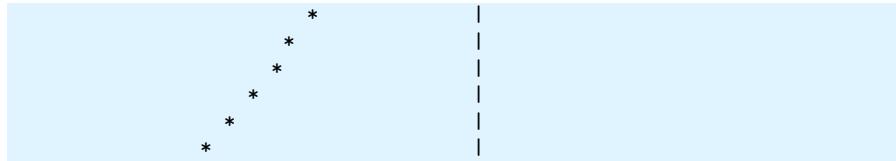
    x = 0
    while True:
        random_step = 1 if np.random.random() > 0.5 else -1
        x = x + dx*random_step
        if x < -1:
            break                # Destination reached!
        print p.plot(0, x)

    # Allow Ctrl+c to abort the simulation
    try:
        time.sleep(0.1)    # Wait for interrupt
    except KeyboardInterrupt:
        print 'Interrupted by Ctrl+c'
        break
```

Observe that we implement an infinite loop, but allow a smooth interrupt of the program by `Ctrl+c` through Python's `KeyboardInterrupt` exception. This is a useful recipe that can be used in many occasions!

The output looks typically like





Positions beyond the limits of the  $x$  axis appear with a value. A file contains the complete ascii plot corresponding to the function `run_random_walk` above.

### 3.7.8 Random walk as a stochastic equation

The (dimensionless) position in a random walk,  $\bar{X}_k$ , can be expressed as a stochastic difference equation:

$$\bar{X}_k = \bar{X}_{k-1} + s, \quad x_0 = 0, \quad (3.117)$$

where  $s$  is a **Bernoulli variable**, taking on the two values  $s = -1$  and  $s = 1$  with equal probability:

$$P(s = 1) = \frac{1}{2}, \quad P(s = -1) = \frac{1}{2}.$$

The  $s$  variable in a step is independent of the  $s$  variable in other steps.

The difference equation expresses essentially the sum of independent Bernoulli variables. Because of the central limit theorem,  $X_k$ , will then be normally distributed with expectation  $kE[s]$  and  $k\text{Var}[s]$ . The expectation and variance of a Bernoulli variable with values  $r = 0$  and  $r = 1$  are  $p$  and  $p(1 - p)$ , respectively. The variable  $s = 2r - 1$  then has expectation  $2E[r] - 1 = 2p - 1 = 0$  and variance  $2^2\text{Var}[r] = 4p(1 - p) = 1$ . The position  $X_k$  is normally distributed with zero expectation and variance  $k$ , as we found in Section 3.7.2.

The central limit theorem tells that as long as  $k$  is not small, the distribution of  $X_k$  remains the same if we replace the Bernoulli variable  $s$  by any other stochastic variable with the same expectation and variance. In particular, may let  $s$  be a standardized Gaussian variable (zero mean, unit variance).

Dividing (3.117) by  $\Delta t$  gives

$$\frac{\bar{X}_k - \bar{X}_{k-1}}{\Delta t} = \frac{1}{\Delta t} s.$$

In the limit  $\Delta t \rightarrow 0$ ,  $s/\Delta t$  approaches a white noise stochastic process ( $s$  is standardized Gaussian variable). With  $\bar{X}(t)$  as the continuous process in the limit  $\Delta t \rightarrow 0$  ( $X_k \rightarrow X(t_k)$ ), we formally get the stochastic differential equation

$$d\bar{X} = dW, \quad (3.118)$$

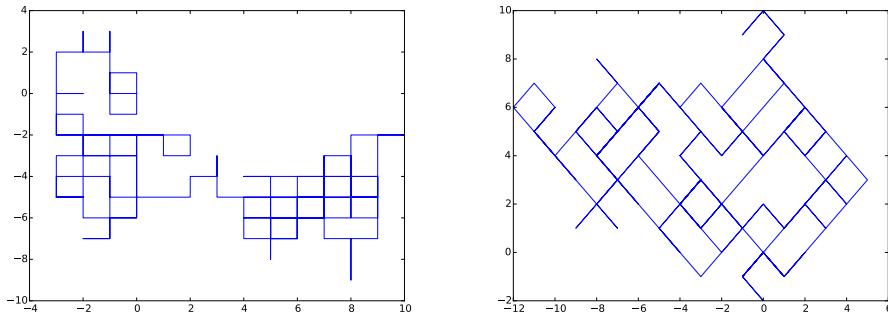
where  $W(t)$  is a [Wiener process](#). Then  $X$  is also a Wiener process. It follows from the stochastic ODE  $dX = dW$  that the probability distribution of  $X$  is given by the [Fokker-Planck equation](#) (3.114). In other words, the key results for random walk we found earlier can alternatively be derived via a stochastic ordinary differential equation and its related Fokker-Planck equation.

### 3.7.9 Random walk in 2D

The most obvious generalization of 1D random walk to two spatial dimensions is to allow movements to the north, east, south, and west, with equal probability  $\frac{1}{4}$ .

```
def random_walk2D(x0, N, p, random=random):
    """2D random walk with 1 particle and N moves: N, E, W, S."""
    # Store position in step k in position[k]
    d = len(x0)
    position = np.zeros((N+1, d))
    position[0,:] = x0
    current_pos = np.array(x0, dtype=float)
    for k in range(N):
        r = random.uniform(0, 1)
        if r <= 0.25:
            current_pos += np.array([0, 1])  # Move north
        elif 0.25 < r <= 0.5:
            current_pos += np.array([1, 0])  # Move east
        elif 0.5 < r <= 0.75:
            current_pos += np.array([0, -1]) # Move south
        else:
            current_pos += np.array([-1, 0]) # Move west
        position[k+1,:] = current_pos
    return position
```

The left plot in Figure 3.21 provides an example on 200 steps with this kind of walk. We may refer to this walk as a walk on a *rectangular mesh* as we move from any spatial mesh point  $(i, j)$  to one of its four neighbors in the rectangular directions:  $(i + 1, j)$ ,  $(i - 1, j)$ ,  $(i, j + 1)$ , or  $(i, j - 1)$ .



**Fig. 3.21** Random walks in 2D with 200 steps: rectangular mesh (left) and diagonal mesh (right).

### 3.7.10 Random walk in any number of space dimensions

From a programming point of view, especially when implementing a random walk in any number of dimensions, it is more natural to consider a walk in the diagonal directions NW, NE, SW, and SE. On a two-dimensional spatial mesh it means that we go from  $(i, j)$  to either  $(i + 1, j + 1)$ ,  $(i - 1, j + 1)$ ,  $(i + 1, j - 1)$ , or  $(i - 1, j - 1)$ . We can with such a *diagonal mesh* (see right plot in Figure 3.21) draw a Bernoulli variable for the step in each spatial direction and trivially write code that works in any number of spatial directions:

```
def random_walkdD(x0, N, p, random=random):
    """Any-D (diagonal) random walk with 1 particle and N moves."""
    # Store position in step k in position[k]
    d = len(x0)
    position = np.zeros((N+1, d))
    position[0,:] = x0
    current_pos = np.array(x0, dtype=float)
    for k in range(N):
        for i in range(d):
            r = random.uniform(0, 1)
            if r <= p:
                current_pos[i] -= 1
            else:
                current_pos[i] += 1
        position[k+1,:] = current_pos
    return position
```

A vectorized version is desired. We follow the ideas from Section 3.7.3, but each step is now a vector in  $d$  spatial dimensions. We therefore need to draw  $Nd$  random numbers in  $\mathbf{r}$ , compute steps in the various directions through `np.where(r <= p, -1, 1)` (each step being  $-1$  or  $1$ ), and then we can reshape this array to an  $N \times d$  array of step *vectors*. Doing an

`np.cumsum` summation along axis 0 will add the vectors, as this demo shows:

```
>>> a = np.arange(6).reshape(3,2)
>>> a
array([[0, 1],
       [2, 3],
       [4, 5]])
>>> np.cumsum(a, axis=0)
array([[ 0,  1],
       [ 2,  4],
       [ 6,  9]])
```

With such summation of step vectors, we get all the positions to be filled in the `position` array:

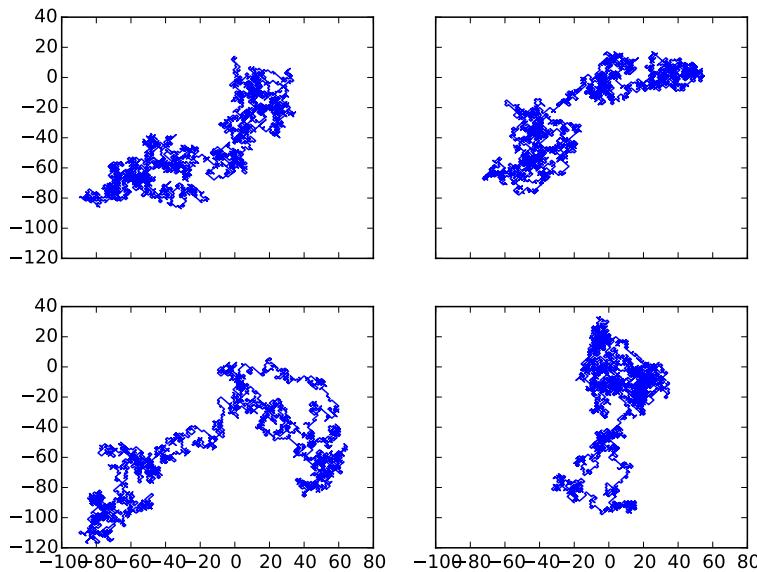
```
def random_walkdD_vec(x0, N, p):
    """Vectorized version of random_walkdD."""
    d = len(x0)
    # Store position in step k in position[k]
    position = np.zeros((N+1,d))
    position[0] = np.array(x0, dtype=float)
    r = np.random.uniform(0, 1, size=N*d)
    steps = np.where(r <= p, -1, 1).reshape(N,d)
    position[1:,:] = x0 + np.cumsum(steps, axis=0)
    return position
```

### 3.7.11 Multiple random walks in any number of space dimensions

As we did in 1D, we extend one single walk to a number of walks (`num_walks` in the code).

```
def random_walksdD(x0, N, p, num_walks=1, num_times=1,
                     random=random):
    """Simulate num_walks random walks from x0 with N steps."""
    d = len(x0)
    position = np.zeros((N+1, d)) # Accumulated positions
    position2 = np.zeros((N+1, d)) # Accumulated positions**2
    # Histogram at num_times selected time points
    pos_hist = np.zeros((num_walks, num_times, d))
    pos_hist_times = [(N//num_times)*i for i in range(num_times)]

    for n in range(num_walks):
        num_times_counter = 0
        current_pos = np.array(x0, dtype=float)
        for k in range(N):
            if k in pos_hist_times:
                pos_hist[n,num_times_counter,:] = current_pos
            r = np.random.uniform(0, 1, size=d)
            steps = np.where(r <= p, -1, 1)
            current_pos += steps
```



**Fig. 3.22** Four random walks with 5000 steps in 2D.

```

        num_times_counter += 1
        # current_pos corresponds to step k+1
        for i in range(d):
            r = random.uniform(0, 1)
            if r <= p:
                current_pos[i] -= 1
            else:
                current_pos[i] += 1
            position [k+1,:] += current_pos
            position2[k+1,:] += current_pos**2
        return position, position2, pos_hist, np.array(pos_hist_times)
    
```

```

def random_walksdD_vec(x0, N, p, num_walks=1, num_times=1):
    """Vectorized version of random_walks1D; no loops."""
    d = len(x0)
    position = np.zeros((N+1, d)) # Accumulated positions
    position2 = np.zeros((N+1, d)) # Accumulated positions**2
    walks = np.zeros((num_walks, N+1, d)) # Positions of each walk
    walks[:,0,:] = x0
    # Histogram at num_times selected time points
    pos_hist = np.zeros((num_walks, num_times, d))
    pos_hist_times = [(N//num_times)*i for i in range(num_times)]
    
```

```

r = np.random.uniform(0, 1, size=N*num_walks*d)
steps = np.where(r <= p, -1, 1).reshape(num_walks, N, d)
walks[:,1:,:,:] = x0 + np.cumsum(steps, axis=1)
position  = np.sum(walks,      axis=0)
position2 = np.sum(walks**2,  axis=0)
pos_hist[:, :, :, :] = walks[:, pos_hist_times, :]
return position, position2, pos_hist, np.array(pos_hist_times)

```

## 3.8 Applications

**hpl 20:** Remaining: comment on boundary conditions for diffusion of substance and heat conduction (include Robin/cooling).

### 3.8.1 Diffusion of a substance

The first process to be considered is a substance that gets transported through a fluid at rest by pure diffusion. We consider an arbitrary volume  $V$  of this fluid, containing the substance with concentration function  $c(\mathbf{x}, t)$ . Physically, we can think of a very small volume with centroid  $\mathbf{x}$  at time  $t$  and assign the ratio of the volume of the substance and the total volume to  $c(\mathbf{x}, t)$ . This means that the mass of the substance in a small volume  $\Delta V$  is approximately  $\varrho c \Delta V$ , where  $\varrho$  is the density of the substance. Consequently, the total mass of the substance inside the volume  $V$  is the sum of all  $\varrho c \Delta V$ , which becomes the volume integral  $\int_V \varrho c dV$ .

Let us reason how the mass of the substance changes and thereby derive a PDE governing the concentration  $c$ . Suppose the substance flows out of  $V$  with a flux  $\mathbf{q}$ . If  $\Delta S$  is a small part of the boundary  $\partial V$  of  $V$ , the volume of the substance flowing out through  $dS$  in a small time interval  $\Delta t$  is  $\varrho \mathbf{q} \cdot \mathbf{n} \Delta t \Delta S$ , where  $\mathbf{n}$  is an outward unit normal to the boundary  $\partial V$ , see Figure 3.23. We realize that only the normal component of  $\mathbf{q}$  is able to transport mass in and out of  $V$ . The total outflow of the mass of the substance in a small time interval  $\Delta t$  becomes the surface integral

$$\int_{\partial V} \varrho \mathbf{q} \cdot \mathbf{n} \Delta t dS.$$

Assuming conservation of mass, this outflow of mass must be balanced by a loss of mass inside the volume. The increase of mass inside the volume, during a small time interval  $\Delta t$ , is

$$\int_V \varrho(c(\mathbf{x}, t + \Delta t) - c(\mathbf{x}, t)) dV,$$

assuming  $\varrho$  is constant, which is reasonable.

Setting the two contributions equal to each other ensures balance of mass inside  $V$ . Dividing by  $\Delta t$  gives

$$\int_V \varrho \frac{c(\mathbf{x}, t + \Delta t) - c(\mathbf{x}, t)}{\Delta t} dV = - \int_{\partial V} \varrho \mathbf{q} \cdot \mathbf{n} dS.$$

Note the minus sign on the right-hand side: the left-hand side expresses loss of mass, while the integral on the right-hand side is the gain of mass.

Now, letting  $\Delta t \rightarrow 0$ , we have

$$\frac{c(\mathbf{x}, t + \Delta t) - c(\mathbf{x}, t)}{\Delta t} \rightarrow \frac{\partial c}{\partial t},$$

so

$$\int_V \varrho \frac{\partial c}{\partial t} dV + \int_{\partial V} \varrho \mathbf{q} \cdot \mathbf{n} dS = 0. \quad (3.119)$$

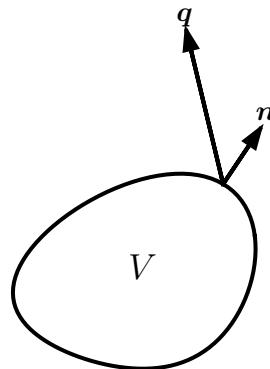
To arrive at a PDE, we express the surface integral as a volume integral using Gauss' divergence theorem:

$$\int_V (\varrho \frac{\partial c}{\partial t} + \nabla \cdot (\varrho \mathbf{q})) dV = 0.$$

Since  $\varrho$  is constant, we can divide by this quantity. If the integral is to vanish for an arbitrary volume  $V$ , the integrand must vanish too, and we get the mass conservation PDE for the substance:

$$\frac{\partial c}{\partial t} + \nabla \cdot \mathbf{q} = 0. \quad (3.120)$$

A fundamental problem is that this is a scalar PDE for four unknowns:  $c$  and the three components of  $\mathbf{q}$ . We therefore need additional equations. Here, Fick's law comes at rescue: it models how the flux  $\mathbf{q}$  of the substance is related to the concentration  $c$ . Diffusion is recognized by mass flowing from regions with high concentration to regions of low concentration.



**Fig. 3.23** An arbitrary volume of a fluid.

This principle suggests that  $\mathbf{q}$  is proportional to the negative gradient of  $c$ :

$$\mathbf{q} = -\alpha \nabla c, \quad (3.121)$$

where  $\alpha$  is an empirically determined constant. The relation (3.121) is known as Fick's law. Inserting (3.121) in (3.120) gives a scalar PDE for the concentration  $c$ :

$$\frac{\partial c}{\partial t} = \alpha \nabla^2 c. \quad (3.122)$$

### 3.8.2 Heat conduction

Heat conduction is a well-known diffusion process. The governing PDE is in this case based on the first law of thermodynamics: the increase in energy of a system is equal to the work done on the system, plus the supplied heat. Here, we shall media at rest and neglect the work done on the system. The principle then reduces to a balance between increase in internal energy and the supplied heat flow by conduction.

Let  $e(\mathbf{x}, t)$  be the *internal energy* per unit mass. The increase of the internal energy in a small volume  $\Delta V$  in a small time interval  $\Delta t$  is then

$$\varrho(e(\mathbf{x}, t + \Delta t) - e(\mathbf{x}, t))\Delta V,$$

where  $\varrho$  is the density of the material subject to heat conduction. In an arbitrary volume  $V$ , as depicted in Figure 3.23, the corresponding increase in internal energy becomes the volume integral

$$\int_V \varrho(e(\mathbf{x}, t + \Delta t) - e(\mathbf{x}, t)) dV.$$

This increase in internal energy is balanced by heat supplied by conduction. Let  $\mathbf{q}$  be the heat flow per time unit. Through the surface  $\partial V$  of  $V$  the following amount of heat flows out of  $V$  during a time interval  $\Delta t$ :

$$\int_{\partial V} \mathbf{q} \cdot \mathbf{n} \Delta t dS.$$

The simplified version of the first law of thermodynamics then states that

$$\int_V \varrho(e(\mathbf{x}, t + \Delta t) - e(\mathbf{x}, t)) dV = - \int_{\partial V} \mathbf{q} \cdot \mathbf{n} \Delta t dS.$$

The minus sign on the right-hand side ensures that the integral there models net *inflow* of heat (since  $\mathbf{n}$  is an outward unit normal,  $\mathbf{q} \cdot \mathbf{n}$  models *outflow*). Dividing by  $\Delta t$  and notifying that

$$\lim_{\Delta t \rightarrow 0} \frac{e(\mathbf{x}, t + \Delta t) - e(\mathbf{x}, t)}{\Delta t} = \frac{\partial e}{\partial t},$$

we get (in the limit  $\Delta t \rightarrow 0$ )

$$\int_V \varrho \frac{\partial e}{\partial t} dV + \int_{\partial V} \mathbf{q} \cdot \mathbf{n} \Delta t dS = 0.$$

This is the integral equation for heat conduction, but we aim at a PDE. The next step is therefore to transform the surface integral to a volume integral via Gauss' divergence theorem. The result is

$$\int_V \left( \varrho \frac{\partial e}{\partial t} + \nabla \cdot \mathbf{q} \right) dV = 0.$$

If this equality is to hold for all volumes  $V$ , the integrand must vanish, and we have the PDE

$$\varrho \frac{\partial e}{\partial t} = -\nabla \cdot \mathbf{q}. \quad (3.123)$$

Sometimes the supplied heat can come from the medium itself. This is the case, for instance, when radioactive rock generates heat. Let us add this effect. If  $f(\mathbf{x}, t)$  is the supplied heat per unit volume per unit time, the heat supplied in a small volume is  $f \Delta t \Delta V$ , and inside an arbitrary volume  $V$  the supplied generated heat becomes

$$\int_V f \Delta t \Delta V.$$

Adding this to the integral statement of the (simplified) first law of thermodynamics, and continuing the derivation, leads to the PDE

$$\varrho \frac{\partial e}{\partial t} = -\nabla \cdot \mathbf{q} + f. \quad (3.124)$$

There are four unknown scalar fields:  $e$  and  $\mathbf{q}$ . Moreover, the temperature  $T$ , which is our primary quantity to compute, does not enter the model yet. We need an additional equation, called the *equation of state*, relating  $e$ ,  $V = 1/\varrho =$ , and  $T$ :  $e = e(V, T)$ . By the chain rule we have

$$\frac{\partial e}{\partial t} = \left. \frac{\partial e}{\partial T} \right|_V \frac{\partial T}{\partial t} + \left. \frac{\partial e}{\partial V} \right|_T \frac{\partial V}{\partial t}.$$

The first coefficient  $\partial e / \partial T$  is called *specific heat capacity at constant volume*, denoted by  $c_v$ :

$$c_v = \left. \frac{\partial e}{\partial T} \right|_V.$$

The specific heat capacity will in general vary with  $T$ , but taking it as a constant is a good approximation in many applications.

The term  $\partial e / \partial V$  models effects due to compressibility and volume expansion. These effects are often small and can be neglected. We shall do so here. Using  $\partial e / \partial t = c_v \partial T / \partial t$  in the PDE gives

$$\varrho c_v \frac{\partial T}{\partial t} = -\nabla \cdot \mathbf{q} + f.$$

We still have four unknown scalar fields ( $T$  and  $\mathbf{q}$ ). To close the system, we need a relation between the heat flux  $\mathbf{q}$  and the temperature  $T$  called Fourier's law:

$$\mathbf{q} = -k \nabla T,$$

which simply states that heat flows from hot to cold areas, along the path of greatest variation. In a solid medium,  $k$  depends on the material of the medium, and in multi-material media one must regard  $k$  as spatially dependent. In a fluid, it is common to assume that  $k$  is a constant. The value of  $k$  reflects how easy heat is conducted through the medium, and  $k$  is named the *coefficient of heat conduction*.

We have now one scalar PDE for the unknown temperature field  $T(\mathbf{x}, t)$ :

$$\varrho c_v \frac{\partial T}{\partial t} = \nabla \cdot (k \nabla T) + f. \quad (3.125)$$

### 3.8.3 Development of flow between two flat plates

Diffusion equations may also arise as simplified versions of other mathematical models, especially in fluid flow. Consider a fluid flowing between two flat, parallel plates. The velocity is uni-directional, say along the  $z$  axis, and depends only on the distance  $x$  from the plates;  $\mathbf{u} = u(x, t)\mathbf{k}$ . The flow is governed by the Navier-Stokes equations,

$$\begin{aligned} \varrho \frac{\partial \mathbf{u}}{\partial t} + \varrho \mathbf{u} \cdot \nabla \mathbf{u} &= -\nabla p + \mu \nabla^2 \mathbf{u} + \varrho \mathbf{f}, \\ \nabla \cdot \mathbf{u} &= 0, \end{aligned}$$

where  $p$  is the pressure field, unknown along with the velocity  $\mathbf{u}$ ,  $\varrho$  is the fluid density,  $\mu$  the dynamic viscosity, and  $\mathbf{f}$  is some external body force. The geometric restrictions of flow between two flat plates puts restrictions on the velocity,  $\mathbf{u} = u(x, t)\mathbf{i}$ , and the  $z$  component of the Navier-Stokes equations collapses to a diffusion equation:

$$\varrho \frac{\partial u}{\partial t} = -\frac{\partial p}{\partial z} + \mu \frac{\partial^2 u}{\partial z^2} + \varrho f_z,$$

if  $f_z$  is the component of  $\mathbf{f}$  in the  $z$  direction.

The boundary conditions are derived from the fact that the fluid sticks to the plates, which means  $\mathbf{u} = 0$  at the plates. Say the location of the plates are  $z = 0$  and  $z = L$ . We then have

$$u(0, t) = u(L, t) = 0.$$

One can easily show that  $\partial p/\partial z$  must be a constant or just a function of time  $t$ . We set  $\partial p/\partial z = -\beta(t)$ . The body force could be a component of gravity, if desired, set as  $f_z = \rho g$ . Switching from  $z$  to  $x$  as independent variable gives a very standard one-dimensional diffusion equation:

$$\rho \frac{\partial u}{\partial t} = \mu \frac{\partial^2 u}{\partial z^2} + \beta(t) + \rho \gamma g, \quad x \in [0, L], \quad t \in (0, T].$$

The boundary conditions are

$$u(0, t) = u(L, t) = 0,$$

while some initial condition

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

must also be prescribed.

The flow is driven by either the pressure gradient  $\beta$  or gravity, or a combination of both. One may also consider one moving plate that drives the fluid. If the plate at  $x = L$  moves with velocity  $U_L(t)$ , we have the adjusted boundary condition

$$u(L, t) = U_L(t).$$

**hpl 21:** Exercises based on this diffusion model.

### 3.8.4 Flow in a straight tube

Now we consider viscous fluid flow in a straight tube with radius  $R$  and rigid walls. The governing equations are the Navier-Stokes equations, but as in Section 3.8.3, it is natural to assume that the velocity is directed along the tube, and that it is axi-symmetric. These assumptions reduced the velocity field to  $\mathbf{u} = u(r, x, t)\mathbf{i}$ , if the  $x$  axis is directed along the tube. From the equation of continuity,  $\nabla \cdot \mathbf{u} = 0$ , we see that  $u$  must be independent of  $x$ . Inserting  $\mathbf{u} = u(r, t)\mathbf{i}$  in the Navier-Stokes equations, expressed in axi-symmetric cylindrical coordinates, results in

$$\rho \frac{\partial u}{\partial t} = \mu \frac{1}{r} \frac{\partial}{\partial r} \left( r \frac{\partial u}{\partial r} \right) + \beta(t) + \rho \gamma g, \quad r \in [0, R], \quad t \in (0, T]. \quad (3.126)$$

Here,  $\beta(t) = -\partial p/\partial x$  is the pressure gradient along the tube. The associated boundary condition is  $u(R, t) = 0$ .

### 3.8.5 Tribology: thin film fluid flow

Thin fluid films are extremely important inside machinery to reduce friction between gliding surfaces. The mathematical model for the fluid motion takes the form of a diffusion problem and is quickly derived here. We consider two solid surfaces whose distance is described by a gap function  $h(x, y)$ . The space between these surfaces is filled with a fluid with dynamic viscosity  $\mu$ . The fluid may move partially because of pressure gradients and partially because the surfaces move. Let  $Ui + Vj$  be the relative velocity of the two surfaces and  $p$  the pressure in the fluid. The mathematical model builds on two principles: 1) conservation of mass, 2) assumption of locally quasi-static flow between flat plates.

The conservation of mass equation reads  $\nabla \cdot \mathbf{u}$ , where  $\mathbf{u}$  is the local fluid velocity. For thin films the detailed variation between the surfaces is not of interest, so  $\nabla \cdot \mathbf{u} = 0$  is integrated (average) in the direction perpendicular to the surfaces. This gives rise to the alternative mass conservation equation

$$\nabla \cdot \mathbf{q} = 0, \quad \mathbf{q} = \int_0^{h(x,y)} \mathbf{u} dz,$$

where  $z$  is the coordinate perpendicular to the surfaces, and  $\mathbf{q}$  is then the volume flux in the fluid gap.

Locally, we may assume that we have steady flow between two flat surfaces, with a pressure gradient and where the lower surface is at rest and the upper moves with velocity  $Ui + Vj$ . The corresponding mathematical problem is actually the limit problem in Section 3.8.3 as  $t \rightarrow \infty$ . The limit problem can be solved analytically, and the local volume flux becomes

$$\mathbf{q}(x, y, z) = \int_0^h \mathbf{u}(x, y, z) dz = -\frac{h^3}{12\mu} \nabla p + \frac{1}{2} U h i + \frac{1}{2} V h j.$$

The idea is to use this expression locally also when the surfaces are not flat, but slowly varying, and if  $U$ ,  $V$ , or  $p$  varies in time, provided

the time variation is sufficiently slow. This is a common quasi-static approximation much used in mathematical modeling.

Inserting the expression for  $\mathbf{q}$  via  $p$ ,  $U$ , and  $V$  in the equation  $\nabla \cdot \mathbf{q} = 0$  gives a diffusion PDE for  $p$ :

$$\nabla \cdot \left( \frac{h^3}{12\mu} \nabla p \right) = \frac{1}{2} \frac{\partial}{\partial x} (hU) + \frac{1}{2} \frac{\partial}{\partial x} (hV). \quad (3.127)$$

The boundary conditions must involve  $p$  or  $\mathbf{q}$  at the boundary.

**hpl 22:** Exercise!

### 3.8.6 Propagation of electrical signals in the brain

One can make a model of how electrical signals are propagated along the neuronal fibers that receive synaptic inputs in the brain. The signal propagation is one-dimensional and can, in the simplest cases, be governed by the [Cable equation](#):

$$c_m \frac{\partial V}{\partial t} = \frac{1}{r_l} \frac{\partial^2 V}{\partial x^2} - \frac{1}{r_m} V \quad (3.128)$$

where  $V(x, t)$  is the voltage to be determined,  $c_m$  is capacitance of the neuronal fiber, while  $r_l$  and  $r_m$  are measures of the resistance. The boundary conditions are often taken as  $V = 0$  at a short circuit or open end,  $\partial V / \partial x = 0$  at a sealed end, or  $\partial V / \partial x \propto V$  where there is an injection of current.

## 3.9 Exercises

### Exercise 3.6: Stabilizing the Crank-Nicolson method by Rannacher time stepping

It is well known that the Crank-Nicolson method may give rise to non-physical oscillations in the solution of diffusion equations if the initial data exhibit jumps (see Section 3.3.6). Rannacher [7] suggested a stabilizing technique consisting of using the Backward Euler scheme for the first two

time steps with step length  $\frac{1}{2}\Delta t$ . One can generalize this idea to taking  $2m$  time steps of size  $\frac{1}{2}\Delta t$  with the Backward Euler method and then continuing with the Crank-Nicolson method, which is of second-order in time. The idea is that the high frequencies of the initial solution are quickly damped out, and the Backward Euler scheme treats these high frequencies correctly. Thereafter, the high frequency content of the solution is gone and the Crank-Nicolson method will do well.

Test this idea for  $m = 1, 2, 3$  on a diffusion problem with a discontinuous initial condition. Measure the convergence rate using the solution (3.45) with the boundary conditions (3.46)-(3.47) for  $t$  values such that the conditions are in the vicinity of  $\pm 1$ . For example,  $t < 5a1.6 \cdot 10^{-2}$  makes the solution diffusion from a step to almost a straight line. The program `diffu_erf_sol.py` shows how to compute the analytical solution.

### Project 3.7: Energy estimates for diffusion problems

This project concerns so-called *energy estimates* for diffusion problems that can be used for qualitative analytical insight and for verification of implementations.

**a)** We start with a 1D homogeneous diffusion equation with zero Dirichlet conditions:

$$u_t = \alpha u_{xx}, \quad x \in \Omega = (0, L), \quad t \in (0, T], \quad (3.129)$$

$$u(0, t) = u(L, t) = 0, \quad t \in (0, T], \quad (3.130)$$

$$u(x, 0) = I(x), \quad x \in [0, L]. \quad (3.131)$$

The energy estimate for this problem reads

$$\|u\|_{L^2} \leq \|I\|_{L^2}, \quad (3.132)$$

where the  $\|\cdot\|_{L^2}$  norm is defined by

$$\|g\|_{L^2} = \sqrt{\int_0^L g^2 dx}. \quad (3.133)$$

The quantity  $\|u\|_{L^2}$  or  $\frac{1}{2}\|u\|_{L^2}$  is known as the *energy* of the solution, although it is not the physical energy of the system. A mathematical tradition has introduced the notion *energy* in this context.

The estimate (3.132) says that the “size of  $u$ ” never exceeds that of the initial condition, or more precisely, it says that the area under the  $u$  curve decreases with time.

To show (3.132), multiply the PDE by  $u$  and integrate from 0 to  $L$ . Use that  $uu_t$  can be expressed as the time derivative of  $u^2$  and that  $u_x xu$  can be integrated by parts to form an integrand  $u_x^2$ . Show that the time derivative of  $\|u\|_{L^2}^2$  must be less than or equal to zero. Integrate this expression and derive (3.132).

**b)** Now we address a slightly different problem,

$$u_t = \alpha u_x x + f(x, t), \quad x \in \Omega = (0, L), \quad t \in (0, T], \quad (3.134)$$

$$u(0, t) = u(L, t) = 0, \quad t \in (0, T], \quad (3.135)$$

$$u(x, 0) = 0, \quad x \in [0, L]. \quad (3.136)$$

The associated energy estimate is

$$\|u\|_{L^2} \leq \|f\|_{L^2}. \quad (3.137)$$

(This result is more difficult to derive.)

Now consider the compound problem with an initial condition  $I(x)$  and a right-hand side  $f(x, t)$ :

$$u_t = \alpha u_x x + f(x, t), \quad x \in \Omega = (0, L), \quad t \in (0, T], \quad (3.138)$$

$$u(0, t) = u(L, t) = 0, \quad t \in (0, T], \quad (3.139)$$

$$u(x, 0) = I(x), \quad x \in [0, L]. \quad (3.140)$$

Show that if  $w_1$  fulfills (3.129)-(3.131) and  $w_2$  fulfills (3.134)-(3.136), then  $u = w_1 + w_2$  is the solution of (3.138)-(3.140). Using the triangle inequality for norms,

$$\|a + b\| \leq \|a\| + \|b\|,$$

show that the energy estimate for (3.138)-(3.140) becomes

$$\|u\|_{L^2} \leq \|I\|_{L^2} + \|f\|_{L^2}. \quad (3.141)$$

**c)** One application of (3.141) is to prove uniqueness of the solution. Suppose  $u_1$  and  $u_2$  both fulfill (3.138)-(3.140). Show that  $u = u_1 - u_2$  then fulfills (3.138)-(3.140) with  $f = 0$  and  $I = 0$ . Use (3.141) to deduce that the energy must be zero for all times and therefore that  $u_1 = u_2$ , which proves that the solution is unique.

**d)** Generalize (3.141) to a 2D/3D diffusion equation  $u_t = \nabla \cdot (\alpha \nabla u)$  for  $x \in \Omega$ .

**Hint.** Use integration by parts in multi dimensions:

$$\int_{\Omega} u \nabla \cdot (\alpha \nabla u) dx = - \int_{\Omega} \alpha \nabla u \cdot \nabla u dx + \int_{\partial\Omega} u \alpha \frac{\partial u}{\partial n},$$

where  $\frac{\partial u}{\partial n} = \mathbf{n} \cdot \nabla u$ ,  $\mathbf{n}$  being the outward unit normal to the boundary  $\partial\Omega$  of the domain  $\Omega$ .

**e)** Now we also consider the multi-dimensional PDE  $u_t = \nabla \cdot (\alpha \nabla u)$ . Integrate both sides over  $\Omega$  and use Gauss' divergence theorem,  $\int_{\Omega} \nabla \cdot \mathbf{q} dx = \int_{\partial\Omega} \mathbf{q} \cdot \mathbf{n} ds$  for a vector field  $\mathbf{q}$ . Show that if we have homogeneous Neumann conditions on the boundary,  $\partial u / \partial n = 0$ , area under the  $u$  surface remains constant in time and

$$\int_{\Omega} u dx = \int_{\Omega} I dx. \quad (3.142)$$

**f)** Establish a code in 1D, 2D, or 3D that can solve a diffusion equation with a source term  $f$ , initial condition  $I$ , and zero Dirichlet or Neumann conditions on the whole boundary.

We can use (3.141) and (3.142) as a partial verification of the code. Choose some functions  $f$  and  $I$  and check that (3.141) is obeyed at any time when zero Dirichlet conditions are used. Iterate over the same  $I$  functions and check that (3.142) is fulfilled when using zero Neumann conditions.

**g)** Make a list of some possible bugs in the code, such as indexing errors in arrays, failure to set the correct boundary conditions, evaluation of a term at a wrong time level, and similar. For each of the bugs, see if the verification tests from the previous subexercise pass or fail. This investigation shows how strong the energy estimates and the estimate (3.142) are for pointing out errors in the implementation.

Filename: `diffu_energy`.

### Exercise 3.8: Splitting methods and preconditioning

In Section 3.6.15, we outlined a class of iterative methods for  $Au = b$  based on splitting  $A$  into  $A = M - N$  and introducing the iteration

$$Mu^k = Nu^k + b.$$

The very simplest splitting is  $M = I$ , where  $I$  is the identity matrix. Show that this choice corresponds to the iteration

$$u^k = u^{k-1} + r^{k-1}, \quad r^{k-1} = b - Au^{k-1}, \quad (3.143)$$

where  $r^{k-1}$  is the residual in the linear system in iteration  $k - 1$ . The formula (3.143) is known as Richardson's iteration. Show that if we apply the simple iteration method (3.143) to the *preconditioned* system  $M^{-1}Au = M^{-1}b$ , we arrive at the Jacobi method by choosing  $M = D$  (the diagonal of  $A$ ) as preconditioner and the SOR method by choosing  $M = \omega^{-1}D + L$  ( $L$  being the lower triangular part of  $A$ ). This equivalence shows that we can apply one iteration of the Jacobi or SOR method as preconditioner.

### Exercise 3.9: Oscillating surface temperature of the earth

Consider a day-and-night or seasonal variation in temperature at the surface of the earth. How deep down in the ground will the surface oscillations reach? For simplicity, we model only the vertical variation along a coordinate  $x$ , where  $x = 0$  at the surface, and  $x$  increases as we go down in the ground. The temperature is governed by the heat equation

$$\varrho c_v \frac{\partial T}{\partial t} = \nabla \cdot (k \nabla T),$$

in some spatial domain  $x \in [0, L]$ , where  $L$  is chosen large enough such that we can assume that  $T$  is approximately constant, independent of the surface oscillations, for  $x > L$ . The parameters  $\varrho$ ,  $c_v$ , and  $k$  are the density, the specific heat capacity at constant volume, and the heat conduction coefficient, respectively.

- a) Derive the mathematical model for computing  $T(x, t)$ . Assume the surface oscillations to be sinusoidal around some mean temperature  $T_m$ . Let  $T = T_m$  initially. At  $x = L$ , assume  $T \approx T_m$ .

**b)** Scale the model in a) assuming  $k$  is constant. Use a time scale  $t_c = \omega^{-1}$  and a length scale  $x_c = \sqrt{2\alpha/\omega}$ , where  $\alpha = k/(\rho c_v)$ . The primary unknown can be scaled as  $\frac{T-T_m}{2A}$ .

Show that the scaled PDE is

$$\frac{\partial u}{\partial \bar{t}} = \frac{1}{2} \frac{\partial^2 u}{\partial \bar{x}^2},$$

with initial condition  $u(\bar{x}, 0) = 0$ , left boundary condition  $u(0, \bar{t}) = \sin(\bar{t})$ , and right boundary condition  $u(\bar{L}, \bar{t}) = 0$ . The bar indicates a dimensionless quantity.

Show that  $u(\bar{x}, \bar{t}) = e^{-\bar{x}} \sin(\bar{x} - \bar{t})$  is a solution that fulfills the PDE and the boundary condition at  $\bar{x} = 0$  (this is the solution we will experience as  $\bar{t} \rightarrow \infty$  and  $\bar{L} \rightarrow \infty$ ). Conclude that an appropriate domain for  $x$  is  $[0, 4]$  if a damping  $e^{-4} \approx 0.18$  is appropriate for implementing  $\bar{u} \approx \text{const}$ ; increasing to  $[0, 6]$  damps  $\bar{u}$  to 0.0025.

**c)** Compute the scaled temperature and make animations comparing two solutions with  $\bar{L} = 4$  and  $\bar{L} = 8$ , respectively (keep  $\Delta x$  the same).

### Exercise 3.10: Oscillating and pulsating flow in tubes

We consider flow in a straight tube with radius  $R$  and straight walls. The flow is driven by a pressure gradient  $\beta(t)$ . The effect of gravity can be neglected. The mathematical problem reads

$$\varrho \frac{\partial u}{\partial t} = \mu \frac{1}{r} \frac{\partial}{\partial r} \left( r \frac{\partial u}{\partial r} \right) + \beta(t), \quad r \in [0, R], \quad t \in (0, T], \quad (3.144)$$

$$u(r, 0) = 0, \quad r \in [0, R], \quad (3.145)$$

$$u(R, 0) = 0, \quad t \in (0, T], \quad (3.146)$$

$$\frac{\partial u}{\partial r}(0, t) = 0, \quad t \in (0, T]. \quad (3.147)$$

We consider two models for  $\beta(t)$ . One plain, sinusoidal oscillation:

$$\beta = A \sin(\omega t), \quad (3.148)$$

and one with periodic pulses,

$$\beta = A \sin^{16}(\omega t), \quad (3.149)$$

Note that both models can be written as  $\beta = A \sin^m(\omega t)$ , with  $m = 1$  and  $m = 16$ , respectively.

**a)** Scale the mathematical model, using the viscous time scale  $\varrho R^2 / \mu$ .

**b)** Implement the scaled model from a), using the unifying  $\theta$  scheme in time and centered differences in space.

**c)** Verify the implementation in b) using a manufactured solution that is quadratic in  $r$  and linear in  $t$ . Make a corresponding test function.

**Hint.** You need to include an extra source term in the equation to allow for such tests. Let the spatial variation be  $1 - r^2$  such that the boundary condition is fulfilled.

**d)** Make animations for  $m = 1, 16$  and  $\alpha = 1, 0.1$ . Choose  $T$  such that the motion has reached a steady state (non-visible changes from period to period in  $u$ ).

**e)** For  $\alpha \gg 1$ , the scaling in a) is not good, because the characteristic time for changes (due to the pressure) is much smaller than the viscous diffusion time scale ( $\alpha$  becomes large). We should in this case base the short time scale on  $1/\omega$ . Scale the model again, and make an animation for  $m = 1, 16$  and  $\alpha = 10$ .

Filename: `axyiyymm_flow`.



Wave (Chapter 2) and diffusion (Chapter 3) equations are solved reliably by finite difference methods. As soon as we add a first-order derivative in space, representing *advective transport* (also known as convective transport), the numerics gets more complicated, and intuitively attractive methods no longer work well. We shall show how and why methods fail and provide remedies. The present chapter builds on basic knowledge about finite difference methods for the diffusion equations, including the analysis by Fourier components, truncation error analysis (Appendix B), and compact difference notation.

## 4.1 One-dimensional time-dependent advection equations

We consider the pure advection model

$$\frac{\partial u}{\partial t} + v \frac{\partial u}{\partial x} = 0, \quad x \in (0, L), \quad t \in (0, T], \quad (4.1)$$

$$u(x, 0) = I(x), \quad x \in (0, L), \quad (4.2)$$

$$u(0, t) = U_0, \quad t \in (0, T]. \quad (4.3)$$

In (4.1),  $v$  is a given parameter, typically reflecting the velocity of transport of a quantity  $u$  with a flow. There is only one boundary condition (4.2) since there is only a first-derivative term in the PDE

(4.1). The information at  $x = 0$  and the initial condition get transported in positive  $x$  direction if  $v > 0$  through the domain.

The solution of (4.1) in an infinite domain (no condition (4.2) at  $x = 0$ ) is

$$u(x, t) = I(x - vt). \quad (4.4)$$

This is also the solution we expect if we let  $I(x)$  be located in the interior of the domain such that  $\lim_{x \rightarrow 0, L} I(x) = 0$  and  $U_0 = 0$ .

### 4.1.1 Simplest scheme: forward in time, centered in space

**Method.** A first attempt to solve a PDE like (4.1) will normally look for a time-discretization scheme that is explicit so we avoid solving systems of linear equations. In space, we anticipate that centered differences are most accurate and therefore best. These two arguments lead us to a Forward Euler scheme in time and centered differences in space:

$$[D_t^+ u + v D_{2x} u = 0]_i^n \quad (4.5)$$

Written out,

$$u^{n+1} = u^n - \frac{1}{2} C (u_{i+1}^n - u_{i-1}^n),$$

with  $C$  as the Courant number

$$C = \frac{v \Delta t}{\Delta x}.$$

**Implementation.** A solver function for our scheme goes as follows.

