

## Lecture notes FYS3150 – Computational Physics, fall 2023

### Introduction

#### Welcome to this course!

- About me:
  - Anders Kvellestad
  - Researcher in the Section for Theoretical Physics
  - Background: Bergen → Oslo → Stockholm → Oslo → London → Oslo
  - Work on exploring new theories in particle physics
  - Keywords: LHC, supersymmetry, dark matter, Higgs, statistics, coding (Python, C++, ...), supercomputers, causing and fixing bugs, responsible for coffee supplies in the theory group
- The teaching team this semester:
  - Even Marius Nordhagen
  - David Richard Shope
  - Ingvild Bergsbak
  - Carl Andreas Lindstrøm
  - Nils Enric Canut Taugbøl
  - Felix Forseth
- In addition, we will have a guest lecture from Anna Kathinka Dalland Evans on how to write scientific reports
- **Question:** Who are you?
  - Study programmes?
  - Level of coding experience?
  - Main motivation for this course?
    - \* Solve those pesky equations
    - \* Learn C++ and other tools
    - \* I just like working with computers
- **Question:** What operating system are you using?
  - Linux?
  - macOS?

- Windows?
- I will need at least two student representatives for the course evaluation
  - You'll join a meeting (~1 hour) with us teachers at the end of the semester, where we discuss what worked and what we can improve in the course
  - If you are willing to do this, just send me an email

## About the course

- Course resources
  - Official UiO course page
  - Our own page, with course material
  - Our Git repo
  - Canvas, for handing in reports
- Teaching language: English
- Programming languages:
  - Main focus on C++
  - Python for data analysis, making plots, etc.
  - Bash for terminal examples, short scripts, etc.
- You *can* use Python for the projects, but we strongly recommend C++
- All lectures and group sessions will assume that you use C++
- This course has been taught by CompPhys guru Morten Hjorth-Jensen for many years
- I took over this course in 2021
- I follow Morten's old course fairly closely, but with a number of personal tweaks from my side
- Course background material: Morten's lecture notes / book draft, available via the UiO course page
  - I will often point you to relevant parts of Morten's notes
  - But our main curriculum is what we discuss in the lectures and as part of the projects
- Course philosophy:
  - Pragmatic, learning by doing (and learning by failing)
  - Will try to focus on concrete examples

- Computational physics is a *huge* and highly active field — this course is just a first introduction
- Lectures:
  - Thursdays and Fridays, 10.15 – 12.00
  - New this year: Will try to lecture in ~30 minute sessions, with two short breaks
  - Lectures are *not* recorded – if you miss a lecture, read the detailed lecture notes
- Group sessions:
  - Also Thursdays and Fridays
  - Probably the most important arena for learning and mastering this course!
  - Four two-hour sessions per week
  - You can come to any group session you want
    - \* Try to avoid all going to the same session
    - \* Be patient with our group teachers — some have taught this course for several years, others are doing it for the first time
  - We start the group sessions already the first week
- Formal requirements
  - Two **problem sets**: Must be **passed**
  - Three **projects**: **Scored** from 0–100
  - Final grade based on weighted average of the project scores. Weighting: 20%, 40%, 40%
- For simplicity, we'll just refer to everything as “projects”, i.e. we'll talk about projects 1–5, and the grade is based on projects 3–5.
- *Tentative* deadlines:
  - Project 1: September 12
  - Project 2: September 26
  - Project 3: October 24
  - Project 4: November 21
  - Project 5: December 12
- Policy on deadlines: **friendly, but strict**
  - Need to be strict on deadlines, both to keep things fair and to keep up with our time schedule
  - There are no second attempts
  - Substantial deadline extensions due to illness require a doctor's note
- **Collaboration is encouraged!**

- We *strongly* encourage you to collaborate in small groups of 2–4 people.
- 3 people per group is ideal
- A group hands in a *joint* project report and code
- By working together *you will learn more*, and we get more time for grading per project report → *more detailed feedback from us*
- Asking questions:
  - **Please ask questions!**
  - Any time during lectures — just cut in and ask
  - For help with your specific project/code:
    - \* Primary forum: group sessions
    - \* Secondary forum: our online discussion forum
  - *Think and try yourself* before you ask for help
  - When writing questions:
    - \* Keep it short and concise
    - \* Have respect for other people’s time (your fellow students, the teachers, ...)
  - *Almost* all questions are welcome: The only type of questions I don’t like are the “questions” that aren’t really questions at all, but just someone trying to show off how much they know.
  - Any personal or procedural issues: Send me an email. (We can also set up a meeting.)
- The broad topics of this course:
  - Learn basic C++, with focus on numerics
  - Matrix operations, eigenvalue problems
  - Solve ordinary and partial differential equations
  - Numerical integration
  - Monte Carlo methods, simulation of stochastic systems
  - Proper presentation of results
  - Debugging :)
- This course is good for your CV (beyond just the grade you get)
  - We’re at a university, so hopefully our main motivation for following/teaching a course should be that learning new stuff is interesting and valuable in itself — that’s at least my main motivation when teaching this!
  - Having said that, after completing this course you can probably also add some new points to your CV:
    - \* Experience with C++
    - \* Experience with the Unix terminal