```
import numpy as np
import matplotlib.pyplot as plt

def solver_FECS(I, U0, v, L, dt, C, T, user_action=None):
    Nt = int(round(T/float(dt)))
    t = np.linspace(0, Nt*dt, Nt+1)    # Mesh points in time
    dx = v*dt/C
    Nx = int(round(L/dx))
    x = np.linspace(0, L, Nx+1)        # Mesh points in space
    # Make sure dx and dt are compatible with x and t
    dx = x[1] - x[0]
    dt = t[1] - t[0]
    C = v*dt/dx

    u = np.zeros(Nx+1)
```

```

u_1 = np.zeros(Nx+1)

# Set initial condition u(x,0) = I(x)
for i in range(0, Nx+1):
    u_1[i] = I(x[i])

if user_action is not None:
    user_action(u_1, x, t, 0)

for n in range(0, Nt):
    # Compute u at inner mesh points
    for i in range(1, Nx):
        u[i] = u_1[i] - 0.5*C*(u_1[i+1] - u_1[i-1])

    # Insert boundary condition
    u[0] = u0

    if user_action is not None:
        user_action(u, x, t, n+1)

    # Switch variables before next step
    u_1, u = u, u_1

```

**Test cases.** The typical solution is  $u$  has the shape of  $I$  and is transported at velocity  $v$  to the right (if  $v > 0$ ). Let us choose a smooth and non-smooth initial condition:

$$u(x, 0) = Ae^{-\frac{1}{2}\left(\frac{x-L/10}{\sigma}\right)^2}, \quad (4.6)$$

$$u(x, 0) = A \cos\left(\frac{5\pi}{L}\left(x - \frac{L}{10}\right)\right), \quad x < \frac{L}{5} \text{ else } 0. \quad (4.7)$$

The parameter  $A$  is the maximum value of the initial condition.

We scale the problem and introduce  $\bar{x} = x/L$  and  $\bar{t} = vt/L$ , which gives

$$\frac{\partial \bar{u}}{\partial \bar{t}} + \frac{\partial \bar{u}}{\partial \bar{x}} = 0.$$

The unknown  $u$  is scaled by the maximum value of the initial condition:  $\bar{u} = u / \max |I(x)|$  such that  $|\bar{u}(\bar{x}, 0)| \in [0, 1]$ . The scaled problem is solved by setting  $v = 1$ ,  $L = 1$ , and  $A = 1$ . From now on we drop the bars.

To run a case and plot the solution, we make the function

```

def run_FECS(case):
    if case == 'gaussian':
        def I(x):
            return np.exp(-0.5*((x-L/10)/sigma)**2)

```

```

        elif case == 'cosinehat':
            def I(x):
                return np.cos(np.pi*5/L*(x - L/10)) if x < L/5 else 0

            L = 1.0
            sigma = 0.02
            legends = []

            def plot(u, x, t, n):
                """Plot every m steps in the same figure."""
                m = 40
                if n % m != 0:
                    return
                print 't=%g, n=%d, u in [%g, %g] w/%d points' % \
                      (t[n], n, u.min(), u.max(), x.size)
                if np.abs(u).max() > 3: # Instability?
                    return
                plt.plot(x, u)
                legends.append('t=%g' % t[n])
                if n > 0:
                    plt.hold('on')

            U0 = 0
            dt = 0.001
            C = 1
            T = 1
            solver(I=I, U0=U0, v=1.0, L=L, dt=dt, C=C, T=T,
                   user_action=plot)
            plt.legend(legends, loc='lower left')
            plt.savefig('tmp.png'); plt.savefig('tmp.pdf')
            plt.axis([0, L, -0.75, 1.1])
            plt.show()

```

**Bug?** Running either of the test cases, the plot becomes a mess, and the printout of  $u$  values in the `plot` function reveals that  $u$  grows very quickly. We may reduce  $\Delta t$  and make it very small, yet the solution just grows. Such behavior points to a bug in the code. However, choosing a coarse mesh and performing a time step by hand calculations produce the same numbers as in the code, so it seems that the implementation is correct. The hypothesis is therefore that the solution is unstable.

### 4.1.2 Analysis of the scheme

We can analyze the finite difference scheme by look at how it treats a Fourier component

$$u_q^n = A^n e^{ikq\Delta x}.$$

The corresponding analytical Fourier component is

$$u = A_e^n e^{ikx}, \quad A_e = e^{-ivk\Delta t} = e^{-iCkx}.$$

In particular,  $|A_e| \leq 1$ .

Inserting the numerical component in the scheme,

$$[D_t^+ A^n e^{ikq\Delta x} + v D_{2x} A^n e^{ikq\Delta x}]_i^n,$$

and making use of (A.25) results in

$$[e^{ikq\Delta x} \left( \frac{A - 1}{\Delta t} + v \frac{1}{\Delta x} i \sin(k\Delta x) \right)]_i^n,$$

which implies

$$A = 1 - iC \sin(k\Delta x).$$

The numerical solution features the formula  $A^n$ . To find out whether  $A^n$  means growth in time, we rewrite  $A$  in polar form:  $A = A_r e^{i\phi}$ , for real numbers  $A_r$  and  $\phi$ , since we then have  $A^n = A_r^n e^{i\phi n}$ . The magnitude of  $A^n$  is  $A_r^n$ . In our case,  $A_r = (1 + C^2 \sin^2(kx))^{1/2} > 1$ , so  $A_r^n$  will increase in time, whereas the exact solution will not. Regardless of  $\Delta t$ , we get unstable numerical solutions.

### 4.1.3 Leapfrog in time, centered differences in space

**Method.** Another explicit scheme is to do a “leapfrog” jump over  $2\Delta t$  in time and combine it with central differences in space:

$$[D_{2t} u + v D_{2x} u] = 0,$$

which results in the updating formula

$$u^{n+1} = u^{n-1} - C(u_{i+1} - u_{i-1}).$$

A special scheme is needed to compute  $u^1$ , but we leave that problem or now.

**Implementation.** We now need to have three time levels and must modify our solver a bit:

```
Nt = int(round(T/float(dt)))
t = np.linspace(0, Nt*dt, Nt+1)    # Mesh points in time
...
```

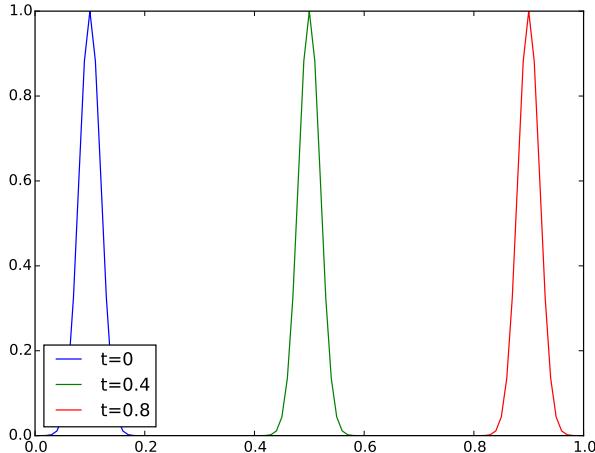
```

u    = np.zeros(Nx+1)
u_1 = np.zeros(Nx+1)
u_2 = np.zeros(Nx+1)
...
for n in range(0, Nt):
    if scheme == 'FECS':
        for i in range(1, Nx):
            u[i] = u_1[i] - 0.5*C*(u_1[i+1] - u_1[i-1])
    elif scheme == 'LFCS':
        if n == 0:
            # Use some scheme for the first step
            for i in range(1, Nx):
                ...
        else:
            for i in range(1, Nx+1):
                u[i] = u_2[i] - C*(u_1[i] - u_1[i-1])

    # Switch variables before next step
    u_2, u_1, u = u_1, u, u_2

```

**Running the test case.** Let us try a coarse mesh such that the smooth Gaussian initial condition is represented by 1 at mesh node 1 and 0 at all other nodes. This triangular initial condition should then be advected to the right. Choosing scaled variables as  $\Delta t = 0.1$ ,  $T = 1$ , and  $C = 1$  gives the plot in Figure 4.1, which is in fact identical to the exact solution (!).



**Fig. 4.1** Leapfrog scheme with  $\Delta t = 0.1$  and  $C = 1$ .

**Analysis.** We can perform a Fourier analysis again. Inserting the numerical Fourier component in the Leapfrog scheme, we get

$$A^2 - i2C \sin(k\Delta x)A - 1 = 0,$$

and

$$A = -iC \sin(k\Delta x) \pm \sqrt{1 - C^2 \sin^2(k\Delta x)}.$$

Rewriting to polar form,  $A = A_r e^{i\phi}$ , we see that  $A_r = 1$ , so the numerical component is not increasing or decreasing in time, which is exactly what we want. However, for  $C > 1$ , the square root can become complex valued, so stability is obtained only as long as  $C \leq 1$ .

**hpl 24:** The complete numerical solution? How are the two roots combined?

We introduce  $p = k\Delta x$ . The amplification factor now reads

$$A = -iC \sin p \pm \sqrt{1 - C^2 \sin^2 p},$$

and is to be compared to the exact amplification factor

$$A_e = e^{-ikv\Delta t} = e^{-ikC\Delta x} = e^{-iCp}.$$

Section 4.1.7 compares many numerical amplification factors with the exact expression.

#### 4.1.4 Upwind differences in space

Since the PDE reflects transport of information along with a flow in positive  $x$  direction  $v > 0$ , it could be natural to go upstream and not downstream in a spatial derivative. That is, we approximate

$$\frac{\partial u}{\partial x}(x_i, t_n) \approx [D_x^- u]_i^n = \frac{u_i^n - u_{i-1}^n}{\Delta x}.$$

This is called an *upwind difference*. This spatial approximation does magic to the scheme we had with Forward Euler in time and centered difference in space. With an upwind difference,

$$[D_t^+ u + v D_x^- u = 0]_i^n, \quad (4.8)$$

written out as

$$u^{n+1} = u_i^n - C(u_i^n - u_{i-1}^n),$$

gives a robust scheme that is stable if  $C \leq 1$ . As with the Leapfrog scheme, it becomes exact if  $C = 1$ . A plot of the solution in case  $C = 1$  is therefore given in Figure 4.1.

The amplification factor can be computed using the formula (A.23),

$$\frac{A - 1}{\Delta t} + \frac{v}{\Delta x} (1 - e^{-ik\Delta x}) = 0,$$

which means

$$A = 1 - C(1 - \cos(p) - i \sin(p)).$$

For  $C < 1$  there is, unfortunately, non-physical damping of discrete Fourier components. This damping can be quite severe.

One can interpret the upwind difference as extra, artificial diffusion in the equation. Solving

$$\frac{\partial u}{\partial t} + v \frac{\partial u}{\partial x} = \nu \frac{\partial^2 u}{\partial x^2},$$

by a forward difference time and centered differences in space,

$$D_t^+ u + v D_{2x} u = \nu D_x D_x u]_i^n,$$

gives actually the upwind scheme (4.8) if  $\nu = v\Delta x/2$ . That is, solving the PDE  $u_t + vu_x = 0$  by centered differences in space and forward difference in time is unsuccessful, but by adding some artificial diffusion  $\nu u_{xx}$ , the method becomes stable.

#### 4.1.5 A Crank-Nicolson discretization in time and centered differences in space

$$A = \frac{1 - \frac{1}{2}iC \sin p}{1 + \frac{1}{2}iC \sin p}.$$

#### 4.1.6 The Lax-Wendroff method

$$A = 1 - iC \sin p - 2C^2 \sin^2(p/2).$$

### 4.1.7 Analysis of dispersion relations

We have developed expressions for  $A(C, p)$  in the exact solution  $u_q^n = A^n e^{ikq\Delta x}$  of the discrete equations. These expressions are valuable for investigating the quality of the numerical solutions. Note that the Fourier component that solves the original PDE problem has no damping and moves with constant velocity  $v$ . There are two basic errors in the numerical Fourier component: there may be damping and the wave velocity may depend on  $C$  and  $p = k\Delta x$ .

The shortest wavelength that can be represented is  $\lambda = 2\Delta x$ . The corresponding  $k$  is  $k = 2\pi/\lambda = \pi/Delta x$ , so  $p = k\Delta x \in (0, \pi]$ .

Given a complex  $A$  as a function of  $C$  and  $p$ , how can we visualize it? The two key ingredients in  $A$  is the magnitude, reflecting damping or growth of the wave, and the angle, closely related to the velocity of the wave. The Fourier component

$$D^n e^{ik(x-ct)}$$

has damping  $D$  and wave velocity  $c$ . Let us express our  $A$  in polar form,  $A = A_r e^{i\phi}$ , and insert this expression in our discrete component  $u_q^n = A^n e^{ikq\Delta x} = A^n e^{ikx}$ :

$$u_q^n = A_r^n e^{i\phi n} e^{ikx} = A_r^n e^{i(kx-n\phi)} = A_r^n e^{i(k(x-ct))},$$

for

$$c = \frac{\phi}{k\Delta t}.$$

Now,

$$k\Delta t = \frac{Ck\Delta x}{v} = \frac{Cp}{v},$$

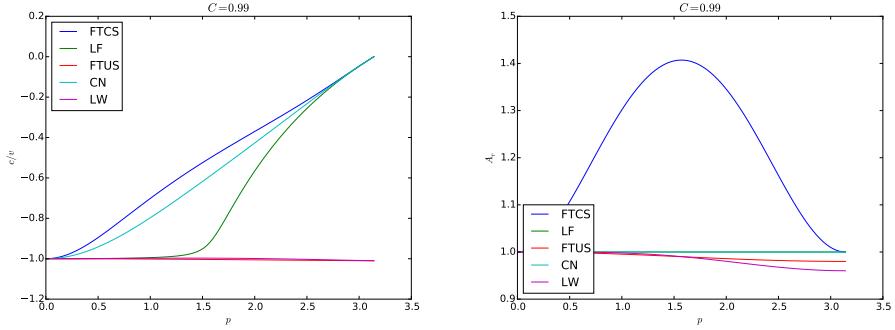
so

$$c = \frac{\phi v}{Cp}.$$

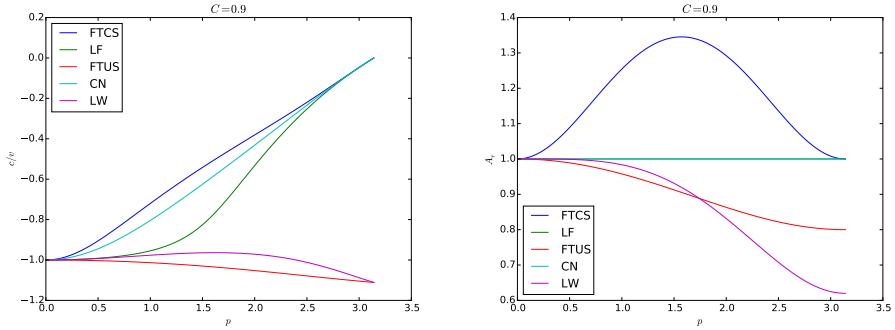
An appropriate dimensionless quantity to plot is  $c/v$ :

$$\frac{c}{v} = \frac{\phi}{Cp}.$$

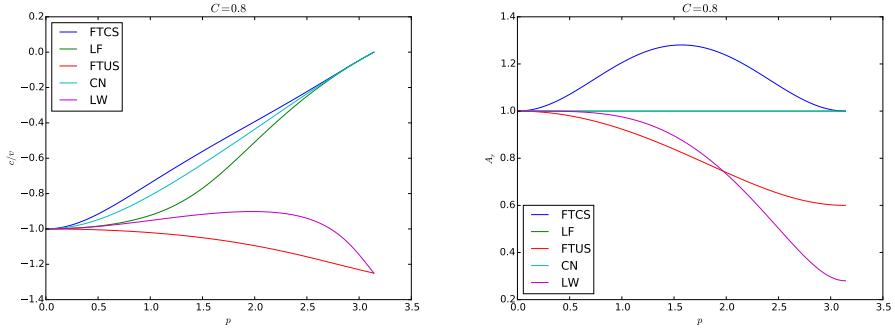
The total damping after some time  $T = n\Delta t$  is reflected by  $A_r(C, p)^n$ . Since normally  $A_r < 1$ , the damping goes like  $A_r^{1/\Delta t}$  and approaches zero



**Fig. 4.2** Dispersion relations for  $C = 0.99$ .

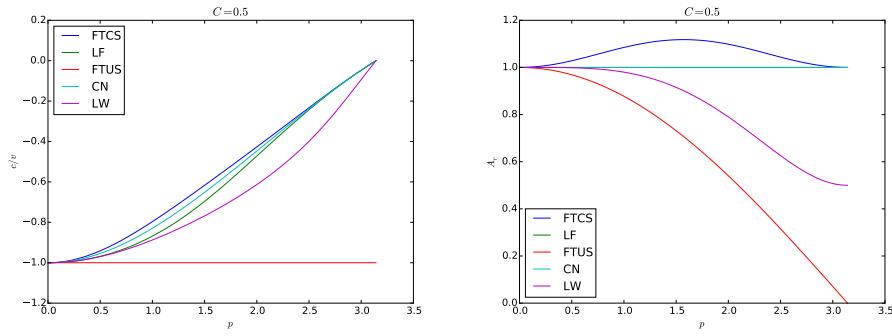


**Fig. 4.3** Dispersion relations for  $C = 0.9$ .

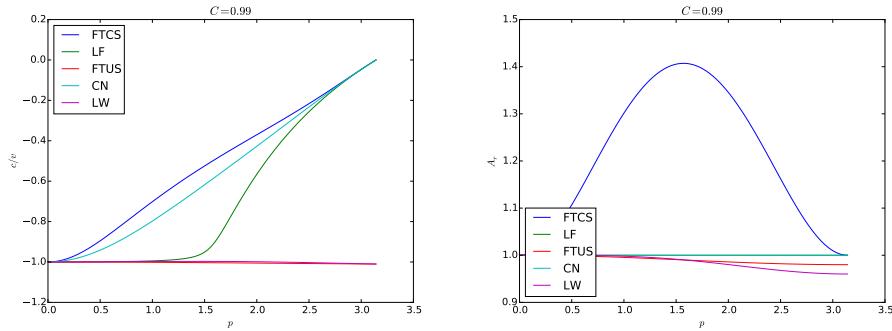


**Fig. 4.4** Dispersion relations for  $C = 0.8$ .

as  $\Delta t \rightarrow 0$ . **hpl 25:** No, how do we explain that reducing  $\Delta t$  reduces the damping, while keeping  $C$  fixed?



**Fig. 4.5** Dispersion relations for  $C = 0.5$ .



**Fig. 4.6** Dispersion relations for  $C = 0.99$ .

## 4.2 One-dimensional stationary advection-diffusion equation

## 4.3 Two-dimensional advection-diffusion equations

## 4.4 Applications of advection equations



## 5.1 Ordinary differential equations

### 5.1.1 The Euler-Cromer scheme on a standard mesh

Consider the fundamental model problem for simple harmonic oscillations,

$$u'' + \omega^2 u = 0, \quad u(0) = I, \quad u'(0) = 0, \quad (5.1)$$

where  $\omega$  is the frequency of the oscillations (the exact solution is  $u(t) = I \cos \omega t$ ). This model can equivalently be formulated as two first-order equations,

$$v' = -\omega^2 u, \quad (5.2)$$

$$u' = v. \quad (5.3)$$

The popular Euler-Cromer scheme for this  $2 \times 2$  system of ODEs applies an explicit forward difference in (5.2) and a backward difference in (5.3):

$$\frac{v^{n+1} - v^n}{\Delta t} = -\omega^2 u^n, \quad (5.4)$$

$$\frac{u^{n+1} - u^n}{\Delta t} = v^{n+1}. \quad (5.5)$$

For a time domain  $[0, T]$ , we have introduced a mesh with points  $0 = t_0 < t_1 < \dots < t_n = T$ . The most common case is a mesh with uniform

spacing  $\Delta t$ :  $t_n = n\Delta t$ . Then  $v^n$  is an approximation to  $v(t)$  at mesh point  $t_n$ , and  $u^n$  is an approximation to  $u(t)$  at the same point. Note that the backward difference in (5.7) leads to an explicit updating formula for  $u^{n+1}$  since  $v^{n+1}$  is already computed:

$$v^{n+1} = v^n - \Delta t \omega^2 u^n, \quad (5.6)$$

$$u^{n+1} = u^n + \Delta t v^{n+1}. \quad (5.7)$$

The Euler-Cromer scheme is equivalent with the standard second-order accurate scheme for (5.1):

$$u^{n+1} = 2u^n - u^{n-1} - \Delta t^2 \omega^2 u^n, \quad n = 1, 2, \dots, \quad (5.8)$$

but for the first time step, the method for (5.1) leads to

$$u^1 = u^0 - \frac{1}{2} \Delta t^2 \omega^2 u^0, \quad (5.9)$$

while Euler-Cromer gives

$$u^1 = u^0 - \Delta t^2 \omega^2 u^0, \quad (5.10)$$

which can be interpreted as a first-order, backward difference approximation of  $u'(0) = 0$  combined with (5.8). At later time steps, however, the alternating use of forward and backward differences in (5.6)-(5.7) leads to a method with error  $\mathcal{O}(\Delta t^2)$ .

### 5.1.2 The Euler-Cromer scheme on a staggered mesh

**hpl 26:** Do the equations in different sequence, first vel, then pos.

The fact that the forward and backward differences used in the Euler-Cromer method yield a second-order accurate method, is not obvious from intuition. A much more intuitive discretization employs solely centered differences and leads to a scheme that is equivalent to the Euler-Cromer scheme. It is in fact fully equivalent to the second-order scheme for (5.1), also for the first time step. This alternative scheme is based on using a *staggered* (or alternating) mesh in time.

In a staggered mesh, the unknowns are sought at different points in the mesh. Specifically,  $u$  is sought at integer time points  $t_n$  and  $v$  is sought at  $t_{n+1/2}$  between two  $u$  points. The unknowns are then  $u^1, v^{3/2}, u^2, v^{5/2}$ , and so on. We typically use the notation  $u^n$  and  $v^{n+\frac{1}{2}}$  for the two unknown

mesh functions. Figure 5.1 presents a graphical sketch of two mesh functions  $u$  and  $v$  on a staggered mesh.

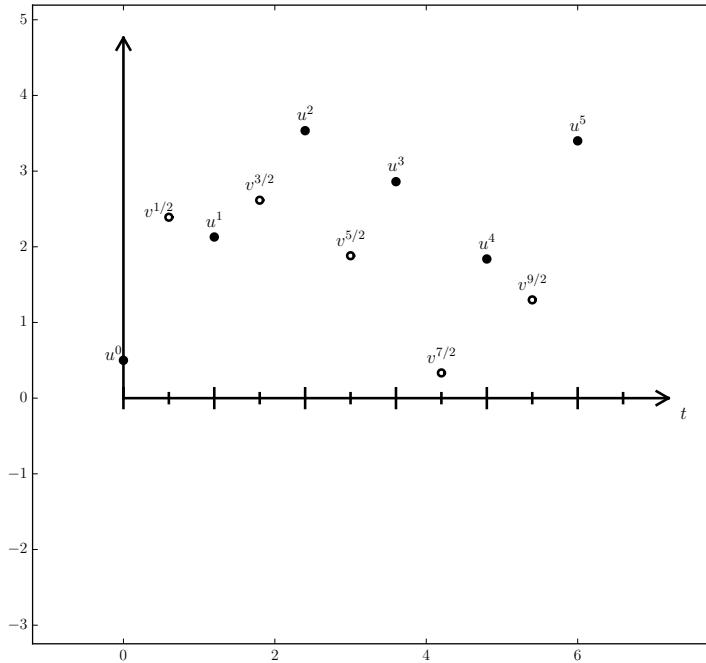


Fig. 5.1 Examples on mesh functions on a staggered mesh in time.

On a staggered mesh it is natural to use centered difference approximations, expressed in operator notation as

$$[D_t u = v]^{n+\frac{1}{2}}, \quad (5.11)$$

$$[D_t v = -\omega u]^{n+1}. \quad (5.12)$$

Writing out the formulas gives

$$u^{n+1} = u^n + \Delta t v^{n+\frac{1}{2}}, \quad (5.13)$$

$$v^{n+\frac{3}{2}} = v^{n+\frac{1}{2}} - \Delta t \omega^2 u^{n+1}. \quad (5.14)$$

Of esthetic reasons one often writes these equations at the previous time level to replace the  $\frac{3}{2}$  by  $\frac{1}{2}$  in the left-most term in the last equation,

$$u^n = u^{n-1} + \Delta t v^{n-\frac{1}{2}}, \quad (5.15)$$

$$v^{n+\frac{1}{2}} = v^{n-\frac{1}{2}} - \Delta t \omega^2 u^n. \quad (5.16)$$

Such a rewrite is only mathematical cosmetics. The important thing is that this centered scheme has exactly the same computational structure as the forward-backward scheme. We shall use the names *forward-backward Euler-Cromer* and *staggered Euler-Cromer* to distinguish the two schemes.

We can eliminate the  $v$  values and get back the centered scheme based on the second-order differential equation, so all these three schemes are equivalent. However, they differ somewhat in the treatment of the initial conditions.

Suppose we have  $u(0) = I$  and  $u'(0) = v(0) = 0$  as mathematical initial conditions. This means  $u^0 = I$  and

$$v(0) \approx \frac{1}{2}(v^{-\frac{1}{2}} + v^{\frac{1}{2}}) = 0, \quad \Rightarrow \quad v^{-\frac{1}{2}} = -v^{\frac{1}{2}}.$$

Using the discretized equation (5.16) for  $n = 0$  yields

$$v^{\frac{1}{2}} = v^{-\frac{1}{2}} - \Delta t \omega^2 I,$$

and eliminating  $v^{-\frac{1}{2}} = -v^{\frac{1}{2}}$  results in  $v^{\frac{1}{2}} = -\frac{1}{2}\Delta t \omega^2 I$  and

$$u^1 = u^0 - \frac{1}{2}\Delta t^2 \omega^2 I,$$

which is exactly the same equation for  $u^1$  as we had in the centered scheme based on the second-order differential equation (and hence corresponds to a centered difference approximation of the initial condition for  $u'(0)$ ). The conclusion is that a staggered mesh is fully equivalent with that scheme, while the forward-backward version gives a slight deviation in the computation of  $u^1$ .

We can redo the derivation of the initial conditions when  $u'(0) = V$ :

$$v(0) \approx \frac{1}{2}(v^{-\frac{1}{2}} + v^{\frac{1}{2}}) = V, \quad \Rightarrow \quad v^{-\frac{1}{2}} = 2V - v^{\frac{1}{2}}.$$

Using this  $v^{-\frac{1}{2}}$  in

$$v^{\frac{1}{2}} = v^{-\frac{1}{2}} - \Delta t \omega^2 I,$$

then gives  $v^{\frac{1}{2}} = V - \frac{1}{2}\Delta t \omega^2 I$ . The general initial conditions are therefore

$$u^0 = I, \quad (5.17)$$

$$v^{\frac{1}{2}} = V - \frac{1}{2} \Delta t \omega^2 I. \quad (5.18)$$

### 5.1.3 Implementation of the scheme on a staggered mesh

The algorithm goes like this:

1. Set the initial values (5.17) and (5.18).
2. For  $n = 1, 2, \dots$ :
  - a. Compute  $u^n$  from (5.15).
  - b. Compute  $v^{n+\frac{1}{2}}$  from (5.16).

**Implementation with integer indices.** Translating the schemes (5.15) and (5.16) to computer code faces the problem of how to store and access  $v^{n+\frac{1}{2}}$ , since arrays only allow integer indices with base 0. We must then introduce a convention:  $v^{1+\frac{1}{2}}$  is stored in  $v[n]$  while  $v^{1-\frac{1}{2}}$  is stored in  $v[n-1]$ . We can then write the algorithm in Python as

```
def solver(I, w, dt, T):
    dt = float(dt)
    Nt = int(round(T/dt))
    u = zeros(Nt+1)
    v = zeros(Nt+1)
    t = linspace(0, Nt*dt, Nt+1) # mesh for u
    t_v = t + dt/2 # mesh for v

    u[0] = I
    v[0] = 0 - 0.5*dt*w**2*u[0]
    for n in range(1, Nt+1):
        u[n] = u[n-1] + dt*v[n-1]
        v[n] = v[n-1] - dt*w**2*u[n]
    return u, t, v, t_v
```

Note that  $u$  and  $v$  are returned together with the mesh points such that the complete mesh function for  $u$  is described by  $u$  and  $t$ , while  $v$  and  $t_v$  represent the mesh function for  $v$ .

**Implementation with half-integer indices.** Some prefer to see a closer relationship between the code and the mathematics for the quantities with half-integer indices. For example, we would like to replace the updating equation for  $v[n]$  by

```
v[n+half] = v[n-half] - dt*w**2*u[n]
```

This is easy to do if we could be sure that `n+half` means `n` and `n-half` means `n-1`. A possible solution is to define `half` as a special object such that an integer plus `half` results in the integer, while an integer minus `half` equals the integer minus 1. A simple Python class may realize the `half` object:

```
class HalfInt:
    def __radd__(self, other):
        return other

    def __rsub__(self, other):
        return other - 1

half = HalfInt()
```

The `__radd__` function is invoked for all expressions `n+half` ("right add" with `self` as `half` and `other` as `n`). Similarly, the `__rsub__` function is invoked for `n-half` and results in `n-1`.

Using the `half` object, we can implement the algorithms in an even more readable way:

```
def solver(I, w, dt, T):
    """
    Solve u'=v, v' = - w**2*u for t in (0,T], u(0)=I and v(0)=0,
    by a central finite difference method with time step dt.
    """
    dt = float(dt)
    Nt = int(round(T/dt))
    u = zeros(Nt+1)
    v = zeros(Nt+1)
    t = linspace(0, Nt*dt, Nt+1) # mesh for u
    t_v = t + dt/2 # mesh for v

    u[0] = I
    v[0+half] = 0 - 0.5*dt*w**2*u[0]
    for n in range(1, Nt+1):
        print n, n+half, n-half
        u[n] = u[n-1] + dt*v[n-half]
        v[n+half] = v[n-half] - dt*w**2*u[n]
    return u, t, v, t_v
```

Verification of this code is easy as we can just compare the computed `u` with the `u` produced by the `solver` function in `vib_undamped.py` (which solves  $u'' + \omega^2 u = 0$  directly). The values should coincide to machine precision since the two numerical methods are mathematically equivalent. We refer to the file `vib_undamped_staggered.py` for the details of a unit test (`test_staggered`) that checks this property.

### 5.1.4 A staggered Euler-Cromer scheme for a generalized model

The more general model for vibration problems,

$$mu'' + f(u') + s(u) = F(t), \quad u(0) = I, \quad u'(0) = V, \quad t \in (0, T], \quad (5.19)$$

can be rewritten as a first-order ODE system

$$u' = v, \quad (5.20)$$

$$v' = m^{-1} (F(t) - f(v) - s(u)). \quad (5.21)$$

It is natural to introduce a staggered mesh (see Section 5.1.2) and seek  $u$  at mesh points  $t_n$  (the numerical value is denoted by  $u^n$ ) and  $v$  between mesh points at  $t_{n+1/2}$  (the numerical value is denoted by  $v^{n+\frac{1}{2}}$ ). A centered difference approximation to (5.20)-(5.21) can then be written in operator notation as

$$[D_t u = v]^{n-\frac{1}{2}}, \quad (5.22)$$

$$[D_t v = m^{-1} (F(t) - f(v) - s(u))]^n. \quad (5.23)$$

Written out,

$$\frac{u^n - u^{n-1}}{\Delta t} = v^{n-\frac{1}{2}}, \quad (5.24)$$

$$\frac{v^{n+\frac{1}{2}} - v^{n-\frac{1}{2}}}{\Delta t} = m^{-1} (F^n - f(v^n) - s(u^n)). \quad (5.25)$$

With linear damping,  $f(v) = bv$ , we can use an arithmetic mean for  $f(v^n)$ :  $f(v^n) \approx \frac{1}{2}(f(v^{n-\frac{1}{2}}) + f(v^{n+\frac{1}{2}}))$ . The system (5.24)-(5.25) can then be solved with respect to the unknowns  $u^n$  and  $v^{n+\frac{1}{2}}$ :

$$u^n = u^{n-1} + \Delta t v^{n-\frac{1}{2}}, \quad (5.26)$$

$$v^{n+\frac{1}{2}} = \left(1 + \frac{b}{2m} \Delta t\right)^{-1} \left(v^{n-\frac{1}{2}} + \Delta t m^{-1} \left(F^n - \frac{1}{2}f(v^{n-\frac{1}{2}}) - s(u^n)\right)\right). \quad (5.27)$$

In case of quadratic damping,  $f(v) = b|v|v$ , we can use a geometric mean:  $f(v^n) \approx b|v^{n-\frac{1}{2}}|v^{n+\frac{1}{2}}$ . Inserting this approximation in (5.24)-(5.25) and solving for the unknowns  $u^n$  and  $v^{n+\frac{1}{2}}$  results in

$$u^n = u^{n-1} + \Delta t v^{n-\frac{1}{2}}, \quad (5.28)$$

$$v^{n+\frac{1}{2}} = \left(1 + \frac{b}{m}|v^{n-\frac{1}{2}}|\Delta t\right)^{-1} \left(v^{n-\frac{1}{2}} + \Delta t m^{-1} (F^n - s(u^n))\right). \quad (5.29)$$

The initial conditions are derived at the end of Section 5.1.2:

$$u^0 = I, \quad (5.30)$$

$$v^{\frac{1}{2}} = V - \frac{1}{2}\Delta t \omega^2 I. \quad (5.31)$$

## 5.2 Exercises

### **Exercise 5.1: Use the forward-backward scheme with quadratic damping**

We consider the generalized model with quadratic damping, expressed as a system of two first-order equations as in Section 5.1.4:

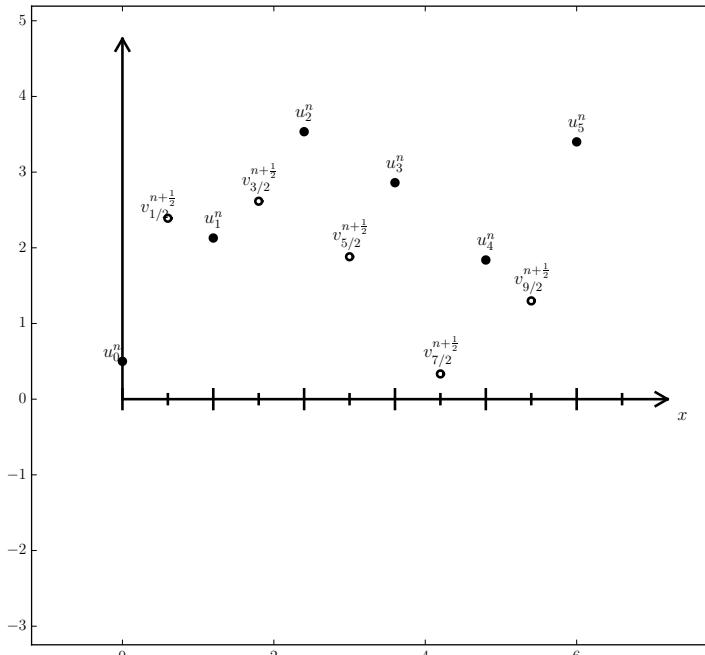
$$\begin{aligned} u' &= v, \\ v' &= \frac{1}{m} (F(t) - \beta|v|v - s(u)). \end{aligned}$$

However, contrary to what is done in Section 5.1.4, we want to apply the idea of a forward-backward discretization:  $u$  is marched forward by a one-sided Forward Euler scheme applied to the first equation, and thereafter  $v$  can be marched forward by a Backward Euler scheme in the second equation, see in Section 1.7. Express the idea in operator notation and write out the scheme. Unfortunately, the backward difference for the  $v$  equation creates a nonlinearity  $|v^{n+1}|v^{n+1}$ . To linearize this nonlinearity, use the known value  $v^n$  inside the absolute value factor, i.e.,  $|v^{n+1}|v^{n+1} \approx |v^n|v^{n+1}$ . Show that the resulting scheme is equivalent to the one in Section 5.1.4 for some time level  $n \geq 1$ .

What we learn from this exercise is that the first-order differences and the linearization trick play together in “the right way” such that

the scheme is as good as when we (in Section 5.1.4) carefully apply centered differences and a geometric mean on a staggered mesh to achieve second-order accuracy. There is a difference in the handling of the initial conditions, though, as explained at the end of Section 1.7. Filename: `vib_gen_bwdamping`.

## 5.3 Partial differential equations



**Fig. 5.2** Examples on mesh functions on a staggered mesh in space.



## 6.1 Introduction of basic concepts

### 6.1.1 Linear versus nonlinear equations

**Algebraic equations.** A linear, scalar, algebraic equation in  $x$  has the form

$$ax + b = 0,$$

for arbitrary real constants  $a$  and  $b$ . The unknown is a number  $x$ . All other algebraic equations, e.g.,  $x^2 + ax + b = 0$ , are nonlinear. The typical feature in a nonlinear algebraic equation is that the unknown appears in products with itself, like  $x^2$  or  $e^x = 1 + x + \frac{1}{2}x^2 + \frac{1}{3!}x^3 + \dots$ .

We know how to solve a linear algebraic equation,  $x = -b/a$ , but there are no general methods for finding the exact solutions of nonlinear algebraic equations, except for very special cases (quadratic equations are a primary example). A nonlinear algebraic equation may have no solution, one solution, or many solutions. The tools for solving nonlinear algebraic equations are *iterative methods*, where we construct a series of linear equations, which we know how to solve, and hope that the solutions of the linear equations converge to a solution of the nonlinear equation we want to solve. Typical methods for nonlinear algebraic equation equations are Newton's method, the Bisection method, and the Secant method.

**Differential equations.** The unknown in a differential equation is a function and not a number. In a linear differential equation, all terms

involving the unknown function are linear in the unknown function or its derivatives. Linear here means that the unknown function, or a derivative of it, is multiplied by a number or a known function. All other differential equations are non-linear.

The easiest way to see if an equation is nonlinear, is to spot nonlinear terms where the unknown function or its derivatives are multiplied by each other. For example, in

$$u'(t) = -a(t)u(t) + b(t),$$

the terms involving the unknown function  $u$  are linear:  $u'$  contains the derivative of the unknown function multiplied by unity, and  $au$  contains the unknown function multiplied by a known function. However,

$$u'(t) = u(t)(1 - u(t)),$$

is nonlinear because of the term  $-u^2$  where the unknown function is multiplied by itself. Also

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = 0,$$

is nonlinear because of the term  $uu_x$  where the unknown function appears in a product with itself or one of its derivatives. (Note here that we use different notations for derivatives:  $u'$  or  $du/dt$  for a function  $u(t)$  of one variable,  $\frac{\partial u}{\partial t}$  or  $u_t$  for a function of more than one variable.)

Another example of a nonlinear equation is

$$u'' + \sin(u) = 0,$$

because  $\sin(u)$  contains products of  $u$  if we expand the function in a Taylor series:

$$\sin(u) = u - \frac{1}{3}u^3 + \dots$$

### Mathematical proof of linearity

To really prove mathematically that some differential equation in an unknown  $u$  is linear, show for each term  $T(u)$  that with  $u = au_1 + bu_2$  for constants  $a$  and  $b$ ,

$$T(au_1 + bu_2) = aT(u_1) + bT(u_2).$$

For example, the term  $T(u) = (\sin^2 t)u'(t)$  is linear because

$$\begin{aligned} T(au_1 + bu_2) &= (\sin^2 t)(au_1(t) + bu_2(t)) \\ &= a(\sin^2 t)u_1(t) + b(\sin^2 t)u_2(t) \\ &= aT(u_1) + bT(u_2). \end{aligned}$$

However,  $T(u) = \sin u$  is nonlinear because

$$T(au_1 + bu_2) = \sin(au_1 + bu_2) \neq a \sin u_1 + b \sin u_2.$$

### 6.1.2 A simple model problem

A series of forthcoming examples will explain how to tackle nonlinear differential equations with various techniques. We start with the (scaled) logistic equation as model problem:

$$u'(t) = u(t)(1 - u(t)). \quad (6.1)$$

This is a nonlinear ordinary differential equation (ODE) which will be solved by different strategies in the following. Depending on the chosen time discretization of (6.1), the mathematical problem to be solved at every time level will either be a linear algebraic equation or a nonlinear algebraic equation. In the former case, the time discretization method transforms the nonlinear ODE into linear subproblems at each time level, and the solution is straightforward to find since linear algebraic equations are easy to solve. However, when the time discretization leads to nonlinear algebraic equations, we cannot (except in very rare cases) solve these without turning to approximate, iterative solution methods.

The next subsections introduce various methods for solving nonlinear differential equations, using (6.1) as model. We shall go through the following set cases:

- explicit time discretization methods (with no need to solve nonlinear algebraic equations)
- implicit Backward Euler discretization, leading to nonlinear algebraic equations solved by

- an exact analytical technique
- Picard iteration based on manual linearization
- a single Picard step
- Newton’s method
- Implicit Crank-Nicolson discretization and linearization via a geometric mean formula

Thereafter, we compare the performance of the various approaches. Despite the simplicity of (6.1), the conclusions reveal typical features of the various methods in much more complicated nonlinear PDE problems.

### 6.1.3 Linearization by explicit time discretization

Time discretization methods are divided into explicit and implicit methods. Explicit methods lead to a closed-form formula for finding new values of the unknowns, while implicit methods give a linear or nonlinear system of equations that couples (all) the unknowns at a new time level. Here we shall demonstrate that explicit methods constitute an efficient way to deal with nonlinear differential equations.

The Forward Euler method is an explicit method. When applied to (6.1), sampled at  $t = t_n$ , it results in

$$\frac{u^{n+1} - u^n}{\Delta t} = u^n(1 - u^n),$$

which is a *linear* algebraic equation for the unknown value  $u^{n+1}$  that we can easily solve:

$$u^{n+1} = u^n + \Delta t u^n(1 - u^n).$$

The nonlinearity in the original equation poses in this case no difficulty in the discrete algebraic equation. Any other explicit scheme in time will also give only linear algebraic equations to solve. For example, a typical 2nd-order Runge-Kutta method for (6.1) leads to the following formulas:

$$\begin{aligned} u^* &= u^n + \Delta t u^n(1 - u^n), \\ u^{n+1} &= u^n + \Delta t \frac{1}{2} (u^n(1 - u^n) + u^*(1 - u^*)) . \end{aligned}$$

The first step is linear in the unknown  $u^*$ . Then  $u^*$  is known in the next step, which is linear in the unknown  $u^{n+1}$ .

### 6.1.4 Exact solution of nonlinear algebraic equations

Switching to a Backward Euler scheme for (6.1),

$$\frac{u^n - u^{n-1}}{\Delta t} = u^n(1 - u^n), \quad (6.2)$$

results in a nonlinear algebraic equation for the unknown value  $u^n$ . The equation is of quadratic type:

$$\Delta t(u^n)^2 + (1 - \Delta t)u^n - u^{n-1} = 0,$$

and may be solved exactly by the well-known formula for such equations. Before we do so, however, we will introduce a shorter, and often cleaner, notation for nonlinear algebraic equations at a given time level. The notation is inspired by the natural notation (i.e., variable names) used in a program, especially in more advanced partial differential equation problems. The unknown in the algebraic equation is denoted by  $u$ , while  $u^{(1)}$  is the value of the unknown at the previous time level (in general,  $u^{(\ell)}$  is the value of the unknown  $\ell$  levels back in time). The notation will be frequently used in later sections. What is meant by  $u$  should be evident from the context:  $u$  may be 1) the exact solution of the ODE/PDE problem, 2) the numerical approximation to the exact solution, or 3) the unknown solution at a certain time level.

The quadratic equation for the unknown  $u^n$  in (6.2) can, with the new notation, be written

$$F(u) = \Delta t u^2 + (1 - \Delta t)u - u^{(1)} = 0. \quad (6.3)$$

The solution is readily found to be

$$u = \frac{1}{2\Delta t} \left( -1 + \Delta t \pm \sqrt{(1 - \Delta t)^2 - 4\Delta t u^{(1)}} \right). \quad (6.4)$$

Now we encounter a fundamental challenge with nonlinear algebraic equations: the equation may have more than one solution. How do we pick the right solution? This is in general a hard problem. In the present simple case, however, we can analyze the roots mathematically and provide an answer. The idea is to expand the roots in a series in  $\Delta t$  and truncate after the linear term since the Backward Euler scheme will introduce an error proportional to  $\Delta t$  anyway. Using `sympy` we find the following Taylor series expansions of the roots:

```
>>> import sympy as sym
```

```
>>> dt, u_1, u = sym.symbols('dt u_1 u')
>>> r1, r2 = sym.solve(dt*u**2 + (1-dt)*u - u_1, u) # find roots
>>> r1
(dt - sqrt(dt**2 + 4*dt*u_1 - 2*dt + 1) - 1)/(2*dt)
>>> r2
(dt + sqrt(dt**2 + 4*dt*u_1 - 2*dt + 1) - 1)/(2*dt)
>>> print r1.series(dt, 0, 2) # 2 terms in dt, around dt=0
-1/dt + 1 - u_1 + dt*(u_1**2 - u_1) + O(dt**2)
>>> print r2.series(dt, 0, 2)
u_1 + dt*(-u_1**2 + u_1) + O(dt**2)
```

We see that the `r1` root, corresponding to a minus sign in front of the square root in (6.4), behaves as  $1/\Delta t$  and will therefore blow up as  $\Delta t \rightarrow 0$ ! Since we know that  $u$  takes on finite values, actually it is less than or equal to 1, only the `r2` root is of relevance in this case: as  $\Delta t \rightarrow 0$ ,  $u \rightarrow u^{(1)}$ , which is the expected result.

For those who are not well experienced with approximating mathematical formulas by series expansion, an alternative method of investigation is simply to compute the limits of the two roots as  $\Delta t \rightarrow 0$  and see if a limit unreasonable:

```
>>> print r1.limit(dt, 0)
-oo
>>> print r2.limit(dt, 0)
u_1
```

### 6.1.5 Linearization

When the time integration of an ODE results in a nonlinear algebraic equation, we must normally find its solution by defining a sequence of linear equations and hope that the solutions of these linear equations converge to the desired solution of the nonlinear algebraic equation. Usually, this means solving the linear equation repeatedly in an iterative fashion. Alternatively, the nonlinear equation can sometimes be approximated by one linear equation, and consequently there is no need for iteration.

Constructing a linear equation from a nonlinear one requires *linearization* of each nonlinear term. This can be done manually as in Picard iteration, or fully algorithmically as in Newton's method. Examples will best illustrate how to linearize nonlinear problems.

### 6.1.6 Picard iteration

Let us write (6.3) in a more compact form

$$F(u) = au^2 + bu + c = 0,$$

with  $a = \Delta t$ ,  $b = 1 - \Delta t$ , and  $c = -u^{(1)}$ . Let  $u^-$  be an available approximation of the unknown  $u$ . Then we can linearize the term  $u^2$  simply by writing  $u^- u$ . The resulting equation,  $\hat{F}(u) = 0$ , is now linear and hence easy to solve:

$$F(u) \approx \hat{F}(u) = au^- u + bu + c = 0.$$

Since the equation  $\hat{F} = 0$  is only approximate, the solution  $u$  does not equal the exact solution  $u_e$  of the exact equation  $F(u_e) = 0$ , but we can hope that  $u$  is closer to  $u_e$  than  $u^-$  is, and hence it makes sense to repeat the procedure, i.e., set  $u^- = u$  and solve  $\hat{F}(u) = 0$  again. There is no guarantee that  $u$  is closer to  $u_e$  than  $u^-$ , but this approach has proven to be effective in a wide range of applications.

The idea of turning a nonlinear equation into a linear one by using an approximation  $u^-$  of  $u$  in nonlinear terms is a widely used approach that goes under many names: *fixed-point iteration*, the method of *successive substitutions*, *nonlinear Richardson iteration*, and *Picard iteration*. We will stick to the latter name.

Picard iteration for solving the nonlinear equation arising from the Backward Euler discretization of the logistic equation can be written as

$$u = -\frac{c}{au^- + b}, \quad u^- \leftarrow u.$$

The  $\leftarrow$  symbol means assignment (we set  $u^-$  equal to the value of  $u$ ). The iteration is started with the value of the unknown at the previous time level:  $u^- = u^{(1)}$ .

Some prefer an explicit iteration counter as superscript in the mathematical notation. Let  $u^k$  be the computed approximation to the solution in iteration  $k$ . In iteration  $k + 1$  we want to solve

$$au^k u^{k+1} + bu^{k+1} + c = 0 \quad \Rightarrow \quad u^{k+1} = -\frac{c}{au^k + b}, \quad k = 0, 1, \dots$$

Since we need to perform the iteration at every time level, the time level counter is often also included:

$$au^{n,k} u^{n,k+1} + bu^{n,k+1} - u^{n-1} = 0 \quad \Rightarrow \quad u^{n,k+1} = \frac{u^n}{au^{n,k} + b}, \quad k = 0, 1, \dots,$$

with the start value  $u^{n,0} = u^{n-1}$  and the final converged value  $u^n = u^{n,k}$  for sufficiently large  $k$ .

However, we will normally apply a mathematical notation in our final formulas that is as close as possible to what we aim to write in a computer code and then it becomes natural to use  $u$  and  $u^-$  instead of  $u^{k+1}$  and  $u^k$  or  $u^{n,k+1}$  and  $u^{n,k}$ .

**Stopping criteria.** The iteration method can typically be terminated when the change in the solution is smaller than a tolerance  $\epsilon_u$ :

$$|u - u^-| \leq \epsilon_u,$$

or when the residual in the equation is sufficiently small ( $< \epsilon_r$ ),

$$|F(u)| = |au^2 + bu + c| < \epsilon_r .$$

**A single Picard iteration.** Instead of iterating until a stopping criterion is fulfilled, one may iterate a specific number of times. Just one Picard iteration is popular as this corresponds to the intuitive idea of approximating a nonlinear term like  $(u^n)^2$  by  $u^{n-1}u^n$ . This follows from the linearization  $u^-u^n$  and the initial choice of  $u^- = u^{n-1}$  at time level  $t_n$ . In other words, a single Picard iteration corresponds to using the solution at the previous time level to linearize nonlinear terms. The resulting discretization becomes (using proper values for  $a$ ,  $b$ , and  $c$ )

$$\frac{u^n - u^{n-1}}{\Delta t} = u^n(1 - u^{n-1}), \quad (6.5)$$

which is a linear algebraic equation in the unknown  $u^n$ , and therefore we can easily solve for  $u^n$ , and there is no need for any alternative notation.

We shall later refer to the strategy of taking one Picard step, or equivalently, linearizing terms with use of the solution at the previous time step, as the *Picard1* method. It is a widely used approach in science and technology, but with some limitations if  $\Delta t$  is not sufficiently small (as will be illustrated later).

### Notice

Equation (6.5) does not correspond to a “pure” finite difference method where the equation is sampled at a point and derivatives replaced by differences (because the  $u^{n-1}$  term on the right-hand side must then be  $u^n$ ). The best interpretation of the scheme (6.5)

is a Backward Euler difference combined with a single (perhaps insufficient) Picard iteration at each time level, with the value at the previous time level as start for the Picard iteration.

### 6.1.7 Linearization by a geometric mean

We consider now a Crank-Nicolson discretization of (6.1). This means that the time derivative is approximated by a centered difference,

$$[D_t u = u(1 - u)]^{n+\frac{1}{2}},$$

written out as

$$\frac{u^{n+1} - u^n}{\Delta t} = u^{n+\frac{1}{2}} - (u^{n+\frac{1}{2}})^2. \quad (6.6)$$

The term  $u^{n+\frac{1}{2}}$  is normally approximated by an arithmetic mean,

$$u^{n+\frac{1}{2}} \approx \frac{1}{2}(u^n + u^{n+1}),$$

such that the scheme involves the unknown function only at the time levels where we actually compute it. The same arithmetic mean applied to the nonlinear term gives

$$(u^{n+\frac{1}{2}})^2 \approx \frac{1}{4}(u^n + u^{n+1})^2,$$

which is nonlinear in the unknown  $u^{n+1}$ . However, using a *geometric mean* for  $(u^{n+\frac{1}{2}})^2$  is a way of linearizing the nonlinear term in (6.6):

$$(u^{n+\frac{1}{2}})^2 \approx u^n u^{n+1}.$$

Using an arithmetic mean on the linear  $u^{n+\frac{1}{2}}$  term in (6.6) and a geometric mean for the second term, results in a linearized equation for the unknown  $u^{n+1}$ :

$$\frac{u^{n+1} - u^n}{\Delta t} = \frac{1}{2}(u^n + u^{n+1}) + u^n u^{n+1},$$

which can readily be solved:

$$u^{n+1} = \frac{1 + \frac{1}{2}\Delta t}{1 + \Delta t u^n - \frac{1}{2}\Delta t} u^n.$$

This scheme can be coded directly, and since there is no nonlinear algebraic equation to iterate over, we skip the simplified notation with  $u$  for  $u^{n+1}$  and  $u^{(1)}$  for  $u^n$ . The technique with using a geometric average is an example of transforming a nonlinear algebraic equation to a linear one, without any need for iterations.

The geometric mean approximation is often very effective for linearizing quadratic nonlinearities. Both the arithmetic and geometric mean approximations have truncation errors of order  $\Delta t^2$  and are therefore compatible with the truncation error  $\mathcal{O}(\Delta t^2)$  of the centered difference approximation for  $u'$  in the Crank-Nicolson method.

Applying the operator notation for the means and finite differences, the linearized Crank-Nicolson scheme for the logistic equation can be compactly expressed as

$$[D_t u = \bar{u}^t + \bar{u}^{2,t,g}]^{n+\frac{1}{2}}.$$

### Remark

If we use an arithmetic instead of a geometric mean for the nonlinear term in (6.6), we end up with a nonlinear term  $(u^{n+1})^2$ . This term can be linearized as  $u^- u^{n+1}$  in a Picard iteration approach and in particular as  $u^n u^{n+1}$  in a Picard1 iteration approach. The latter gives a scheme almost identical to the one arising from a geometric mean (the difference in  $u^{n+1}$  being  $\frac{1}{4}\Delta t u^n (u^{n+1} - u^n) \approx \frac{1}{4}\Delta t^2 u' u$ , i.e., a difference of  $\mathcal{O}(\Delta t^2)$ ).

### 6.1.8 Newton's method

The Backward Euler scheme (6.2) for the logistic equation leads to a nonlinear algebraic equation (6.3). Now we write any nonlinear algebraic equation in the general and compact form

$$F(u) = 0.$$

Newton's method linearizes this equation by approximating  $F(u)$  by its Taylor series expansion around a computed value  $u^-$  and keeping only the linear part:

$$\begin{aligned} F(u) &= F(u^-) + F'(u^-)(u - u^-) + \frac{1}{2}F''(u^-)(u - u^-)^2 + \dots \\ &\approx F(u^-) + F'(u^-)(u - u^-) = \hat{F}(u). \end{aligned}$$

The linear equation  $\hat{F}(u) = 0$  has the solution

$$u = u^- - \frac{F(u^-)}{F'(u^-)}.$$

Expressed with an iteration index in the unknown, Newton's method takes on the more familiar mathematical form

$$u^{k+1} = u^k - \frac{F(u^k)}{F'(u^k)}, \quad k = 0, 1, \dots$$

It can be shown that the error in iteration  $k + 1$  of Newton's method is proportional to the square of the error in iteration  $k$ , a result referred to as *quadratic convergence*. This means that for small errors the method converges very fast, and in particular much faster than Picard iteration and other iteration methods. (The proof of this result is found in most textbooks on numerical analysis.) However, the quadratic convergence appears only if  $u^k$  is sufficiently close to the solution. Further away from the solution the method can easily converge very slowly or diverge. The reader is encouraged to do Exercise 6.3 to get a better understanding for the behavior of the method.

Application of Newton's method to the logistic equation discretized by the Backward Euler method is straightforward as we have

$$F(u) = au^2 + bu + c, \quad a = \Delta t, \quad b = 1 - \Delta t, \quad c = -u^{(1)},$$

and then

$$F'(u) = 2au + b.$$

The iteration method becomes

$$u = u^- + \frac{a(u^-)^2 + bu^- + c}{2au^- + b}, \quad u^- \leftarrow u. \quad (6.7)$$

At each time level, we start the iteration by setting  $u^- = u^{(1)}$ . Stopping criteria as listed for the Picard iteration can be used also for Newton's method.

An alternative mathematical form, where we write out  $a$ ,  $b$ , and  $c$ , and use a time level counter  $n$  and an iteration counter  $k$ , takes the form

$$u^{n,k+1} = u^{n,k} + \frac{\Delta t(u^{n,k})^2 + (1 - \Delta t)u^{n,k} - u^{n-1}}{2\Delta t u^{n,k} + 1 - \Delta t}, \quad u^{n,0} = u^{n-1}, \quad k = 0, 1, \dots \quad (6.8)$$

A program implementation is much closer to (6.7) than to (6.8), but the latter is better aligned with the established mathematical notation used in the literature.

### 6.1.9 Relaxation

One iteration in Newton's method or Picard iteration consists of solving a linear problem  $\hat{F}(u) = 0$ . Sometimes convergence problems arise because the new solution  $u$  of  $\hat{F}(u) = 0$  is “too far away” from the previously computed solution  $u^-$ . A remedy is to introduce a relaxation, meaning that we first solve  $\hat{F}(u^*) = 0$  for a suggested value  $u^*$  and then we take  $u$  as a weighted mean of what we had,  $u^-$ , and what our linearized equation  $\hat{F} = 0$  suggests,  $u^*$ :

$$u = \omega u^* + (1 - \omega)u^-.$$

The parameter  $\omega$  is known as a *relaxation parameter*, and a choice  $\omega < 1$  may prevent divergent iterations.

Relaxation in Newton's method can be directly incorporated in the basic iteration formula:

$$u = u^- - \omega \frac{F(u^-)}{F'(u^-)}. \quad (6.9)$$

### 6.1.10 Implementation and experiments

The program `logistic.py` contains implementations of all the methods described above. Below is an extract of the file showing how the Picard and Newton methods are implemented for a Backward Euler discretization of the logistic equation.

```

def BE_logistic(u0, dt, Nt, choice='Picard',
                eps_r=1E-3, omega=1, max_iter=1000):
    if choice == 'Picard1':
        choice = 'Picard'
        max_iter = 1

    u = np.zeros(Nt+1)
    iterations = []
    u[0] = u0
    for n in range(1, Nt+1):
        a = dt
        b = 1 - dt
        c = -u[n-1]

        if choice == 'Picard':

            def F(u):
                return a*u**2 + b*u + c

            u_ = u[n-1]
            k = 0
            while abs(F(u_)) > eps_r and k < max_iter:
                u_ = omega*(-c/(a*u_ + b)) + (1-omega)*u_
                k += 1
            u[n] = u_
            iterations.append(k)

        elif choice == 'Newton':

            def F(u):
                return a*u**2 + b*u + c

            def dF(u):
                return 2*a*u + b

            u_ = u[n-1]
            k = 0
            while abs(F(u_)) > eps_r and k < max_iter:
                u_ = u_ - F(u_)/dF(u_)
                k += 1
            u[n] = u_
            iterations.append(k)

    return u, iterations

```

The Crank-Nicolson method utilizing a linearization based on the geometric mean gives a simpler algorithm:

```

def CN_logistic(u0, dt, Nt):
    u = np.zeros(Nt+1)
    u[0] = u0
    for n in range(0, Nt):
        u[n+1] = (1 + 0.5*dt)/(1 + dt*u[n] - 0.5*dt)*u[n]
    return u

```

We may run experiments with the model problem (6.1) and the different strategies for dealing with nonlinearities as described above. For a quite coarse time resolution,  $\Delta t = 0.9$ , use of a tolerance  $\epsilon_r = 0.1$  in the stopping criterion introduces an iteration error, especially in the Picard iterations, that is visibly much larger than the time discretization error due to a large  $\Delta t$ . This is illustrated by comparing the upper two plots in Figure 6.1. The one to the right has a stricter tolerance  $\epsilon = 10^{-3}$ , which leads to all the curves corresponding to Picard and Newton iteration to be on top of each other (and no changes can be visually observed by reducing  $\epsilon_r$  further). The reason why Newton's method does much better than Picard iteration in the upper left plot is that Newton's method with one step comes far below the  $\epsilon_r$  tolerance, while the Picard iteration needs on average 7 iterations to bring the residual down to  $\epsilon_r = 10^{-1}$ , which gives insufficient accuracy in the solution of the nonlinear equation. It is obvious that the Picard1 method gives significant errors in addition to the time discretization unless the time step is as small as in the lower right plot.

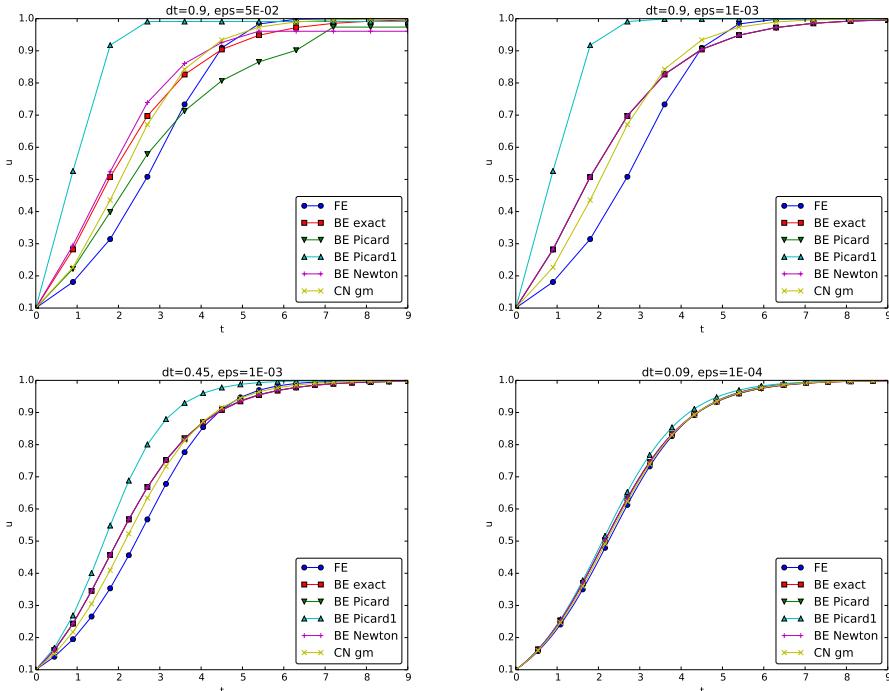
The *BE exact* curve corresponds to using the exact solution of the quadratic equation at each time level, so this curve is only affected by the Backward Euler time discretization. The *CN gm* curve corresponds to the theoretically more accurate Crank-Nicolson discretization, combined with a geometric mean for linearization. This curve appears more accurate, especially if we take the plot in the lower right with a small  $\Delta t$  and an appropriately small  $\epsilon_r$  value as the exact curve.

When it comes to the need for iterations, Figure 6.2 displays the number of iterations required at each time level for Newton's method and Picard iteration. The smaller  $\Delta t$  is, the better starting value we have for the iteration, and the faster the convergence is. With  $\Delta t = 0.9$  Picard iteration requires on average 32 iterations per time step, but this number is dramatically reduced as  $\Delta t$  is reduced.

However, introducing relaxation and a parameter  $\omega = 0.8$  immediately reduces the average of 32 to 7, indicating that for the large  $\Delta t = 0.9$ , Picard iteration takes too long steps. An approximately optimal value for  $\omega$  in this case is 0.5, which results in an average of only 2 iterations! An even more dramatic impact of  $\omega$  appears when  $\Delta t = 1$ : Picard iteration does not converge in 1000 iterations, but  $\omega = 0.5$  again brings the average number of iterations down to 2.

**hpl 27:** Is this remark really relevant now? Compare with text.

**Remark.** The simple Crank-Nicolson method with a geometric mean for the quadratic nonlinearity gives visually more accurate solutions than



**Fig. 6.1** Impact of solution strategy and time step length on the solution.

the Backward Euler discretization. Even with a tolerance of  $\epsilon_r = 10^{-3}$ , all the methods for treating the nonlinearities in the Backward Euler discretization give graphs that cannot be distinguished. So for accuracy in this problem, the time discretization is much more crucial than  $\epsilon_r$ . Ideally, one should estimate the error in the time discretization, as the solution progresses, and set  $\epsilon_r$  accordingly.

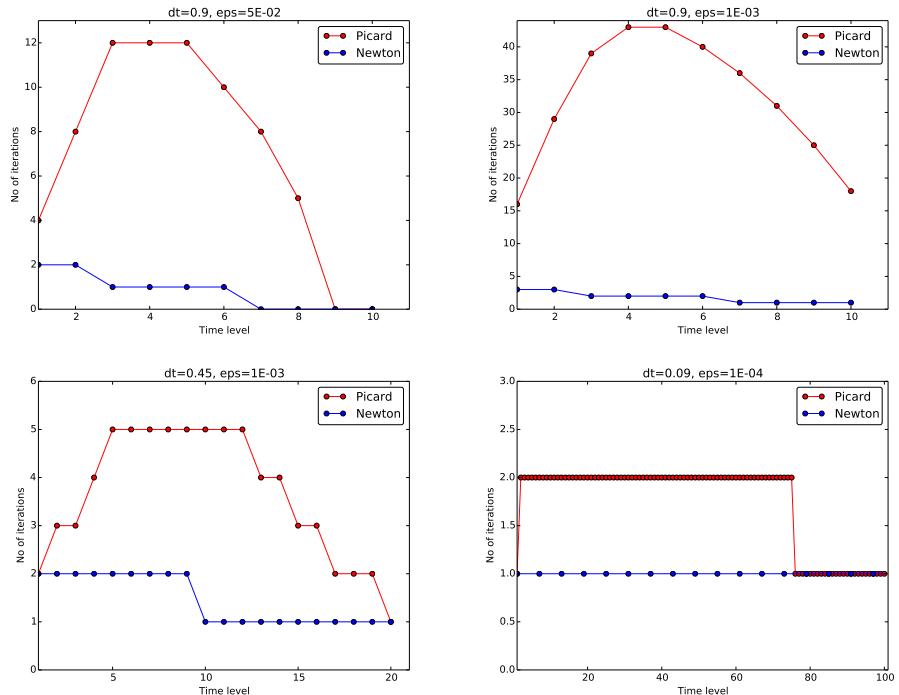
### 6.1.11 Generalization to a general nonlinear ODE

Let us see how the various methods in the previous sections can be applied to the more generic model

$$u' = f(u, t), \quad (6.10)$$

where  $f$  is a nonlinear function of  $u$ .

**Explicit time discretization.** Explicit ODE methods like the Forward Euler scheme, Runge-Kutta methods, Adams-Bashforth methods all



**Fig. 6.2** Comparison of the number of iterations at various time levels for Picard and Newton iteration.

evaluate  $f$  at time levels where  $u$  is already computed, so nonlinearities in  $f$  do not pose any difficulties.

**Backward Euler discretization.** Approximating  $u'$  by a backward difference leads to a Backward Euler scheme, which can be written as

$$F(u^n) = u^n - \Delta t f(u^n, t_n) - u^{n-1} = 0,$$

or alternatively

$$F(u) = u - \Delta t f(u, t_n) - u^{(1)} = 0.$$

A simple Picard iteration, not knowing anything about the nonlinear structure of  $f$ , must approximate  $f(u, t_n)$  by  $f(u^-, t_n)$ :

$$\hat{F}(u) = u - \Delta t f(u^-, t_n) - u^{(1)}.$$

The iteration starts with  $u^- = u^{(1)}$  and proceeds with repeating

$$u^* = \Delta t f(u^-, t_n) + u^{(1)}, \quad u = \omega u^* + (1 - \omega)u^-, \quad u^- \leftarrow u,$$

until a stopping criterion is fulfilled.

### Explicit vs implicit treatment of nonlinear terms

Evaluating  $f$  for a known  $u^-$  is referred to as *explicit* treatment of  $f$ , while if  $f(u, t)$  has some structure, say  $f(u, t) = u^3$ , parts of  $f$  can involve the known  $u$ , as in the manual linearization like  $(u^-)^2u$ , and then the treatment of  $f$  is “more implicit” and “less explicit”. This terminology is inspired by time discretization of  $u' = f(u, t)$ , where evaluating  $f$  for known  $u$  values gives explicit schemes, while treating  $f$  or parts of  $f$  implicitly, makes  $f$  contribute to the unknown terms in the equation at the new time level.

Explicit treatment of  $f$  usually means stricter conditions on  $\Delta t$  to achieve stability of time discretization schemes. The same applies to iteration techniques for nonlinear algebraic equations: the “less” we linearize  $f$  (i.e., the more we keep of  $u$  in the original formula), the faster the convergence may be.

We may say that  $f(u, t) = u^3$  is treated explicitly if we evaluate  $f$  as  $(u^-)^3$ , partially implicit if we linearize as  $(u^-)^2u$  and fully implicit if we represent  $f$  by  $u^3$ . (Of course, the fully implicit representation will require further linearization, but with  $f(u, t) = u^2$  a fully implicit treatment is possible if the resulting quadratic equation is solved with a formula.)

For the ODE  $u' = -u^3$  with  $f(u, t) = -u^3$  and coarse time resolution  $\Delta t = 0.4$ , Picard iteration with  $(u^-)^2u$  requires 8 iterations with  $\epsilon_r = 10^{-3}$  for the first time step, while  $(u^-)^3$  leads to 22 iterations. After about 10 time steps both approaches are down to about 2 iterations per time step, but this example shows a potential of treating  $f$  more implicitly.

A trick to treat  $f$  implicitly in Picard iteration is to evaluate it as  $f(u^-, t)u/u^-$ . For a polynomial  $f$ ,  $f(u, t) = u^m$ , this corresponds to  $(u^-)^m u/u^- = (u^-)^{m-1}u$ . Sometimes this more implicit treatment has no effect, as with  $f(u, t) = \exp(-u)$  and  $f(u, t) = \ln(1 + u)$ , but with  $f(u, t) = \sin(2(u + 1))$ , the  $f(u^-, t)u/u^-$  trick leads to 7, 9, and 11 iterations during the first three steps, while  $f(u^-, t)$  demands 17, 21, and 20 iterations. (Experiments can be done with the code [ODE\\_Picard\\_tricks.py](#).)

Newton's method applied to a Backward Euler discretization of  $u' = f(u, t)$  requires the computation of the derivative

$$F'(u) = 1 - \Delta t \frac{\partial f}{\partial u}(u, t_n).$$

Starting with the solution at the previous time level,  $u^- = u^{(1)}$ , we can just use the standard formula

$$u = u^- - \omega \frac{F(u^-)}{F'(u^-)} = u^- - \omega \frac{u^- - \Delta t f(u^-, t_n) - u^{(1)}}{1 - \Delta t \frac{\partial f}{\partial u}(u^-, t_n)}. \quad (6.11)$$

**Crank-Nicolson discretization.** The standard Crank-Nicolson scheme with arithmetic mean approximation of  $f$  takes the form

$$\frac{u^{n+1} - u^n}{\Delta t} = \frac{1}{2}(f(u^{n+1}, t_{n+1}) + f(u^n, t_n)).$$

We can write the scheme as a nonlinear algebraic equation

$$F(u) = u - u^{(1)} - \Delta t \frac{1}{2} f(u, t_{n+1}) - \Delta t \frac{1}{2} f(u^{(1)}, t_n) = 0. \quad (6.12)$$

A Picard iteration scheme must in general employ the linearization

$$\hat{F}(u) = u - u^{(1)} - \Delta t \frac{1}{2} f(u^-, t_{n+1}) - \Delta t \frac{1}{2} f(u^{(1)}, t_n),$$

while Newton's method can apply the general formula (6.11) with  $F(u)$  given in (6.12) and

$$F'(u) = 1 - \frac{1}{2} \Delta t \frac{\partial f}{\partial u}(u, t_{n+1}).$$

### 6.1.12 Systems of ODEs

We may write a system of ODEs

$$\begin{aligned}\frac{d}{dt}u_0(t) &= f_0(u_0(t), u_1(t), \dots, u_N(t), t), \\ \frac{d}{dt}u_1(t) &= f_1(u_0(t), u_1(t), \dots, u_N(t), t), \\ &\vdots \\ \frac{d}{dt}u_m(t) &= f_m(u_0(t), u_1(t), \dots, u_N(t), t),\end{aligned}$$

as

$$u' = f(u, t), \quad u(0) = U_0, \quad (6.13)$$

if we interpret  $u$  as a vector  $u = (u_0(t), u_1(t), \dots, u_N(t))$  and  $f$  as a vector function with components  $(f_0(u, t), f_1(u, t), \dots, f_N(u, t))$ .

Most solution methods for scalar ODEs, including the Forward and Backward Euler schemes and the Crank-Nicolson method, generalize in a straightforward way to systems of ODEs simply by using vector arithmetics instead of scalar arithmetics, which corresponds to applying the scalar scheme to each component of the system. For example, here is a backward difference scheme applied to each component,

$$\begin{aligned}\frac{u_0^n - u_0^{n-1}}{\Delta t} &= f_0(u^n, t_n), \\ \frac{u_1^n - u_1^{n-1}}{\Delta t} &= f_1(u^n, t_n), \\ &\vdots \\ \frac{u_N^n - u_N^{n-1}}{\Delta t} &= f_N(u^n, t_n),\end{aligned}$$

which can be written more compactly in vector form as

$$\frac{u^n - u^{n-1}}{\Delta t} = f(u^n, t_n).$$

This is a *system of algebraic equations*,

$$u^n - \Delta t f(u^n, t_n) - u^{n-1} = 0,$$

or written out

$$u_0^n - \Delta t f_0(u^n, t_n) - u_0^{n-1} = 0,$$

⋮

$$u_N^n - \Delta t f_N(u^n, t_n) - u_N^{n-1} = 0.$$

**Example.** We shall address the  $2 \times 2$  ODE system for oscillations of a pendulum subject to gravity and air drag. The system can be written as

$$\dot{\omega} = -\sin \theta - \beta \omega |\omega|, \quad (6.14)$$

$$\dot{\theta} = \omega, \quad (6.15)$$

where  $\beta$  is a dimensionless parameter (this is the scaled, dimensionless version of the original, physical model). The unknown components of the system are the angle  $\theta(t)$  and the angular velocity  $\omega(t)$ . We introduce  $u_0 = \omega$  and  $u_1 = \theta$ , which leads to

$$\begin{aligned} u'_0 &= f_0(u, t) = -\sin u_1 - \beta u_0 |u_0|, \\ u'_1 &= f_1(u, t) = u_0. \end{aligned}$$

A Crank-Nicolson scheme reads

$$\begin{aligned} \frac{u_0^{n+1} - u_0^n}{\Delta t} &= -\sin u_1^{n+\frac{1}{2}} - \beta u_0^{n+\frac{1}{2}} |u_0^{n+\frac{1}{2}}| \\ &\approx -\sin \left( \frac{1}{2}(u_1^{n+1} + u_1^n) \right) - \beta \frac{1}{4}(u_0^{n+1} + u_0^n) |u_0^{n+1} + u_0^n|, \end{aligned} \quad (6.16)$$

$$\frac{u_1^{n+1} - u_1^n}{\Delta t} = u_0^{n+\frac{1}{2}} \approx \frac{1}{2}(u_0^{n+1} + u_0^n). \quad (6.17)$$

This is a *coupled system* of two nonlinear algebraic equations in two unknowns  $u_0^{n+1}$  and  $u_1^{n+1}$ .

Using the notation  $u_0$  and  $u_1$  for the unknowns  $u_0^{n+1}$  and  $u_1^{n+1}$  in this system, writing  $u_0^{(1)}$  and  $u_1^{(1)}$  for the previous values  $u_0^n$  and  $u_1^n$ , multiplying by  $\Delta t$  and moving the terms to the left-hand sides, gives

$$u_0 - u_0^{(1)} + \Delta t \sin \left( \frac{1}{2}(u_1 + u_1^{(1)}) \right) + \frac{1}{4} \Delta t \beta(u_0 + u_0^{(1)}) |u_0 + u_0^{(1)}| = 0, \quad (6.18)$$

$$u_1 - u_1^{(1)} - \frac{1}{2} \Delta t (u_0 + u_0^{(1)}) = 0. \quad (6.19)$$

Obviously, we have a need for solving systems of nonlinear algebraic equations, which is the topic of the next section.

## 6.2 Systems of nonlinear algebraic equations

Implicit time discretization methods for a system of ODEs, or a PDE, lead to *systems* of nonlinear algebraic equations, written compactly as

$$F(u) = 0,$$

where  $u$  is a vector of unknowns  $u = (u_0, \dots, u_N)$ , and  $F$  is a vector function:  $F = (F_0, \dots, F_N)$ . The system at the end of Section 6.1.12 fits this notation with  $N = 2$ ,  $F_0(u)$  given by the left-hand side of (6.18), while  $F_1(u)$  is the left-hand side of (6.19).

Sometimes the equation system has a special structure because of the underlying problem, e.g.,

$$A(u)u = b(u),$$

with  $A(u)$  as an  $(N+1) \times (N+1)$  matrix function of  $u$  and  $b$  as a vector function:  $b = (b_0, \dots, b_N)$ .

We shall next explain how Picard iteration and Newton's method can be applied to systems like  $F(u) = 0$  and  $A(u)u = b(u)$ . The exposition has a focus on ideas and practical computations. More theoretical considerations, including quite general results on convergence properties of these methods, can be found in Kelley [3].

### 6.2.1 Picard iteration

We cannot apply Picard iteration to nonlinear equations unless there is some special structure. For the commonly arising case  $A(u)u = b(u)$  we can linearize the product  $A(u)u$  to  $A(u^-)u$  and  $b(u)$  as  $b(u^-)$ . That is,

we use the most previously computed approximation in  $A$  and  $b$  to arrive at a *linear system* for  $u$ :

$$A(u^-)u = b(u^-).$$

A relaxed iteration takes the form

$$A(u^-)u^* = b(u^-), \quad u = \omega u^* + (1 - \omega)u^-.$$

In other words, we solve a system of nonlinear algebraic equations as a sequence of linear systems.

### Algorithm for relaxed Picard iteration

Given  $A(u)u = b(u)$  and an initial guess  $u^-$ , iterate until convergence:

1. solve  $A(u^-)u^* = b(u^-)$  with respect to  $u^*$
2.  $u = \omega u^* + (1 - \omega)u^-$
3.  $u^- \leftarrow u$

“Until convergence” means that the iteration is stopped when the change in the unknown,  $\|u - u^-\|$ , or the residual  $\|A(u)u - b\|$ , is sufficiently small, see Section 6.2.3 for more details.

### 6.2.2 Newton’s method

The natural starting point for Newton’s method is the general nonlinear vector equation  $F(u) = 0$ . As for a scalar equation, the idea is to approximate  $F$  around a known value  $u^-$  by a linear function  $\hat{F}$ , calculated from the first two terms of a Taylor expansion of  $F$ . In the multi-variate case these two terms become

$$F(u^-) + J(u^-) \cdot (u - u^-),$$

where  $J$  is the *Jacobian* of  $F$ , defined by

$$J_{i,j} = \frac{\partial F_i}{\partial u_j}.$$

So, the original nonlinear system is approximated by

$$\hat{F}(u) = F(u^-) + J(u^-) \cdot (u - u^-) = 0,$$

which is linear in  $u$  and can be solved in a two-step procedure: first solve  $J\delta u = -F(u^-)$  with respect to the vector  $\delta u$  and then update  $u = u^- + \delta u$ . A relaxation parameter can easily be incorporated:

$$u = \omega(u^- + \delta u) + (1 - \omega)u^- = u^- + \omega\delta u.$$

### Algorithm for Newton's method

Given  $F(u) = 0$  and an initial guess  $u^-$ , iterate until convergence:

1. solve  $J\delta u = -F(u^-)$  with respect to  $\delta u$
2.  $u = u^- + \omega\delta u$
3.  $u^- \leftarrow u$

For the special system with structure  $A(u)u = b(u)$ ,

$$F_i = \sum_k A_{i,k}(u)u_k - b_i(u),$$

one gets

$$J_{i,j} = \sum_k \frac{\partial A_{i,k}}{\partial u_j} u_k + A_{i,j} - \frac{\partial b_i}{\partial u_j}. \quad (6.20)$$

We realize that the Jacobian needed in Newton's method consists of  $A(u^-)$  as in the Picard iteration plus two additional terms arising from the differentiation. Using the notation  $A'(u)$  for  $\partial A / \partial u$  (a quantity with three indices:  $\partial A_{i,k} / \partial u_j$ ), and  $b'(u)$  for  $\partial b / \partial u$  (a quantity with two indices:  $\partial b_i / \partial u_j$ ), we can write the linear system to be solved as

$$(A + A'u + b')\delta u = -Au + b,$$

or

$$(A(u^-) + A'(u^-)u^- + b'(u^-))\delta u = -A(u^-)u^- + b(u^-).$$

Rearranging the terms demonstrates the difference from the system solved in each Picard iteration:

$$\underbrace{A(u^-)(u^- + \delta u) - b(u^-)}_{\text{Picard system}} + \gamma(A'(u^-)u^- + b'(u^-))\delta u = 0.$$

Here we have inserted a parameter  $\gamma$  such that  $\gamma = 0$  gives the Picard system and  $\gamma = 1$  gives the Newton system. Such a parameter can be handy in software to easily switch between the methods.

### Combined algorithm for Picard and Newton iteration

Given  $A(u)$ ,  $b(u)$ , and an initial guess  $u^-$ , iterate until convergence:

1. solve  $(A + \gamma(A'(u^-)u^- + b'(u^-)))\delta u = -A(u^-)u^- + b(u^-)$  with respect to  $\delta u$
2.  $u = u^- + \omega\delta u$
3.  $u^- \leftarrow u$

$\gamma = 1$  gives a Newton method while  $\gamma = 0$  corresponds to Picard iteration.

### 6.2.3 Stopping criteria

Let  $\|\cdot\|$  be the standard Euclidean vector norm. Four termination criteria are much in use:

- Absolute change in solution:  $\|u - u^-\| \leq \epsilon_u$
- Relative change in solution:  $\|u - u^-\| \leq \epsilon_u \|u_0\|$ , where  $u_0$  denotes the start value of  $u^-$  in the iteration
- Absolute residual:  $\|F(u)\| \leq \epsilon_r$
- Relative residual:  $\|F(u)\| \leq \epsilon_r \|F(u_0)\|$

To prevent divergent iterations to run forever, one terminates the iterations when the current number of iterations  $k$  exceeds a maximum value  $k_{\max}$ .

The relative criteria are most used since they are not sensitive to the characteristic size of  $u$ . Nevertheless, the relative criteria can be misleading when the initial start value for the iteration is very close to the solution, since an unnecessary reduction in the error measure is enforced. In such cases the absolute criteria work better. It is common

to combine the absolute and relative measures of the size of the residual, as in

$$\|F(u)\| \leq \epsilon_{rr}\|F(u_0)\| + \epsilon_{ra}, \quad (6.21)$$

where  $\epsilon_{rr}$  is the tolerance in the relative criterion and  $\epsilon_{ra}$  is the tolerance in the absolute criterion. With a very good initial guess for the iteration (typically the solution of a differential equation at the previous time level), the term  $\|F(u_0)\|$  is small and  $\epsilon_{ra}$  is the dominating tolerance. Otherwise,  $\epsilon_{rr}\|F(u_0)\|$  and the relative criterion dominates.

With the change in solution as criterion we can formulate a combined absolute and relative measure of the change in the solution:

$$\|\delta u\| \leq \epsilon_{ur}\|u_0\| + \epsilon_{ua}, \quad (6.22)$$

The ultimate termination criterion, combining the residual and the change in solution with a test on the maximum number of iterations, can be expressed as

$$\|F(u)\| \leq \epsilon_{rr}\|F(u_0)\| + \epsilon_{ra} \quad \text{or} \quad \|\delta u\| \leq \epsilon_{ur}\|u_0\| + \epsilon_{ua} \quad \text{or} \quad k > k_{\max}. \quad (6.23)$$

#### 6.2.4 Example: A nonlinear ODE model from epidemiology

The simplest model spreading of a disease, such as a flu, takes the form of a  $2 \times 2$  ODE system

$$S' = -\beta SI, \quad (6.24)$$

$$I' = \beta SI - \nu I, \quad (6.25)$$

where  $S(t)$  is the number of people who can get ill (susceptibles) and  $I(t)$  is the number of people who are ill (infected). The constants  $\beta > 0$  and  $\nu > 0$  must be given along with initial conditions  $S(0)$  and  $I(0)$ .

**Implicit time discretization.** A Crank-Nicolson scheme leads to a  $2 \times 2$  system of nonlinear algebraic equations in the unknowns  $S^{n+1}$  and  $I^{n+1}$ :

$$\frac{S^{n+1} - S^n}{\Delta t} = -\beta[SI]^{n+\frac{1}{2}} \approx -\frac{\beta}{2}(S^n I^n + S^{n+1} I^{n+1}), \quad (6.26)$$

$$\frac{I^{n+1} - I^n}{\Delta t} = \beta[SI]^{n+\frac{1}{2}} - \nu I^{n+\frac{1}{2}} \approx \frac{\beta}{2}(S^n I^n + S^{n+1} I^{n+1}) - \frac{\nu}{2}(I^n + I^{n+1}). \quad (6.27)$$

Introducing  $S$  for  $S^{n+1}$ ,  $S^{(1)}$  for  $S^n$ ,  $I$  for  $I^{n+1}$ ,  $I^{(1)}$  for  $I^n$ , we can rewrite the system as

$$F_S(S, I) = S - S^{(1)} + \frac{1}{2}\Delta t\beta(S^{(1)}I^{(1)} + SI) = 0, \quad (6.28)$$

$$F_I(S, I) = I - I^{(1)} - \frac{1}{2}\Delta t\beta(S^{(1)}I^{(1)} + SI) + \frac{1}{2}\Delta t\nu(I^{(1)} + I) = 0. \quad (6.29)$$

**A Picard iteration.** We assume that we have approximations  $S^-$  and  $I^-$  to  $S$  and  $I$ . A way of linearizing the only nonlinear term  $SI$  is to write  $I^-S$  in the  $F_S = 0$  equation and  $S^-I$  in the  $F_I = 0$  equation, which also *decouples* the equations. Solving the resulting linear equations with respect to the unknowns  $S$  and  $I$  gives

$$S = \frac{S^{(1)} - \frac{1}{2}\Delta t\beta S^{(1)}I^{(1)}}{1 + \frac{1}{2}\Delta t\beta I^-},$$

$$I = \frac{I^{(1)} + \frac{1}{2}\Delta t\beta S^{(1)}I^{(1)} - \frac{1}{2}\Delta t\nu I^{(1)}}{1 - \frac{1}{2}\Delta t\beta S^- + \frac{1}{2}\Delta t\nu}.$$

Before a new iteration, we must update  $S^- \leftarrow S$  and  $I^- \leftarrow I$ .

**Newton's method.** The nonlinear system (6.28)-(6.29) can be written as  $F(u) = 0$  with  $F = (F_S, F_I)$  and  $u = (S, I)$ . The Jacobian becomes

$$J = \begin{pmatrix} \frac{\partial}{\partial S} F_S & \frac{\partial}{\partial I} F_S \\ \frac{\partial}{\partial S} F_I & \frac{\partial}{\partial I} F_I \end{pmatrix} = \begin{pmatrix} 1 + \frac{1}{2}\Delta t\beta I & \frac{1}{2}\Delta t\beta S \\ -\frac{1}{2}\Delta t\beta I & 1 - \frac{1}{2}\Delta t\beta S + \frac{1}{2}\Delta t\nu \end{pmatrix}.$$

The Newton system  $J(u^-)\delta u = -F(u^-)$  to be solved in each iteration is then

$$\begin{pmatrix} 1 + \frac{1}{2}\Delta t\beta I^- & \frac{1}{2}\Delta t\beta S^- \\ -\frac{1}{2}\Delta t\beta I^- & 1 - \frac{1}{2}\Delta t\beta S^- + \frac{1}{2}\Delta t\nu \end{pmatrix} \begin{pmatrix} \delta S \\ \delta I \end{pmatrix} = \begin{pmatrix} S^- - S^{(1)} + \frac{1}{2}\Delta t\beta(S^{(1)}I^{(1)} + S^-I^-) \\ I^- - I^{(1)} - \frac{1}{2}\Delta t\beta(S^{(1)}I^{(1)} + S^-I^-) + \frac{1}{2}\Delta t\nu(I^{(1)} + I^-) \end{pmatrix}$$

**Remark.** For this particular system of ODEs, explicit time integration methods work very well. Even a Forward Euler scheme is fine, but (as also experienced more generally) the 4-th order Runge-Kutta method is an excellent balance between high accuracy, high efficiency, and simplicity.

## 6.3 Linearization at the differential equation level

The attention is now turned to nonlinear partial differential equations (PDEs) and application of the techniques explained above for ODEs. The model problem is a nonlinear diffusion equation for  $u(\mathbf{x}, t)$ :

$$\frac{\partial u}{\partial t} = \nabla \cdot (\alpha(u)\nabla u) + f(u), \quad \mathbf{x} \in \Omega, \quad t \in (0, T], \quad (6.30)$$

$$-\alpha(u) \frac{\partial u}{\partial n} = g, \quad \mathbf{x} \in \partial\Omega_N, \quad t \in (0, T], \quad (6.31)$$

$$u = u_0, \quad \mathbf{x} \in \partial\Omega_D, \quad t \in (0, T]. \quad (6.32)$$

In the present section, our aim is to discretize this problem in time and then present techniques for linearizing the time-discrete PDE problem “at the PDE level” such that we transform the nonlinear stationary PDE problem at each time level into a sequence of linear PDE problems, which can be solved using any method for linear PDEs. This strategy avoids the solution of systems of nonlinear algebraic equations. In Section 6.4 we shall take the opposite (and more common) approach: discretize the nonlinear problem in time and space first, and then solve the resulting nonlinear algebraic equations at each time level by the methods of Section 6.2. Very often, the two approaches are mathematically identical, so there is no preference from a computational efficiency point of view. The details of the ideas sketched above will hopefully become clear through the forthcoming examples.

### 6.3.1 Explicit time integration

The nonlinearities in the PDE are trivial to deal with if we choose an explicit time integration method for (6.30), such as the Forward Euler method:

$$[D_t^+ u = \nabla \cdot (\alpha(u) \nabla u) + f(u)]^n,$$

or written out,

$$\frac{u^{n+1} - u^n}{\Delta t} = \nabla \cdot (\alpha(u^n) \nabla u^n) + f(u^n),$$

which is a linear equation in the unknown  $u^{n+1}$  with solution

$$u^{n+1} = u^n + \Delta t \nabla \cdot (\alpha(u^n) \nabla u^n) + \Delta t f(u^n).$$

The disadvantage with this discretization is the strict stability criterion  $\Delta t \leq h^2/(6 \max \alpha)$  for the case  $f = 0$  and a standard 2nd-order finite difference discretization in 3D space with mesh cell sizes  $h = \Delta x = \Delta y = \Delta z$ .

### 6.3.2 Backward Euler scheme and Picard iteration

A Backward Euler scheme for (6.30) reads

$$[D_t^- u = \nabla \cdot (\alpha(u) \nabla u) + f(u)]^n.$$

Written out,

$$\frac{u^n - u^{n-1}}{\Delta t} = \nabla \cdot (\alpha(u^n) \nabla u^n) + f(u^n). \quad (6.33)$$

This is a nonlinear PDE for the unknown function  $u^n(\mathbf{x})$ . Such a PDE can be viewed as a time-independent PDE where  $u^{n-1}(\mathbf{x})$  is a known function.

We introduce a Picard iteration with  $k$  as iteration counter. A typical linearization of the  $\nabla \cdot (\alpha(u^n) \nabla u^n)$  term in iteration  $k + 1$  is to use the previously computed  $u^{n,k}$  approximation in the diffusion coefficient:  $\alpha(u^{n,k})$ . The nonlinear source term is treated similarly:  $f(u^{n,k})$ . The unknown function  $u^{n,k+1}$  then fulfills the linear PDE

$$\frac{u^{n,k+1} - u^{n-1}}{\Delta t} = \nabla \cdot (\alpha(u^{n,k}) \nabla u^{n,k+1}) + f(u^{n,k}). \quad (6.34)$$

The initial guess for the Picard iteration at this time level can be taken as the solution at the previous time level:  $u^{n,0} = u^{n-1}$ .

We can alternatively apply the implementation-friendly notation where  $u$  corresponds to the unknown we want to solve for, i.e.,  $u^{n,k+1}$  above, and  $u^-$  is the most recently computed value,  $u^{n,k}$  above. Moreover,  $u^{(1)}$  denotes the unknown function at the previous time level,  $u^{n-1}$  above. The PDE to be solved in a Picard iteration then looks like

$$\frac{u - u^{(1)}}{\Delta t} = \nabla \cdot (\alpha(u^-) \nabla u) + f(u^-). \quad (6.35)$$

At the beginning of the iteration we start with the value from the previous time level:  $u^- = u^{(1)}$ , and after each iteration,  $u^-$  is updated to  $u$ .

### Remark on notation

The previous derivations of the numerical scheme for time discretizations of PDEs have, strictly speaking, a somewhat sloppy notation, but it is much used and convenient to read. A more precise notation must distinguish clearly between the exact solution of the PDE problem, here denoted  $u_e(\mathbf{x}, t)$ , and the exact solution of the spatial problem, arising after time discretization at each time level, where (6.33) is an example. The latter is here represented as  $u^n(\mathbf{x})$  and is an approximation to  $u_e(\mathbf{x}, t_n)$ . Then we have another approximation  $u^{n,k}(\mathbf{x})$  to  $u^n(\mathbf{x})$  when solving the nonlinear PDE problem for  $u^n$  by iteration methods, as in (6.34).

In our notation,  $u$  is a synonym for  $u^{n,k+1}$  and  $u^{(1)}$  is a synonym for  $u^{n-1}$ , inspired by what are natural variable names in a code. We will usually state the PDE problem in terms of  $u$  and quickly redefine the symbol  $u$  to mean the numerical approximation, while  $u_e$  is not explicitly introduced unless we need to talk about the exact solution and the approximate solution at the same time.

### 6.3.3 Backward Euler scheme and Newton's method

At time level  $n$ , we have to solve the stationary PDE (6.33). In the previous section, we saw how this can be done with Picard iterations. Another alternative is to apply the idea of Newton's method in a clever way. Normally, Newton's method is defined for systems of *algebraic*

*equations*, but the idea of the method can be applied at the PDE level too.

**Linearization via Taylor expansions.** Let  $u^{n,k}$  be an approximation to the unknown  $u^n$ . We seek a better approximation on the form

$$u^n = u^{n,k} + \delta u. \quad (6.36)$$

The idea is to insert (6.36) in (6.33), Taylor expand the nonlinearities and keep only the terms that are linear in  $\delta u$  (which makes (6.36) an approximation for  $u^n$ ). Then we can solve a linear PDE for the correction  $\delta u$  and use (6.36) to find a new approximation

$$u^{n,k+1} = u^{n,k} + \delta u$$

to  $u^n$ . Repeating this procedure gives a sequence  $u^{n,k+1}$ ,  $k = 0, 1, \dots$  that hopefully converges to the goal  $u^n$ .

Let us carry out all the mathematical details for the nonlinear diffusion PDE discretized by the Backward Euler method. Inserting (6.36) in (6.33) gives

$$\frac{u^{n,k} + \delta u - u^{n-1}}{\Delta t} = \nabla \cdot (\alpha(u^{n,k} + \delta u) \nabla(u^{n,k} + \delta u)) + f(u^{n,k} + \delta u). \quad (6.37)$$

We can Taylor expand  $\alpha(u^{n,k} + \delta u)$  and  $f(u^{n,k} + \delta u)$ :

$$\begin{aligned} \alpha(u^{n,k} + \delta u) &= \alpha(u^{n,k}) + \frac{d\alpha}{du}(u^{n,k})\delta u + \mathcal{O}(\delta u^2) \approx \alpha(u^{n,k}) + \alpha'(u^{n,k})\delta u, \\ f(u^{n,k} + \delta u) &= f(u^{n,k}) + \frac{df}{du}(u^{n,k})\delta u + \mathcal{O}(\delta u^2) \approx f(u^{n,k}) + f'(u^{n,k})\delta u. \end{aligned}$$

Inserting the linear approximations of  $\alpha$  and  $f$  in (6.37) results in

$$\begin{aligned} \frac{u^{n,k} + \delta u - u^{n-1}}{\Delta t} &= \nabla \cdot (\alpha(u^{n,k}) \nabla u^{n,k}) + f(u^{n,k}) + \\ &\quad \nabla \cdot (\alpha(u^{n,k}) \nabla \delta u) + \nabla \cdot (\alpha'(u^{n,k}) \delta u \nabla u^{n,k}) + \\ &\quad \nabla \cdot (\alpha'(u^{n,k}) \delta u \nabla \delta u) + f'(u^{n,k}) \delta u. \end{aligned} \quad (6.38)$$

The term  $\alpha'(u^{n,k}) \delta u \nabla \delta u$  is of order  $\delta u^2$  and therefore omitted since we expect the correction  $\delta u$  to be small ( $\delta u \gg \delta u^2$ ). Reorganizing the equation gives a PDE for  $\delta u$  that we can write in short form as

$$\delta F(\delta u; u^{n,k}) = -F(u^{n,k}),$$

where

$$F(u^{n,k}) = \frac{u^{n,k} - u^{n-1}}{\Delta t} - \nabla \cdot (\alpha(u^{n,k}) \nabla u^{n,k}) + f(u^{n,k}), \quad (6.39)$$

$$\begin{aligned} \delta F(\delta u; u^{n,k}) &= -\frac{1}{\Delta t} \delta u + \nabla \cdot (\alpha(u^{n,k}) \nabla \delta u) + \\ &\quad \nabla \cdot (\alpha'(u^{n,k}) \delta u \nabla u^{n,k}) + f'(u^{n,k}) \delta u. \end{aligned} \quad (6.40)$$

Note that  $\delta F$  is a linear function of  $\delta u$ , and  $F$  contains only terms that are known, such that the PDE for  $\delta u$  is indeed linear.

### Observations

The notational form  $\delta F = -F$  resembles the Newton system  $J\delta u = -F$  for systems of algebraic equations, with  $\delta F$  as  $J\delta u$ . The unknown vector in a linear system of algebraic equations enters the system as a linear operator in terms of a matrix-vector product ( $J\delta u$ ), while at the PDE level we have a linear differential operator instead ( $\delta F$ ).

**Similarity with Picard iteration.** We can rewrite the PDE for  $\delta u$  in a slightly different way too if we define  $u^{n,k} + \delta u$  as  $u^{n,k+1}$ .

$$\begin{aligned} \frac{u^{n,k+1} - u^{n-1}}{\Delta t} &= \nabla \cdot (\alpha(u^{n,k}) \nabla u^{n,k+1}) + f(u^{n,k}) \\ &\quad + \nabla \cdot (\alpha'(u^{n,k}) \delta u \nabla u^{n,k}) + f'(u^{n,k}) \delta u. \end{aligned} \quad (6.41)$$

Note that the first line is the same PDE as arise in the Picard iteration, while the remaining terms arise from the differentiations that are an inherent ingredient in Newton's method.

**Implementation.** For coding we want to introduce  $u$  for  $u^n$ ,  $u^-$  for  $u^{n,k}$  and  $u^{(1)}$  for  $u^{n-1}$ . The formulas for  $F$  and  $\delta F$  are then more clearly written as

$$F(u^-) = \frac{u^- - u^{(1)}}{\Delta t} - \nabla \cdot (\alpha(u^-) \nabla u^-) + f(u^-), \quad (6.42)$$

$$\begin{aligned} \delta F(\delta u; u^-) = & -\frac{1}{\Delta t} \delta u + \nabla \cdot (\alpha(u^-) \nabla \delta u) + \\ & \nabla \cdot (\alpha'(u^-) \delta u \nabla u^-) + f'(u^-) \delta u. \end{aligned} \quad (6.43)$$

The form that orders the PDE as the Picard iteration terms plus the Newton method's derivative terms becomes

$$\begin{aligned} \frac{u - u^{(1)}}{\Delta t} = & \nabla \cdot (\alpha(u^-) \nabla u) + f(u^-) + \\ & \gamma (\nabla \cdot (\alpha'(u^-) (u - u^-) \nabla u^-) + f'(u^-) (u - u^-)). \end{aligned} \quad (6.44)$$

The Picard and full Newton versions correspond to  $\gamma = 0$  and  $\gamma = 1$ , respectively.

**Derivation with alternative notation.** Some may prefer to derive the linearized PDE for  $\delta u$  using the more compact notation. We start with inserting  $u^n = u^- + \delta u$  to get

$$\frac{u^- + \delta u - u^{n-1}}{\Delta t} = \nabla \cdot (\alpha(u^- + \delta u) \nabla (u^- + \delta u)) + f(u^- + \delta u).$$

Taylor expanding,

$$\begin{aligned} \alpha(u^- + \delta u) & \approx \alpha(u^-) + \alpha'(u^-) \delta u, \\ f(u^- + \delta u) & \approx f(u^-) + f'(u^-) \delta u, \end{aligned}$$

and inserting these expressions gives a less cluttered PDE for  $\delta u$ :

$$\begin{aligned} \frac{u^- + \delta u - u^{n-1}}{\Delta t} = & \nabla \cdot (\alpha(u^-) \nabla u^-) + f(u^-) + \\ & \nabla \cdot (\alpha(u^-) \nabla \delta u) + \nabla \cdot (\alpha'(u^-) \delta u \nabla u^-) + \\ & \nabla \cdot (\alpha'(u^-) \delta u \nabla \delta u) + f'(u^-) \delta u. \end{aligned}$$

### 6.3.4 Crank-Nicolson discretization

A Crank-Nicolson discretization of (6.30) applies a centered difference at  $t_{n+\frac{1}{2}}$ :

$$[D_t u = \nabla \cdot (\alpha(u) \nabla u) + f(u)]^{n+\frac{1}{2}}.$$

The standard technique is to apply an arithmetic average for quantities defined between two mesh points, e.g.,

$$u^{n+\frac{1}{2}} \approx \frac{1}{2}(u^n + u^{n+1}).$$

However, with nonlinear terms we have many choices of formulating an arithmetic mean:

$$[f(u)]^{n+\frac{1}{2}} \approx f\left(\frac{1}{2}(u^n + u^{n+1})\right) = [f(\bar{u}^t)]^{n+\frac{1}{2}}, \quad (6.45)$$

$$[f(u)]^{n+\frac{1}{2}} \approx \frac{1}{2}(f(u^n) + f(u^{n+1})) = [\overline{f(u)}^t]^{n+\frac{1}{2}}, \quad (6.46)$$

$$[\alpha(u) \nabla u]^{n+\frac{1}{2}} \approx \alpha\left(\frac{1}{2}(u^n + u^{n+1})\right) \nabla\left(\frac{1}{2}(u^n + u^{n+1})\right) = [\alpha(\bar{u}^t) \nabla \bar{u}^t]^{n+\frac{1}{2}}, \quad (6.47)$$

$$[\alpha(u) \nabla u]^{n+\frac{1}{2}} \approx \frac{1}{2}(\alpha(u^n) + \alpha(u^{n+1})) \nabla\left(\frac{1}{2}(u^n + u^{n+1})\right) = [\overline{\alpha(u)}^t \nabla \bar{u}^t]^{n+\frac{1}{2}}, \quad (6.48)$$

$$[\alpha(u) \nabla u]^{n+\frac{1}{2}} \approx \frac{1}{2}(\alpha(u^n) \nabla u^n + \alpha(u^{n+1}) \nabla u^{n+1}) = [\overline{\alpha(u) \nabla u}^t]^{n+\frac{1}{2}}. \quad (6.49)$$

A big question is whether there are significant differences in accuracy between taking the products of arithmetic means or taking the arithmetic mean of products. Exercise 6.6 investigates this question, and the answer is that the approximation is  $\mathcal{O}(\Delta t^2)$  in both cases.

## 6.4 Discretization of 1D stationary nonlinear differential equations

Section 6.3 presented methods for linearizing time-discrete PDEs directly prior to discretization in space. We can alternatively carry out the discretization in space of the time-discrete nonlinear PDE problem and get a system of nonlinear algebraic equations, which can be solved by Picard iteration or Newton's method as presented in Section 6.2. This latter approach will now be described in detail.

We shall work with the 1D problem

$$-(\alpha(u)u')' + au = f(u), \quad x \in (0, L), \quad \alpha(u(0))u'(0) = C, \quad u(L) = D. \quad (6.50)$$

The problem (6.50) arises from the stationary limit of a diffusion equation,

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left( \alpha(u) \frac{\partial u}{\partial x} \right) - au + f(u), \quad (6.51)$$

as  $t \rightarrow \infty$  and  $\partial u / \partial t \rightarrow 0$ . Alternatively, the problem (6.50) arises at each time level from implicit time discretization of (6.51). For example, a Backward Euler scheme for (6.51) leads to

$$\frac{u^n - u^{n-1}}{\Delta t} = \frac{d}{dx} \left( \alpha(u^n) \frac{du^n}{dx} \right) - au^n + f(u^n). \quad (6.52)$$

Introducing  $u(x)$  for  $u^n(x)$ ,  $u^{(1)}$  for  $u^{n-1}$ , and defining  $f(u)$  in (6.50) to be  $f(u)$  in (6.52) plus  $u^{n-1}/\Delta t$ , gives (6.50) with  $a = 1/\Delta t$ .

### 6.4.1 Finite difference discretization

The nonlinearity in the differential equation (6.50) poses no more difficulty than a variable coefficient, as in the term  $(\alpha(x)u')'$ . We can therefore use a standard finite difference approach to discretizing the Laplace term with a variable coefficient:

$$[-D_x \alpha D_x u + au = f]_i.$$

Writing this out for a uniform mesh with points  $x_i = i\Delta x$ ,  $i = 0, \dots, N_x$ , leads to

$$-\frac{1}{\Delta x^2} \left( \alpha_{i+\frac{1}{2}} (u_{i+1} - u_i) - \alpha_{i-\frac{1}{2}} (u_i - u_{i-1}) \right) + au_i = f(u_i). \quad (6.53)$$

This equation is valid at all the mesh points  $i = 0, 1, \dots, N_x - 1$ . At  $i = N_x$  we have the Dirichlet condition  $u_i = 0$ . The only difference from the case with  $(\alpha(x)u')'$  and  $f(x)$  is that now  $\alpha$  and  $f$  are functions of  $u$  and not only on  $x$ :  $(\alpha(u(x))u')'$  and  $f(u(x))$ .

The quantity  $\alpha_{i+\frac{1}{2}}$ , evaluated between two mesh points, needs a comment. Since  $\alpha$  depends on  $u$  and  $u$  is only known at the mesh points, we need to express  $\alpha_{i+\frac{1}{2}}$  in terms of  $u_i$  and  $u_{i+1}$ . For this purpose we use

an arithmetic mean, although a harmonic mean is also common in this context if  $\alpha$  features large jumps. There are two choices of arithmetic means:

$$\alpha_{i+\frac{1}{2}} \approx \alpha\left(\frac{1}{2}(u_i + u_{i+1})\right) = [\alpha(\bar{u}^x)]^{i+\frac{1}{2}}, \quad (6.54)$$

$$\alpha_{i+\frac{1}{2}} \approx \frac{1}{2}(\alpha(u_i) + \alpha(u_{i+1})) = [\overline{\alpha(u)}^x]^{i+\frac{1}{2}} \quad (6.55)$$

Equation (6.53) with the latter approximation then looks like

$$\begin{aligned} & -\frac{1}{2\Delta x^2} ((\alpha(u_i) + \alpha(u_{i+1}))(u_{i+1} - u_i) - (\alpha(u_{i-1}) + \alpha(u_i))(u_i - u_{i-1})) \\ & + au_i = f(u_i), \end{aligned} \quad (6.56)$$

or written more compactly,

$$[-D_x \overline{\alpha}^x D_x u + au = f]_i.$$

At mesh point  $i = 0$  we have the boundary condition  $\alpha(u)u' = C$ , which is discretized by

$$[\alpha(u)D_{2x}u = C]_0,$$

meaning

$$\alpha(u_0) \frac{u_1 - u_{-1}}{2\Delta x} = C. \quad (6.57)$$

The fictitious value  $u_{-1}$  can be eliminated with the aid of (6.56) for  $i = 0$ . Formally, (6.56) should be solved with respect to  $u_{i-1}$  and that value (for  $i = 0$ ) should be inserted in (6.57), but it is algebraically much easier to do it the other way around. Alternatively, one can use a ghost cell  $[-\Delta x, 0]$  and update the  $u_{-1}$  value in the ghost cell according to (6.57) after every Picard or Newton iteration. Such an approach means that we use a known  $u_{-1}$  value in (6.56) from the previous iteration.

## 6.4.2 Solution of algebraic equations

**The structure of the equation system.** The nonlinear algebraic equations (6.56) are of the form  $A(u)u = b(u)$  with

$$\begin{aligned} A_{i,i} &= \frac{1}{2\Delta x^2}(\alpha(u_{i-1}) + 2\alpha(u_i)\alpha(u_{i+1})) + a, \\ A_{i,i-1} &= -\frac{1}{2\Delta x^2}(\alpha(u_{i-1}) + \alpha(u_i)), \\ A_{i,i+1} &= -\frac{1}{2\Delta x^2}(\alpha(u_i) + \alpha(u_{i+1})), \\ b_i &= f(u_i). \end{aligned}$$

The matrix  $A(u)$  is tridiagonal:  $A_{i,j} = 0$  for  $j > i + 1$  and  $j < i - 1$ .

The above expressions are valid for internal mesh points  $1 \leq i \leq N_x - 1$ . For  $i = 0$  we need to express  $u_{i-1} = u_{-1}$  in terms of  $u_1$  using (6.57):

$$u_{-1} = u_1 - \frac{2\Delta x}{\alpha(u_0)}C. \quad (6.58)$$

This value must be inserted in  $A_{0,0}$ . The expression for  $A_{i,i+1}$  applies for  $i = 0$ , and  $A_{i,i-1}$  does not enter the system when  $i = 0$ .

Regarding the last equation, its form depends on whether we include the Dirichlet condition  $u(L) = D$ , meaning  $u_{N_x} = D$ , in the nonlinear algebraic equation system or not. Suppose we choose  $(u_0, u_1, \dots, u_{N_x-1})$  as unknowns, later referred to as *systems without Dirichlet conditions*. The last equation corresponds to  $i = N_x - 1$ . It involves the boundary value  $u_{N_x}$ , which is substituted by  $D$ . If the unknown vector includes the boundary value,  $(u_0, u_1, \dots, u_{N_x})$ , later referred to as *system including Dirichlet conditions*, the equation for  $i = N_x - 1$  just involves the unknown  $u_{N_x}$ , and the final equation becomes  $u_{N_x} = D$ , corresponding to  $A_{i,i} = 1$  and  $b_i = D$  for  $i = N_x$ .

**Picard iteration.** The obvious Picard iteration scheme is to use previously computed values of  $u_i$  in  $A(u)$  and  $b(u)$ , as described more in detail in Section 6.2. With the notation  $u^-$  for the most recently computed value of  $u$ , we have the system  $F(u) \approx \hat{F}(u) = A(u^-)u - b(u^-)$ , with  $F = (F_0, F_1, \dots, F_m)$ ,  $u = (u_0, u_1, \dots, u_m)$ . The index  $m$  is  $N_x$  if the system includes the Dirichlet condition as a separate equation and  $N_x - 1$  otherwise. The matrix  $A(u^-)$  is tridiagonal, so the solution procedure is to fill a tridiagonal matrix data structure and the right-hand side vector with the right numbers and call a Gaussian elimination routine for tridiagonal linear systems.

**Mesh with two cells.** It helps on the understanding of the details to write out all the mathematics in a specific case with a small mesh, say just two cells ( $N_x = 2$ ). We use  $u_i^-$  for the  $i$ -th component in  $u^-$ .

The starting point is the basic expressions for the nonlinear equations at mesh point  $i = 0$  and  $i = 1$  are

$$A_{0,-1}u_{-1} + A_{0,0}u_0 + A_{0,1}u_1 = b_0, \quad (6.59)$$

$$A_{1,0}u_0 + A_{1,1}u_1 + A_{1,2}u_2 = b_1. \quad (6.60)$$

Equation (6.59) written out reads

$$\begin{aligned} & \frac{1}{2\Delta x^2}(-(\alpha(u_{-1}) + \alpha(u_0))u_{-1} + \\ & (\alpha(u_{-1}) + 2\alpha(u_0) + \alpha(u_1))u_0 - \\ & (\alpha(u_0) + \alpha(u_1)))u_1 + au_0 = f(u_0). \end{aligned}$$

We must then replace  $u_{-1}$  by (6.58). With Picard iteration we get

$$\begin{aligned} & \frac{1}{2\Delta x^2}(-(\alpha(u_{-1}^-) + 2\alpha(u_0^-) + \alpha(u_1^-))u_1 + \\ & (\alpha(u_{-1}^-) + 2\alpha(u_0^-) + \alpha(u_1^-))u_0 + au_0 \\ & = f(u_0^-) - \frac{1}{\alpha(u_0^-)\Delta x}(\alpha(u_{-1}^-) + \alpha(u_0^-))C, \end{aligned}$$

where

$$u_{-1}^- = u_1^- - \frac{2\Delta x}{\alpha(u_0^-)}C.$$

Equation (6.60) contains the unknown  $u_2$  for which we have a Dirichlet condition. In case we omit the condition as a separate equation, (6.60) with Picard iteration becomes

$$\begin{aligned} & \frac{1}{2\Delta x^2}(-(\alpha(u_0^-) + \alpha(u_1^-))u_0 + \\ & (\alpha(u_0^-) + 2\alpha(u_1^-) + \alpha(u_2^-))u_1 - \\ & (\alpha(u_1^-) + \alpha(u_2^-)))u_2 + au_1 = f(u_1^-). \end{aligned}$$

We must now move the  $u_2$  term to the right-hand side and replace all occurrences of  $u_2$  by  $D$ :

$$\begin{aligned} & \frac{1}{2\Delta x^2}(-(\alpha(u_0^-) + \alpha(u_1^-))u_0 + \\ & (\alpha(u_0^-) + 2\alpha(u_1^-) + \alpha(D))u_1 + au_1 \\ & = f(u_1^-) + \frac{1}{2\Delta x^2}(\alpha(u_1^-) + \alpha(D))D. \end{aligned}$$

The two equations can be written as a  $2 \times 2$  system:

$$\begin{pmatrix} B_{0,0} & B_{0,1} \\ B_{1,0} & B_{1,1} \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \end{pmatrix} = \begin{pmatrix} d_0 \\ d_1 \end{pmatrix},$$

where

$$B_{0,0} = \frac{1}{2\Delta x^2}(\alpha(u_{-1}^-) + 2\alpha(u_0^-) + \alpha(u_1^-)) + a \quad (6.61)$$

$$B_{0,1} = -\frac{1}{2\Delta x^2}(\alpha(u_{-1}^-) + 2\alpha(u_0^-) + \alpha(u_1^-)), \quad (6.62)$$

$$B_{1,0} = -\frac{1}{2\Delta x^2}(\alpha(u_0^-) + \alpha(u_1^-)), \quad (6.63)$$

$$B_{1,1} = \frac{1}{2\Delta x^2}(\alpha(u_0^-) + 2\alpha(u_1^-) + \alpha(D)) + a, \quad (6.64)$$

$$d_0 = f(u_0^-) - \frac{1}{\alpha(u_0^-)\Delta x}(\alpha(u_{-1}^-) + \alpha(u_0^-))C, \quad (6.65)$$

$$d_1 = f(u_1^-) + \frac{1}{2\Delta x^2}(\alpha(u_1^-) + \alpha(D))D. \quad (6.66)$$

The system with the Dirichlet condition becomes

$$\begin{pmatrix} B_{0,0} & B_{0,1} & 0 \\ B_{1,0} & B_{1,1} & B_{1,2} \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} u_0 \\ u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} d_0 \\ d_1 \\ D \end{pmatrix},$$

with

$$B_{1,1} = \frac{1}{2\Delta x^2}(\alpha(u_0^-) + 2\alpha(u_1^-) + \alpha(u_2^-)) + a, \quad (6.67)$$

$$B_{1,2} = -\frac{1}{2\Delta x^2}(\alpha(u_1^-) + \alpha(u_2^-)), \quad (6.68)$$

$$d_1 = f(u_1^-). \quad (6.69)$$

Other entries are as in the  $2 \times 2$  system.

**Newton's method.** The Jacobian must be derived in order to use Newton's method. Here it means that we need to differentiate  $F(u) = A(u)u - b(u)$  with respect to the unknown parameters  $u_0, u_1, \dots, u_m$  ( $m = N_x$  or  $m = N_x - 1$ , depending on whether the Dirichlet condition is included in the nonlinear system  $F(u) = 0$  or not). Nonlinear equation number  $i$  has the structure

$$F_i = A_{i,i-1}(u_{i-1}, u_i)u_{i-1} + A_{i,i}(u_{i-1}, u_i, u_{i+1})u_i + A_{i,i+1}(u_i, u_{i+1})u_{i+1} - b_i(u_i).$$

Computing the Jacobian requires careful differentiation. For example,

$$\begin{aligned} \frac{\partial}{\partial u_i} (A_{i,i}(u_{i-1}, u_i, u_{i+1})u_i) &= \frac{\partial A_{i,i}}{\partial u_i} u_i + A_{i,i} \frac{\partial u_i}{\partial u_i} \\ &= \frac{\partial}{\partial u_i} \left( \frac{1}{2\Delta x^2} (\alpha(u_{i-1}) + 2\alpha(u_i) + \alpha(u_{i+1})) + a \right) u_i + \\ &\quad \frac{1}{2\Delta x^2} (\alpha(u_{i-1}) + 2\alpha(u_i) + \alpha(u_{i+1})) + a \\ &= \frac{1}{2\Delta x^2} (2\alpha'(u_i)u_i + \alpha(u_{i-1}) + 2\alpha(u_i) + \alpha(u_{i+1})) + a. \end{aligned}$$

The complete Jacobian becomes

$$\begin{aligned} J_{i,i} &= \frac{\partial F_i}{\partial u_i} = \frac{\partial A_{i,i-1}}{\partial u_i} u_{i-1} + \frac{\partial A_{i,i}}{\partial u_i} u_i + A_{i,i} + \frac{\partial A_{i,i+1}}{\partial u_i} u_{i+1} - \frac{\partial b_i}{\partial u_i} \\ &= \frac{1}{2\Delta x^2} (-\alpha'(u_i)u_{i-1} + 2\alpha'(u_i)u_i + \alpha(u_{i-1}) + 2\alpha(u_i) + \alpha(u_{i+1})) + \\ &\quad a - \frac{1}{2\Delta x^2} \alpha'(u_i)u_{i+1} - b'(u_i), \\ J_{i,i-1} &= \frac{\partial F_i}{\partial u_{i-1}} = \frac{\partial A_{i,i-1}}{\partial u_{i-1}} u_{i-1} + A_{i-1,i} + \frac{\partial A_{i,i}}{\partial u_{i-1}} u_i - \frac{\partial b_i}{\partial u_{i-1}} \\ &= \frac{1}{2\Delta x^2} (-\alpha'(u_{i-1})u_{i-1} - (\alpha(u_{i-1}) + \alpha(u_i)) + \alpha'(u_{i-1})u_i), \\ J_{i,i+1} &= \frac{\partial A_{i,i+1}}{\partial u_{i-1}} u_{i+1} + A_{i+1,i} + \frac{\partial A_{i,i}}{\partial u_{i+1}} u_i - \frac{\partial b_i}{\partial u_{i+1}} \\ &= \frac{1}{2\Delta x^2} (-\alpha'(u_{i+1})u_{i+1} - (\alpha(u_i) + \alpha(u_{i+1})) + \alpha'(u_{i+1})u_i). \end{aligned}$$

The explicit expression for nonlinear equation number  $i$ ,  $F_i(u_0, u_1, \dots)$ , arises from moving the  $f(u_i)$  term in (6.56) to the left-hand side:

$$\begin{aligned} F_i = -\frac{1}{2\Delta x^2} & ((\alpha(u_i) + \alpha(u_{i+1}))(u_{i+1} - u_i) - (\alpha(u_{i-1}) + \alpha(u_i))(u_i - u_{i-1})) \\ & + au_i - f(u_i) = 0. \end{aligned} \quad (6.70)$$

At the boundary point  $i = 0$ ,  $u_{-1}$  must be replaced using the formula (6.58). When the Dirichlet condition at  $i = N_x$  is not a part of the equation system, the last equation  $F_m = 0$  for  $m = N_x - 1$  involves the quantity  $u_{N_x-1}$  which must be replaced by  $D$ . If  $u_{N_x}$  is treated as an unknown in the system, the last equation  $F_m = 0$  has  $m = N_x$  and reads

$$F_{N_x}(u_0, \dots, u_{N_x}) = u_{N_x} - D = 0.$$

Similar replacement of  $u_{-1}$  and  $u_{N_x}$  must be done in the Jacobian for the first and last row. When  $u_{N_x}$  is included as an unknown, the last row in the Jacobian must help implement the condition  $\delta u_{N_x} = 0$ , since we assume that  $u$  contains the right Dirichlet value at the beginning of the iteration ( $u_{N_x} = D$ ), and then the Newton update should be zero for  $i = 0$ , i.e.,  $\delta u_{N_x} = 0$ . This also forces the right-hand side to be  $b_i = 0$ ,  $i = N_x$ .

We have seen, and can see from the present example, that the linear system in Newton's method contains all the terms present in the system that arises in the Picard iteration method. The extra terms in Newton's method can be multiplied by a factor such that it is easy to program one linear system and set this factor to 0 or 1 to generate the Picard or Newton system.

## 6.5 Multi-dimensional PDE problems

The fundamental ideas in the derivation of  $F_i$  and  $J_{i,j}$  in the 1D model problem are easily generalized to multi-dimensional problems. Nevertheless, the expressions involved are slightly different, with derivatives in  $x$  replaced by  $\nabla$ , so we present some examples below in detail.

### 6.5.1 Finite difference discretization

A typical diffusion equation

$$u_t = \nabla \cdot (\alpha(u) \nabla u) + f(u),$$

can be discretized by (e.g.) a Backward Euler scheme, which in 2D can be written

$$[D_t^- u = D_x \overline{\alpha(u)}^x D_x u + D_y \overline{\alpha(u)}^y D_y u + f(u)]_{i,j}^n.$$

We do not dive into the details of handling boundary conditions now. Dirichlet and Neumann conditions are handled as in a corresponding linear, variable-coefficient diffusion problems.

Writing the scheme out, putting the unknown values on the left-hand side and known values on the right-hand side, and introducing  $\Delta x = \Delta y = h$  to save some writing, one gets

$$\begin{aligned} u_{i,j}^n - \frac{\Delta t}{h^2} & \left( \frac{1}{2}(\alpha(u_{i,j}^n) + \alpha(u_{i+1,j}^n))(u_{i+1,j}^n - u_{i,j}^n) \right. \\ & - \frac{1}{2}(\alpha(u_{i-1,j}^n) + \alpha(u_{i,j}^n))(u_{i,j}^n - u_{i-1,j}^n) \\ & + \frac{1}{2}(\alpha(u_{i,j}^n) + \alpha(u_{i,j+1}^n))(u_{i,j+1}^n - u_{i,j}^n) \\ & \left. - \frac{1}{2}(\alpha(u_{i,j-1}^n) + \alpha(u_{i,j}^n))(u_{i,j}^n - u_{i-1,j-1}^n) \right) - \Delta t f(u_{i,j}^n) = u_{i,j}^{n-1} \end{aligned}$$

This defines a nonlinear algebraic system on the form  $A(u)u = b(u)$ .

**Picard iteration.** The most recently computed values  $u^-$  of  $u^n$  can be used in  $\alpha$  and  $f$  for a Picard iteration, or equivalently, we solve  $A(u^-)u = b(u^-)$ . The result is a linear system of the same type as arising from  $u_t = \nabla \cdot (\alpha(\mathbf{x}) \nabla u) + f(\mathbf{x}, t)$ .

The Picard iteration scheme can also be expressed in operator notation:

$$[D_t^- u = D_x \overline{\alpha(u^-)}^x D_x u + D_y \overline{\alpha(u^-)}^y D_y u + f(u^-)]_{i,j}^n.$$

**Newton's method.** As always, Newton's method is technically more involved than Picard iteration. We first define the nonlinear algebraic equations to be solved, drop the superscript  $n$  (use  $u$  for  $u^n$ ), and introduce  $u^{(1)}$  for  $u^{n-1}$ :

$$\begin{aligned}
F_{i,j} = u_{i,j} - \frac{\Delta t}{h^2} ( & \\
& \frac{1}{2}(\alpha(u_{i,j}) + \alpha(u_{i+1,j}))(u_{i+1,j} - u_{i,j}) - \\
& \frac{1}{2}(\alpha(u_{i-1,j}) + \alpha(u_{i,j}))(u_{i,j} - u_{i-1,j}) + \\
& \frac{1}{2}(\alpha(u_{i,j}) + \alpha(u_{i,j+1}))(u_{i,j+1} - u_{i,j}) - \\
& \frac{1}{2}(\alpha(u_{i,j-1}) + \alpha(u_{i,j}))(u_{i,j} - u_{i-1,j-1})) - \Delta t f(u_{i,j}) - u_{i,j}^{(1)} = 0 .
\end{aligned}$$

It is convenient to work with two indices  $i$  and  $j$  in 2D finite difference discretizations, but it complicates the derivation of the Jacobian, which then gets four indices. (Make sure you really understand the 1D version of this problem as treated in Section 6.4.1.) The left-hand expression of an equation  $F_{i,j} = 0$  is to be differentiated with respect to each of the unknowns  $u_{r,s}$  (recall that this is short notation for  $u_{r,s}^n$ ),  $r \in \mathcal{I}_x$ ,  $s \in \mathcal{I}_y$ :

$$J_{i,j,r,s} = \frac{\partial F_{i,j}}{\partial u_{r,s}} .$$

The Newton system to be solved in each iteration can be written as

$$\sum_{r \in \mathcal{I}_x} \sum_{s \in \mathcal{I}_y} J_{i,j,r,s} \delta u_{r,s} = -F_{i,j}, \quad i \in \mathcal{I}_x, \quad j \in \mathcal{I}_y .$$

Given  $i$  and  $j$ , only a few  $r$  and  $s$  indices give nonzero contribution to the Jacobian since  $F_{i,j}$  contains  $u_{i \pm 1,j}$ ,  $u_{i,j \pm 1}$ , and  $u_{i,j}$ . This means that  $J_{i,j,r,s}$  has nonzero contributions only if  $r = i \pm 1$ ,  $s = j \pm 1$ , as well as  $r = i$  and  $s = j$ . The corresponding terms in  $J_{i,j,r,s}$  are  $J_{i,j,i-1,j}$ ,  $J_{i,j,i+1,j}$ ,  $J_{i,j,i,j-1}$ ,  $J_{i,j,i,j+1}$ , and  $J_{i,j,i,j}$ . Therefore, the left-hand side of the Newton system,  $\sum_r \sum_s J_{i,j,r,s} \delta u_{r,s}$  collapses to

$$\begin{aligned}
J_{i,j,r,s} \delta u_{r,s} = & J_{i,j,i,j} \delta u_{i,j} + J_{i,j,i-1,j} \delta u_{i-1,j} + J_{i,j,i+1,j} \delta u_{i+1,j} + J_{i,j,i,j-1} \delta u_{i,j-1} \\
& + J_{i,j,i,j+1} \delta u_{i,j+1}
\end{aligned}$$

The specific derivatives become

$$\begin{aligned}
J_{i,j,i-1,j} &= \frac{\partial F_{i,j}}{\partial u_{i-1,j}} \\
&= \frac{\Delta t}{h^2} (\alpha'(u_{i-1,j})(u_{i,j} - u_{i-1,j}) + \alpha(u_{i-1,j})(-1)), \\
J_{i,j,i+1,j} &= \frac{\partial F_{i,j}}{\partial u_{i+1,j}} \\
&= \frac{\Delta t}{h^2} (-\alpha'(u_{i+1,j})(u_{i+1,j} - u_{i,j}) - \alpha(u_{i-1,j})), \\
J_{i,j,i,j-1} &= \frac{\partial F_{i,j}}{\partial u_{i,j-1}} \\
&= \frac{\Delta t}{h^2} (\alpha'(u_{i,j-1})(u_{i,j} - u_{i,j-1}) + \alpha(u_{i,j-1})(-1)), \\
J_{i,j,i,j+1} &= \frac{\partial F_{i,j}}{\partial u_{i,j+1}} \\
&= \frac{\Delta t}{h^2} (-\alpha'(u_{i,j+1})(u_{i,j+1} - u_{i,j}) - \alpha(u_{i,j-1})). 
\end{aligned}$$

The  $J_{i,j,i,j}$  entry has a few more terms and is left as an exercise. Inserting the most recent approximation  $u^-$  for  $u$  in the  $J$  and  $F$  formulas and then forming  $J\delta u = -F$  gives the linear system to be solved in each Newton iteration. Boundary conditions will affect the formulas when any of the indices coincide with a boundary value of an index.

### 6.5.2 Continuation methods

Picard iteration or Newton's method may diverge when solving PDEs with severe nonlinearities. Relaxation with  $\omega < 1$  may help, but in highly nonlinear problems it can be necessary to introduce a *continuation parameter*  $\Lambda$  in the problem:  $\Lambda = 0$  gives a version of the problem that is easy to solve, while  $\Lambda = 1$  is the target problem. The idea is then to increase  $\Lambda$  in steps,  $\Lambda_0 = 0, \Lambda_1 < \dots < \Lambda_n = 1$ , and use the solution from the problem with  $\Lambda_{i-1}$  as initial guess for the iterations in the problem corresponding to  $\Lambda_i$ .

The continuation method is easiest to understand through an example. Suppose we intend to solve

$$-\nabla \cdot (||\nabla u||^q \nabla u) = f,$$

which is an equation modeling the flow of a non-Newtonian fluid through a channel or pipe. For  $q = 0$  we have the Poisson equation (corresponding to

a Newtonian fluid) and the problem is linear. A typical value for pseudo-plastic fluids may be  $q_n = -0.8$ . We can introduce the continuation parameter  $\Lambda \in [0, 1]$  such that  $q = q_n \Lambda$ . Let  $\{\Lambda_\ell\}_{\ell=0}^n$  be the sequence of  $\Lambda$  values in  $[0, 1]$ , with corresponding  $q$  values  $\{q_\ell\}_{\ell=0}^n$ . We can then solve a sequence of problems

$$-\nabla \cdot (||\nabla u^\ell||_\ell^q \nabla u^\ell) = f, \quad \ell = 0, \dots, n,$$

where the initial guess for iterating on  $u^\ell$  is the previously computed solution  $u^{\ell-1}$ . If a particular  $\Lambda_\ell$  leads to convergence problems, one may try a smaller increase in  $\Lambda$ :  $\Lambda_* = \frac{1}{2}(\Lambda_{\ell-1} + \Lambda_\ell)$ , and repeat halving the step in  $\Lambda$  until convergence is reestablished.

## 6.6 Exercises

### Problem 6.1: Determine if equations are nonlinear or not

Classify each term in the following equations as linear or nonlinear. Assume that  $u$ ,  $\mathbf{u}$ , and  $p$  are unknown functions and that all other symbols are known quantities.

1.  $mu'' + \beta|u'|u' + cu = F(t)$
2.  $u_t = \alpha u_{xx}$
3.  $u_{tt} = c^2 \nabla^2 u$
4.  $u_t = \nabla \cdot (\alpha(u) \nabla u) + f(x, y)$
5.  $u_t + f(u)_x = 0$
6.  $\mathbf{u}_t + \mathbf{u} \cdot \nabla \mathbf{u} = -\nabla p + r \nabla^2 \mathbf{u}$ ,  $\nabla \cdot \mathbf{u} = 0$  ( $\mathbf{u}$  is a vector field)
7.  $u' = f(u, t)$
8.  $\nabla^2 u = \lambda e^u$

Filename: `nonlinear_vs_linear.`

### Exercise 6.2: Derive and investigate a generalized logistic model

The logistic model for population growth is derived by assuming a nonlinear growth rate,

$$u' = a(u)u, \quad u(0) = I, \tag{6.71}$$

and the logistic model arises from the simplest possible choice of  $a(u)$ :  $r(u) = \varrho(1 - u/M)$ , where  $M$  is the maximum value of  $u$  that the environment can sustain, and  $\varrho$  is the growth under unlimited access to resources (as in the beginning when  $u$  is small). The idea is that  $a(u) \sim \varrho$  when  $u$  is small and that  $a(t) \rightarrow 0$  as  $u \rightarrow M$ .

An  $a(u)$  that generalizes the linear choice is the polynomial form

$$a(u) = \varrho(1 - u/M)^p, \quad (6.72)$$

where  $p > 0$  is some real number.

- a)** Formulate a Forward Euler, Backward Euler, and a Crank-Nicolson scheme for (6.71).

**Hint.** Use a geometric mean approximation in the Crank-Nicolson scheme:  $[a(u)u]^{n+1/2} \approx a(u^n)u^{n+1}$ .

- b)** Formulate Picard and Newton iteration for the Backward Euler scheme in a).

- c)** Implement the numerical solution methods from a) and b). Use `logistic.py` to compare the case  $p = 1$  and the choice (6.72).

- d)** Implement unit tests that check the asymptotic limit of the solutions:  $u \rightarrow M$  as  $t \rightarrow \infty$ .

**Hint.** You need to experiment to find what “infinite time” is (increases substantially with  $p$ ) and what the appropriate tolerance is for testing the asymptotic limit.

- e)** Perform experiments with Newton and Picard iteration for the model (6.72). See how sensitive the number of iterations is to  $\Delta t$  and  $p$ .

Filename: `logistic_p`.

### Problem 6.3: Experience the behavior of Newton’s method

The program `Newton_demo.py` illustrates graphically each step in Newton’s method and is run like

---

Terminal> python Newton_demo.py f dfdx x0 xmin xmax	Terminal
---	----------

---

Use this program to investigate potential problems with Newton’s method when solving  $e^{-0.5x^2} \cos(\pi x) = 0$ . Try a starting point  $x_0 = 0.8$  and  $x_0 = 0.85$  and watch the different behavior. Just run

---

	Terminal	
	Terminal> python Newton_demo.py '0.2 + exp(-0.5*x**2)*cos(pi*x), \	
	'-x*exp(-x**2)*cos(pi*x) - pi*exp(-x**2)*sin(pi*x)' \	
	0.85 -3 3	

---

and repeat with 0.85 replaced by 0.8.

### Problem 6.4: Compute the Jacobian of a $2 \times 2$ system

Write up the system (6.18)-(6.19) in the form  $F(u) = 0$ ,  $F = (F_0, F_1)$ ,  $u = (u_0, u_1)$ , and compute the Jacobian  $J_{i,j} = \partial F_i / \partial u_j$ .

### Problem 6.5: Solve nonlinear equations arising from a vibration ODE

Consider a nonlinear vibration problem

$$mu'' + bu'|u'| + s(u) = F(t), \quad (6.73)$$

where  $m > 0$  is a constant,  $b \geq 0$  is a constant,  $s(u)$  a possibly nonlinear function of  $u$ , and  $F(t)$  is a prescribed function. Such models arise from Newton's second law of motion in mechanical vibration problems where  $s(u)$  is a spring or restoring force,  $mu''$  is mass times acceleration, and  $bu'|u'|$  models water or air drag.

- a)** Rewrite the equation for  $u$  as a system of two first-order ODEs, and discretize this system by a Crank-Nicolson (centered difference) method. With  $v = u'$ , we get a nonlinear term  $v^{n+\frac{1}{2}}|v^{n+\frac{1}{2}}|$ . Use a geometric average for  $v^{n+\frac{1}{2}}$ .
  - b)** Formulate a Picard iteration method to solve the system of nonlinear algebraic equations.
  - c)** Explain how to apply Newton's method to solve the nonlinear equations at each time level. Derive expressions for the Jacobian and the right-hand side in each Newton iteration.
- Filename: `nonlin_vib.`

### Exercise 6.6: Find the truncation error of arithmetic mean of products

In Section 6.3.4 we introduce alternative arithmetic means of a product. Say the product is  $P(t)Q(t)$  evaluated at  $t = t_{n+\frac{1}{2}}$ . The exact value is

$$[PQ]^{n+\frac{1}{2}} = P^{n+\frac{1}{2}}Q^{n+\frac{1}{2}}$$

There are two obvious candidates for evaluating  $[PQ]^{n+\frac{1}{2}}$  as a mean of values of  $P$  and  $Q$  at  $t_n$  and  $t_{n+1}$ . Either we can take the arithmetic mean of each factor  $P$  and  $Q$ ,

$$[PQ]^{n+\frac{1}{2}} \approx \frac{1}{2}(P^n + P^{n+1})\frac{1}{2}(Q^n + Q^{n+1}), \quad (6.74)$$

or we can take the arithmetic mean of the product  $PQ$ :

$$[PQ]^{n+\frac{1}{2}} \approx \frac{1}{2}(P^nQ^n + P^{n+1}Q^{n+1}). \quad (6.75)$$

The arithmetic average of  $P(t_{n+\frac{1}{2}})$  is  $\mathcal{O}(\Delta t^2)$ :

$$P(t_{n+\frac{1}{2}}) = \frac{1}{2}(P^n + P^{n+1}) + \mathcal{O}(\Delta t^2).$$

A fundamental question is whether (6.74) and (6.75) have different orders of accuracy in  $\Delta t = t_{n+1} - t_n$ . To investigate this question, expand quantities at  $t_{n+1}$  and  $t_n$  in Taylor series around  $t_{n+\frac{1}{2}}$ , and subtract the true value  $[PQ]^{n+\frac{1}{2}}$  from the approximations (6.74) and (6.75) to see what the order of the error terms are.

**Hint.** You may explore `sympy` for carrying out the tedious calculations. A general Taylor series expansion of  $P(t + \frac{1}{2}\Delta t)$  around  $t$  involving just a general function  $P(t)$  can be created as follows:

```
>>> from sympy import *
>>> t, dt = symbols('t dt')
>>> P = symbols('P', cls=Function)
>>> P(t).series(t, 0, 4)
P(0) + t*Subs(Derivative(P(_x), _x), (_x,), (0,)) +
t**2*Subs(Derivative(P(_x), _x, _x), (_x,), (0,))/2 +
t**3*Subs(Derivative(P(_x), _x, _x, _x), (_x,), (0,))/6 + O(t**4)
>>> P_p = P(t).series(t, 0, 4).subs(t, dt/2)
>>> P_p
P(0) + dt*Subs(Derivative(P(_x), _x), (_x,), (0,))/2 +
dt**2*Subs(Derivative(P(_x), _x, _x), (_x,), (0,))/8 +
dt**3*Subs(Derivative(P(_x), _x, _x, _x), (_x,), (0,))/48 + O(dt**4)
```

The error of the arithmetic mean,  $\frac{1}{2}(P(-\frac{1}{2}\Delta t) + P(-\frac{1}{2}\Delta t))$  for  $t = 0$  is then