- \* Experience with git and GitHub
- \* Experience with writing technical reports in LaTeX
- \* ...

## The most useful advice you'll get all year

- Something you don't understand?
  - *Read and think*
  - *Discuss* with your fellow students
  - *Ask us*
- Code isn't working?
  - Don't just try stuff at random!
    - \* This rarely works, and when it does you typically still can't trust the results...
  - *Read the documentation* for the command/tool you are using
  - *Search online for the error message*, after removing things that are specific to your code (variable names, file names, etc.)
    - \* *Read the explanations* you find, don't just copy code
  - Try to isolate and reproduce the problem in a small, separate example code. (*A minimal working example.*)
  - Read the course pages on debugging
  - We'll also discuss debugging in the lectures
- How you present your results **really matters**
  - Quality of language
  - Quality of figures
  - Layout
  - Report structure
  - Referencing
  - Code comments and documentation

**Figure 1:** Presentation quality matters

- Spend time with pen and paper before you start coding
  - Make a rough sketch of program parts and flow
  - Sketch your program with code comments first, then start filling in the code
  - Make a sketch of discretisations, to avoid mistakes with indices

$\begin{array}{cccccc} | & | & | & | & | & | \\ x_0 & x_1 & x_2 & x_3 & x_4 & x_5 \end{array}$

- Boundary cond. at  $x_0$  and  $x_5$
- 6 elements in x array
- X range is split in 5 steps

**Figure 2:** Sketch discretisations

- Make sure you understand the quantities you present in plots and tables
  - Makes it much easier to spot mistakes
  - Pay attention to units!
  - *Tip:* Always set axis ranges manually
- Read the report template we provide, plus the example student reports
- And read the *Checklist for reports* page on our webpage, to avoid many common mistakes

## Plagiarism

- Plagiarism is **very serious**

- Have seen a few cases in the past
- Can have very serious consequences, e.g. losing the right to study
- You must:
  - Write your own text — never copy text from others (unless it is marked a direct quote)
  - Write your own code, unless it's code we have provided to help
  - Always acknowledge contributions from others
  - Properly cite articles, books, webpages, ...
    - \* We'll discuss this more in detail when you start writing project reports

## Use of ChatGPT and related tools

- ChatGPT and other AI-based *large language models* (LLMs) can, like any new tool, be used in wise ways and not-so-wise ways
- The policy on LLM use in our course is as follows:
  - If you use an LLM in your work, you need to add to your report a description of what you used the LLM for. This would be part of the *Methods* section of your report. (We'll discuss report writing in detail later in the course.)
  - In the report template we will probably add a dedicated subsection called e.g. *Tools*, where you mention the key tools you have used, and what you have used them for, e.g. sentences like "All figures in this report have been made using the Python package `matplotlib`."
  - If we see that you have used an LLM in ways that you haven't described in the report, this will lead to a lower score, analogous to what happens if you don't provide proper references, or just have a very incomplete description of the methods you've used.
  - **Important:** When you hand in a report or code, you take full responsibility for all the content. That is, you can never put the blame for anything on an LLM model.
- Some advice:
  - Don't use LLMs like ChatGPT as search engines. An LLM is *not* a new, cool way of searching the web. There is no database, no in-built checks for correctness of content, etc.
    - \* So you should *not* use LLM output as a reference for a statement in your report.
  - For you to be able to judge the quality, correctness and appropriateness of some LLM output, you first need to actually build up your own expertise. That is, you need to
    - \* study the given scientific topic

- \* know/learn how to write good texts
- \* know/learn how a given coding language works
- \* ...
- The best way to learn these things is to sit down and do them yourself, mostly from scratch
- *Once* you have built up the necessary expertise, LLMs can become a useful tool for some tasks
- Examples of tasks where an LLM may be useful in this course:
  - \* Help with debugging code problems
  - \* Help with suggesting language improvements (to text that you have already drafted)
- My main advice: Don't use LLMs too much!
  - \* Learning how to use LLMs is itself a useful skill
  - \* But overuse will probably reduce your learning outcome in this course!
  - \* The most “painful” moments in your work – when you work through the math yourself, when you try to formulate a correct and good sentence for your report, when you systematically go through your code to find that one strange bug, or when you think carefully about whether a given result makes sense – these are the moments when you actually learn the most!