```
>>> P_m = P(t).series(t, 0, 4).subs(t, -dt/2)
>>> mean = Rational(1,2)*(P_m + P_p)
>>> error = simplify(expand(mean) - P(0))
>>> error
dt**2*Subs(Derivative(P(_x), _x, _x), (_x,), (0,))/8 + O(dt**4)
```

Use these examples to investigate the error of (6.74) and (6.75) for  $n = 0$ . (Choosing  $n = 0$  is necessary for not making the expressions too complicated for `sympy`, but there is of course no lack of generality by using  $n = 0$  rather than an arbitrary  $n$  - the main point is the product and addition of Taylor series.)

Filename: `product_arith_mean`.

### Problem 6.7: Newton's method for linear problems

Suppose we have a linear system  $F(u) = Au - b = 0$ . Apply Newton's method to this system, and show that the method converges in one iteration. Filename: `Newton_linear`.

### Exercise 6.8: Discretize a 1D problem with a nonlinear coefficient

We consider the problem

$$((1+u^2)u')' = 1, \quad x \in (0, 1), \quad u(0) = u(1) = 0. \quad (6.76)$$

Discretize (6.76) by a centered finite difference method on a uniform mesh. Filename: `nonlin_1D_coeff_discretize`.

### Exercise 6.9: Linearize a 1D problem with a nonlinear coefficient

We have a two-point boundary value problem

$$((1+u^2)u')' = 1, \quad x \in (0, 1), \quad u(0) = u(1) = 0. \quad (6.77)$$

**a)** Construct a Picard iteration method for (6.77) without discretizing in space.

- b)** Apply Newton's method to (6.77) without discretizing in space.
- c)** Discretize (6.77) by a centered finite difference scheme. Construct a Picard method for the resulting system of nonlinear algebraic equations.
- d)** Discretize (6.77) by a centered finite difference scheme. Define the system of nonlinear algebraic equations, calculate the Jacobian, and set up Newton's method for solving the system.  
Filename: `nonlin_1D_coeff_linearize`.

### Problem 6.10: Finite differences for the 1D Bratu problem

We address the so-called Bratu problem

$$u'' + \lambda e^u = 0, \quad x \in (0, 1), \quad u(0) = u(1) = 0, \quad (6.78)$$

where  $\lambda$  is a given parameter and  $u$  is a function of  $x$ . This is a widely used model problem for studying numerical methods for nonlinear differential equations. The problem (6.78) has an exact solution

$$u_e(x) = -2 \ln \left( \frac{\cosh((x - \frac{1}{2})\theta/2)}{\cosh(\theta/4)} \right),$$

where  $\theta$  solves

$$\theta = \sqrt{2\lambda} \cosh(\theta/4).$$

There are two solutions of (6.78) for  $0 < \lambda < \lambda_c$  and no solution for  $\lambda > \lambda_c$ . For  $\lambda = \lambda_c$  there is one unique solution. The critical value  $\lambda_c$  solves

$$1 = \sqrt{2\lambda_c} \frac{1}{4} \sinh(\theta(\lambda_c)/4).$$

A numerical value is  $\lambda_c = 3.513830719$ .

- a)** Discretize (6.78) by a centered finite difference method.
- b)** Set up the nonlinear equations  $F_i(u_0, u_1, \dots, u_{N_x}) = 0$  from a). Calculate the associated Jacobian.
- c)** Implement a solver that can compute  $u(x)$  using Newton's method. Plot the error as a function of  $x$  in each iteration.

- d)** Investigate whether Newton's method gives second-order convergence by computing  $\|u_e - u\|/\|u_e - u^-\|^2$  in each iteration, where  $u$  is solution in the current iteration and  $u^-$  is the solution in the previous iteration. Filename: `nonlin_1D_Bratu_fd`.

### Exercise 6.11: Discretize a nonlinear 1D heat conduction PDE by finite differences

We address the 1D heat conduction PDE

$$\varrho c(T)T_t = (k(T)T_x)_x,$$

for  $x \in [0, L]$ , where  $\varrho$  is the density of the solid material,  $c(T)$  is the heat capacity,  $T$  is the temperature, and  $k(T)$  is the heat conduction coefficient.  $T(x, 0) = I(x)$ , and ends are subject to a cooling law:

$$k(T)T_x|_{x=0} = h(T)(T - T_s), \quad -k(T)T_x|_{x=L} = h(T)(T - T_s),$$

where  $h(T)$  is a heat transfer coefficient and  $T_s$  is the given surrounding temperature.

- a)** Discretize this PDE in time using either a Backward Euler or Crank-Nicolson scheme.
- b)** Formulate a Picard iteration method for the time-discrete problem (i.e., an iteration method before discretizing in space).
- c)** Formulate a Newton method for the time-discrete problem in b).
- d)** Discretize the PDE by a finite difference method in space. Derive the matrix and right-hand side of a Picard iteration method applied to the space-time discretized PDE.
- e)** Derive the matrix and right-hand side of a Newton method applied to the discretized PDE in d).

Filename: `nonlin_1D_heat_FD`.

### Exercise 6.12: Differentiate a highly nonlinear term

The operator  $\nabla \cdot (\alpha(u)\nabla u)$  with  $\alpha(u) = |\nabla u|^q$  appears in several physical problems, especially flow of Non-Newtonian fluids. The expression  $|\nabla u|$  is defined as the Euclidean norm of a vector:  $|\nabla u|^2 = \nabla u \cdot \nabla u$ . In a

Newton method one has to carry out the differentiation  $\partial\alpha(u)/\partial c_j$ , for  $u = \sum_k c_k \psi_k$ . Show that

$$\frac{\partial}{\partial u_j} |\nabla u|^q = q |\nabla u|^{q-2} \nabla u \cdot \nabla \psi_j.$$

Filename: `nonlin_differentiate`.

### Exercise 6.13: Crank-Nicolson for a nonlinear 3D diffusion equation

Redo Section 6.5.1 when a Crank-Nicolson scheme is used to discretize the equations in time and the problem is formulated for three spatial dimensions.

**Hint.** Express the Jacobian as  $J_{i,j,k,r,s,t} = \partial F_{i,j,k} / \partial u_{r,s,t}$  and observe, as in the 2D case, that  $J_{i,j,k,r,s,t}$  is very sparse:  $J_{i,j,k,r,s,t} \neq 0$  only for  $r = i \pm 1$ ,  $s = j \pm 1$ , and  $t = k \pm 1$  as well as  $r = i$ ,  $s = j$ , and  $t = k$ .

Filename: `nonlin_heat_FD_CN_2D`.

### Exercise 6.14: Find the sparsity of the Jacobian

Consider a typical nonlinear Laplace term like  $\nabla \cdot \alpha(u) \nabla u$  discretized by centered finite differences. Explain why the Jacobian corresponding to this term has the same sparsity pattern as the matrix associated with the corresponding linear term  $\alpha \nabla^2 u$ .

**Hint.** Set up the unknowns that enter the difference equation at a point  $(i, j)$  in 2D or  $(i, j, k)$  in 3D, and identify the nonzero entries of the Jacobian that can arise from such a type of difference equation.

Filename: `nonlin_sparsity_Jacobian`.

### Problem 6.15: Investigate a 1D problem with a continuation method

Flow of a pseudo-plastic power-law fluid between two flat plates can be modeled by

$$\frac{d}{dx} \left( \mu_0 \left| \frac{du}{dx} \right|^{n-1} \frac{du}{dx} \right) = -\beta, \quad u'(0) = 0, \quad u(H) = 0,$$

where  $\beta > 0$  and  $\mu_0 > 0$  are constants. A target value of  $n$  may be  $n = 0.2$ .

- a)** Formulate a Picard iteration method directly for the differential equation problem.
- b)** Perform a finite difference discretization of the problem in each Picard iteration. Implement a solver that can compute  $u$  on a mesh. Verify that the solver gives an exact solution for  $n = 1$  on a uniform mesh regardless of the cell size.
- c)** Given a sequence of decreasing  $n$  values, solve the problem for each  $n$  using the solution for the previous  $n$  as initial guess for the Picard iteration. This is called a continuation method. Experiment with  $n = (1, 0.6, 0.2)$  and  $n = (1, 0.9, 0.8, \dots, 0.2)$  and make a table of the number of Picard iterations versus  $n$ .
- d)** Derive a Newton method at the differential equation level and discretize the resulting linear equations in each Newton iteration with the finite difference method.
- e)** Investigate if Newton's method has better convergence properties than Picard iteration, both in combination with a continuation method.

## Useful formulas

A

### A.1 Finite difference operator notation

$$u'(t_n) \approx [D_t u]^n = \frac{u^{n+\frac{1}{2}} - u^{n-\frac{1}{2}}}{\Delta t} \quad (\text{A.1})$$

$$u'(t_n) \approx [D_{2t} u]^n = \frac{u^{n+1} - u^{n-1}}{2\Delta t} \quad (\text{A.2})$$

$$u'(t_n) = [D_t^- u]^n = \frac{u^n - u^{n-1}}{\Delta t} \quad (\text{A.3})$$

$$u'(t_n) \approx [D_t^+ u]^n = \frac{u^{n+1} - u^n}{\Delta t} \quad (\text{A.4})$$

$$u'(t_{n+\theta}) = [\bar{D}_t u]^{n+\theta} = \frac{u^{n+1} - u^n}{\Delta t} \quad (\text{A.5})$$

$$u'(t_n) \approx [D_t^{2-} u]^n = \frac{3u^n - 4u^{n-1} + u^{n-2}}{2\Delta t} \quad (\text{A.6})$$

$$u''(t_n) \approx [D_t D_t u]^n = \frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} \quad (\text{A.7})$$

$$u(t_{n+\frac{1}{2}}) \approx [\bar{u}^t]^{n+\frac{1}{2}} = \frac{1}{2}(u^{n+1} + u^n) \quad (\text{A.8})$$

$$u(t_{n+\frac{1}{2}})^2 \approx [\bar{u}^{2,t,g}]^{n+\frac{1}{2}} = u^{n+1} u^n \quad (\text{A.9})$$

$$u(t_{n+\frac{1}{2}}) \approx [\bar{u}^{t,h}]^{n+\frac{1}{2}} = \frac{2}{\frac{1}{u^{n+1}} + \frac{1}{u^n}} \quad (\text{A.10})$$

$$u(t_{n+\theta}) \approx [\bar{u}^{t,\theta}]^{n+\theta} = \theta u^{n+1} + (1 - \theta) u^n, \quad (\text{A.11})$$

$$t_{n+\theta} = \theta t_{n+1} + (1 - \theta) t_{n-1} \quad (\text{A.12})$$

## A.2 Truncation errors of finite difference approximations

$$\begin{aligned} u'_e(t_n) &= [D_t u_e]^n + R^n = \frac{u_e^{n+\frac{1}{2}} - u_e^{n-\frac{1}{2}}}{\Delta t} + R^n, \\ R^n &= -\frac{1}{24} u'''_e(t_n) \Delta t^2 + \mathcal{O}(\Delta t^4) \end{aligned} \quad (\text{A.13})$$

$$\begin{aligned} u'_e(t_n) &= [D_{2t} u_e]^n + R^n = \frac{u_e^{n+1} - u_e^{n-1}}{2\Delta t} + R^n, \\ R^n &= -\frac{1}{6} u'''_e(t_n) \Delta t^2 + \mathcal{O}(\Delta t^4) \end{aligned} \quad (\text{A.14})$$

$$\begin{aligned} u'_e(t_n) &= [D_t^- u_e]^n + R^n = \frac{u_e^n - u_e^{n-1}}{\Delta t} + R^n, \\ R^n &= -\frac{1}{2} u''_e(t_n) \Delta t + \mathcal{O}(\Delta t^2) \end{aligned} \quad (\text{A.15})$$

$$\begin{aligned} u'_e(t_n) &= [D_t^+ u_e]^n + R^n = \frac{u_e^{n+1} - u_e^n}{\Delta t} + R^n, \\ R^n &= -\frac{1}{2} u''_e(t_n) \Delta t + \mathcal{O}(\Delta t^2) \end{aligned} \quad (\text{A.16})$$

$$\begin{aligned} u'_e(t_{n+\theta}) &= [\bar{D}_t u_e]^{n+\theta} + R^{n+\theta} = \frac{u_e^{n+1} - u_e^n}{\Delta t} + R^{n+\theta}, \\ R^{n+\theta} &= -\frac{1}{2}(1-2\theta) u''_e(t_{n+\theta}) \Delta t + \frac{1}{6}((1-\theta)^3 - \theta^3) u'''_e(t_{n+\theta}) \Delta t^2 + \\ &\quad \mathcal{O}(\Delta t^3) \end{aligned} \quad (\text{A.17})$$

$$\begin{aligned} u'_e(t_n) &= [D_t^{2-} u_e]^n + R^n = \frac{3u_e^n - 4u_e^{n-1} + u_e^{n-2}}{2\Delta t} + R^n, \\ R^n &= \frac{1}{3} u'''_e(t_n) \Delta t^2 + \mathcal{O}(\Delta t^3) \end{aligned} \quad (\text{A.18})$$

$$\begin{aligned} u''_e(t_n) &= [D_t D_t u_e]^n + R^n = \frac{u_e^{n+1} - 2u_e^n + u_e^{n-1}}{\Delta t^2} + R^n, \\ R^n &= -\frac{1}{12} u''''_e(t_n) \Delta t^2 + \mathcal{O}(\Delta t^4) \end{aligned} \quad (\text{A.19})$$

$$\begin{aligned} u_e(t_{n+\theta}) &= [\bar{u}_e^{t,\theta}]^{n+\theta} + R^{n+\theta} = \theta u_e^{n+1} + (1-\theta) u_e^n + R^{n+\theta}, \\ R^{n+\theta} &= -\frac{1}{2} u''_e(t_{n+\theta}) \Delta t^2 \theta(1-\theta) + \mathcal{O}(\Delta t^3). \end{aligned} \quad (\text{A.20})$$

### A.3 Finite differences of exponential functions

**Complex exponentials.** Let  $u^n = \exp(i\omega n \Delta t) = e^{i\omega t_n}$ .

$$[D_t D_t u]^n = u^n \frac{2}{\Delta t} (\cos \omega \Delta t - 1) = -\frac{4}{\Delta t} \sin^2 \left( \frac{\omega \Delta t}{2} \right), \quad (\text{A.21})$$

$$[D_t^+ u]^n = u^n \frac{1}{\Delta t} (\exp(i\omega \Delta t) - 1), \quad (\text{A.22})$$

$$[D_t^- u]^n = u^n \frac{1}{\Delta t} (1 - \exp(-i\omega \Delta t)), \quad (\text{A.23})$$

$$[D_t u]^n = u^n \frac{2}{\Delta t} i \sin \left( \frac{\omega \Delta t}{2} \right), \quad (\text{A.24})$$

$$[D_{2t} u]^n = u^n \frac{1}{\Delta t} i \sin(\omega \Delta t). \quad (\text{A.25})$$

**Real exponentials.** Let  $u^n = \exp(\omega n \Delta t) = e^{\omega t_n}$ .

$$[D_t D_t u]^n = u^n \frac{2}{\Delta t} (\cos \omega \Delta t - 1) = -\frac{4}{\Delta t} \sin^2 \left( \frac{\omega \Delta t}{2} \right), \quad (\text{A.26})$$

$$[D_t^+ u]^n = u^n \frac{1}{\Delta t} (\exp(i\omega \Delta t) - 1), \quad (\text{A.27})$$

$$[D_t^- u]^n = u^n \frac{1}{\Delta t} (1 - \exp(-i\omega \Delta t)), \quad (\text{A.28})$$

$$[D_t u]^n = u^n \frac{2}{\Delta t} i \sin \left( \frac{\omega \Delta t}{2} \right), \quad (\text{A.29})$$

$$[D_{2t} u]^n = u^n \frac{1}{\Delta t} i \sin(\omega \Delta t). \quad (\text{A.30})$$

#### A.4 Finite differences of $t^n$

The following results are useful when checking if a polynomial term in a solution fulfills the discrete equation for the numerical method.

$$[D_t^+ t]^n = 1, \quad (\text{A.31})$$

$$[D_t^- t]^n = 1, \quad (\text{A.32})$$

$$[D_t t]^n = 1, \quad (\text{A.33})$$

$$[D_{2t} t]^n = 1, \quad (\text{A.34})$$

$$[D_t D_t t]^n = 0. \quad (\text{A.35})$$

The next formulas concern the action of difference operators on a  $t^2$  term.

$$[D_t^+ t^2]^n = (2n + 1)\Delta t, \quad (\text{A.36})$$

$$[D_t^- t^2]^n = (2n - 1)\Delta t, \quad (\text{A.37})$$

$$[D_t t^2]^n = 2n\Delta t, \quad (\text{A.38})$$

$$[D_{2t} t^2]^n = 2n\Delta t, \quad (\text{A.39})$$

$$[D_t D_t t^2]^n = 2, \quad (\text{A.40})$$

Finally, we present formulas for a  $t^3$  term: **These must be controlled against lib.py**. Use  $t_n$  instead of  $n\Delta t$ ??

$$[D_t^+ t^3]^n = 3(n\Delta t)^2 + 3n\Delta t^2 + \Delta t^2, \quad (\text{A.41})$$

$$[D_t^- t^3]^n = 3(n\Delta t)^2 - 3n\Delta t^2 + \Delta t^2, \quad (\text{A.42})$$

$$[D_t t^3]^n = 3(n\Delta t)^2 + \frac{1}{4}\Delta t^2, \quad (\text{A.43})$$

$$[D_{2t} t^3]^n = 3(n\Delta t)^2 + \Delta t^2, \quad (\text{A.44})$$

$$[D_t D_t t^3]^n = 6n\Delta t, \quad (\text{A.45})$$

### A.4.1 Software

Application of finite difference operators to polynomials and exponential functions, resulting in the formulas above, can easily be computed by some `sympy` code:

```
from sympy import *
t, dt, n, w = symbols('t dt n w', real=True)

# Finite difference operators

def D_t_forward(u):
    return (u(t + dt) - u(t))/dt

def D_t_backward(u):
    return (u(t) - u(t-dt))/dt

def D_t_centered(u):
    return (u(t + dt/2) - u(t-dt/2))/dt

def D_2t_centered(u):
    return (u(t + dt) - u(t-dt))/(2*dt)

def D_tD_t(u):
```

```

        return (u(t + dt) - 2*u(t) + u(t-dt))/(dt**2)

op_list = [D_t_forward, D_t_backward,
           D_t_centered, D_2t_centered, D_t_D_t]

def ft1(t):
    return t

def ft2(t):
    return t**2

def ft3(t):
    return t**3

def f_expiwt(t):
    return exp(I*w*t)

def f_expwt(t):
    return exp(w*t)

func_list = [ft1, ft2, ft3, f_expiwt, f_expwt]

```

To see the results, one can now make a simple loop like

```

for func in func_list:
    for op in op_list:
        f = func
        e = op(f)
        e = simplify(expand(e))
        print e
        if func in [f_expiwt, f_expwt]:
            e = e/f(t)
        e = e.subs(t, n*dt)
        print expand(e)
        print factor(simplify(expand(e)))

```

## Summary

Truncation error analysis provides a widely applicable framework for analyzing the accuracy of finite difference schemes. This type of analysis can also be used for finite element and finite volume methods if the discrete equations are written in finite difference form. The result of the analysis is an asymptotic estimate of the error in the scheme on the form  $Ch^r$ , where  $h$  is a discretization parameter ( $\Delta t$ ,  $\Delta x$ , etc.),  $r$  is a number, known as the convergence rate, and  $C$  is a constant, typically dependent on the derivatives of the exact solution.

Knowing  $r$  gives understanding of the accuracy of the scheme. But maybe even more important, a powerful verification method for computer codes is to check that the empirically observed convergence rates in experiments coincide with the theoretical value of  $r$  found from truncation error analysis.

The analysis can be carried out by hand, by symbolic software, and also numerically. All three methods will be illustrated. From examining the symbolic expressions of the truncation error we can add correction terms to the differential equations in order to increase the numerical accuracy.

In general, the term truncation error refers to the discrepancy that arises from performing a finite number of steps to approximate a process with infinitely many steps. The term is used in a number of contexts,

including truncation of infinite series, finite precision arithmetic, finite differences, and differential equations. We shall be concerned with computing truncation errors arising in finite difference formulas and in finite difference discretizations of differential equations.

## B.1 Overview of truncation error analysis

### B.1.1 Abstract problem setting

Consider an abstract differential equation

$$\mathcal{L}(u) = 0,$$

where  $\mathcal{L}(u)$  is some formula involving the unknown  $u$  and its derivatives. One example is  $\mathcal{L}(u) = u'(t) + a(t)u(t) - b(t)$ , where  $a$  and  $b$  are constants or functions of time. We can discretize the differential equation and obtain a corresponding discrete model, here written as

$$\mathcal{L}_\Delta(u) = 0.$$

The solution  $u$  of this equation is the *numerical solution*. To distinguish the numerical solution from the exact solution of the differential equation problem, we denote the latter by  $u_e$  and write the differential equation and its discrete counterpart as

$$\begin{aligned}\mathcal{L}(u_e) &= 0, \\ \mathcal{L}_\Delta(u) &= 0.\end{aligned}$$

Initial and/or boundary conditions can usually be left out of the truncation error analysis and are omitted in the following.

The numerical solution  $u$  is, in a finite difference method, computed at a collection of mesh points. The discrete equations represented by the abstract equation  $\mathcal{L}_\Delta(u) = 0$  are usually algebraic equations involving  $u$  at some neighboring mesh points.

### B.1.2 Error measures

A key issue is how accurate the numerical solution is. The ultimate way of addressing this issue would be to compute the error  $u_e - u$  at the mesh

points. This is usually extremely demanding. In very simplified problem settings we may, however, manage to derive formulas for the numerical solution  $u$ , and therefore closed form expressions for the error  $u_e - u$ . Such special cases can provide considerable insight regarding accuracy and stability, but the results are established for special problems.

The error  $u_e - u$  can be computed empirically in special cases where we know  $u_e$ . Such cases can be constructed by the method of manufactured solutions, where we choose some exact solution  $u_e = v$  and fit a source term  $f$  in the governing differential equation  $\mathcal{L}(u_e) = f$  such that  $u_e = v$  is a solution (i.e.,  $f = \mathcal{L}(v)$ ). Assuming an error model of the form  $Ch^r$ , where  $h$  is the discretization parameter, such as  $\Delta t$  or  $\Delta x$ , one can estimate the convergence rate  $r$ . This is a widely applicable procedure, but the validity of the results is, strictly speaking, tied to the chosen test problems.

Another error measure arises by asking to what extent the exact solution  $u_e$  fits the discrete equations. Clearly,  $u_e$  is in general not a solution of  $\mathcal{L}_\Delta(u) = 0$ , but we can define the residual

$$R = \mathcal{L}_\Delta(u_e),$$

and investigate how close  $R$  is to zero. A small  $R$  means intuitively that the discrete equations are close to the differential equation, and then we are tempted to think that  $u^n$  must also be close to  $u_e(t_n)$ .

The residual  $R$  is known as the truncation error of the finite difference scheme  $\mathcal{L}_\Delta(u) = 0$ . It appears that the truncation error is relatively straightforward to compute by hand or symbolic software *without specializing the differential equation and the discrete model to a special case*. The resulting  $R$  is found as a power series in the discretization parameters. The leading-order terms in the series provide an asymptotic measure of the accuracy of the numerical solution method (as the discretization parameters tend to zero). An advantage of truncation error analysis, compared to empirical estimation of convergence rates, or detailed analysis of a special problem with a mathematical expression for the numerical solution, is that the truncation error analysis reveals the accuracy of the various building blocks in the numerical method and how each building block impacts the overall accuracy. The analysis can therefore be used to detect building blocks with lower accuracy than the others.

Knowing the truncation error or other error measures is important for verification of programs by empirically establishing convergence rates.

The forthcoming text will provide many examples on how to compute truncation errors for finite difference discretizations of ODEs and PDEs.

## B.2 Truncation errors in finite difference formulas

The accuracy of a finite difference formula is a fundamental issue when discretizing differential equations. We shall first go through a particular example in detail and thereafter list the truncation error in the most common finite difference approximation formulas.

### B.2.1 Example: The backward difference for $u'(t)$

Consider a backward finite difference approximation of the first-order derivative  $u'$ :

$$[D_t^- u]^n = \frac{u^n - u^{n-1}}{\Delta t} \approx u'(t_n). \quad (\text{B.1})$$

Here,  $u^n$  means the value of some function  $u(t)$  at a point  $t_n$ , and  $[D_t^- u]^n$  is the *discrete derivative* of  $u(t)$  at  $t = t_n$ . The discrete derivative computed by a finite difference is not exactly equal to the derivative  $u'(t_n)$ . The error in the approximation is

$$R^n = [D_t^- u]^n - u'(t_n). \quad (\text{B.2})$$

The common way of calculating  $R^n$  is to

1. expand  $u(t)$  in a Taylor series around the point where the derivative is evaluated, here  $t_n$ ,
2. insert this Taylor series in (B.2), and
3. collect terms that cancel and simplify the expression.

The result is an expression for  $R^n$  in terms of a power series in  $\Delta t$ . The error  $R^n$  is commonly referred to as the *truncation error* of the finite difference formula.

The Taylor series formula often found in calculus books takes the form

$$f(x + h) = \sum_{i=0}^{\infty} \frac{1}{i!} \frac{d^i f}{dx^i}(x) h^i.$$

In our application, we expand the Taylor series around the point where the finite difference formula approximates the derivative. The Taylor series of  $u^n$  at  $t_n$  is simply  $u(t_n)$ , while the Taylor series of  $u^{n-1}$  at  $t_n$  must employ the general formula,

$$\begin{aligned} u(t_{n-1}) &= u(t - \Delta t) = \sum_{i=0}^{\infty} \frac{1}{i!} \frac{d^i u}{dt^i}(t_n)(-\Delta t)^i \\ &= u(t_n) - u'(t_n)\Delta t + \frac{1}{2}u''(t_n)\Delta t^2 + \mathcal{O}(\Delta t^3), \end{aligned}$$

where  $\mathcal{O}(\Delta t^3)$  means a power-series in  $\Delta t$  where the lowest power is  $\Delta t^3$ . We assume that  $\Delta t$  is small such that  $\Delta t^p \gg \Delta t^q$  if  $p$  is smaller than  $q$ . The details of higher-order terms in  $\Delta t$  are therefore not of much interest. Inserting the Taylor series above in the left-hand side of (B.2) gives rise to some algebra:

$$\begin{aligned} [D_t^- u]^n - u'(t_n) &= \frac{u(t_n) - u(t_{n-1})}{\Delta t} - u'(t_n) \\ &= \frac{u(t_n) - (u(t_n) - u'(t_n)\Delta t + \frac{1}{2}u''(t_n)\Delta t^2 + \mathcal{O}(\Delta t^3))}{\Delta t} - u'(t_n) \\ &= -\frac{1}{2}u''(t_n)\Delta t + \mathcal{O}(\Delta t^2)), \end{aligned}$$

which is, according to (B.2), the truncation error:

$$R^n = -\frac{1}{2}u''(t_n)\Delta t + \mathcal{O}(\Delta t^2). \quad (\text{B.3})$$

The dominating term for small  $\Delta t$  is  $-\frac{1}{2}u''(t_n)\Delta t$ , which is proportional to  $\Delta t$ , and we say that the truncation error is of *first order* in  $\Delta t$ .

### B.2.2 Example: The forward difference for $u'(t)$

We can analyze the approximation error in the forward difference

$$u'(t_n) \approx [D_t^+ u]^n = \frac{u^{n+1} - u^n}{\Delta t},$$

by writing

$$R^n = [D_t^+ u]^n - u'(t_n),$$

and expanding  $u^{n+1}$  in a Taylor series around  $t_n$ ,

$$u(t_{n+1}) = u(t_n) + u'(t_n)\Delta t + \frac{1}{2}u''(t_n)\Delta t^2 + \mathcal{O}(\Delta t^3).$$

The result becomes

$$R = \frac{1}{2}u''(t_n)\Delta t + \mathcal{O}(\Delta t^2),$$

showing that also the forward difference is of first order.

### B.2.3 Example: The central difference for $u'(t)$

For the central difference approximation,

$$u'(t_n) \approx [D_t u]^n, \quad [D_t u]^n = \frac{u^{n+\frac{1}{2}} - u^{n-\frac{1}{2}}}{\Delta t},$$

we write

$$R^n = [D_t u]^n - u'(t_n),$$

and expand  $u(t_{n+\frac{1}{2}})$  and  $u(t_{n-1/2})$  in Taylor series around the point  $t_n$  where the derivative is evaluated. We have

$$\begin{aligned} u(t_{n+\frac{1}{2}}) &= u(t_n) + u'(t_n)\frac{1}{2}\Delta t + \frac{1}{2}u''(t_n)(\frac{1}{2}\Delta t)^2 + \\ &\quad \frac{1}{6}u'''(t_n)(\frac{1}{2}\Delta t)^3 + \frac{1}{24}u''''(t_n)(\frac{1}{2}\Delta t)^4 + \\ &\quad \frac{1}{120}u'''''(t_n)(\frac{1}{2}\Delta t)^5 + \mathcal{O}(\Delta t^6), \\ u(t_{n-1/2}) &= u(t_n) - u'(t_n)\frac{1}{2}\Delta t + \frac{1}{2}u''(t_n)(\frac{1}{2}\Delta t)^2 - \\ &\quad \frac{1}{6}u'''(t_n)(\frac{1}{2}\Delta t)^3 + \frac{1}{24}u''''(t_n)(\frac{1}{2}\Delta t)^4 - \\ &\quad \frac{1}{120}u'''''(t_n)(\frac{1}{2}\Delta t)^5 + \mathcal{O}(\Delta t^6). \end{aligned}$$

Now,

$$u(t_{n+\frac{1}{2}}) - u(t_{n-1/2}) = u'(t_n)\Delta t + \frac{1}{24}u'''(t_n)\Delta t^3 + \frac{1}{960}u''''(t_n)\Delta t^5 + \mathcal{O}(\Delta t^7).$$

By collecting terms in  $[D_t u]^n - u'(t_n)$  we find the truncation error to be

$$R^n = \frac{1}{24}u'''(t_n)\Delta t^2 + \mathcal{O}(\Delta t^4), \quad (\text{B.4})$$

with only even powers of  $\Delta t$ . Since  $R \sim \Delta t^2$  we say the centered difference is of *second order* in  $\Delta t$ .

### B.2.4 Overview of leading-order error terms in finite difference formulas

Here we list the leading-order terms of the truncation errors associated with several common finite difference formulas for the first and second derivatives.

$$[D_t u]^n = \frac{u^{n+\frac{1}{2}} - u^{n-\frac{1}{2}}}{\Delta t} = u'(t_n) + R^n, \quad (\text{B.5})$$

$$R^n = \frac{1}{24} u'''(t_n) \Delta t^2 + \mathcal{O}(\Delta t^4) \quad (\text{B.6})$$

$$[D_{2t} u]^n = \frac{u^{n+1} - u^{n-1}}{2\Delta t} = u'(t_n) + R^n, \quad (\text{B.7})$$

$$R^n = \frac{1}{6} u'''(t_n) \Delta t^2 + \mathcal{O}(\Delta t^4) \quad (\text{B.8})$$

$$[D_t^- u]^n = \frac{u^n - u^{n-1}}{\Delta t} = u'(t_n) + R^n, \quad (\text{B.9})$$

$$R^n = -\frac{1}{2} u''(t_n) \Delta t + \mathcal{O}(\Delta t^2) \quad (\text{B.10})$$

$$[D_t^+ u]^n = \frac{u^{n+1} - u^n}{\Delta t} = u'(t_n) + R^n, \quad (\text{B.11})$$

$$R^n = \frac{1}{2} u''(t_n) \Delta t + \mathcal{O}(\Delta t^2) \quad (\text{B.12})$$

$$[D_t u]^{n+\theta} = \frac{u^{n+1} - u^n}{\Delta t} = u'(t_{n+\theta}) + R^{n+\theta}, \quad (\text{B.13})$$

$$R^{n+\theta} = \frac{1}{2}(1-2\theta)u''(t_{n+\theta})\Delta t - \frac{1}{6}((1-\theta)^3 - \theta^3)u'''(t_{n+\theta})\Delta t^2 + \mathcal{O}(\Delta t^3) \quad (\text{B.14})$$

$$[D_t^{2-} u]^n = \frac{3u^n - 4u^{n-1} + u^{n-2}}{2\Delta t} = u'(t_n) + R^n, \quad (\text{B.15})$$

$$R^n = -\frac{1}{3} u'''(t_n) \Delta t^2 + \mathcal{O}(\Delta t^3) \quad (\text{B.16})$$

$$[D_t D_t u]^n = \frac{u^{n+1} - 2u^n + u^{n-1}}{\Delta t^2} = u''(t_n) + R^n, \quad (\text{B.17})$$

$$R^n = \frac{1}{12} u'''(t_n) \Delta t^2 + \mathcal{O}(\Delta t^4) \quad (\text{B.18})$$

It will also be convenient to have the truncation errors for various means or averages. The weighted arithmetic mean leads to

$$[\bar{u}^{t,\theta}]^{n+\theta} = \theta u^{n+1} + (1 - \theta)u^n = u(t_{n+\theta}) + R^{n+\theta}, \quad (\text{B.19})$$

$$R^{n+\theta} = \frac{1}{2}u''(t_{n+\theta})\Delta t^2\theta(1 - \theta) + \mathcal{O}(\Delta t^3). \quad (\text{B.20})$$

The standard arithmetic mean follows from this formula when  $\theta = 1/2$ . Expressed at point  $t_n$  we get

$$[\bar{u}^t]^n = \frac{1}{2}(u^{n-\frac{1}{2}} + u^{n+\frac{1}{2}}) = u(t_n) + R^n, \quad (\text{B.21})$$

$$R^n = \frac{1}{8}u''(t_n)\Delta t^2 + \frac{1}{384}u'''(t_n)\Delta t^4 + \mathcal{O}(\Delta t^6). \quad (\text{B.22})$$

The geometric mean also has an error  $\mathcal{O}(\Delta t^2)$ :

$$[\bar{u}^{2^t,g}]^n = u^{n-\frac{1}{2}}u^{n+\frac{1}{2}} = (u^n)^2 + R^n, \quad (\text{B.23})$$

$$R^n = -\frac{1}{4}u'(t_n)^2\Delta t^2 + \frac{1}{4}u(t_n)u''(t_n)\Delta t^2 + \mathcal{O}(\Delta t^4). \quad (\text{B.24})$$

The harmonic mean is also second-order accurate:

$$[\bar{u}^{t,h}]^n = u^n = \frac{2}{\frac{1}{u^{n-\frac{1}{2}}} + \frac{1}{u^{n+\frac{1}{2}}}} + R^{n+\frac{1}{2}}, \quad (\text{B.25})$$

$$R^n = -\frac{u'(t_n)^2}{4u(t_n)}\Delta t^2 + \frac{1}{8}u''(t_n)\Delta t^2. \quad (\text{B.26})$$

### B.2.5 Software for computing truncation errors

We can use `sympy` to aid calculations with Taylor series. The derivatives can be defined as symbols, say `D3f` for the 3rd derivative of some function  $f$ . A truncated Taylor series can then be written as `f + D1f*h + D2f*h**2/2`. The following class takes some symbol `f` for the function in question and makes a list of symbols for the derivatives. The `__call__` method computes the symbolic form of the series truncated at `num_terms` terms.

```

import sympy as sym

class TaylorSeries:
    """Class for symbolic Taylor series."""
    def __init__(self, f, num_terms=4):
        self.f = f
        self.N = num_terms
        # Introduce symbols for the derivatives
        self.df = [f]
        for i in range(1, self.N+1):
            self.df.append(sym.Symbol('D%d' % (i, f.name)))

    def __call__(self, h):
        """Return the truncated Taylor series at x+h."""
        terms = self.f
        for i in range(1, self.N+1):
            terms += sym.Rational(1, sym.factorial(i))*self.df[i]*h**i
        return terms

```

We may, for example, use this class to compute the truncation error of the Forward Euler finite difference formula:

```

>>> from truncation_errors import TaylorSeries
>>> from sympy import *
>>> u, dt = symbols('u dt')
>>> u_Taylor = TaylorSeries(u, 4)
>>> u_Taylor(dt)
D1u*dt + D2u*dt**2/2 + D3u*dt**3/6 + D4u*dt**4/24 + u
>>> FE = (u_Taylor(dt) - u)/dt
>>> FE
(D1u*dt + D2u*dt**2/2 + D3u*dt**3/6 + D4u*dt**4/24)/dt
>>> simplify(FE)
D1u + D2u*dt/2 + D3u*dt**2/6 + D4u*dt**3/24

```

The truncation error consists of the terms after the first one ( $u'$ ).

The module file `trunc/truncation_errors.py` contains another class `DiffOp` with symbolic expressions for most of the truncation errors listed in the previous section. For example:

```

>>> from truncation_errors import DiffOp
>>> from sympy import *
>>> u = Symbol('u')
>>> diffop = DiffOp(u, independent_variable='t')
>>> diffop['geometric_mean']
-D1u**2*dt**2/4 - D1u*D3u*dt**4/48 + D2u**2*dt**4/64 + ...
>>> diffop['Dtm']
D1u + D2u*dt/2 + D3u*dt**2/6 + D4u*dt**3/24
>>> >>> diffop.operator_names()
['geometric_mean', 'harmonic_mean', 'Dtm', 'D2t', 'DtDt',
 'weighted_arithmetic_mean', 'Dtp', 'Dt']

```

The indexing of `diffop` applies names that correspond to the operators:  $D_{tp}$  for  $D_t^+$ ,  $D_{tm}$  for  $D_t^-$ ,  $Dt$  for  $D_t$ ,  $D2t$  for  $D_{2t}$ ,  $DtDt$  for  $D_tD_t$ .

### B.3 Truncation errors in exponential decay ODE

We shall now compute the truncation error of a finite difference scheme for a differential equation. Our first problem involves the following the linear ODE modeling exponential decay,

$$u'(t) = -au(t). \quad (\text{B.27})$$

#### B.3.1 Truncation error of the Forward Euler scheme

We begin with the Forward Euler scheme for discretizing (B.27):

$$[D_t^+ u = -au]^n. \quad (\text{B.28})$$

The idea behind the truncation error computation is to insert the exact solution  $u_e$  of the differential equation problem (B.27) in the discrete equations (B.28) and find the residual that arises because  $u_e$  does not solve the discrete equations. Instead,  $u_e$  solves the discrete equations with a residual  $R^n$ :

$$[D_t^+ u_e + au_e = R]^n. \quad (\text{B.29})$$

From (B.11)-(B.12) it follows that

$$[D_t^+ u_e]^n = u'_e(t_n) + \frac{1}{2}u''_e(t_n)\Delta t + \mathcal{O}(\Delta t^2),$$

which inserted in (B.29) results in

$$u'_e(t_n) + \frac{1}{2}u''_e(t_n)\Delta t + \mathcal{O}(\Delta t^2) + au_e(t_n) = R^n.$$

Now,  $u'_e(t_n) + au_e^n = 0$  since  $u_e$  solves the differential equation. The remaining terms constitute the residual:

$$R^n = \frac{1}{2}u''_e(t_n)\Delta t + \mathcal{O}(\Delta t^2). \quad (\text{B.30})$$

This is the truncation error  $R^n$  of the Forward Euler scheme.

Because  $R^n$  is proportional to  $\Delta t$ , we say that the Forward Euler scheme is of first order in  $\Delta t$ . However, the truncation error is just one error measure, and it is not equal to the true error  $u_e^n - u^n$ . For this simple model problem we can compute a range of different error measures

for the Forward Euler scheme, including the true error  $u_e^n - u^n$ , and all of them have dominating terms proportional to  $\Delta t$ .

### B.3.2 Truncation error of the Crank-Nicolson scheme

For the Crank-Nicolson scheme,

$$[D_t u = -au]^{n+\frac{1}{2}}, \quad (\text{B.31})$$

we compute the truncation error by inserting the exact solution of the ODE and adding a residual  $R$ ,

$$[D_t u_e + a\bar{u}_e^t = R]^{n+\frac{1}{2}}. \quad (\text{B.32})$$

The term  $[D_t u_e]^{n+\frac{1}{2}}$  is easily computed from (B.5)-(B.6) by replacing  $n$  with  $n + \frac{1}{2}$  in the formula,

$$[D_t u_e]^{n+\frac{1}{2}} = u'_e(t_{n+\frac{1}{2}}) + \frac{1}{24} u'''_e(t_{n+\frac{1}{2}}) \Delta t^2 + \mathcal{O}(\Delta t^4).$$

The arithmetic mean is related to  $u(t_{n+\frac{1}{2}})$  by (B.21)-(B.22) so

$$[a\bar{u}_e^t]^{n+\frac{1}{2}} = u_e(t_{n+\frac{1}{2}}) + \frac{1}{8} u''_e(t_n) \Delta t^2 + \mathcal{O}(\Delta t^4).$$

Inserting these expressions in (B.32) and observing that  $u'_e(t_{n+\frac{1}{2}}) + a u_e^{n+\frac{1}{2}} = 0$ , because  $u_e(t)$  solves the ODE  $u'(t) = -au(t)$  at any point  $t$ , we find that

$$R^{n+\frac{1}{2}} = \left( \frac{1}{24} u'''_e(t_{n+\frac{1}{2}}) + \frac{1}{8} u''_e(t_n) \right) \Delta t^2 + \mathcal{O}(\Delta t^4) \quad (\text{B.33})$$

Here, the truncation error is of second order because the leading term in  $R$  is proportional to  $\Delta t^2$ .

At this point it is wise to redo some of the computations above to establish the truncation error of the Backward Euler scheme, see Exercise B.7.

### B.3.3 Truncation error of the $\theta$ -rule

We may also compute the truncation error of the  $\theta$ -rule,

$$[\bar{D}_t u = -a\bar{u}^{t,\theta}]^{n+\theta}.$$

Our computational task is to find  $R^{n+\theta}$  in

$$[\bar{D}_t u_e + a \bar{u}_e^{t,\theta} = R]^{n+\theta}.$$

From (B.13)-(B.14) and (B.19)-(B.20) we get expressions for the terms with  $u_e$ . Using that  $u'_e(t_{n+\theta}) + au_e(t_{n+\theta}) = 0$ , we end up with

$$\begin{aligned} R^{n+\theta} = & (\frac{1}{2} - \theta)u''_e(t_{n+\theta})\Delta t + \frac{1}{2}\theta(1-\theta)u''_e(t_{n+\theta})\Delta t^2 + \\ & \frac{1}{2}(\theta^2 - \theta + 3)u'''_e(t_{n+\theta})\Delta t^2 + \mathcal{O}(\Delta t^3) \end{aligned} \quad (\text{B.34})$$

For  $\theta = 1/2$  the first-order term vanishes and the scheme is of second order, while for  $\theta \neq 1/2$  we only have a first-order scheme.

### B.3.4 Using symbolic software

The previously mentioned `truncation_error` module can be used to automate the Taylor series expansions and the process of collecting terms. Here is an example on possible use:

```
from truncation_error import DiffOp
from sympy import *

def decay():
    u, a = symbols('u a')
    diffop = DiffOp(u, independent_variable='t',
                    num_terms_Taylor_series=3)
    D1u = diffop.D(1)    # symbol for du/dt
    ODE = D1u + a*u      # define ODE

    # Define schemes
    FE = diffop['Dtp'] + a*u
    CN = diffop['Dt'] + a*u
    BE = diffop['Dtm'] + a*u
    theta = diffop['barDt'] + a*diffop['weighted_arithmetic_mean']
    theta = sm.simplify(sm.expand(theta))
    # Residuals (truncation errors)
    R = {'FE': FE-ODE, 'BE': BE-ODE, 'CN': CN-ODE,
          'theta': theta-ODE}
    return R
```

The returned dictionary becomes

```
decay: {
    'BE': D2u*dt/2 + D3u*dt**2/6,
    'FE': -D2u*dt/2 + D3u*dt**2/6,
```

```

'CN': D3u*dt**2/24,
'theta': -D2u*a*dt**2*theta**2/2 + D2u*a*dt**2*theta/2 -
          D2u*dt*theta + D2u*dt/2 + D3u*a*dt**3*theta**3/3 -
          D3u*a*dt**3*theta**2/2 + D3u*a*dt**3*theta/6 +
          D3u*dt**2*theta**2/2 - D3u*dt**2*theta/2 + D3u*dt**2/6,
}

```

The results are in correspondence with our hand-derived expressions.

### B.3.5 Empirical verification of the truncation error

The task of this section is to demonstrate how we can compute the truncation error  $R$  numerically. For example, the truncation error of the Forward Euler scheme applied to the decay ODE  $u' = -ua$  is

$$R^n = [D_t^+ u_e + au_e]^n. \quad (\text{B.35})$$

If we happen to know the exact solution  $u_e(t)$ , we can easily evaluate  $R^n$  from the above formula.

To estimate how  $R$  varies with the discretization parameter  $\Delta t$ , which has been our focus in the previous mathematical derivations, we first make the assumption that  $R = C\Delta t^r$  for appropriate constants  $C$  and  $r$  and small enough  $\Delta t$ . The rate  $r$  can be estimated from a series of experiments where  $\Delta t$  is varied. Suppose we have  $m$  experiments  $(\Delta t_i, R_i)$ ,  $i = 0, \dots, m-1$ . For two consecutive experiments  $(\Delta t_{i-1}, R_{i-1})$  and  $(\Delta t_i, R_i)$ , a corresponding  $r_{i-1}$  can be estimated by

$$r_{i-1} = \frac{\ln(R_{i-1}/R_i)}{\ln(\Delta t_{i-1}/\Delta t_i)}, \quad (\text{B.36})$$

for  $i = 1, \dots, m-1$ . Note that the truncation error  $R_i$  varies through the mesh, so (B.36) is to be applied pointwise. A complicating issue is that  $R_i$  and  $R_{i-1}$  refer to different meshes. Pointwise comparisons of the truncation error at a certain point in all meshes therefore requires any computed  $R$  to be restricted to the *coarsest mesh* and that all finer meshes contain all the points in the coarsest mesh. Suppose we have  $N_0$  intervals in the coarsest mesh. Inserting a superscript  $n$  in (B.36), where  $n$  counts mesh points in the coarsest mesh,  $n = 0, \dots, N_0$ , leads to the formula

$$r_{i-1}^n = \frac{\ln(R_{i-1}^n/R_i^n)}{\ln(\Delta t_{i-1}/\Delta t_i)}. \quad (\text{B.37})$$

Experiments are most conveniently defined by  $N_0$  and a number of refinements  $m$ . Suppose each mesh have twice as many cells  $N_i$  as the previous one:

$$N_i = 2^i N_0, \quad \Delta t_i = T N_i^{-1},$$

where  $[0, T]$  is the total time interval for the computations. Suppose the computed  $R_i$  values on the mesh with  $N_i$  intervals are stored in an array  $\mathbf{R[i]}$  ( $\mathbf{R}$  being a list of arrays, one for each mesh). Restricting this  $R_i$  function to the coarsest mesh means extracting every  $N_i/N_0$  point and is done as follows:

```
stride = N[i]/N_0
R[i] = R[i][::stride]
```

The quantity  $\mathbf{R[i][n]}$  now corresponds to  $R_i^n$ .

In addition to estimating  $r$  for the pointwise values of  $R = C\Delta t^r$ , we may also consider an integrated quantity on mesh  $i$ ,

$$R_{I,i} = \left( \Delta t_i \sum_{n=0}^{N_i} (R_i^n)^2 \right)^{\frac{1}{2}} \approx \int_0^T R_i(t) dt. \quad (\text{B.38})$$

The sequence  $R_{I,i}$ ,  $i = 0, \dots, m - 1$ , is also expected to behave as  $C\Delta t^r$ , with the same  $r$  as for the pointwise quantity  $R$ , as  $\Delta t \rightarrow 0$ .

The function below computes the  $R_i$  and  $R_{I,i}$  quantities, plots them and compares with the theoretically derived truncation error ( $\mathbf{R_a}$ ) if available.

```
import numpy as np
import scitools.std as plt

def estimate(truncation_error, T, N_0, m, makeplot=True):
    """
    Compute the truncation error in a problem with one independent
    variable, using m meshes, and estimate the convergence
    rate of the truncation error.

    The user-supplied function truncation_error(dt, N) computes
    the truncation error on a uniform mesh with N intervals of
    length dt:::

    R, t, R_a = truncation_error(dt, N)

    where R holds the truncation error at points in the array t,
    and R_a are the corresponding theoretical truncation error
    values (None if not available).

    The truncation_error function is run on a series of meshes
    with 2**i*N_0 intervals, i=0,1,...,m-1.
    """

    # Create a vector of time points
    t = np.linspace(0, T, 1000)

    # Create a list of meshes
    N = [2**i * N_0 for i in range(m)]
    R = []
    R_a = []

    for i in range(m):
        dt = T / N[i]
        R.append(truncation_error(dt, N[i]))
        R_a.append(truncation_error(dt, N[i], True))

    # Plot the results
    if makeplot:
        plt.figure()
        plt.loglog(N, R, 'o')
        plt.loglog(N, R_a, 'x')
        plt.title('Truncation error vs. number of intervals')
        plt.xlabel('Number of intervals (N)')
        plt.ylabel('Truncation error')
        plt.legend(['Numerical', 'Theoretical'])

    return R, t, R_a
```

```

The values of R and R_a are restricted to the coarsest mesh.
and based on these data, the convergence rate of R (pointwise)
and time-integrated R can be estimated empirically.
"""
N = [2**i*N_0 for i in range(m)]

R_I = np.zeros(m) # time-integrated R values on various meshes
R = [None]*m      # time series of R restricted to coarsest mesh
R_a = [None]*m    # time series of R_a restricted to coarsest mesh
dt = np.zeros(m)
legends_R = []; legends_R_a = [] # all legends of curves

for i in range(m):
    dt[i] = T/float(N[i])
    R[i], t, R_a[i] = truncation_error(dt[i], N[i])

    R_I[i] = np.sqrt(dt[i]*np.sum(R[i]**2))

    if i == 0:
        t_coarse = t           # the coarsest mesh

    stride = N[i]/N_0
    R[i] = R[i][::stride]      # restrict to coarsest mesh
    R_a[i] = R_a[i][::stride]

    if makeplot:
        plt.figure(1)
        plt.plot(t_coarse, R[i], log='y')
        legends_R.append('N=%d' % N[i])
        plt.hold('on')

        plt.figure(2)
        plt.plot(t_coarse, R_a[i] - R[i], log='y')
        plt.hold('on')
        legends_R_a.append('N=%d' % N[i])

if makeplot:
    plt.figure(1)
    plt.xlabel('time')
    plt.ylabel('pointwise truncation error')
    plt.legend(legends_R)
    plt.savefig('R_series.png')
    plt.savefig('R_series.pdf')
    plt.figure(2)
    plt.xlabel('time')
    plt.ylabel('pointwise error in estimated truncation error')
    plt.legend(legends_R_a)
    plt.savefig('R_error.png')
    plt.savefig('R_error.pdf')

# Convergence rates
r_R_I = convergence_rates(dt, R_I)
print 'R integrated in time; r:',
print ' '.join(['%.1f' % r for r in r_R_I])

```

```
R = np.array(R) # two-dim. numpy array
r_R = [convergence_rates(dt, R[:,n])[-1]
       for n in range(len(t_coarse))]
```

The first `makeplot` block demonstrates how to build up two figures in parallel, using `plt.figure(i)` to create and switch to figure number  $i$ . Figure numbers start at 1. A logarithmic scale is used on the  $y$  axis since we expect that  $R$  as a function of time (or mesh points) is exponential. The reason is that the theoretical estimate (B.30) contains  $u_e''$ , which for the present model goes like  $e^{-at}$ . Taking the logarithm makes a straight line.

The code follows closely the previously stated mathematical formulas, but the statements for computing the convergence rates might deserve an explanation. The generic help function `convergence_rate(h, E)` computes and returns  $r_{i-1}$ ,  $i = 1, \dots, m-1$  from (B.37), given  $\Delta t_i$  in  $h$  and  $R_i^n$  in  $E$ :

```
def convergence_rates(h, E):
    from math import log
    r = [log(E[i]/E[i-1])/log(h[i]/h[i-1])
         for i in range(1, len(h))]
    return r
```

Calling `r_R_I = convergence_rates(dt, R_I)` computes the sequence of rates  $r_0, r_1, \dots, r_{m-2}$  for the model  $R_I \sim \Delta t^r$ , while the statements

```
R = np.array(R) # two-dim. numpy array
r_R = [convergence_rates(dt, R[:,n])[-1]
       for n in range(len(t_coarse))]
```

compute the final rate  $r_{m-2}$  for  $R^n \sim \Delta t^r$  at each mesh point  $t_n$  in the coarsest mesh. This latter computation deserves more explanation. Since  $R[i][n]$  holds the estimated truncation error  $R_i^n$  on mesh  $i$ , at point  $t_n$  in the coarsest mesh,  $R[:,n]$  picks out the sequence  $R_i^n$  for  $i = 0, \dots, m-1$ . The `convergence_rate` function computes the rates at  $t_n$ , and by indexing `[-1]` on the returned array from `convergence_rate`, we pick the rate  $r_{m-2}$ , which we believe is the best estimation since it is based on the two finest meshes.

The `estimate` function is available in a module `trunc_emir.py`. Let us apply this function to estimate the truncation error of the Forward Euler scheme. We need a function `decay_FE(dt, N)` that can compute (B.35) at the points in a mesh with time step  $dt$  and  $N$  intervals:

```
import numpy as np
import trunc_emir
```

```

def decay_FE(dt, N):
    dt = float(dt)
    t = np.linspace(0, N*dt, N+1)
    u_e = I*np.exp(-a*t) # exact solution, I and a are global
    u = u_e # naming convention when writing up the scheme
    R = np.zeros(N)

    for n in range(0, N):
        R[n] = (u[n+1] - u[n])/dt + a*u[n]

    # Theoretical expression for the truncation error
    R_a = 0.5*I*(-a)**2*np.exp(-a*t)*dt

    return R, t[:-1], R_a[:-1]

if __name__ == '__main__':
    I = 1; a = 2 # global variables needed in decay_FE
    trunc_empir.estimate(decay_FE, T=2.5, N_0=6, m=4, makeplot=True)

```

The estimated rates for the integrated truncation error  $R_I$  become 1.1, 1.0, and 1.0 for this sequence of four meshes. All the rates for  $R^n$ , computed as  $r_R$ , are also very close to 1 at all mesh points. The agreement between the theoretical formula (B.30) and the computed quantity (ref(B.35)) is very good, as illustrated in Figures B.1 and B.2. The program `trunc_decay_FE.py` was used to perform the simulations and it can easily be modified to test other schemes (see also Exercise B.7).

### B.3.6 Increasing the accuracy by adding correction terms

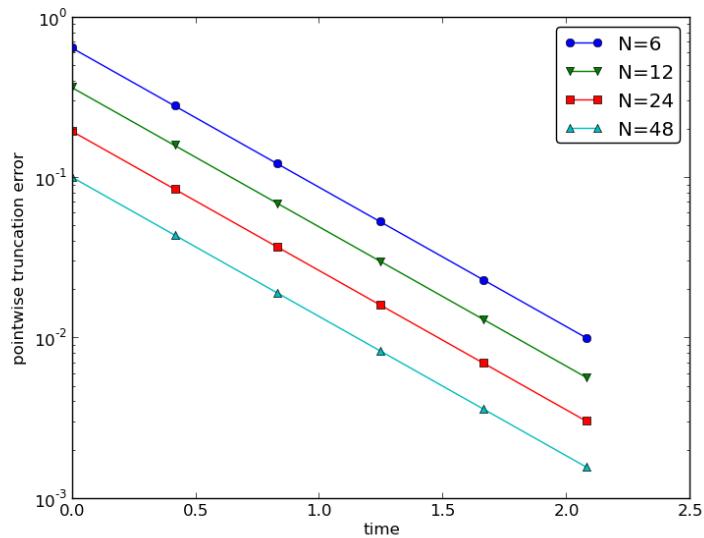
Now we ask the question: can we add terms in the differential equation that can help increase the order of the truncation error? To be precise, let us revisit the Forward Euler scheme for  $u' = -au$ , insert the exact solution  $u_e$ , include a residual  $R$ , but also include new terms  $C$ :

$$[D_t^+ u_e + au_e = C + R]^n. \quad (\text{B.39})$$

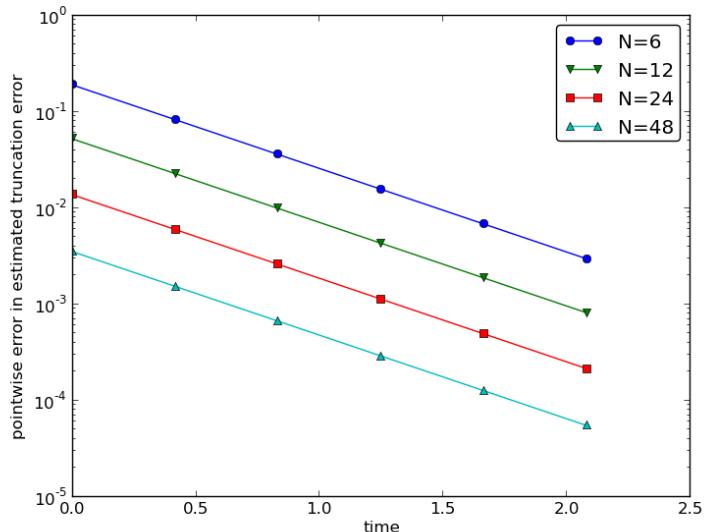
Inserting the Taylor expansions for  $[D_t^+ u_e]^n$  and keeping terms up to 3rd order in  $\Delta t$  gives the equation

$$\frac{1}{2}u_e''(t_n)\Delta t - \frac{1}{6}u_e'''(t_n)\Delta t^2 + \frac{1}{24}u_e''''(t_n)\Delta t^3 + \mathcal{O}(\Delta t^4) = C^n + R^n.$$

Can we find  $C^n$  such that  $R^n$  is  $\mathcal{O}(\Delta t^2)$ ? Yes, by setting



**Fig. B.1** Estimated truncation error at mesh points for different meshes.



**Fig. B.2** Difference between theoretical and estimated truncation error at mesh points for different meshes.

$$C^n = \frac{1}{2} u_e''(t_n) \Delta t,$$

we manage to cancel the first-order term and

$$R^n = \frac{1}{6}u_e'''(t_n)\Delta t^2 + \mathcal{O}(\Delta t^3).$$

The correction term  $C^n$  introduces  $\frac{1}{2}\Delta t u''$  in the discrete equation, and we have to get rid of the derivative  $u''$ . One idea is to approximate  $u''$  by a second-order accurate finite difference formula,  $u'' \approx (u^{n+1} - 2u^n + u^{n-1})/\Delta t^2$ , but this introduces an additional time level with  $u^{n-1}$ . Another approach is to rewrite  $u''$  in terms of  $u'$  or  $u$  using the ODE:

$$u' = -au \quad \Rightarrow \quad u'' = -au' = -a(-au) = a^2u.$$

This means that we can simply set  $C^n = \frac{1}{2}a^2\Delta t u^n$ . We can then either solve the discrete equation

$$[D_t^+ u = -au + \frac{1}{2}a^2\Delta t u]^n, \quad (\text{B.40})$$

or we can equivalently discretize the perturbed ODE

$$u' = -\hat{a}u, \quad \hat{a} = a(1 - \frac{1}{2}a\Delta t), \quad (\text{B.41})$$

by a Forward Euler method. That is, we replace the original coefficient  $a$  by the perturbed coefficient  $\hat{a}$ . Observe that  $\hat{a} \rightarrow a$  as  $\Delta t \rightarrow 0$ .

The Forward Euler method applied to (B.41) results in

$$[D_t^+ u = -a(1 - \frac{1}{2}a\Delta t)u]^n.$$

We can control our computations and verify that the truncation error of the scheme above is indeed  $\mathcal{O}(\Delta t^2)$ .

Another way of revealing the fact that the perturbed ODE leads to a more accurate solution is to look at the amplification factor. Our scheme can be written as

$$u^{n+1} = Au^n, \quad A = 1 - \hat{a}\Delta t = 1 - p + \frac{1}{2}p^2, \quad p = a\Delta t,$$

The amplification factor  $A$  as a function of  $p = a\Delta t$  is seen to be the first three terms of the Taylor series for the exact amplification factor  $e^{-p}$ . The Forward Euler scheme for  $u = -au$  gives only the first two terms  $1 - p$  of the Taylor series for  $e^{-p}$ . That is, using  $\hat{a}$  increases the order of the accuracy in the amplification factor.

Instead of replacing  $u''$  by  $a^2u$ , we use the relation  $u'' = -au'$  and add a term  $-\frac{1}{2}a\Delta tu'$  in the ODE:

$$u' = -au - \frac{1}{2}a\Delta tu' \Rightarrow \left(1 + \frac{1}{2}a\Delta t\right)u' = -au.$$

Using a Forward Euler method results in

$$\left(1 + \frac{1}{2}a\Delta t\right) \frac{u^{n+1} - u^n}{\Delta t} = -au^n,$$

which after some algebra can be written as

$$u^{n+1} = \frac{1 - \frac{1}{2}a\Delta t}{1 + \frac{1}{2}a\Delta t} u^n.$$

This is the same formula as the one arising from a Crank-Nicolson scheme applied to  $u' = -au$ ! It now recommended to do Exercise B.7 and repeat the above steps to see what kind of correction term is needed in the Backward Euler scheme to make it second order.

The Crank-Nicolson scheme is a bit more challenging to analyze, but the ideas and techniques are the same. The discrete equation reads

$$[D_t u = -au]^{n+\frac{1}{2}},$$

and the truncation error is defined through

$$[D_t u_e + a\bar{u}_e^{-t} = C + R]^{n+\frac{1}{2}},$$

where we have added a correction term. We need to Taylor expand both the discrete derivative and the arithmetic mean with aid of (B.5)-(B.6) and (B.21)-(B.22), respectively. The result is

$$\frac{1}{24}u'''_e(t_{n+\frac{1}{2}})\Delta t^2 + \mathcal{O}(\Delta t^4) + \frac{a}{8}u''_e(t_{n+\frac{1}{2}})\Delta t^2 + \mathcal{O}(\Delta t^4) = C^{n+\frac{1}{2}} + R^{n+\frac{1}{2}}.$$

The goal now is to make  $C^{n+\frac{1}{2}}$  cancel the  $\Delta t^2$  terms:

$$C^{n+\frac{1}{2}} = \frac{1}{24}u'''_e(t_{n+\frac{1}{2}})\Delta t^2 + \frac{a}{8}u''_e(t_n)\Delta t^2.$$

Using  $u' = -au$ , we have that  $u'' = a^2u$ , and we find that  $u''' = -a^3u$ . We can therefore solve the perturbed ODE problem

$$u' = -\hat{a}u, \quad \hat{a} = a(1 - \frac{1}{12}a^2 \Delta t^2),$$

by the Crank-Nicolson scheme and obtain a method that is of fourth order in  $\Delta t$ . Exercise B.7 encourages you to implement these correction terms and calculate empirical convergence rates to verify that higher-order accuracy is indeed obtained in real computations.

### B.3.7 Extension to variable coefficients

Let us address the decay ODE with variable coefficients,

$$u'(t) = -a(t)u(t) + b(t),$$

discretized by the Forward Euler scheme,

$$[D_t^+ u = -au + b]^n. \quad (\text{B.42})$$

The truncation error  $R$  is as always found by inserting the exact solution  $u_e(t)$  in the discrete scheme:

$$[D_t^+ u_e + au_e - b = R]^n. \quad (\text{B.43})$$

Using (B.11)-(B.12),

$$u'_e(t_n) - \frac{1}{2}u''_e(t_n)\Delta t + \mathcal{O}(\Delta t^2) + a(t_n)u_e(t_n) - b(t_n) = R^n.$$

Because of the ODE,

$$u'_e(t_n) + a(t_n)u_e(t_n) - b(t_n) = 0,$$

so we are left with the result

$$R^n = -\frac{1}{2}u''_e(t_n)\Delta t + \mathcal{O}(\Delta t^2). \quad (\text{B.44})$$

We see that the variable coefficients do not pose any additional difficulties in this case. Exercise B.7 takes the analysis above one step further to the Crank-Nicolson scheme.

### B.3.8 Exact solutions of the finite difference equations

Having a mathematical expression for the numerical solution is very valuable in program verification since we then know the exact numbers that the program should produce. Looking at the various formulas for the truncation errors in (B.5)-(B.6) and (B.25)-(B.26) in Section B.2.4, we see that all but two of the  $R$  expressions contains a second or higher order derivative of  $u_e$ . The exceptions are the geometric and harmonic means where the truncation error involves  $u'_e$  and even  $u_e$  in case of the harmonic mean. So, apart from these two means, choosing  $u_e$  to be a linear function of  $t$ ,  $u_e = ct + d$  for constants  $c$  and  $d$ , will make the truncation error vanish since  $u''_e = 0$ . Consequently, the truncation error of a finite difference scheme will be zero since the various approximations used will all be exact. This means that the linear solution is an exact solution of the discrete equations.

In a particular differential equation problem, the reasoning above can be used to determine if we expect a linear  $u_e$  to fulfill the discrete equations. To actually prove that this is true, we can either compute the truncation error and see that it vanishes, or we can simply insert  $u_e(t) = ct + d$  in the scheme and see that it fulfills the equations. The latter method is usually the simplest. It will often be necessary to add some source term to the ODE in order to allow a linear solution.

Many ODEs are discretized by centered differences. From Section B.2.4 we see that all the centered difference formulas have truncation errors involving  $u'''_e$  or higher-order derivatives. A quadratic solution, e.g.,  $u_e(t) = t^2 + ct + d$ , will then make the truncation errors vanish. This observation can be used to test if a quadratic solution will fulfill the discrete equations. Note that a quadratic solution will not obey the equations for a Crank-Nicolson scheme for  $u' = -au + b$  because the approximation applies an arithmetic mean, which involves a truncation error with  $u''_e$ .

### B.3.9 Computing truncation errors in nonlinear problems

The general nonlinear ODE

$$u' = f(u, t), \quad (\text{B.45})$$

can be solved by a Crank-Nicolson scheme

$$[D_t u = \bar{f}^t]^{n+\frac{1}{2}}. \quad (\text{B.46})$$

The truncation error is as always defined as the residual arising when inserting the exact solution  $u_e$  in the scheme:

$$[D_t u_e - \bar{f}^t = R]^{n+\frac{1}{2}}. \quad (\text{B.47})$$

Using (B.21)-(B.22) for  $\bar{f}^t$  results in

$$\begin{aligned} [\bar{f}^t]^{n+\frac{1}{2}} &= \frac{1}{2}(f(u_e^n, t_n) + f(u_e^{n+1}, t_{n+1})) \\ &= f(u_e^{n+\frac{1}{2}}, t_{n+\frac{1}{2}}) + \frac{1}{8}u_e''(t_{n+\frac{1}{2}})\Delta t^2 + \mathcal{O}(\Delta t^4). \end{aligned}$$

With (B.5)-(B.6) the discrete equations (B.47) lead to

$$u'_e(t_{n+\frac{1}{2}}) + \frac{1}{24}u'''_e(t_{n+\frac{1}{2}})\Delta t^2 - f(u_e^{n+\frac{1}{2}}, t_{n+\frac{1}{2}}) - \frac{1}{8}u''_e(t_{n+\frac{1}{2}})\Delta t^2 + \mathcal{O}(\Delta t^4) = R^{n+\frac{1}{2}}.$$

Since  $u'_e(t_{n+\frac{1}{2}}) - f(u_e^{n+\frac{1}{2}}, t_{n+\frac{1}{2}}) = 0$ , the truncation error becomes

$$R^{n+\frac{1}{2}} = \left( \frac{1}{24}u'''_e(t_{n+\frac{1}{2}}) - \frac{1}{8}u''_e(t_{n+\frac{1}{2}}) \right) \Delta t^2.$$

The computational techniques worked well even for this nonlinear ODE.

## B.4 Truncation errors in vibration ODEs

### B.4.1 Linear model without damping

The next example on computing the truncation error involves the following ODE for vibration problems:

$$u''(t) + \omega^2 u(t) = 0. \quad (\text{B.48})$$

Here,  $\omega$  is a given constant.

**The truncation error of a centered finite difference scheme.** Using a standard, second-ordered, central difference for the second-order derivative in time, we have the scheme

$$[D_t D_t u + \omega^2 u = 0]^n. \quad (\text{B.49})$$

Inserting the exact solution  $u_e$  in this equation and adding a residual  $R$  so that  $u_e$  can fulfill the equation results in

$$[D_t D_t u_e + \omega^2 u_e = R]^n. \quad (\text{B.50})$$

To calculate the truncation error  $R^n$ , we use (B.17)-(B.18), i.e.,

$$[D_t D_t u_e]^n = u_e''(t_n) + \frac{1}{12} u_e'''(t_n) \Delta t^2,$$

and the fact that  $u_e''(t) + \omega^2 u_e(t) = 0$ . The result is

$$R^n = \frac{1}{12} u_e'''(t_n) \Delta t^2 + \mathcal{O}(\Delta t^4). \quad (\text{B.51})$$

**The truncation error of approximating  $u'(0)$ .** The initial conditions for (B.48) are  $u(0) = I$  and  $u'(0) = V$ . The latter involves a finite difference approximation. The standard choice

$$[D_{2t} u = V]^0,$$

where  $u^{-1}$  is eliminated with the aid of the discretized ODE for  $n = 0$ , involves a centered difference with an  $\mathcal{O}(\Delta t^2)$  truncation error given by (B.7)-(B.8). The simpler choice

$$[D_t^+ u = V]^0,$$

is based on a forward difference with a truncation error  $\mathcal{O}(\Delta t)$ . A central question is if this initial error will impact the order of the scheme throughout the simulation. Exercise B.7 asks you to perform an experiment to investigate this question.

**Truncation error of the equation for the first step.** We have shown that the truncation error of the difference used to approximate the initial condition  $u'(0) = 0$  is  $\mathcal{O}(\Delta t^2)$ , but can also investigate the difference equation used for the first step. In a truncation error setting, the right way to view this equation is not to use the initial condition  $[D_{2t} u = V]^0$  to express  $u^{-1} = u^1 - 2\Delta t V$  in order to eliminate  $u^{-1}$  from the discretized differential equation, but the other way around: the fundamental equation is the discretized initial condition  $[D_{2t} u = V]^0$  and we use the discretized ODE  $[D_t D_t + \omega^2 u = 0]^0$  to eliminate  $u^{-1}$  in the discretized initial condition. From  $[D_t D_t + \omega^2 u = 0]^0$  we have

$$u^{-1} = 2u^0 - u^1 - \Delta t^2 \omega^2 u^0,$$

which inserted in  $[D_{2t}u = V]^0$  gives

$$\frac{u^1 - u^0}{\Delta t} + \frac{1}{2}\omega^2 \Delta t u^0 = V. \quad (\text{B.52})$$

The first term can be recognized as a forward difference such that the equation can be written in operator notation as

$$[D_t^+ u + \frac{1}{2}\omega^2 \Delta t u = V]^0.$$

The truncation error is defined as

$$[D_t^+ u_e + \frac{1}{2}\omega^2 \Delta t u_e - V = R]^0.$$

Using (B.11)-(B.12) with one more term in the Taylor series, we get that

$$u'_e(0) + \frac{1}{2}u''_e(0)\Delta t + \frac{1}{6}u'''_e(0)\Delta t^2 + \mathcal{O}(\Delta t^3) + \frac{1}{2}\omega^2 \Delta t u_e(0) - V = R^n.$$

Now,  $u'_e(0) = V$  and  $u''_e(0) = -\omega^2 u_e(0)$  so we get

$$R^n = \frac{1}{6}u'''_e(0)\Delta t^2 + \mathcal{O}(\Delta t^3).$$

There is another way of analyzing the discrete initial condition, because eliminating  $u^{-1}$  via the discretized ODE can be expressed as

$$[D_{2t}u + \Delta t(D_t D_t u - \omega^2 u) = V]^0. \quad (\text{B.53})$$

Writing out (B.53) shows that the equation is equivalent to (B.52). The truncation error is defined by

$$[D_{2t}u_e + \Delta t(D_t D_t u_e - \omega^2 u_e) = V + R]^0.$$

Replacing the difference via (B.7)-(B.8) and (B.17)-(B.18), as well as using  $u'_e(0) = V$  and  $u''_e(0) = -\omega^2 u_e(0)$ , gives

$$R^n = \frac{1}{6}u'''_e(0)\Delta t^2 + \mathcal{O}(\Delta t^3).$$

**Computing correction terms.** The idea of using correction terms to increase the order of  $R^n$  can be applied as described in Section B.3.6. We look at

$$[D_t D_t u_e + \omega^2 u_e = C + R]^n,$$

and observe that  $C^n$  must be chosen to cancel the  $\Delta t^2$  term in  $R^n$ . That is,

$$C^n = \frac{1}{12} u_e'''(t_n) \Delta t^2.$$

To get rid of the 4th-order derivative we can use the differential equation:  $u'' = -\omega^2 u$ , which implies  $u'''' = \omega^4 u$ . Adding the correction term to the ODE results in

$$u'' + \omega^2 \left(1 - \frac{1}{12} \omega^2 \Delta t^2\right) u = 0. \quad (\text{B.54})$$

Solving this equation by the standard scheme

$$[D_t D_t u + \omega^2 \left(1 - \frac{1}{12} \omega^2 \Delta t^2\right) u = 0]^n,$$

will result in a scheme with truncation error  $\mathcal{O}(\Delta t^4)$ .

We can use another set of arguments to justify that (B.54) leads to a higher-order method. Mathematical analysis of the scheme (B.49) reveals that the numerical frequency  $\tilde{\omega}$  is (approximately as  $\Delta t \rightarrow 0$ )

$$\tilde{\omega} = \omega \left(1 + \frac{1}{24} \omega^2 \Delta t^2\right).$$

One can therefore attempt to replace  $\omega$  in the ODE by a slightly smaller  $\omega$  since the numerics will make it larger:

$$[u'' + (\omega \left(1 - \frac{1}{24} \omega^2 \Delta t^2\right))^2 u = 0].$$

Expanding the squared term and omitting the higher-order term  $\Delta t^4$  gives exactly the ODE (B.54). Experiments show that  $u^n$  is computed to 4th order in  $\Delta t$ .

## B.4.2 Model with damping and nonlinearity

The model (B.48) can be extended to include damping  $\beta u'$ , a nonlinear restoring (spring) force  $s(u)$ , and some known excitation force  $F(t)$ :

$$m u'' + \beta u' + s(u) = F(t). \quad (\text{B.55})$$

The coefficient  $m$  usually represents the mass of the system. This governing equation can be discretized by centered differences:

$$[mD_t D_t u + \beta D_{2t} u + s(u) = F]^n. \quad (\text{B.56})$$

The exact solution  $u_e$  fulfills the discrete equations with a residual term:

$$[mD_t D_t u_e + \beta D_{2t} u_e + s(u_e) = F + R]^n. \quad (\text{B.57})$$

Using (B.17)-(B.18) and (B.7)-(B.8) we get

$$\begin{aligned} [mD_t D_t u_e + \beta D_{2t} u_e]^n &= mu_e''(t_n) + \beta u_e'(t_n) + \\ &\quad \left( \frac{m}{12} u_e'''(t_n) + \frac{\beta}{6} u_e''(t_n) \right) \Delta t^2 + \mathcal{O}(\Delta t^4) \end{aligned}$$

Combining this with the previous equation, we can collect the terms

$$mu_e''(t_n) + \beta u_e'(t_n) + \omega^2 u_e(t_n) + s(u_e(t_n)) - F^n,$$

and set this sum to zero because  $u_e$  solves the differential equation. We are left with the truncation error

$$R^n = \left( \frac{m}{12} u_e'''(t_n) + \frac{\beta}{6} u_e''(t_n) \right) \Delta t^2 + \mathcal{O}(\Delta t^4), \quad (\text{B.58})$$

so the scheme is of second order.

According to (B.58), we can add correction terms

$$C^n = \left( \frac{m}{12} u_e'''(t_n) + \frac{\beta}{6} u_e''(t_n) \right) \Delta t^2,$$

to the right-hand side of the ODE to obtain a fourth-order scheme. However, expressing  $u'''$  and  $u''$  in terms of lower-order derivatives is now harder because the differential equation is more complicated:

$$\begin{aligned} u''' &= \frac{1}{m} (F' - \beta u'' - s'(u)u'), \\ u'''' &= \frac{1}{m} (F'' - \beta u''' - s''(u)(u')^2 - s'(u)u''), \\ &= \frac{1}{m} (F'' - \beta \frac{1}{m} (F' - \beta u'' - s'(u)u') - s''(u)(u')^2 - s'(u)u''). \end{aligned}$$

It is not impossible to discretize the resulting modified ODE, but it is up to debate whether correction terms are feasible and the way to go. Computing with a smaller  $\Delta t$  is usually always possible in these problems to achieve the desired accuracy.

### B.4.3 Extension to quadratic damping

Instead of the linear damping term  $\beta u'$  in (B.55) we now consider quadratic damping  $\beta|u'|u'$ :

$$mu'' + \beta|u'|u' + s(u) = F(t). \quad (\text{B.59})$$

A centered difference for  $u'$  gives rise to a nonlinearity, which can be linearized using a geometric mean:  $||u'|u'|^n \approx ||u'|^{n-\frac{1}{2}}||u'|^{n+\frac{1}{2}}$ . The resulting scheme becomes

$$[mD_tD_tu]^n + \beta|[D_tu]^{n-\frac{1}{2}}|[D_tu]^{n+\frac{1}{2}} + s(u^n) = F^n. \quad (\text{B.60})$$

The truncation error is defined through

$$[mD_tD_tu_e]^n + \beta|[D_tu_e]^{n-\frac{1}{2}}|[D_tu_e]^{n+\frac{1}{2}} + s(u_e^n) - F^n = R^n. \quad (\text{B.61})$$

We start with expressing the truncation error of the geometric mean. According to (B.23)-(B.24), **hpl 30: Should it be  $u'_e$  instead of  $u'$  below?**

$$\begin{aligned} |[D_tu_e]^{n-\frac{1}{2}}|[D_tu_e]^{n+\frac{1}{2}} &= |[D_tu_e|D_tu_e]|^n - \frac{1}{4}u'_e(t_n)^2\Delta t^2 + \\ &\quad \frac{1}{4}u_e(t_n)u''_e(t_n)\Delta t^2 + \mathcal{O}(\Delta t^4). \end{aligned}$$

Using (B.5)-(B.6) for the  $D_tu_e$  factors results in

$$|[D_tu_e|D_tu_e]|^n = |u'_e + \frac{1}{24}u'''_e(t_n)\Delta t^2 + \mathcal{O}(\Delta t^4)|(|u'_e + \frac{1}{24}u'''_e(t_n)\Delta t^2 + \mathcal{O}(\Delta t^4))$$

We can remove the absolute value since it essentially gives a factor 1 or -1 only. Calculating the product, we have the leading-order terms

$$[D_tu_eD_tu_e]^n = (u'_e(t_n))^2 + \frac{1}{12}u_e(t_n)u''_e(t_n)\Delta t^2 + \mathcal{O}(\Delta t^4).$$

With

$$m[D_tD_tu_e]^n = mu''_e(t_n) + \frac{m}{12}u'''_e(t_n)\Delta t^2 + \mathcal{O}(\Delta t^4),$$

and using the differential equation on the form  $mu'' + \beta(u')^2 + s(u) = F$ , we end up with

$$R^n = \left( \frac{m}{12} u_e'''(t_n) + \frac{\beta}{12} u_e(t_n) u_e'''(t_n) \right) \Delta t^2 + \mathcal{O}(\Delta t^4).$$

This result demonstrates that we have second-order accuracy also with quadratic damping. The key elements that lead to the second-order accuracy is that the difference approximations are  $\mathcal{O}(\Delta t^2)$  and the geometric mean approximation is also  $\mathcal{O}(\Delta t^2)$ .

#### B.4.4 The general model formulated as first-order ODEs

The second-order model (B.59) can be formulated as a first-order system,

$$u' = v, \quad (\text{B.62})$$

$$v' = \frac{1}{m} (F(t) - \beta|v|v - s(u)). \quad (\text{B.63})$$

The system (B.62)-(B.62) can be solved either by a forward-backward scheme or a centered scheme on a staggered mesh.

**hpl 31:** The next section is wrong. Use Euler-Cromer so we get away with explicit update of  $v$  first, then  $u$ . If we cannot really explain who terms cancel to get second-order accuracy, should we not just focus on the staggered scheme? Everything below must be rewritten.

**The forward-backward scheme.** The discretization is based on the idea of stepping (B.62) forward in time and then using a backward difference in (B.63) with the recently computed (and therefore known)  $u$ :

$$[D_t^+ u = v]^n, \quad (\text{B.64})$$

$$[D_t^- v = \frac{1}{m} (F(t) - \beta|v|v - s(u))]^{n+1}. \quad (\text{B.65})$$

The term  $|v|v$  gives rise to a nonlinearity  $|v^{n+1}|v^{n+1}$ , which can be linearized as  $|v^n|v^{n+1}$ :

$$[D_t^+ u = v]^n, \quad (\text{B.66})$$

$$[D_t^- v]^{n+1} = \frac{1}{m} (F(t_{n+1}) - \beta|v^n|v^{n+1} - s(u^{n+1})). \quad (\text{B.67})$$

Each ODE will have a truncation error when inserting the exact solutions  $u_e$  and  $v_e$  in (B.64)-(B.65):

$$[D_t^+ u_e = v_e + R_u]^n, \quad (\text{B.68})$$

$$[D_t^- v_e]^{n+1} = \frac{1}{m}(F(t_{n+1}) - \beta|v_e(t_n)|v_e(t_{n+1}) - s(u_e(t_{n+1}))) + R_v^{n+1}. \quad (\text{B.69})$$

Application of (B.11)-(B.12) and (B.9)-(B.10) in (B.68) and (B.69), respectively, gives

$$u'_e(t_n) + \frac{1}{2}u''_e(t_n)\Delta t + \mathcal{O}(\Delta t^2) = v_e(t_n) + R_u^n, \quad (\text{B.70})$$

$$v'_e(t_{n+1}) - \frac{1}{2}v''_e(t_{n+1})\Delta t + \mathcal{O}(\Delta t^2) = \frac{1}{m}(F(t_{n+1}) - \beta|v_e(t_n)|v_e(t_{n+1}) + s(u_e(t_{n+1}))) + R_v^n. \quad (\text{B.71})$$

Since  $u'_e = v_e$ , (B.70) gives

$$R_u^n = \frac{1}{2}u''_e(t_n)\Delta t + \mathcal{O}(\Delta t^2).$$

In (B.71) we can collect the terms that constitute the ODE, but the damping term has the wrong form. Let us drop the absolute value in the damping term for simplicity. Adding a subtracting the right form  $v^{n+1}v^{n+1}$  helps:

$$v'_e(t_{n+1}) - \frac{1}{m}(F(t_{n+1}) - \beta v_e(t_{n+1})v_e(t_{n+1}) + s(u_e(t_{n+1}))) + (\beta v_e(t_n)v_e(t_{n+1}) - \beta v_e(t_{n+1})v_e(t_{n+1}))),$$

which reduces to

$$\begin{aligned} \frac{\beta}{m}v_e(t_{n+1}(v_e(t_n) - v_e(t_{n+1})) &= \frac{\beta}{m}v_e(t_{n+1}[D_t^- v_e]^{n+1}\Delta t \\ &= \frac{\beta}{m}v_e(t_{n+1}(v'_e(t_{n+1})\Delta t + -\frac{1}{2}v'''_e(t_{n+1})\Delta t^+\mathcal{O}(\Delta t^3))). \end{aligned}$$

We end with  $R_u^n$  and  $R_v^{n+1}$  as  $\mathcal{O}(\Delta t)$ , simply because all the building blocks in the schemes (the forward and backward differences and the

linearization trick) are only first-order accurate. However, this analysis is misleading: the building blocks play together in a way that makes the scheme second-order accurate. This is shown by considering an alternative, yet equivalent, formulation of the above scheme.

**A centered scheme on a staggered mesh.** We now introduce a staggered mesh where we seek  $u$  at mesh points  $t_n$  and  $v$  at points  $t_{n+\frac{1}{2}}$  in between the  $u$  points. The staggered mesh makes it easy to formulate centered differences in the system (B.62)-(B.62):

$$[D_t u = v]^{n-\frac{1}{2}}, \quad (\text{B.72})$$

$$[D_t v = \frac{1}{m}(F(t) - \beta|v|v - s(u))]^n. \quad (\text{B.73})$$

The term  $|v^n|v^n$  causes trouble since  $v^n$  is not computed, only  $v^{n-\frac{1}{2}}$  and  $v^{n+\frac{1}{2}}$ . Using geometric mean, we can express  $|v^n|v^n$  in terms of known quantities:  $|v^n|v^n \approx |v^{n-\frac{1}{2}}|v^{n+\frac{1}{2}}$ . We then have

$$[D_t u]^{n-\frac{1}{2}} = v^{n-\frac{1}{2}}, \quad (\text{B.74})$$

$$[D_t v]^n = \frac{1}{m}(F(t_n) - \beta|v^{n-\frac{1}{2}}|v^{n+\frac{1}{2}} - s(u^n)). \quad (\text{B.75})$$

The truncation error in each equation fulfills

$$\begin{aligned} [D_t u_e]^{n-\frac{1}{2}} &= v_e(t_{n-\frac{1}{2}}) + R_u^{n-\frac{1}{2}}, \\ [D_t v_e]^n &= \frac{1}{m}(F(t_n) - \beta|v_e(t_{n-\frac{1}{2}})|v_e(t_{n+\frac{1}{2}}) - s(u^n)) + R_v^n. \end{aligned}$$

The truncation error of the centered differences is given by (B.5)-(B.6), and the geometric mean approximation analysis can be taken from (B.23)-(B.24). These results lead to

$$u'_e(t_{n-\frac{1}{2}}) + \frac{1}{24}u'''_e(t_{n-\frac{1}{2}})\Delta t^2 + \mathcal{O}(\Delta t^4) = v_e(t_{n-\frac{1}{2}}) + R_u^{n-\frac{1}{2}},$$

and

$$v'_e(t_n) = \frac{1}{m}(F(t_n) - \beta|v_e(t_n)|v_e(t_n) + \mathcal{O}(\Delta t^2) - s(u^n)) + R_v^n.$$

The ODEs fulfilled by  $u_e$  and  $v_e$  are evident in these equations, and we achieve second-order accuracy for the truncation error in both equations:

$$R_u^{n-\frac{1}{2}} = \mathcal{O}(\Delta t^2), \quad R_v^n = \mathcal{O}(\Delta t^2).$$

Comparing (B.74)-(B.75) with (B.66)-(B.67), we can hopefully realize that these schemes are equivalent (which becomes clear when we implement both). The obvious advantage with the staggered mesh approach is that we can all the way use second-order accurate building blocks and in this way convince ourselves that the resulting scheme has an error of  $\mathcal{O}(\Delta t^2)$ .

## B.5 Truncation errors in wave equations

### B.5.1 Linear wave equation in 1D

The standard, linear wave equation in 1D for a function  $u(x, t)$  reads

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2} + f(x, t), \quad x \in (0, L), \quad t \in (0, T], \quad (\text{B.76})$$

where  $c$  is the constant wave velocity of the physical medium in  $[0, L]$ . The equation can also be more compactly written as

$$u_{tt} = c^2 u_{xx} + f, \quad x \in (0, L), \quad t \in (0, T], \quad (\text{B.77})$$

Centered, second-order finite differences are a natural choice for discretizing the derivatives, leading to

$$[D_t D_t u]_i^n = c^2 [D_x D_x u]_i^n + f_i^n. \quad (\text{B.78})$$

Inserting the exact solution  $u_e(x, t)$  in (B.78) makes this function fulfill the equation if we add the term  $R$ :

$$[D_t D_t u_e]_i^n = c^2 [D_x D_x u_e]_i^n + f_i^n + R_i^n \quad (\text{B.79})$$

Our purpose is to calculate the truncation error  $R$ . From (B.17)-(B.18) we have that

$$[D_t D_t u_e]_i^n = u_{e,tt}(x_i, t_n) + \frac{1}{12} u_{e,tttt}(x_i, t_n) \Delta t^2 + \mathcal{O}(\Delta t^4),$$

when we use a notation taking into account that  $u_e$  is a function of two variables and that derivatives must be partial derivatives. The notation  $u_{e,tt}$  means  $\partial^2 u_e / \partial t^2$ .

The same formula may also be applied to the  $x$ -derivative term:

$$[D_x D_x u_e]_i^n = u_{e,xx}(x_i, t_n) + \frac{1}{12} u_{e,xxxx}(x_i, t_n) \Delta x^2 + \mathcal{O}(\Delta x^4),$$

Equation (B.79) now becomes

$$u_{e,tt} + \frac{1}{12} u_{e,tttt}(x_i, t_n) \Delta t^2 = c^2 u_{e,xx} + c^2 \frac{1}{12} u_{e,xxxx}(x_i, t_n) \Delta x^2 + f(x_i, t_n) + \mathcal{O}(\Delta t^4, \Delta x^4) + R_i^n.$$

Because  $u_e$  fulfills the partial differential equation (PDE) (B.77), the first, third, and fifth term cancel out, and we are left with

$$R_i^n = \frac{1}{12} u_{e,tttt}(x_i, t_n) \Delta t^2 - c^2 \frac{1}{12} u_{e,xxxx}(x_i, t_n) \Delta x^2 + \mathcal{O}(\Delta t^4, \Delta x^4), \quad (\text{B.80})$$

showing that the scheme (B.78) is of second order in the time and space mesh spacing.

## B.5.2 Finding correction terms

Can we add correction terms to the PDE and increase the order of  $R_i^n$  in (B.80)? The starting point is

$$[D_t D_t u_e = c^2 D_x D_x u_e + f + C + R]_i^n \quad (\text{B.81})$$

From the previous analysis we simply get (B.80) again, but now with  $C$ :

$$R_i^n + C_i^n = \frac{1}{12} u_{e,tttt}(x_i, t_n) \Delta t^2 - c^2 \frac{1}{12} u_{e,xxxx}(x_i, t_n) \Delta x^2 + \mathcal{O}(\Delta t^4, \Delta x^4). \quad (\text{B.82})$$

The idea is to let  $C_i^n$  cancel the  $\Delta t^2$  and  $\Delta x^2$  terms to make  $R_i^n = \mathcal{O}(\Delta t^4, \Delta x^4)$ :

$$C_i^n = \frac{1}{12} u_{e,tttt}(x_i, t_n) \Delta t^2 - c^2 \frac{1}{12} u_{e,xxxx}(x_i, t_n) \Delta x^2.$$

Essentially, it means that we add a new term

$$C = \frac{1}{12} (u_{tttt} \Delta t^2 - c^2 u_{xxxx} \Delta x^2),$$

to the right-hand side of the PDE. We must either discretize these 4th-order derivatives directly or rewrite them in terms of lower-order derivatives with the aid of the PDE. The latter approach is more feasible. From the PDE we have the operator equality

$$\frac{\partial^2}{\partial t^2} = c^2 \frac{\partial^2}{\partial x^2},$$

so

$$u_{tttt} = c^2 u_{xxtt}, \quad u_{xxxx} = c^{-2} u_{ttxx}.$$

Assuming  $u$  is smooth enough, so that  $u_{xxtt} = u_{ttxx}$ , these relations lead to

$$C = \frac{1}{12} ((c^2 \Delta t^2 - \Delta x^2) u_{xx})_{tt}.$$

A natural discretization is

$$C_i^n = \frac{1}{12} ((c^2 \Delta t^2 - \Delta x^2) [D_x D_x D_t D_t u]_i^n).$$

Writing out  $[D_x D_x D_t D_t u]_i^n$  as  $[D_x D_x (D_t D_t u)]_i^n$  gives

$$\frac{1}{\Delta t^2} \left( \frac{u_{i+1}^{n+1} - 2u_i^n + u_{i-1}^{n-1}}{\Delta x^2} - 2 \right. \\ \left. \frac{u_i^{n+1} - 2u_i^n + u_i^{n-1}}{\Delta x^2} + \frac{u_{i-1}^{n+1} - 2u_{i-1}^n + u_{i-1}^{n-1}}{\Delta x^2} \right)$$

Now the unknown values  $u_{i+1}^{n+1}$ ,  $u_i^{n+1}$ , and  $u_{i-1}^{n+1}$  are *coupled*, and we must solve a tridiagonal system to find them. This is in principle straightforward, but it results in an implicit finite difference scheme, while we had a convenient explicit scheme without the correction terms.

### B.5.3 Extension to variable coefficients

Now we address the variable coefficient version of the linear 1D wave equation,

$$\frac{\partial^2 u}{\partial t^2} = \frac{\partial}{\partial x} \left( \lambda(x) \frac{\partial u}{\partial x} \right),$$

or written more compactly as

$$u_{tt} = (\lambda u_x)_x. \quad (\text{B.83})$$

The discrete counterpart to this equation, using arithmetic mean for  $\lambda$  and centered differences, reads

$$[D_t D_t u = D_x \bar{\lambda}^x D_x u]_i^n. \quad (\text{B.84})$$

The truncation error is the residual  $R$  in the equation

$$[D_t D_t u_e = D_x \bar{\lambda}^x D_x u_e + R]_i^n. \quad (\text{B.85})$$

The difficulty with (B.85) is how to compute the truncation error of the term  $[D_x \bar{\lambda}^x D_x u_e]_i^n$ .

We start by writing out the outer operator:

$$[D_x \bar{\lambda}^x D_x u_e]_i^n = \frac{1}{\Delta x} \left( [\bar{\lambda}^x D_x u_e]_{i+\frac{1}{2}}^n - [\bar{\lambda}^x D_x u_e]_{i-\frac{1}{2}}^n \right). \quad (\text{B.86})$$

With the aid of (B.5)-(B.6) and (B.21)-(B.22) we have

$$\begin{aligned} [D_x u_e]_{i+\frac{1}{2}}^n &= u_{e,x}(x_{i+\frac{1}{2}}, t_n) + \frac{1}{24} u_{e,xxx}(x_{i+\frac{1}{2}}, t_n) \Delta x^2 + \mathcal{O}(\Delta x^4), \\ [\bar{\lambda}^x]_{i+\frac{1}{2}} &= \lambda(x_{i+\frac{1}{2}}) + \frac{1}{8} \lambda''(x_{i+\frac{1}{2}}) \Delta x^2 + \mathcal{O}(\Delta x^4), \\ [\bar{\lambda}^x D_x u_e]_{i+\frac{1}{2}}^n &= (\lambda(x_{i+\frac{1}{2}}) + \frac{1}{8} \lambda''(x_{i+\frac{1}{2}}) \Delta x^2 + \mathcal{O}(\Delta x^4)) \times \\ &\quad (u_{e,x}(x_{i+\frac{1}{2}}, t_n) + \frac{1}{24} u_{e,xxx}(x_{i+\frac{1}{2}}, t_n) \Delta x^2 + \mathcal{O}(\Delta x^4)) \\ &= \lambda(x_{i+\frac{1}{2}}) u_{e,x}(x_{i+\frac{1}{2}}, t_n) + \lambda(x_{i+\frac{1}{2}}) \frac{1}{24} u_{e,xxx}(x_{i+\frac{1}{2}}, t_n) \Delta x^2 + \\ &\quad u_{e,x}(x_{i+\frac{1}{2}}) \frac{1}{8} \lambda''(x_{i+\frac{1}{2}}) \Delta x^2 + \mathcal{O}(\Delta x^4) \\ &= [\lambda u_{e,x}]_{i+\frac{1}{2}}^n + G_{i+\frac{1}{2}}^n \Delta x^2 + \mathcal{O}(\Delta x^4), \end{aligned}$$

where we have introduced the short form

$$G_{i+\frac{1}{2}}^n = \left( \frac{1}{24} u_{\text{e},xxx}(x_{i+\frac{1}{2}}, t_n) \lambda((x_{i+\frac{1}{2}}) + u_{\text{e},x}(x_{i+\frac{1}{2}}, t_n) \frac{1}{8} \lambda''(x_{i+\frac{1}{2}})) \Delta x^2 \right).$$

Similarly, we find that

$$[\bar{\lambda}^x D_x u_{\text{e}}]_{i-\frac{1}{2}}^n = [\lambda u_{\text{e},x}]_{i-\frac{1}{2}}^n + G_{i-\frac{1}{2}}^n \Delta x^2 + \mathcal{O}(\Delta x^4).$$

Inserting these expressions in the outer operator (B.86) results in

$$\begin{aligned} [D_x \bar{\lambda}^x D_x u_{\text{e}}]_i^n &= \frac{1}{\Delta x} ([\bar{\lambda}^x D_x u_{\text{e}}]_{i+\frac{1}{2}}^n - [\bar{\lambda}^x D_x u_{\text{e}}]_{i-\frac{1}{2}}^n) \\ &= \frac{1}{\Delta x} ([\lambda u_{\text{e},x}]_{i+\frac{1}{2}}^n + G_{i+\frac{1}{2}}^n \Delta x^2 - [\lambda u_{\text{e},x}]_{i-\frac{1}{2}}^n - G_{i-\frac{1}{2}}^n \Delta x^2 + \mathcal{O}(\Delta x^4)) \\ &= [D_x \lambda u_{\text{e},x}]_i^n + [D_x G]_i^n \Delta x^2 + \mathcal{O}(\Delta x^4). \end{aligned}$$

The reason for  $\mathcal{O}(\Delta x^4)$  in the remainder is that there are coefficients in front of this term, say  $H \Delta x^4$ , and the subtraction and division by  $\Delta x$  results in  $[D_x H]_i^n \Delta x^4$ .

We can now use (B.5)-(B.6) to express the  $D_x$  operator in  $[D_x \lambda u_{\text{e},x}]_i^n$  as a derivative and a truncation error:

$$[D_x \lambda u_{\text{e},x}]_i^n = \frac{\partial}{\partial x} \lambda(x_i) u_{\text{e},x}(x_i, t_n) + \frac{1}{24} (\lambda u_{\text{e},x})_{xxx}(x_i, t_n) \Delta x^2 + \mathcal{O}(\Delta x^4).$$

Expressions like  $[D_x G]_i^n \Delta x^2$  can be treated in an identical way,

$$[D_x G]_i^n \Delta x^2 = G_x(x_i, t_n) \Delta x^2 + \frac{1}{24} G_{xxx}(x_i, t_n) \Delta x^4 + \mathcal{O}(\Delta x^4).$$

There will be a number of terms with the  $\Delta x^2$  factor. We lump these now into  $\mathcal{O}(\Delta x^2)$ . The result of the truncation error analysis of the spatial derivative is therefore summarized as

$$[D_x \bar{\lambda}^x D_x u_{\text{e}}]_i^n = \frac{\partial}{\partial x} \lambda(x_i) u_{\text{e},x}(x_i, t_n) + \mathcal{O}(\Delta x^2).$$

After having treated the  $[D_t D_t u_{\text{e}}]_i^n$  term as well, we achieve

$$R_i^n = \mathcal{O}(\Delta x^2) + \frac{1}{12} u_{\text{e},tttt}(x_i, t_n) \Delta t^2.$$

The main conclusion is that the scheme is of second-order in time and space also in this variable coefficient case. The key ingredients for second

order are the centered differences and the arithmetic mean for  $\lambda$ : all those building blocks feature second-order accuracy.

### B.5.4 1D wave equation on a staggered mesh

**hpl 32:** Write out.

### B.5.5 Linear wave equation in 2D/3D

The two-dimensional extension of (B.76) takes the form

$$\frac{\partial^2 u}{\partial t^2} = c^2 \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) + f(x, y, t), \quad (x, y) \in (0, L) \times (0, H), \quad t \in (0, T], \quad (\text{B.87})$$

where now  $c(x, y)$  is the constant wave velocity of the physical medium  $[0, L] \times [0, H]$ . In the compact notation, the PDE (B.87) can be written

$$u_{tt} = c^2(u_{xx} + u_{yy}) + f(x, y, t), \quad (x, y) \in (0, L) \times (0, H), \quad t \in (0, T], \quad (\text{B.88})$$

in 2D, while the 3D version reads

$$u_{tt} = c^2(u_{xx} + u_{yy} + u_{zz}) + f(x, y, z, t), \quad (\text{B.89})$$

for  $(x, y, z) \in (0, L) \times (0, H) \times (0, B)$  and  $t \in (0, T]$ .

Approximating the second-order derivatives by the standard formulas (B.17)-(B.18) yields the scheme

$$[D_t D_t u = c^2(D_x D_x u + D_y D_y u) + f]_{i,j,k}^n. \quad (\text{B.90})$$

The truncation error is found from

$$[D_t D_t u_e = c^2(D_x D_x u_e + D_y D_y u_e) + f + R]^n. \quad (\text{B.91})$$

The calculations from the 1D case can be repeated to the terms in the  $y$  and  $z$  directions. Collecting terms that fulfill the PDE, we end up with

$$R_{i,j,k}^n = \left[ \frac{1}{12} u_{e,tttt} \Delta t^2 - c^2 \frac{1}{12} (u_{e,xxxx} \Delta x^2 + u_{e,yyyy} \Delta x^2 + u_{e,zzzz} \Delta z^2) \right]_{i,j,k}^n + \mathcal{O}(\Delta t^4, \Delta x^4, \Delta y^4, \Delta z^4). \quad (\text{B.92})$$

## B.6 Truncation errors in diffusion equations

### B.6.1 Linear diffusion equation in 1D

The standard, linear, 1D diffusion equation takes the form

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2} + f(x, t), \quad x \in (0, L), \quad t \in (0, T], \quad (\text{B.93})$$

where  $\alpha > 0$  is the constant diffusion coefficient. A more compact form of the diffusion equation is  $u_t = \alpha u_{xx} + f$ .

The spatial derivative in the diffusion equation,  $\alpha u_{xx}$ , is commonly discretized as  $[D_x D_x u]_i^n$ . The time-derivative, however, can be treated by a variety of methods.

**The Forward Euler scheme in time.** Let us start with the simple Forward Euler scheme:

$$[D_t^+ u = \alpha D_x D_x u + f]^n.$$

The truncation error arises as the residual  $R$  when inserting the exact solution  $u_e$  in the discrete equations:

$$[D_t^+ u_e = \alpha D_x D_x u_e + f + R_i^n].$$

Now, using (B.11)-(B.12) and (B.17)-(B.18), we can transform the difference operators to derivatives:

$$\begin{aligned} u_{e,t}(x_i, t_n) + \frac{1}{2} u_{e,tt}(t_n) \Delta t + \mathcal{O}(\Delta t^2) &= \alpha u_{e,xx}(x_i, t_n) + \\ \frac{\alpha}{12} u_{e,xxxx}(x_i, t_n) \Delta x^2 + \mathcal{O}(\Delta x^4) &+ f(x_i, t_n) + R_i^n. \end{aligned}$$

The terms  $u_{e,t}(x_i, t_n) - \alpha u_{e,xx}(x_i, t_n) - f(x_i, t_n)$  vanish because  $u_e$  solves the PDE. The truncation error then becomes

$$R_i^n = \frac{1}{2} u_{e,tt}(t_n) \Delta t + \mathcal{O}(\Delta t^2) - \frac{\alpha}{12} u_{e,xxxx}(x_i, t_n) \Delta x^2 + \mathcal{O}(\Delta x^4).$$

**The Crank-Nicolson scheme in time.** The Crank-Nicolson method consists of using a centered difference for  $u_t$  and an arithmetic average of the  $u_{xx}$  term:

$$[D_t u]_i^{n+\frac{1}{2}} = \alpha \frac{1}{2} ([D_x D_x u]_i^n + [D_x D_x u]_i^{n+1}) + f_i^{n+\frac{1}{2}}.$$

The equation for the truncation error is

$$[D_t u_e]_i^{n+\frac{1}{2}} = \alpha \frac{1}{2} ([D_x D_x u_e]_i^n + [D_x D_x u_e]_i^{n+1}) + f_i^{n+\frac{1}{2}} + R_i^{n+\frac{1}{2}}.$$

To find the truncation error, we start by expressing the arithmetic average in terms of values at time  $t_{n+\frac{1}{2}}$ . According to (B.21)-(B.22),

$$\frac{1}{2} ([D_x D_x u_e]_i^n + [D_x D_x u_e]_i^{n+1}) = [D_x D_x u_e]_i^{n+\frac{1}{2}} + \frac{1}{8} [D_x D_x u_{e,tt}]_i^{n+\frac{1}{2}} \Delta t^2 + \mathcal{O}(\Delta t^4).$$

With (B.17)-(B.18) we can express the difference operator  $D_x D_x u$  in terms of a derivative:

$$[D_x D_x u_e]_i^{n+\frac{1}{2}} = u_{e,xx}(x_i, t_{n+\frac{1}{2}}) + \frac{1}{12} u_{e,xxxx}(x_i, t_{n+\frac{1}{2}}) \Delta x^2 + \mathcal{O}(\Delta x^4).$$

The error term from the arithmetic mean is similarly expanded,

$$\frac{1}{8} [D_x D_x u_{e,tt}]_i^{n+\frac{1}{2}} \Delta t^2 = \frac{1}{8} u_{e,txx}(x_i, t_{n+\frac{1}{2}}) \Delta t^2 + \mathcal{O}(\Delta t^2 \Delta x^2)$$

The time derivative is analyzed using (B.5)-(B.6):

$$[D_t u]_i^{n+\frac{1}{2}} = u_{e,t}(x_i, t_{n+\frac{1}{2}}) + \frac{1}{24} u_{e,ttt}(x_i, t_{n+\frac{1}{2}}) \Delta t^2 + \mathcal{O}(\Delta t^4).$$

Summing up all the contributions and notifying that

$$u_{e,t}(x_i, t_{n+\frac{1}{2}}) = \alpha u_{e,xx}(x_i, t_{n+\frac{1}{2}}) + f(x_i, t_{n+\frac{1}{2}}),$$

the truncation error is given by

$$R_i^{n+\frac{1}{2}} = \frac{1}{8}u_{e,xx}(x_i, t_{n+\frac{1}{2}})\Delta t^2 + \frac{1}{12}u_{e,xxxx}(x_i, t_{n+\frac{1}{2}})\Delta x^2 + \\ \frac{1}{24}u_{e,ttt}(x_i, t_{n+\frac{1}{2}})\Delta t^2 + \mathcal{O}(\Delta x^4) + \mathcal{O}(\Delta t^4) + \mathcal{O}(\Delta t^2\Delta x^2)$$

### B.6.2 Linear diffusion equation in 2D/3D

### B.6.3 A nonlinear diffusion equation in 2D

## B.7 Exercises

### Exercise B.1: Truncation error of a weighted mean

Derive the truncation error of the weighted mean in (B.19)-(B.20).

**Hint.** Expand  $u_e^{n+1}$  and  $u_e^n$  around  $t_{n+\theta}$ .

Filename: `trunc_weighted_mean`.

### Exercise B.2: Simulate the error of a weighted mean

We consider the weighted mean

$$u_e(t_n) \approx \theta u_e^{n+1} + (1 - \theta)u_e^n.$$

Choose some specific function for  $u_e(t)$  and compute the error in this approximation for a sequence of decreasing  $\Delta t = t_{n+1} - t_n$  and for  $\theta = 0, 0.25, 0.5, 0.75, 1$ . Assuming that the error equals  $C\Delta t^r$ , for some constants  $C$  and  $r$ , compute  $r$  for the two smallest  $\Delta t$  values for each choice of  $\theta$  and compare with the truncation error (B.19)-(B.20). Filename: `trunc_theta_avg`.

### Exercise B.3: Verify a truncation error formula

Set up a numerical experiment as explained in Section B.3.5 for verifying the formulas (B.15)-(B.16). Filename: `trunc_backward_2level`.

### Exercise B.4: Truncation error of the Backward Euler scheme

Derive the truncation error of the Backward Euler scheme for the decay ODE  $u' = -au$  with constant  $a$ . Extend the analysis to cover the variable-coefficient case  $u' = -a(t)u + b(t)$ . Filename: `trunc_decay_BE`.

**Exercise B.5: Empirical estimation of truncation errors**

Use the ideas and tools from Section B.3.5 to estimate the rate of the truncation error of the Backward Euler and Crank-Nicolson schemes applied to the exponential decay model  $u' = -au$ ,  $u(0) = I$ .

**Hint.** In the Backward Euler scheme, the truncation error can be estimated at mesh points  $n = 1, \dots, N$ , while the truncation error must be estimated at midpoints  $t_{n+\frac{1}{2}}$ ,  $n = 0, \dots, N - 1$  for the Crank-Nicolson scheme. The `truncation_error(dt, N)` function to be supplied to the `estimate` function needs to carefully implement these details and return the right `t` array such that `t[i]` is the time point corresponding to the quantities `R[i]` and `R_a[i]`.

Filename: `trunc_decay_BNCN`.

**Exercise B.6: Correction term for a Backward Euler scheme**

Consider the model  $u' = -au$ ,  $u(0) = I$ . Use the ideas of Section B.3.6 to add a correction term to the ODE such that the Backward Euler scheme applied to the perturbed ODE problem is of second order in  $\Delta t$ . Find the amplification factor. Filename: `trunc_decay_BE_corr`.

**Exercise B.7: Verify the effect of correction terms**

The program `decay_convrate.py` solves  $u' = -au$ ,  $u(0) = I$ , by the  $\theta$ -rule and computes convergence rates. Copy this file and adjust  $a$  in the `solver` function such that it incorporates correction terms. Run the program to verify that the error from the Forward and Backward Euler schemes with perturbed  $a$  is  $\mathcal{O}(\Delta t^2)$ , while the error arising from the Crank-Nicolson scheme with perturbed  $a$  is  $\mathcal{O}(\Delta t^4)$ . Filename: `trunc_decay_corr_verify`.

**Exercise B.8: Truncation error of the Crank-Nicolson scheme**

The variable-coefficient ODE  $u' = -a(t)u + b(t)$  can be discretized in two different ways by the Crank-Nicolson scheme, depending on whether we use averages for  $a$  and  $b$  or compute them at the midpoint  $t_{n+\frac{1}{2}}$ :

$$[D_t u = -a\bar{u}^t + b]^{n+\frac{1}{2}}, \quad (\text{B.94})$$

$$[D_t u = \overline{-au + b}^t]^{n+\frac{1}{2}}. \quad (\text{B.95})$$

Compute the truncation error in both cases. Filename: `trunc_decay_CN_vc`.

### Exercise B.9: Truncation error of $u' = f(u, t)$

Consider the general nonlinear first-order scalar ODE

$$u'(t) = f(u(t), t).$$

Show that the truncation error in the Forward Euler scheme,

$$[D_t^+ u = f(u, t)]^n,$$

and in the Backward Euler scheme,

$$[D_t^- u = f(u, t)]^n,$$

both are of first order, regardless of what  $f$  is.

Showing the order of the truncation error in the Crank-Nicolson scheme,

$$[D_t u = f(u, t)]^{n+\frac{1}{2}},$$

is somewhat more involved: Taylor expand  $u_e^n$ ,  $u_e^{n+1}$ ,  $f(u_e^n, t_n)$ , and  $f(u_e^{n+1}, t_{n+1})$  around  $t_{n+\frac{1}{2}}$ , and use that

$$\frac{df}{dt} = \frac{\partial f}{\partial u} u' + \frac{\partial f}{\partial t}.$$

Check that the derived truncation error is consistent with previous results for the case  $f(u, t) = -au$ . Filename: `trunc_nonlinear_ODE`.

### Exercise B.10: Truncation error of $[D_t D_t u]^n$

Derive the truncation error of the finite difference approximation (B.17)-(B.18) to the second-order derivative. Filename: `trunc_d2u`.

### Exercise B.11: Investigate the impact of approximating $u'(0)$

Section B.4.1 describes two ways of discretizing the initial condition  $u'(0) = V$  for a vibration model  $u'' + \omega^2 u = 0$ : a centered difference  $[D_{2t}u = V]^0$  or a forward difference  $[D_t^+ u = V]^0$ . The program `vib_undamped.py` solves  $u'' + \omega^2 u = 0$  with  $[D_{2t}u = 0]^0$  and features a function `convergence_rates` for computing the order of the error in the numerical solution. Modify this program such that it applies the forward difference  $[D_t^+ u = 0]^0$  and report how this simpler and more convenient approximation impacts the overall convergence rate of the scheme. Filename: `trunc_vib_ic_fw`.

### Exercise B.12: Investigate the accuracy of a simplified scheme

Consider the ODE

$$mu'' + \beta|u'|u' + s(u) = F(t).$$

The term  $|u'|u'$  quickly gives rise to nonlinearities and complicates the scheme. Why not simply apply a backward difference to this term such that it only involves known values? That is, we propose to solve

$$[mD_tD_t u + \beta|D_t^- u|D_t^- u + s(u) = F]^n.$$

Drop the absolute value for simplicity and find the truncation error of the scheme. Perform numerical experiments with the scheme and compare with the one based on centered differences. Can you illustrate the accuracy loss visually in real computations, or is the asymptotic analysis here mainly of theoretical interest? Filename: `trunc_vib_bw_damping`.



# Software engineering; wave equation model

C

## C.1 A 1D wave equation simulator

### C.1.1 Mathematical model

Let  $u_t$ ,  $u_{tt}$ ,  $u_x$ ,  $u_{xx}$  denote derivatives of  $u$  with respect to the subscript, i.e.,  $u_{tt}$  is a second-order time derivative and  $u_x$  is a first-order space derivative. The initial-boundary value problem implemented in the `wave1D_dn_vc.py` code is

$$u_{tt} = (q(x)u_x)_x + f(x, t), \quad x \in (0, L), \quad t \in (0, T] \quad (\text{C.1})$$

$$u(x, 0) = I(x), \quad x \in [0, L] \quad (\text{C.2})$$

$$u_t(x, 0) = V(t), \quad x \in [0, L] \quad (\text{C.3})$$

$$u(0, t) = U_0(t) \text{ or } u_x(0, t) = 0, \quad t \in (0, T] \quad (\text{C.4})$$

$$u(L, t) = U_L(t) \text{ or } u_x(L, t) = 0, \quad t \in (0, T] \quad (\text{C.5})$$

We allow variable wave velocity  $c^2(x) = q(x)$ , and Dirichlet or homogeneous Neumann conditions at the boundaries.

### C.1.2 Numerical discretization

The PDE is discretized by second-order finite differences in time and space, with arithmetic mean for the variable coefficient

$$[D_t D_t u = \varrho^{-1} D_x \bar{q}^x D_x u + f]_i^n . \quad (\text{C.6})$$

The Neumann boundary conditions are discretized by

$$[D_{2x} u]_i^n = 0,$$

at a boundary point  $i$ . The details of how the numerical scheme is worked out are described in Sections 2.6 and 2.7.

### C.1.3 A solver function

The general initial-boundary value problem (C.1)-(C.5) solved by finite difference methods can be implemented in the following `solver` function (taken from the file `wave1D_dn_vc.py`). This function builds on simpler versions described in Sections 2.3, 2.4 2.6, and 2.7. There are several quite advanced constructs that will be commented upon later.

```
def solver(I, V, f, c, U_0, U_L, L, dt, C, T,
           user_action=None, version='scalar',
           stability_safety_factor=1.0):
    """Solve u_tt=(c^2*u_x)_x + f on (0,L)x(0,T]."""
    Nt = int(round(T/dt))
    t = np.linspace(0, Nt*dt, Nt+1)      # Mesh points in time

    # Find max(c) using a fake mesh and adapt dx to C and dt
    if isinstance(c, (float,int)):
        c_max = c
    elif callable(c):
        c_max = max([c(x_) for x_ in np.linspace(0, L, 101)])
    dx = dt*c_max/(stability_safety_factor*C)
    Nx = int(round(L/dx))
    x = np.linspace(0, L, Nx+1)          # Mesh points in space
    # Make sure dx and dt are compatible with x and t
    dx = x[1] - x[0]
    dt = t[1] - t[0]

    # Treat c(x) as array
    if isinstance(c, (float,int)):
        c = np.zeros(x.shape) + c
    elif callable(c):
        # Call c(x) and fill array c
        c_ = np.zeros(x.shape)
        for i in range(Nx+1):
            c_[i] = c(x[i])
        c = c_

    q = c**2
    C2 = (dt/dx)**2; dt2 = dt*dt      # Help variables in the scheme
```

```

# Wrap user-given f, I, V, U_0, U_L if None or 0
if f is None or f == 0:
    f = (lambda x, t: 0) if version == 'scalar' else \
        lambda x, t: np.zeros(x.shape)
if I is None or I == 0:
    I = (lambda x: 0) if version == 'scalar' else \
        lambda x: np.zeros(x.shape)
if V is None or V == 0:
    V = (lambda x: 0) if version == 'scalar' else \
        lambda x: np.zeros(x.shape)
if U_0 is not None:
    if isinstance(U_0, (float,int)) and U_0 == 0:
        U_0 = lambda t: 0
if U_L is not None:
    if isinstance(U_L, (float,int)) and U_L == 0:
        U_L = lambda t: 0

# Make hash of all input data
import hashlib, inspect
data = inspect.getsource(I) + ' ' + inspect.getsource(V) + \
    ' ' + inspect.getsource(f) + ' ' + str(c) + ' ' + \
    ('None' if U_0 is None else inspect.getsource(U_0)) + \
    ('None' if U_L is None else inspect.getsource(U_L)) + \
    ' ' + str(L) + str(dt) + ' ' + str(C) + ' ' + str(T) + \
    ' ' + str(stability_safety_factor)
hashed_input = hashlib.sha1(data).hexdigest()
if os.path.isfile('.' + hashed_input + '_archive.npz'):
    # Simulation is already run
    return -1, hashed_input

u    = np.zeros(Nx+1)    # Solution array at new time level
u_1 = np.zeros(Nx+1)    # Solution at 1 time level back
u_2 = np.zeros(Nx+1)    # Solution at 2 time levels back

import time; t0 = time.clock() # CPU time measurement

Ix = range(0, Nx+1)
It = range(0, Nt+1)

# Load initial condition into u_1
for i in range(0,Nx+1):
    u_1[i] = I(x[i])

if user_action is not None:
    user_action(u_1, x, t, 0)

# Special formula for the first step
for i in Ix[1:-1]:
    u[i] = u_1[i] + dt*V(x[i]) + \
        0.5*C2*(0.5*(q[i] + q[i+1])*(u_1[i+1] - u_1[i]) - \
        0.5*(q[i] + q[i-1])*(u_1[i] - u_1[i-1])) + \
        0.5*dt2*f(x[i], t[0])

i = Ix[0]

```

```

if U_0 is None:
    # Set boundary values (x=0: i-1 -> i+1 since u[i-1]=u[i+1]
    # when du/dn = 0, on x=L: i+1 -> i-1 since u[i+1]=u[i-1])
    ip1 = i+1
    im1 = ip1 # i-1 -> i+1
    u[i] = u_1[i] + dt*V(x[i]) + \
        0.5*C2*(0.5*(q[i] + q[ip1])*(u_1[ip1] - u_1[i]) - \
        0.5*(q[i] + q[im1])*(u_1[i] - u_1[im1])) + \
        0.5*dt2*f(x[i], t[0])
else:
    u[i] = U_0(dt)

i = Ix[-1]
if U_L is None:
    im1 = i-1
    ip1 = im1 # i+1 -> i-1
    u[i] = u_1[i] + dt*V(x[i]) + \
        0.5*C2*(0.5*(q[i] + q[ip1])*(u_1[ip1] - u_1[i]) - \
        0.5*(q[i] + q[im1])*(u_1[i] - u_1[im1])) + \
        0.5*dt2*f(x[i], t[0])
else:
    u[i] = U_L(dt)

if user_action is not None:
    user_action(u, x, t, 1)

# Update data structures for next step
#u_2[:] = u_1; u_1[:] = u # safe, but slower
u_2, u_1, u = u_1, u, u_2

for n in It[1:-1]:
    # Update all inner points
    if version == 'scalar':
        for i in Ix[1:-1]:
            u[i] = - u_2[i] + 2*u_1[i] + \
                C2*(0.5*(q[i] + q[i+1])*(u_1[i+1] - u_1[i]) - \
                0.5*(q[i] + q[i-1])*(u_1[i] - u_1[i-1])) + \
                dt2*f(x[i], t[n])

    elif version == 'vectorized':
        u[1:-1] = - u_2[1:-1] + 2*u_1[1:-1] + \
            C2*(0.5*(q[1:-1] + q[2:])*(u_1[2:] - u_1[1:-1]) - \
            0.5*(q[1:-1] + q[:-2])*(u_1[1:-1] - u_1[:-2])) + \
            dt2*f(x[1:-1], t[n])
    else:
        raise ValueError('version=%s' % version)

    # Insert boundary conditions
    i = Ix[0]
    if U_0 is None:
        # Set boundary values
        # x=0: i-1 -> i+1 since u[i-1]=u[i+1] when du/dn=0
        # x=L: i+1 -> i-1 since u[i+1]=u[i-1] when du/dn=0
        ip1 = i+1

```

```

im1 = ip1
u[i] = - u_2[i] + 2*u_1[i] + \
        C2*(0.5*(q[i] + q[ip1])*(u_1[ip1] - u_1[i]) - \
              0.5*(q[i] + q[im1])*(u_1[i] - u_1[im1])) + \
        dt2*f(x[i], t[n])
else:
    u[i] = U_O(t[n+1])

i = Ix[-1]
if U_L is None:
    im1 = i-1
    ip1 = im1
    u[i] = - u_2[i] + 2*u_1[i] + \
            C2*(0.5*(q[i] + q[ip1])*(u_1[ip1] - u_1[i]) - \
                  0.5*(q[i] + q[im1])*(u_1[i] - u_1[im1])) + \
            dt2*f(x[i], t[n])
else:
    u[i] = U_L(t[n+1])

if user_action is not None:
    if user_action(u, x, t, n+1):
        break

# Update data structures for next step
#u_2[:] = u_1; u_1[:] = u # safe, but slower
u_2, u_1, u = u_1, u, u_2

# Important to correct the mathematically wrong u=u_2 above
# before returning u
u = u_1
cpu_time = time.clock() - t0
return cpu_time, hashed_input

```

Or maybe copy section by section...?

## C.2 Saving large arrays in files

Numerical simulations produce large arrays as results and the software needs to store these arrays on disk. Several methods are available in Python. We recommend to use tailored solutions for large arrays and not standard file storage tools such as `pickle` (`cPickle` for speed in Python version 2) and `shelve`, because the tailored solutions have been optimized for array data and are hence much more faster than the standard tools.

### C.2.1 Using savez to store arrays in files

**Storing individual arrays.** The `numpy.storez` function can store a set of arrays to a named file in a zip archive. An associated function `numpy.load` can be used to read the file later. Basically, we call `numpy.storez(filename, **kwargs)`, where `kwargs` is a dictionary containing array names as keys and the corresponding array objects as values. Very often, the solution at a time point is given a natural name where the name of the variable and the time level counter are combined, e.g., `u11` or `v39`. Suppose `n` is the time level counter and we have two solution arrays, `u` and `v`, that we want to save to a zip archive. The appropriate code is

```
import numpy as np
u_name = 'u%04d' % n    # array name
v_name = 'v%04d' % n    # array name
kwargs = {u_name: u, v_name: v}    # keyword args for savez
fname = '.mydata%04d.dat' % n
np.savez(fname, **kwargs)
if n == 0:                # store x once
    np.savez('.mydata_x.dat', x=x)
```

Since the name of the array must be given as a keyword argument to `savez`, and the name must be constructed as shown, it becomes a little tricky to do the call, but with a dictionary `kwargs` and `**kwargs`, which sends each key-value pair as individual keyword arguments, the task gets accomplished.

**Merging zip archives.** Each separate call to `np.savez` creates a new file (zip archive) with extension `.npz`. It is very convenient to collect all results in one archive instead. This can be done by merging all the individual `.npz` files into a single zip archive:

```
def merge_zip_archives(individual_archives, archive_name):
    """
    Merge individual zip archives made with numpy.savez into
    one archive with name archive_name.
    The individual archives can be given as a list of names
    or as a Unix wild card filename expression for glob.glob.
    The result of this function is that all the individual
    archives are deleted and the new single archive made.
    """
    import zipfile
    archive = zipfile.ZipFile(
        archive_name, 'w', zipfile.ZIP_DEFLATED,
        allowZip64=True)
    if isinstance(individual_archives, (list,tuple)):
        filenames = individual_archives
    elif isinstance(individual_archives, str):
```

```

filenames = glob.glob(individual_archives)

# Open each archive and write to the common archive
for filename in filenames:
    f = zipfile.ZipFile(filename, 'r',
                        zipfile.ZIP_DEFLATED)
    for name in f.namelist():
        data = f.open(name, 'r')
        # Save under name without .npy
        archive.writestr(name[:-4], data.read())
    f.close()
    os.remove(filename)
archive.close()

```

Here we remark that `savez` automatically adds the `.npz` extension to the names of the arrays we store. We do not want this extension in the final archive.

**Reading arrays from zip archives.** Archives created by `savez` or the merged archive we describe above with name of the form `myarchive.npz`, can be conveniently read by the `numpy.load` function:

```

import numpy as np
array_names = np.load('myarchive.npz')
for array_name in array_names:
    # array_names[array_name] is the array itself
    # e.g. plot(array_names['t'], array_names[array_name])

```

## C.2.2 Using joblib to store arrays in files

The Python package `joblib` has nice functionality for efficient storage of arrays on disk. The following class applies this functionality so that one can save an array, or in fact any Python data structure (e.g., a dictionary of arrays), to disk under a certain name. Later, we can retrieve the object by use of its name. The name of the directory under which the arrays are stored by `joblib` can be given by the user.

```

class Storage(object):
    """
    Store large data structures (e.g. numpy arrays) efficiently
    using joblib.

    Use:

    >>> from Storage import Storage
    >>> storage = Storage(cachedir='tmp_u01', verbose=1)
    >>> import numpy as np
    >>> a = np.linspace(0, 1, 100000) # large array

```

```

>>> b = np.linspace(0, 1, 100000) # large array
>>> storage.save('a', a)
>>> storage.save('b', b)
>>> # later
>>> a = storage.retrieve('a')
>>> b = storage.retrieve('b')
"""

def __init__(self, cachedir='tmp', verbose=1):
    """
    Parameters
    -----
    cachedir: str
        Name of directory where objects are stored in files.
    verbose: bool, int
        Let joblib and this class speak when storing files
        to disk.
    """
    import joblib
    self.memory = joblib.Memory(cachedir=cachedir,
                                 verbose=verbose)
    self.verbose = verbose
    self.retrieve = self.memory.cache(
        self.retrieve, ignore=['data'])
    self.save = self.retrieve

    def retrieve(self, name, data=None):
        if self.verbose > 0:
            print 'joblib save of', name
        return data

```

The `retrieve` and `save` functions, which do the work, seem quite magic. The idea is that `joblib` looks at the `name` parameter and saves the return value `data` to disk if the `name` parameter has not been used in a previous call. Otherwise, if `name` is already registered, `joblib` fetches the `data` object from file and returns it (this is example of a memoize function, see Section ??in [5]).

### C.2.3 Using a hash to create a file or directory name

The user of array storage techniques like those outlined in Sections C.2.2 and C.2.1 demand the user to assign a name for the file(s) or directory where the solution is to be stored. Ideally, this name should reflect parameters in the problem such that one can recognize an already run simulation. One technique is to make a hash string out of the input data. A hash string is a 40-character long hexadecimal string that uniquely reflects another potentially much longer string. (You may be used to hash

strings from the Git version control system: every committed version of the files in Git is recognized by a hash string.)

Suppose you have some input data in the form of functions, `numpy` arrays, and other objects. To turn these input data into a string, we may grab the source code of the functions, use a very efficient hash method for potentially large arrays, and simply convert all other objects via `str` to a string representation. The final string, merging all input data, is then converted to an SHA1 hash string such that we represent the input with a 40-character long string.

```
def myfunction(func1, func2, array1, array2, obj1, obj2):
    # Convert arguments to hash
    import inspect, joblib, hashlib
    data = (inspect.getsource(func1),
            inspect.getsource(func2),
            joblib.hash(array1),
            joblib.hash(array2),
            str(obj1),
            str(obj2))
    hash_input = hashlib.sha1(data).hexdigest()
```

It is wise to use `joblib.hash` and not try to do a `str(array1)`, since that string can be *very* long, and `joblib.hash` is more efficient than `hashlib` to turn these data into a hash.

**Remark: turning function objects into their source code is unreliable!**

The idea of turning a function object into a string via its source code may look smart, but is not a completely reliable solution. Suppose we have some function

```
x0 = 0.1
f = lambda x: 0 if x <= x0 else 1
```

The source code will be `f = lambda x: 0 if x <= x0 else 1`, so if the calling code changes the value of `x0` (which `f` remembers - it is a closure), the source remains unchanged, the hash is the same, and the change in input data is unnoticed. Consequently, the technique above must be used with care. The user can always just remove the stored files in disk and thereby force a recomputation (provided the software applies to hash to test if a zip archive or `joblib` subdirectory exists and if so avoids recomputation).

### C.3 Software for the 1D wave equation

We use `numpy.storez` to store the solution at each time level on disk. Such actions must be taken care of outside the `solver` function, more precisely in the `user_action` function that is called at every time level.

We have in the `wave1D_dn_vc.py` code implemented the `user_action` callback function as a class `PlotAndStoreSolution` with a `__call__(self, x, t, t, n)` method for the `user_action` function. Basically, `__call__` stores and plots the solution. The storage makes use of the `numpy.savez` function for saving a set of arrays to a zip archive. Here, in this callback function, we want to save one array, `u`. Since there will be many such arrays, we introduce the array names '`u%04d' % n` and closely related filenames. The usage of `numpy.savez` in `__call__` goes like this:

```
from numpy import savez
name = 'u%04d' % n    # array name
kwargs = {name: u}    # keyword args for savez
fname = '.' + self.filename + '_' + name + '.dat'
self.t.append(t[n])  # store corresponding time value
savez(fname, **kwargs)
if n == 0:            # store x once
    savez('.' + self.filename + '_x.dat', x=x)
```

For example, if `n` is 10 and `self.filename` is `tmp`, the above call to `savez` becomes `savez('.tmp_u0010.dat', u0010=u)`. The actual filename becomes `.tmp_u0010.dat.npz`. The actual array name becomes `u0010.npy`.

Each `savez` call results in a file, so after the simulation we have one file per time level. Each file produced by `savez` is a zip archive. It makes sense to merge all the files into one. This is done in the `close_file` method in the `PlotAndStoreSolution` class. The code goes as follows.

```
class PlotAndStoreSolution:
    ...
    def close_file(self, hashed_input):
        """
        Merge all files from savez calls into one archive.
        hashed_input is a string reflecting input data
        for this simulation (made by solver).
        """
        if self.filename is not None:
            # Save all the time points where solutions are saved
            savez('.' + self.filename + '_t.dat',
                  t=array(self.t, dtype=float))
            # Merge all savez files to one zip archive
            archive_name = '.' + hashed_input + '_archive.npz'
```

```

filenames = glob.glob('.' + self.filename + '*.dat.npz')
merge_zip_archives(filenames, archive_name)

```

We use various `ZipFile` functionality to extract the content of the individual files (each with name `filename`) and write it to the merged archive (`archive`). There is only one array in each individual file (`filename`) so strictly speaking, there is no need for the loop `for name in f.namelist()` (as `f.namelist()` returns a list of length 1). However, in other applications where we compute more arrays at each time level, `savez` will store all these and then there is need for iterating over `f.namelist()`.

Instead of merging the archives written by `savez` we could make an alternative implementation that writes all our arrays into one archive. This is the subject of Exercise C.9.

### C.3.1 Making hash strings from input data

The `hashed_input` argument, used to name the resulting archive file with all solutions, is supposed to be a hash reflecting all import parameters in the problem such that this simulation has a unique name. The `hashed_input` string is made in the `solver` function, using the `hashlib` and `inspect` modules, based on the arguments to `solver`:

```

# Make hash of all input data
import hashlib, inspect
data = inspect.getsource(I) + '_' + inspect.getsource(V) + \
      '_' + inspect.getsource(f) + '_' + str(c) + '_' + \
      ('None' if U_0 is None else inspect.getsource(U_0)) + \
      ('None' if U_L is None else inspect.getsource(U_L)) + \
      '_' + str(L) + str(dt) + '_' + str(C) + '_' + str(T) + \
      '_' + str(stability_safety_factor)
hashed_input = hashlib.sha1(data).hexdigest()

```

NOTE: All this is now explained!

To get the source code of a function `f` as a string, we use `inspect.getsource(f)`. All input, functions as well as variables, is then merged to a string `data`, and then `hashlib.sha1` makes a unique, much shorter (40 characters long), fixed-length string out of `data` that we can use in the archive filename.

#### Remark

Note that the construction of the `data` string is not fool proof: if, e.g.,  $I$  is a formula with parameters and the parameters change, the source code is still the same and `data` and hence the hash remains unaltered. The implementation must therefore be used with care!

### C.3.2 Avoiding rerunning previously run cases

If the archive file whose name is based on `hashed_input` already exists, the simulation with the current set of parameters has been done before and one can avoid redoing the work. The `solver` function returns the CPU time and `hashed_input`, and a negative CPU time means that no simulation was run. In that case we should not call the `close_file` method above (otherwise we overwrite the archive with just the `self.t` array). The typical usage goes like

```
action = PlotAndStoreSolution(...)
dt = (L/Nx)/C # choose the stability limit with given Nx
cpu, hashed_input = solver(
    I=lambda x: ...,
    V=0, f=0, c=1, U_0=lambda t: 0, U_L=None, L=1,
    dt=dt, C=C, T=T,
    user_action=action, version='vectorized',
    stability_safety_factor=1)
action.make_movie_file()
if cpu > 0: # did we generate new data?
    action.close_file(hashed_input)
```

### C.3.3 Verification

Exact solutions of the numerical equations are always attractive for verification purposes since the software should reproduce such solutions to machine precision. With Dirichlet boundary conditions we can construct a function that is linear in  $t$  and quadratic in  $x$  that is an exact solution of the scheme, while with Neumann conditions we are left with testing just a constant solution (see comments in Section 2.6.5).

A more general method for verification is to check the convergence rates.

**hpl 33:** Do convergence rates here!

## C.4 Programming the solver with classes

Many who know about class programming prefer to organize their software in terms of classes. We can easily port our function-based code in ... to a class version.

We will create a class `Problem` to hold the physical parameters of the problem and a class `Solver` to hold the numerical parameters and the solver function. In addition, it is convenient to collect the arrays that describe the mesh in a special `Mesh` class and make a class `Function` for a mesh function (mesh point values and its mesh).

### C.4.1 Class Problem

### C.4.2 Class Mesh

The `Mesh` class can be made valid for a space-time mesh in any number of space dimensions. To make versatile, the constructor accepts either a tuple/list of number of cells in each spatial dimension or a tuple/list of cell spacings. In addition, we need the size of the hypercube mesh as a tuple/list of 2-tuples with lower and upper limits of the mesh coordinates in each direction. For 1D meshes it is more natural to just write the number of cells or the cell size and not wrap it in a list. We also need the time interval from `t0` to `T`. Giving no spatial discretization information implies a time mesh only, and vice versa. The `Mesh` class with documentation and a doc test should now be self-explanatory:

```
import numpy as np

class Mesh(object):
    """
        Holds data structures for a uniform mesh on a hypercube in
        space, plus a uniform mesh in time.

    ===== =====
    Argument          Explanation
    ===== =====
    L      List of 2-lists of min and max coordinates
          in each spatial direction.
    T      Final time in time mesh.
    Nt     Number of cells in time mesh.
    dt     Time step. Either Nt or dt must be given.
    N      List of number of cells in the spatial directions.
    d      List of cell sizes in the spatial directions.
          Either N or d must be given.
    ===== =====
```

Users can access all the parameters mentioned above, plus “`x[i]`” and “`t`” for the coordinates in direction “`i`” and the time coordinates, respectively.

Examples:

```
>>> from UniformFDMesh import Mesh
>>>
>>> # Simple space mesh
>>> m = Mesh(L=[0,1], N=4)
>>> print m.dump()
space: [0,1] N=4 d=0.25
>>>
>>> # Simple time mesh
>>> m = Mesh(T=4, dt=0.5)
>>> print m.dump()
time: [0,4] Nt=8 dt=0.5
>>>
>>> # 2D space mesh
>>> m = Mesh(L=[[0,1], [-1,1]], d=[0.5, 1])
>>> print m.dump()
space: [0,1]x[-1,1] N=2x2 d=0.5,1
>>>
>>> # 2D space mesh and time mesh
>>> m = Mesh(L=[[0,1], [-1,1]], d=[0.5, 1], Nt=10, T=3)
>>> print m.dump()
space: [0,1]x[-1,1] N=2x2 d=0.5,1 time: [0,3] Nt=10 dt=0.3

"""
def __init__(self,
             L=None, T=None, t0=0,
             N=None, d=None,
             Nt=None, dt=None):
    if N is None and d is None:
        # No spatial mesh
        if Nt is None and dt is None:
            raise ValueError(
                'Mesh constructor: either Nt or dt must be given')
        if T is None:
            raise ValueError(
                'Mesh constructor: T must be given')
    if Nt is None and dt is None:
        if N is None and d is None:
            raise ValueError(
                'Mesh constructor: either N or d must be given')
        if L is None:
            raise ValueError(
                'Mesh constructor: L must be given')

    # Allow 1D interface without nested lists with one element
    if L is not None and isinstance(L[0], (float,int)):
        # Only an interval was given
        L = [L]
```

```
if N is not None and isinstance(N, (float,int)):
    N = [N]
if d is not None and isinstance(d, (float,int)):
    d = [d]

# Set all attributes to None
self.x = None
self.t = None
self.Nt = None
self.dt = None
self.N = None
self.d = None
self.t0 = t0

if N is None and d is not None and L is not None:
    self.L = L
    if len(d) != len(L):
        raise ValueError(
            'd has different size (no of space dim.) from '
            'L: %d vs %d', len(d), len(L))
    self.d = d
    self.N = [int(round(float(self.L[i][1] -
                                self.L[i][0])/d[i]))]
    for i in range(len(d))

if d is None and N is not None and L is not None:
    self.L = L
    if len(N) != len(L):
        raise ValueError(
            'N has different size (no of space dim.) from '
            'L: %d vs %d', len(N), len(L))
    self.N = N
    self.d = [float(self.L[i][1] - self.L[i][0])/N[i]
              for i in range(len(N))]

if Nt is None and dt is not None and T is not None:
    self.T = T
    self.dt = dt
    self.Nt = int(round(T/dt))
if dt is None and Nt is not None and T is not None:
    self.T = T
    self.Nt = Nt
    self.dt = T/float(Nt)

if self.N is not None:
    self.x = [np.linspace(
        self.L[i][0], self.L[i][1], self.N[i]+1)
        for i in range(len(self.L))]
if Nt is not None:
    self.t = np.linspace(self.t0, self.T, self.Nt+1)

def get_num_space_dim(self):
    return len(self.d) if self.d is not None else 0

def has_space(self):
```

```

    return self.d is not None

def has_time(self):
    return self.dt is not None

def dump(self):
    s = ''
    if self.has_space():
        s += 'space: ' + \
            'x'.join(['[%g,%g]' % (self.L[i][0], self.L[i][1])
                      for i in range(len(self.L))]) + ' N='
        s += 'x'.join([str(Ni) for Ni in self.N]) + ' d='
        s += ','.join([str(di) for di in self.d])
    if self.has_space() and self.has_time():
        s += ','
    if self.has_time():
        s += 'time: ' + '[%g,%g]' % (self.t0, self.T) + \
            ' Nt=%g' % self.Nt + ' dt=%g' % self.dt
    return s

```

### We rely on attribute access - not get/set functions!

Java programmers in particular are used to get/set functions in classes to access internal data. In Python, we usually apply direct access of the attribute, such as `m.N[i]` if `m` is a `Mesh` object. A widely used convention is to do this as long as access to an attribute does not require additional code. In that case, one applies a property construction. The original interface remains the same after a property is introduced (in contrast to Java), so user will not notice a change to properties.

The only argument against direct attribute access in class `Mesh` is that the attributes are read-only so we could avoid offering a set function. Instead, we rely on the user that she does not assign new values to the attributes.

### C.4.3 Class Function

A class `Function` is handy to hold a mesh and corresponding values for a scalar or vector function over the mesh. Since we may have a time or space mesh, or a combined time and space mesh, with one or more components in the function, some if tests are needed for allocating the right array sizes. To help the user, an `indices` attribute with the name

of the indices in the final array  $\mathbf{u}$  for the function values is made. The examples in the doc string should explain the functionality.

```
class Function(object):
    """
    A scalar or vector function over a mesh (of class Mesh).

    ===== Argument Explanation =====
    mesh      Class Mesh object: spatial and/or temporal mesh.
    num_comp  Number of components in function (1 for scalar).
    space_only True if the function is defined on the space mesh
                 only (to save space). False if function has values
                 in space and time.
    =====

    The indexing of "u", which holds the mesh point values of the
    function, depends on whether we have a space and/or time mesh.
```

Examples:

```
>>> from UniformFDMesh import Mesh, Function
>>>
>>> # Simple space mesh
>>> m = Mesh(L=[0,1], N=4)
>>> print m.dump()
space: [0,1] N=4 d=0.25
>>> f = Function(m)
>>> f.indices
['x0']
>>> f.u.shape
(5,)
>>> f.u[4]  # space point 4
0.0
>>>
>>> # Simple time mesh for two components
>>> m = Mesh(T=4, dt=0.5)
>>> print m.dump()
time: [0,4] Nt=8 dt=0.5
>>> f = Function(m, num_comp=2)
>>> f.indices
['time', 'component']
>>> f.u.shape
(9, 2)
>>> f.u[3,1]  # time point 3, comp=1 (2nd comp.)
0.0
>>>
>>> # 2D space mesh
>>> m = Mesh(L=[[0,1], [-1,1]], d=[0.5, 1])
>>> print m.dump()
space: [0,1]x[-1,1] N=2x2 d=0.5,1
>>> f = Function(m)
```

```

>>> f.indices
['x0', 'x1']
>>> f.u.shape
(3, 3)
>>> f.u[1,2] # space point (1,2)
0.0
>>>
>>> # 2D space mesh and time mesh
>>> m = Mesh(L=[[0,1],[-1,1]], d=[0.5,1], Nt=10, T=3)
>>> print m.dump()
space: [0,1]x[-1,1] N=2x2 d=0.5,1 time: [0,3] Nt=10 dt=0.3
>>> f = Function(m, num_comp=2, space_only=False)
>>> f.indices
['time', 'x0', 'x1', 'component']
>>> f.u.shape
(11, 3, 3, 2)
>>> f.u[2,1,2,0] # time step 2, space point (1,2), comp=0
0.0
>>> # Function with space data only
>>> f = Function(m, num_comp=1, space_only=True)
>>> f.indices
['x0', 'x1']
>>> f.u.shape
(3, 3)
>>> f.u[1,2] # space point (1,2)
0.0
"""
def __init__(self, mesh, num_comp=1, space_only=True):
    self.mesh = mesh
    self.num_comp = num_comp
    self.indices = []

# Create array(s) to store mesh point values
if (self.mesh.has_space() and not self.mesh.has_time()) or \
    (self.mesh.has_space() and self.mesh.has_time() and \
     space_only):
    # Space mesh only
    if num_comp == 1:
        self.u = np.zeros(
            [self.mesh.N[i] + 1
             for i in range(len(self.mesh.N))])
        self.indices = [
            'x'+str(i) for i in range(len(self.mesh.N))]
    else:
        self.u = np.zeros(
            [self.mesh.N[i] + 1
             for i in range(len(self.mesh.N))] +
            [num_comp])
        self.indices = [
            'x'+str(i)
            for i in range(len(self.mesh.N))] +\
            ['component']
if not self.mesh.has_space() and self.mesh.has_time():
    # Time mesh only

```

```

if num_comp == 1:
    self.u = np.zeros(self.mesh.Nt+1)
    self.indices = ['time']
else:
    # Need num_comp entries per time step
    self.u = np.zeros((self.mesh.Nt+1, num_comp))
    self.indices = ['time', 'component']
if self.mesh.has_space() and self.mesh.has_time() \
and not space_only:
    # Space-time mesh
    size = [self.mesh.Nt+1] + \
            [self.mesh.N[i]+1
                for i in range(len(self.mesh.N))]
    if num_comp > 1:
        self.indices = ['time'] + \
                    ['x'+str(i)
                        for i in range(len(self.mesh.N))] +\
                    ['component']
    size += [num_comp]
else:
    self.indices = ['time'] + ['x'+str(i)
                                for i in range(len(self.mesh.N))]
self.u = np.zeros(size)

```

#### C.4.4 Class Solver

With the `Mesh` and `Function` classes in place, we can rewrite the `solver` function, but we put it as a method in class `Solver`:

**hpl 34:** Rewrite solver!

## C.5 Migrating loops to Cython

We now consider the `wave2D_u0.py` code for solving the 2D linear wave equation with constant wave velocity and homogeneous Dirichlet boundary conditions  $u = 0$ . We shall in the present chapter extend this code with computational modules written in other languages than Python. This extended version is called `wave2D_u0_adv.py`.

The `wave2D_u0.py` file contains a `solver` function, which calls an `advance_*` function to advance the numerical scheme one level forward in time. The function `advance_scalar` applies standard Python loops to implement the scheme, while `advance_vectorized` performs corresponding vectorized arithmetics with array slices. The statements of this solver are explained in Section 2.12, in particular Sections 2.12.1 and 2.12.2.

Although vectorization can bring down the CPU time dramatically compared with scalar code, there is still some factor 5-10 to win in these types of applications by implementing the finite difference scheme in compiled code, typically in Fortran, C, or C++. This can quite easily be done by adding a little extra code to our program. Cython is an extension of Python that offers the easiest way to nail our Python loops in the scalar code down to machine code and achieve the efficiency of C.

Cython can be viewed as an extended Python language where variables are declared with types and where functions are marked to be implemented in C. Migrating Python code to Cython is done by copying the desired code segments to functions (or classes) and placing them in one or more separate files with extension .pyx.

### C.5.1 Declaring variables and annotating the code

Our starting point is the plain `advance_scalar` function for a scalar implementation of the updating algorithm for new values  $u_{i,j}^{n+1}$ :

```
def advance_scalar(u, u_1, u_2, f, x, y, t, n, Cx2, Cy2, dt2,
                   V=None, step1=False):
    Ix = range(0, u.shape[0]); Iy = range(0, u.shape[1])
    if step1:
        dt = sqrt(dt2) # save
        Cx2 = 0.5*Cx2; Cy2 = 0.5*Cy2; dt2 = 0.5*dt2 # redefine
        D1 = 1; D2 = 0
    else:
        D1 = 2; D2 = 1
    for i in Ix[1:-1]:
        for j in Iy[1:-1]:
            u_xx = u_1[i-1,j] - 2*u_1[i,j] + u_1[i+1,j]
            u_yy = u_1[i,j-1] - 2*u_1[i,j] + u_1[i,j+1]
            u[i,j] = D1*u_1[i,j] - D2*u_2[i,j] +
                      Cx2*u_xx + Cy2*u_yy + dt2*f(x[i], y[j], t[n])
            if step1:
                u[i,j] += dt*V(x[i], y[j])
    # Boundary condition u=0
    j = Iy[0]
    for i in Ix: u[i,j] = 0
    j = Iy[-1]
    for i in Ix: u[i,j] = 0
    i = Ix[0]
    for j in Iy: u[i,j] = 0
    i = Ix[-1]
    for j in Iy: u[i,j] = 0
    return u
```

We simply take a copy of this function and put it in a file `wave2D_u0_1loop_cy.pyx`. The relevant Cython implementation arises from declaring variables with types and adding some important annotations to speed up array computing in Cython. Let us first list the complete code in the `.pyx` file:

```
import numpy as np
cimport numpy as np
cimport cython
ctypedef np.float64_t DT      # data type

@cython.boundscheck(False)    # turn off array bounds check
@cython.wraparound(False)     # turn off negative indices (u[-1,-1])
cpdef advance(
    np.ndarray[DT, ndim=2, mode='c'] u,
    np.ndarray[DT, ndim=2, mode='c'] u_1,
    np.ndarray[DT, ndim=2, mode='c'] u_2,
    np.ndarray[DT, ndim=2, mode='c'] f,
    double Cx2, double Cy2, double dt2):

    cdef:
        int Ix_start = 0
        int Iy_start = 0
        int Ix_end = u.shape[0]-1
        int Iy_end = u.shape[1]-1
        int i, j
        double u_xx, u_yy

    for i in range(Ix_start+1, Ix_end):
        for j in range(Iy_start+1, Iy_end):
            u_xx = u_1[i-1,j] - 2*u_1[i,j] + u_1[i+1,j]
            u_yy = u_1[i,j-1] - 2*u_1[i,j] + u_1[i,j+1]
            u[i,j] = 2*u_1[i,j] - u_2[i,j] + \
                      Cx2*u_xx + Cy2*u_yy + dt2*f[i,j]
    # Boundary condition u=0
    j = Iy_start
    for i in range(Ix_start, Ix_end+1): u[i,j] = 0
    j = Iy_end
    for i in range(Ix_start, Ix_end+1): u[i,j] = 0
    i = Ix_start
    for j in range(Iy_start, Iy_end+1): u[i,j] = 0
    i = Ix_end
    for j in range(Iy_start, Iy_end+1): u[i,j] = 0
    return u
```

This example may act as a recipe on how to transform array-intensive code with loops into Cython.

1. Variables are declared with types: for example, `double v` in the argument list instead of just `v`, and `cdef double v` for a variable `v` in the body of the function. A Python `float` object is declared as

- `double` for translation to C by Cython, while an `int` object is declared by `int`.
2. Arrays need a comprehensive type declaration involving
    - the type `np.ndarray`,
    - the data type of the elements, here 64-bit floats, abbreviated as `DT` through `ctypedef np.float64_t DT` (instead of `DT` we could use the full name of the data type: `np.float64_t`, which is a Cython-defined type),
    - the dimensions of the array, here `ndim=2` and `ndim=1`,
    - specification of contiguous memory for the array (`mode='c'`).
  3. Functions declared with `cpdef` are translated to C but are also accessible from Python.
  4. In addition to the standard `numpy` import we also need a special Cython import of `numpy`: `cimport numpy as np`, to appear *after* the standard import.
  5. By default, array indices are checked to be within their legal limits. To speed up the code one should turn off this feature for a specific function by placing `@cython.boundscheck(False)` above the function header.
  6. Also by default, array indices can be negative (counting from the end), but this feature has a performance penalty and is therefore here turned off by writing `@cython.wraparound(False)` right above the function header.
  7. The use of index sets `Ix` and `Iy` in the scalar code cannot be successfully translated to C. One reason is that constructions like `Ix[1:-1]` involve negative indices, and these are now turned off. Another reason is that Cython loops must take the form `for i in xrange` or `for i in range` for being translated into efficient C loops. We have therefore introduced `Ix_start` as `Ix[0]` and `Ix_end` as `Ix[-1]` to hold the start and end of the values of index *i*. Similar variables are introduced for the *j* index. A loop `for i in Ix` is with these new variables written as `for i in range(Ix_start, Ix_end+1)`.

### Array declaration syntax in Cython

We have used the syntax `np.ndarray[DT, ndim=2, mode='c']` to declare `numpy` arrays in Cython. There is a simpler, alternative syntax, employing [typed memory views](#), where the declaration looks like `double [:,:]`. However, the full support for this functionality

is not yet ready, and in this text we use the full array declaration syntax.

### C.5.2 Visual inspection of the C translation

Cython can visually explain how successfully it translated a code from Python to C. The command

```
Terminal> cython -a wave2D_u0_loop_cy.pyx
```

produces an HTML file `wave2D_u0_loop_cy.html`, which can be loaded into a web browser to illustrate which lines of the code have been translated to C. Figure C.1 shows the illustrated code. Yellow lines indicate the lines that Cython did not manage to translate to efficient C code and that remain in Python. For the present code we see that Cython is able to translate all the loops with array computing to C, which is our primary goal.

```
Raw output: wave2D_u0_loop_cy.c
1: import numpy as np
2: cimport numpy as np
3: cimport cython
4: ctypedef np.float64_t DT    # data type
5:
6: #cython.boundscheck(False) # turn off array bounds check
7: #cython.wraparound(False) # turn off negative indices (u[-1,-1])
8: cpdef void f(np.ndarray[DT, ndim=2, mode="c"] u,
9:             np.ndarray[DT, ndim=2, mode="c"] u1,
10:            np.ndarray[DT, ndim=2, mode="c"] u2,
11:            np.ndarray[DT, ndim=2, mode="c"] u3,
12:            np.ndarray[DT, ndim=2, mode="c"] u4,
13:            double Cx2, double Cy2, double dt2):
14:
15:     cdef int Ix_start = 0
16:     cdef int Ix_end = 0
17:     cdef int Ix_end = u.shape[0]-1
18:     cdef int Iy_end = u.shape[1]-1
19:     cdef int i, j
20:     cdef double u_xx, u_yy
21:
22:     for i in range(Ix_start+1, Ix_end):
23:         for j in range(Iy_start+1, Iy_end):
24:             u_xx = 2*u[i-1,j-1] + 2*u[i-1,j] + u[i-1,j+1]
25:             u_yy = u[i-1,j-1] + 2*u[i-1,j] + u[i-1,j+1]
26:             u[i,j] = 2*u[i-1,j] - u[2i,j] + \
27:                      Cx2*u_xx + Cy2*u_yy + dt2*f(i,j)
28:
29:     # Boundary condition u=0
29:     j = Iy_start
30:     for i in range(Ix_start, Ix_end+1): u[i,j] = 0
31:     i = Ix_start
32:     for j in range(Iy_start, Iy_end+1): u[i,j] = 0
33:     i = Ix_start
34:     for j in range(Iy_start, Iy_end+1): u[i,j] = 0
35:     i = Ix_end
36:     for j in range(Iy_start, Iy_end+1): u[i,j] = 0
37:     return u
```

**Fig. C.1** Visual illustration of Cython's ability to translate Python to C.

You can also inspect the generated C code directly, as it appears in the file `wave2D_u0_loop_cy.c`. Nevertheless, understanding this C code requires some familiarity with writing Python extension modules in C by hand. Deep down in the file we can see in detail how the compute-intensive statements have been translated into some complex C code that

is quite different from what a human would write (at least if a direct correspondence to the mathematical notation was intended).

### C.5.3 Building the extension module

Cython code must be translated to C, compiled, and linked to form what is known in the Python world as a *C extension module*. This is usually done by making a `setup.py` script, which is the standard way of building and installing Python software. For an extension module arising from Cython code, the following `setup.py` script is all we need to build and install the module:

```
from distutils.core import setup
from distutils.extension import Extension
from Cython.Distutils import build_ext

cymodule = 'wave2D_u0_loop_cy'
setup(
    name=cymodule
    ext_modules=[Extension(cymodule, [cymodule + '.pyx'],)],
    cmdclass={'build_ext': build_ext},
)
```

We run the script by

---

Terminal

Terminal> python setup.py build\_ext --inplace

---

The `-inplace` option makes the extension module available in the current directory as the file `wave2D_u0_loop_cy.so`. This file acts as a normal Python module that can be imported and inspected:

```
>>> import wave2D_u0_loop_cy
>>> dir(wave2D_u0_loop_cy)
['__builtins__', '__doc__', '__file__', '__name__',
 '__package__', '__test__', 'advance', 'np']
```

The important output from the `dir` function is our Cython function `advance` (the module also features the imported `numpy` module under the name `np` as well as many standard Python objects with double underscores in their names).

The `setup.py` file makes use of the `distutils` package in Python and Cython's extension of this package. These tools know how Python was built on the computer and will use compatible compiler(s) and options when building other code in Cython, C, or C++. Quite some experience

with building large program systems is needed to do the build process manually, so using a `setup.py` script is strongly recommended.

### Simplified build of a Cython module

When there is no need to link the C code with special libraries, Cython offers a shortcut for generating and importing the extension module:

```
import pyximport; pyximport.install()
```

This makes the `setup.py` script redundant. However, in the `wave2D_u0_adv.py` code we do not use `pyximport` and require an explicit build process of this and many other modules.

### C.5.4 Calling the Cython function from Python

The `wave2D_u0_loop_cy` module contains our `advance` function, which we now may call from the Python program for the wave equation:

```
import wave2D_u0_loop_cy
advance = wave2D_u0_loop_cy.advance
...
for n in It[1:-1]:                      # time loop
    f_a[:, :] = f(xv, yv, t[n])          # precompute, size as u
    u = advance(u, u_1, u_2, f_a, x, y, t, Cx2, Cy2, dt2)
```

**Efficiency.** For a mesh consisting of  $120 \times 120$  cells, the scalar Python code require 1370 CPU time units, the vectorized version requires 5.5, while the Cython version requires only 1! For a smaller mesh with  $60 \times 60$  cells Cython is about 1000 times faster than the scalar Python code, and the vectorized version is about 6 times slower than the Cython version.

## C.6 Migrating loops to Fortran

Instead of relying on Cython's (excellent) ability to translate Python to C, we can invoke a compiled language directly and write the loops ourselves. Let us start with Fortran 77, because this is a language with more convenient array handling than C (or plain C++), because we can use the same multi-dimensional indices in the Fortran code as in

the `numpy` arrays in the Python code, while in C these arrays are one-dimensional and requires us to reduce multi-dimensional indices to a single index.

### C.6.1 The Fortran subroutine

We write a Fortran subroutine `advance` in a file `wave2D_u0_loop_f77.f` for implementing the updating formula (2.117) and setting the solution to zero at the boundaries:

```

subroutine advance(u, u_1, u_2, f, Cx2, Cy2, dt2, Nx, Ny)
integer Nx, Ny
real*8 u(0:Nx,0:Ny), u_1(0:Nx,0:Ny), u_2(0:Nx,0:Ny)
real*8 f(0:Nx,0:Ny), Cx2, Cy2, dt2
integer i, j
real*8 u_xx, u_yy
Cf2py intent(in, out) u

C     Scheme at interior points
do j = 1, Ny-1
    do i = 1, Nx-1
        u_xx = u_1(i-1,j) - 2*u_1(i,j) + u_1(i+1,j)
        u_yy = u_1(i,j-1) - 2*u_1(i,j) + u_1(i,j+1)
        u(i,j) = 2*u_1(i,j) - u_2(i,j) + Cx2*u_xx + Cy2*u_yy +
        & dt2*f(i,j)
    end do
end do

C     Boundary conditions
j = 0
do i = 0, Nx
    u(i,j) = 0
end do
j = Ny
do i = 0, Nx
    u(i,j) = 0
end do
i = 0
do j = 0, Ny
    u(i,j) = 0
end do
i = Nx
do j = 0, Ny
    u(i,j) = 0
end do
return
end

```

This code is plain Fortran 77, except for the special `Cf2py` comment line, which here specifies that `u` is both an input argument *and* an object to

be returned from the `advance` routine. Or more precisely, Fortran is not able return an array from a function, but we need a *wrapper code* in C for the Fortran subroutine to enable calling it from Python, and from this wrapper code one can return `u` to the calling Python code.

### Tip: Return all computed objects to the calling code

It is not strictly necessary to return `u` to the calling Python code since the `advance` function will modify the elements of `u`, but the convention in Python is to get all output from a function as returned values. That is, the right way of calling the above Fortran subroutine from Python is

```
u = advance(u, u_1, u_2, f, Cx2, Cy2, dt2)
```

The less encouraged style, which works and resembles the way the Fortran subroutine is called from Fortran, reads

```
advance(u, u_1, u_2, f, Cx2, Cy2, dt2)
```

## C.6.2 Building the Fortran module with `f2py`

The nice feature of writing loops in Fortran is that, without much effort, the tool `f2py` can produce a C extension module such that we can call the Fortran version of `advance` from Python. The necessary commands to run are

---

Terminal
----------

```
Terminal> f2py -m wave2D_u0_loop_f77 -h wave2D_u0_loop_f77.pyf \
           --overwrite-signature wave2D_u0_loop_f77.f
Terminal> f2py -c wave2D_u0_loop_f77.pyf --build-dir build_f77 \
           -DF2PY_REPORT_ON_ARRAY_COPY=1 wave2D_u0_loop_f77.f
```

---

The first command asks `f2py` to interpret the Fortran code and make a Fortran 90 specification of the extension module in the file `wave2D_u0_loop_f77.pyf`. The second command makes `f2py` generate all necessary wrapper code, compile our Fortran file and the wrapper code, and finally build the module. The build process takes place in the specified subdirectory `build_f77` so that files can be inspected if something goes wrong. The option `-DF2PY_REPORT_ON_ARRAY_COPY=1` makes `f2py` write a message for every array that is copied in the commu-

nication between Fortran and Python, which is very useful for avoiding unnecessary array copying (see below). The name of the module file is `wave2D_u0_loop_f77.so`, and this file can be imported and inspected as any other Python module:

```
>>> import wave2D_u0_loop_f77
>>> dir(wave2D_u0_loop_f77)
['__doc__', '__file__', '__name__', '__package__',
 '__version__', 'advance']
>>> print wave2D_u0_loop_f77.__doc__
This module 'wave2D_u0_loop_f77' is auto-generated with f2py....
Functions:
  u = advance(u,u_1,u_2,f,cx2,cy2,dt2,
              nx=(shape(u,0)-1),ny=(shape(u,1)-1))
```

### Examine the doc strings!

Printing the doc strings of the module and its functions is extremely important after having created a module with `f2py`. The reason is that `f2py` makes Python interfaces to the Fortran functions that are different from how the functions are declared in the Fortran code (!). The rationale for this behavior is that `f2py` creates *Pythonic* interfaces such that Fortran routines can be called in the same way as one calls Python functions. Output data from Python functions is always returned to the calling code, but this is technically impossible in Fortran. Also, arrays in Python are passed to Python functions without their dimensions because that information is packed with the array data in the array objects. This is not possible in Fortran, however. Therefore, `f2py` removes array dimensions from the argument list, and `f2py` makes it possible to return objects back to Python.

Let us follow the advice of examining the doc strings and take a close look at the documentation `f2py` has generated for our Fortran `advance` subroutine:

```
>>> print wave2D_u0_loop_f77.advance.__doc__
This module 'wave2D_u0_loop_f77' is auto-generated with f2py
Functions:
  u = advance(u,u_1,u_2,f,cx2,cy2,dt2,
              nx=(shape(u,0)-1),ny=(shape(u,1)-1))

  advance - Function signature:
    u = advance(u,u_1,u_2,f,cx2,cy2,dt2,[nx,ny])
Required arguments:
```

```

u : input rank-2 array('d') with bounds (nx + 1,ny + 1)
u_1 : input rank-2 array('d') with bounds (nx + 1,ny + 1)
u_2 : input rank-2 array('d') with bounds (nx + 1,ny + 1)
f : input rank-2 array('d') with bounds (nx + 1,ny + 1)
cx2 : input float
cy2 : input float
dt2 : input float
Optional arguments:
  nx := (shape(u,0)-1) input int
  ny := (shape(u,1)-1) input int
Return objects:
  u : rank-2 array('d') with bounds (nx + 1,ny + 1)

```

Here we see that the `nx` and `ny` parameters declared in Fortran are optional arguments that can be omitted when calling `advance` from Python.

We strongly recommend to print out the documentation of *every* Fortran function to be called from Python and make sure the call syntax is exactly as listed in the documentation.

### C.6.3 How to avoid array copying

Multi-dimensional arrays are stored as a stream of numbers in memory. For a two-dimensional array consisting of rows and columns there are two ways of creating such a stream: *row-major ordering*, which means that rows are stored consecutively in memory, or *column-major ordering*, which means that the columns are stored one after each other. All programming languages inherited from C, including Python, apply the row-major ordering, but Fortran uses column-major storage. Thinking of a two-dimensional array in Python or C as a matrix, it means that Fortran works with the transposed matrix.

Fortunately, `f2py` creates extra code so that accessing `u(i,j)` in the Fortran subroutine corresponds to the element `u[i,j]` in the underlying `numpy` array (without the extra code, `u(i,j)` in Fortran would access `u[j,i]` in the `numpy` array). Technically, `f2py` takes a copy of our `numpy` array and reorders the data before sending the array to Fortran. Such copying can be costly. For 2D wave simulations on a  $60 \times 60$  grid the overhead of copying is a factor of 5, which means that almost the whole performance gain of Fortran over vectorized `numpy` code is lost!

To avoid having `f2py` to copy arrays with C storage to the corresponding Fortran storage, we declare the arrays with Fortran storage:

```

order = 'Fortran' if version == 'f77' else 'C'
u    = zeros((Nx+1,Ny+1), order=order)  # solution array

```

```
u_1 = zeros((Nx+1,Ny+1), order=order) # solution at t-dt
u_2 = zeros((Nx+1,Ny+1), order=order) # solution at t-2*dt
```

In the compile and build step of using f2py, it is recommended to add an extra option for making f2py report on array copying:

---

Terminal	
----------	--

```
Terminal> f2py -c wave2D_u0_loop_f77.pyf --build-dir build_f77 \
-DF2PY_REPORT_ON_ARRAY_COPY=1 wave2D_u0_loop_f77.f
```

---

It can sometimes be a challenge to track down which array that causes a copying. There are two principal reasons for copying array data: either the array does not have Fortran storage or the element types do not match those declared in the Fortran code. The latter cause is usually effectively eliminated by using `real*8` data in the Fortran code and `float64` (the default `float` type in `numpy`) in the arrays on the Python side. The former reason is more common, and to check whether an array before a Fortran call has the right storage one can print the result of `isfortran(a)`, which is `True` if the array `a` has Fortran storage.

Let us look at an example where we face problems with array storage. A typical problem in the `wave2D_u0.py` code is to set

```
f_a = f(xv, yv, t[n])
```

before the call to the Fortran `advance` routine. This computation creates a new array with C storage. An undesired copy of `f_a` will be produced when sending `f_a` to a Fortran routine. There are two remedies, either direct insertion of data in an array with Fortran storage,

```
f_a = zeros((Nx+1, Ny+1), order='Fortran')
...
f_a[:, :] = f(xv, yv, t[n])
```

or remaking the `f(xv, yv, t[n])` array,

```
f_a = asarray(f(xv, yv, t[n]), order='Fortran')
```

The former remedy is most efficient if the `asarray` operation is to be performed a large number of times.

**Efficiency.** The efficiency of this Fortran code is very similar to the Cython code. There is usually nothing more to gain, from a computational efficiency point of view, by implementing the *complete* Python program in Fortran or C. That will just be a lot more code for all administering work that is needed in scientific software, especially if we extend our sample program `wave2D_u0.py` to handle a real scientific problem. Then

only a small portion will consist of loops with intensive array calculations. These can be migrated to Cython or Fortran as explained, while the rest of the programming can be more conveniently done in Python.

## C.7 Migrating loops to C via Cython

The computationally intensive loops can alternatively be implemented in C code. Just as Fortran calls for care regarding the storage of two-dimensional arrays, working with two-dimensional arrays in C is a bit tricky. The reason is that `numpy` arrays are viewed as one-dimensional arrays when transferred to C, while C programmers will think of `u`, `u_1`, and `u_2` as two dimensional arrays and index them like `u[i] [j]`. The C code must declare `u` as `double* u` and translate an index pair `[i] [j]` to a corresponding single index when `u` is viewed as one-dimensional. This translation requires knowledge of how the numbers in `u` are stored in memory.

### C.7.1 Translating index pairs to single indices

Two-dimensional `numpy` arrays with the default C storage are stored row by row. In general, multi-dimensional arrays with C storage are stored such that the last index has the fastest variation, then the next last index, and so on, ending up with the slowest variation in the first index. For a two-dimensional `u` declared as `zeros((Nx+1,Ny+1))` in Python, the individual elements are stored in the following order:

```
u[0,0], u[0,1], u[0,2], ..., u[0,Ny], u[1,0], u[1,1], ...,
u[1,Ny], u[2,0], ..., u[Nx,0], u[Nx,1], ..., u[Nx, Ny]
```

Viewing `u` as one-dimensional, the index pair  $(i, j)$  translates to  $i(N_y + 1) + j$ . So, where a C programmer would naturally write an index `u[i] [j]`, the indexing must read `u[i*(Ny+1) + j]`. This is tedious to write, so it can be handy to define a C macro,

```
#define idx(i,j) (i)*(Ny+1) + j
```

so that we can write `u[idx(i,j)]`, which reads much better and is easier to debug.

### Be careful with macro definitions

Macros just perform simple text substitutions: `idx(hello,world)` is expanded to `(hello)*(Ny+1) + world`. The parenthesis in `(i)` are essential - using the natural mathematical formula `i*(Ny+1) + j` in the macro definition, `idx(i-1,j)` would expand to `i-1*(Ny+1) + j`, which is the wrong formula. Macros are handy, but requires careful use. In C++, inline functions are safer and replace the need for macros.

### C.7.2 The complete C code

The C version of our function `advance` can be coded as follows.

```
#define idx(i,j) (i)*(Ny+1) + j

void advance(double* u, double* u_1, double* u_2, double* f,
            double Cx2, double Cy2, double dt2, int Nx, int Ny)
{
    int i, j;
    double u_xx, u_yy;
    /* Scheme at interior points */
    for (i=1; i<=Nx-1; i++) {
        for (j=1; j<=Ny-1; j++) {
            u_xx = u_1[idx(i-1,j)] - 2*u_1[idx(i,j)] + u_1[idx(i+1,j)];
            u_yy = u_1[idx(i,j-1)] - 2*u_1[idx(i,j)] + u_1[idx(i,j+1)];
            u[idx(i,j)] = 2*u_1[idx(i,j)] - u_2[idx(i,j)] +
                Cx2*u_xx + Cy2*u_yy + dt2*f[idx(i,j)];
        }
    }
    /* Boundary conditions */
    j = 0; for (i=0; i<=Nx; i++) u[idx(i,j)] = 0;
    j = Ny; for (i=0; i<=Nx; i++) u[idx(i,j)] = 0;
    i = 0; for (j=0; j<=Ny; j++) u[idx(i,j)] = 0;
    i = Nx; for (j=0; j<=Ny; j++) u[idx(i,j)] = 0;
}
```

### C.7.3 The Cython interface file

All the code above appears in a file `wave2D_u0_loop_c.c`. We need to compile this file together with C wrapper code such that `advance` can be called from Python. Cython can be used to generate appropriate wrapper code. The relevant Cython code for interfacing C is placed in a file with

extension .pyx. Here this file, called `wave2D_u0_loop_c_cy.pyx`, looks like

```
import numpy as np
cimport numpy as np
cimport cython

cdef extern from "wave2D_u0_loop_c.h":
    void advance(double* u, double* u_1, double* u_2, double* f,
                double Cx2, double Cy2, double dt2,
                int Nx, int Ny)

@cython.boundscheck(False)
@cython.wraparound(False)
def advance_cwrap(
    np.ndarray[double, ndim=2, mode='c'] u,
    np.ndarray[double, ndim=2, mode='c'] u_1,
    np.ndarray[double, ndim=2, mode='c'] u_2,
    np.ndarray[double, ndim=2, mode='c'] f,
    double Cx2, double Cy2, double dt2):
    advance(&u[0,0], &u_1[0,0], &u_2[0,0], &f[0,0],
           Cx2, Cy2, dt2,
           u.shape[0]-1, u.shape[1]-1)
    return u
```

We first declare the C functions to be interfaced. These must also appear in a C header file, `wave2D_u0_loop_c.h`,

```
extern void advance(double* u, double* u_1, double* u_2, double* f,
                    double Cx2, double Cy2, double dt2,
                    int Nx, int Ny);
```

The next step is to write a Cython function with Python objects as arguments. The name `advance` is already used for the C function so the function to be called from Python is named `advance_cwrap`. The contents of this function is simply a call to the `advance` version in C. To this end, the right information from the Python objects must be passed on as arguments to `advance`. Arrays are sent with their C pointers to the first element, obtained in Cython as `&u[0,0]` (the `&` takes the address of a C variable). The `Nx` and `Ny` arguments in `advance` are easily obtained from the shape of the `numpy` array `u`. Finally, `u` must be returned such that we can set `u = advance(...)` in Python.

### C.7.4 Building the extension module

It remains to build the extension module. An appropriate `setup.py` file is

```

from distutils.core import setup
from distutils.extension import Extension
from Cython.Distutils import build_ext

sources = ['wave2D_u0_loop_c.c', 'wave2D_u0_loop_c_pyx.pyx']
module = 'wave2D_u0_loop_c_pyx'
setup(
    name=module,
    ext_modules=[Extension(module, sources,
                           libraries=[], # C libs to link with
                           )],
    cmdclass={'build_ext': build_ext},
)

```

All we need to specify is the .c file(s) and the .pyx interface file. Cython is automatically run to generate the necessary wrapper code. Files are then compiled and linked to an extension module residing in the file `wave2D_u0_loop_c_pyx.so`. Here is a session with running `setup.py` and examining the resulting module in Python

---

Terminal> python setup.py build_ext --inplace	Terminal
---	----------

```

Terminal> python setup.py build_ext --inplace
Terminal> python
>>> import wave2D_u0_loop_c_pyx as m
>>> dir(m)
['__builtins__', '__doc__', '__file__', '__name__', '__package__',
 '__test__', 'advance_cwrap', 'np']

```

---

The call to the C version of `advance` can go like this in Python:

```

import wave2D_u0_loop_c_pyx
advance = wave2D_u0_loop_c_pyx.advance_cwrap
...
f_a[:, :] = f(xv, yv, t[n])
u = advance(u, u_1, u_2, f_a, Cx2, Cy2, dt2)

```

**Efficiency.** In this example, the C and Fortran code runs at the same speed, and there are no significant differences in the efficiency of the wrapper code. The overhead implied by the wrapper code is negligible as long as we do not work with very small meshes and consequently little numerical work in the `advance` function.

## C.8 Migrating loops to C via f2py

An alternative to using Cython for interfacing C code is to apply `f2py`. The C code is the same, just the details of specifying how it is to be

called from Python differ. The `f2py` tool requires the call specification to be a Fortran 90 module defined in a `.pyf` file. This file was automatically generated when we interfaced a Fortran subroutine. With a C function we need to write this module ourselves, or we can use a trick and let `f2py` generate it for us. The trick consists in writing the signature of the C function with Fortran syntax and place it in a Fortran file, here `wave2D_u0_loop_c_f2py_signature.f`:

```

subroutine advance(u, u_1, u_2, f, Cx2, Cy2, dt2, Nx, Ny)
Cf2py intent(c) advance
integer Nx, Ny, N
real*8 u(0:Nx,0:Ny), u_1(0:Nx,0:Ny), u_2(0:Nx,0:Ny)
real*8 f(0:Nx, 0:Ny), Cx2, Cy2, dt2
Cf2py intent(in, out) u
Cf2py intent(c) u, u_1, u_2, f, Cx2, Cy2, dt2, Nx, Ny
return
end
```

Note that we need a special `f2py` instruction, through a `Cf2py` comment line, to specify that all the function arguments are C variables. We also need to tell that the function is actually in C: `intent(c) advance`.

Since `f2py` is just concerned with the function signature and not the complete contents of the function body, it can easily generate the Fortran 90 module specification based solely on the signature above:

---

Terminal
Terminal> f2py -m wave2D\_u0\_loop\_c\_f2py \
-h wave2D\_u0\_loop\_c\_f2py.pyf --overwrite-signature \
wave2D\_u0\_loop\_c\_f2py\_signature.f

---

The compile and build step is as for the Fortran code, except that we list C files instead of Fortran files:

---

Terminal
Terminal> f2py -c wave2D\_u0\_loop\_c\_f2py.pyf \
--build-dir tmp\_build\_c \
-DF2PY\_REPORT\_ON\_ARRAY\_COPY=1 wave2D\_u0\_loop\_c.c

---

As when interfacing Fortran code with `f2py`, we need to print out the doc string to see the exact call syntax from the Python side. This doc string is identical for the C and Fortran versions of `advance`.

### C.8.1 Migrating loops to C++ via f2py

C++ is a much more versatile language than C or Fortran and has over the last two decades become very popular for numerical computing.

Many will therefore prefer to migrate compute-intensive Python code to C++. This is, in principle, easy: just write the desired C++ code and use some tool for interfacing it from Python. A tool like [SWIG](#) can interpret the C++ code and generate interfaces for a wide range of languages, including Python, Perl, Ruby, and Java. However, SWIG is a comprehensive tool with a correspondingly steep learning curve. Alternative tools, such as [Boost Python](#), [SIP](#), and [Shiboken](#) are similarly comprehensive. Simpler tools include [PyBindGen](#),

A technically much easier way of interfacing C++ code is to drop the possibility to use C++ classes directly from Python, but instead make a C interface to the C++ code. The C interface can be handled by [f2py](#) as shown in the example with pure C code. Such a solution means that classes in Python and C++ cannot be mixed and that only primitive data types like numbers, strings, and arrays can be transferred between Python and C++. Actually, this is often a very good solution because it forces the C++ code to work on array data, which usually gives faster code than if fancy data structures with classes are used. The arrays coming from Python, and looking like plain C/C++ arrays, can be efficiently wrapped in more user-friendly C++ array classes in the C++ code, if desired.

## C.9 Exercises

### Exercise C.1: Make an improved `numpy.savez` function

The `numpy.savez` function can save multiple arrays to a zip archive. Unfortunately, if we want to use `savez` in time-dependent problems and call it multiple times (once per time level), each call leads to a separate zip archive. It is more convenient to have all arrays in one archive, which can be read by `numpy.load`. Section C.2 provides a recipe for merging all the individual zip archives into one archive. An alternative is to write a new `savez` function that allows multiple calls and storage into the same archive prior to a final `close` method to close the archive and make it ready for reading. Implement such an improved `savez` function as a class `Savez`.

The class should pass the following unit test:

```
def test_Savez():
    import tempfile, os
    tmp = 'tmp_testarchive'
```

```
database = Savez(tmp)
for i in range(4):
    array = np.linspace(0, 5+i, 3)
    kwargs = {'myarray_%02d' % i: array}
    database.savez(**kwargs)
database.close()

database = np.load(tmp+'.npz')

expected = {
    'myarray_00': np.array([ 0.,  2.5,  5. ]),
    'myarray_01': np.array([ 0.,  3.,  6.]),
    'myarray_02': np.array([ 0.,  3.5,  7.]),
    'myarray_03': np.array([ 0.,  4.,  8.]),
}
for name in database:
    computed = database[name]
    diff = np.abs(expected[name] - computed).max()
    assert diff < 1E-13
database.close()
os.remove(tmp+'.npz')
```

**Hint.** Study the source code for function `savez` (or more precisely, function `_savez`).

Filename: `Savez`.



## References

- [1] R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. Van der Vorst. *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*. SIAM, second edition, 1994. [http://www.netlib.org/linalg/html\\_templates/Templates.html](http://www.netlib.org/linalg/html_templates/Templates.html).
- [2] C. Greif and U. M. Ascher. *A First Course in Numerical Methods*. Computational Science and Engineering. SIAM, 2011.
- [3] C. T. Kelley. *Iterative Methods for Linear and Nonlinear Equations*. SIAM, 1995.
- [4] H. P. Langtangen. *Finite Difference Computing with Exponential Decay Models*. Springer, 2015. <http://tinyurl.com/nclmcng/web>.
- [5] H. P. Langtangen and G. K. Pedersen. *Scaling of Differential Equations*. 2015. <http://tinyurl.com/qfjgxmfp/web>.
- [6] R. LeVeque. *Finite Difference Methods for Ordinary and Partial Differential Equations: Steady-State and Time-Dependent Problems*. SIAM, 2007.
- [7] R. Rannacher. Finite element solution of diffusion problems with irregular data. *Numerische Mathematik*, 43:309–327, 1984.
- [8] Y. Saad. *Iterative Methods for Sparse Linear Systems*. SIAM, second edition, 2003. [http://www-users.cs.umn.edu/~saad/IterMethBook\\_2ndEd.pdf](http://www-users.cs.umn.edu/~saad/IterMethBook_2ndEd.pdf).
- [9] L. N. Trefethen. *Trefethen's index cards - Forty years of notes about People, Words and Mathematics*. World Scientific, 2011.



# Index

- alternating mesh, 362
- amplification factor, 249
- animation, 13
- `argparse` (Python module), 56
- `ArgumentParser` (Python class), 56
  - arithmetic mean, 141
  - array computing, 117
  - array slices, 117
  - averaging
    - arithmetic, 141
    - geometric, 52, 141
    - harmonic, 141
- boundary condition
  - open (radiation), 155
- boundary conditions
  - Dirichlet, 129
  - Neumann, 129
  - periodic, 157
- C extension module, 496
- C/Python array storage, 501
- callback function, 107
- centered difference, 3
- Cholesky factorization, 303
- closure, 112
- column-major ordering, 501
- continuation method, 413, 421
- correction terms, 445
- Courant number, 167
- Cython, 491
- `cython -a` (Python-C translation in HTML), 495
- decay ODE, 438
- declaration of variables in Cython, 493
- diffusion equation, 1D, 219
- diffusion limit of random walk, 315
- Dirichlet conditions, 129
- discrete Fourier transform, 164
- `distutils`, 496
- DOF (degree of freedom), 67
- energy estimates (diffusion), 342
- energy principle, 39
- error

- global, 26
- explicit discretization methods, 221
- finite differences
  - backward, 432
  - centered, 3, 434
  - forward, 433
- fixed-point iteration, 376
- Flash (video format), 13
- Fokker-Planck equation, 329
- forced vibrations, 50
- Fortran array storage, 501
- Fortran subroutine, 498
- Forward Euler scheme, 221
- forward-backward Euler-Cromer scheme, 44
- Fourier series, 164
- Fourier transform, 164
- frequency (of oscillations), 2
- Gauss-Seidel method, 297
- geometric mean, 52, 141
- harmonic average, 141
- heat equation, 1D, 219
- homogeneous Dirichlet conditions, 129
- homogeneous Neumann conditions, 129
- HTML5 video tag, 14
- Hz (unit), 2
- implicit discretization methods, 234
- index set notation, 132, 183
- interrupt a program by Ctrl+c, 327
- Jacobi method, 290
- lambda function (Python), 121
- linearization, 376
- explicit time integration, 374
- fixed-point iteration, 376
- Picard iteration, 376
- successive substitutions, 376
- LU factorization, 303
- making movies, 13
- mechanical energy, 39
- mechanical vibrations, 2
- mesh
  - finite differences, 2, 94
- mesh function, 2, 95
- MP4 (video format), 13
- Neumann conditions, 129
- nonlinear restoring force, 50
- nonlinear spring, 50
- nose, 7, 109
- Ogg (video format), 13
- open boundary condition, 155
- oscillations, 2
- parallelism, 117
- period (of oscillations), 2
- periodic boundary conditions, 157
- phase plane plot, 35
- Picard iteration, 376
- Plotter** class (SciTools), 327
- preconditioning, 306, 345
- pytest, 7, 109
- radiation condition, 155
- random walk, 309
- red-black numbering, 299
- relaxation (nonlinear equations), 382
- resonance, 88
- Richardson iteration, 345
- row-major ordering, 501
- scalar code, 117
- scitools movie** command, 15

- `scitools.avplotter`, 327  
seed (random numbers), 313  
`setup.py`, 496  
single Picard iteration technique,  
    378  
slice, 117  
SOR method, 298  
sparse matrix, 274  
stability criterion, 27, 168  
staggered Euler-Cromer scheme,  
    362  
staggered mesh, 362  
stationary solution, 219  
stencil  
    1D wave equation, 95  
    Neumann boundary, 130  
stochastic difference equation, 328  
stochastic ODE, 329  
stopping criteria (nonlinear prob-  
    lems), 378, 394  
successive substitutions, 376
- test function, 7, 109  
truncation error  
    Backward Euler scheme, 432  
    correction terms, 445  
    Crank-Nicolson scheme, 434  
    Forward Euler scheme, 433  
    general, 430  
    table of formulas, 435
- unit testing, 7, 109  
upwind difference, 355
- vectorization, 117, 312  
verification, 313, 450  
    convergence rates, 8  
    hand calculations, 7  
    polynomial solution, 109  
    polynomial solutions, 8
- vibration ODE, 2  
video formats, 13  
wave equation  
    1D, 93  
    1D, analytical properties, 162  
    1D, exact numerical solution,  
        166  
    1D, finite difference method,  
        94  
    1D, implementation, 106  
    1D, stability, 168  
    2D, implementation, 181
- waves  
    on a string, 93  
WebM (video format), 13  
Wiener process, 329  
wrapper code, 498