## Is it safe to use ChatGPT for your task?

Aleksandr Tiulkanov | January 19, 2023



**Figure 3:** Example considerations to make before using ChatGPT or similar tools. Flowchart by A. Tiulkanov, included in the UNESCO report *ChatGPT and Artificial Intelligence in higher education*.

## In-lecture code discussion #1

- We have two pages with coding resources
  - [anderkve.github.io/FYS3150](https://anderkve.github.io/FYS3150)
  - [github.com/anderkve/FYS3150/tree/master/code\\_examples](https://github.com/anderkve/FYS3150/tree/master/code_examples)
- All code examples I discuss in the lectures can be found in one of these places
- Long code examples, e.g. example programs involving multiple files, are typically found in the `code_examples` directory of our Git repo.
- Make sure to explore these pages on your own! There's lots of help and hints to be found there!
  - In the first group sessions, spend some time going through the different introductory material on [anderkve.github.io/FYS3150](https://anderkve.github.io/FYS3150) before you start on project 1.
- Now let's introduce C++!
  - (Note that we won't have time in the lectures to talk about all C++ details you need for the projects.)
  - Intro:  
[anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/intro](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/intro)
  - Hello World:  
[anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/hello\\_world](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/hello_world)
  - Compiling and linking:  
[anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/compiling\\_and\\_linking\\_take\\_1](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/compiling_and_linking_take_1)
  - Source files and header files:  
[anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/source\\_files\\_and\\_header\\_files](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/source_files_and_header_files)
  - Code structure:  
[anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/code\\_structure](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/code_structure)
    - \* See also this example:  
[github.com/anderkve/FYS3150/tree/master/code\\_examples/code\\_structure/example\\_1](https://github.com/anderkve/FYS3150/tree/master/code_examples/code_structure/example_1)
  - Compilation and linking example with multiple files:  
[github.com/anderkve/FYS3150/tree/master/code\\_examples/compilation\\_linking/example\\_1](https://github.com/anderkve/FYS3150/tree/master/code_examples/compilation_linking/example_1)
  - *Strongly typed* languages (e.g. C++) vs *weakly typed* languages (e.g. Python).
    - \* [anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/variables](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/variables)

- Write to file:

[anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/write\\_to\\_file](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/write_to_file)

- \* Also, remember that in cases with small output, simply *redirecting* terminal output into a file can be an easy and quick way to store output to a file – see the *Write terminal output to file* section of [anderkve.github.io/FYS3150/book/using\\_the\\_terminal/basics](https://anderkve.github.io/FYS3150/book/using_the_terminal/basics)

## Topics in project 1

Some things are covered in the lectures, other things via examples on the webpage

- Discretisation of a continuous problem, in this case the following boundary value problem (BVP):

$$\begin{aligned}-\frac{d^2 u}{dx^2} &= f(x) \\ x &\in [0, 1] \\ u(0) &= 0 \\ u(1) &= 0\end{aligned}$$

- Mathematical approx. to second derivative (suitable for discretisation)
- Connection between a BVP and a standard matrix equation ( $\mathbf{A}\vec{x} = \vec{b}$ ), and approaches to solve this
  - Gaussian elimination
  - LU decomposition
- Errors!
  - Truncation error (purely math)
  - Numerical roundoff error (can't represent numbers with infinite precision on computers)
    - \*  $\rightarrow$  *loss of numerical precision*
- Counting floating-point operations (FLOPs)
- Coding:
  - Working with arrays/vectors and matrices
  - Input/output (nicely formatted output)
  - Timing the code
  - Compilation and linking, basic code design

## Discretisation of continuous functions

- Computers can't represent all possible numbers (finite range and "resolution")  
→ Need to discretise!
- Take some function  $u(x)$ , with  $x \in [x_{\min}, x_{\max}]$ . ( $u(x)$  might e.g. be the solution of our diff. eq. in project 1.)
- $u$  and  $x$  are *continuous* quantities



Figure 4: Continuous function

- Discretised representation



Figure 5: Discretised representation

*Tip:* When testing and debugging your code or trying to understand your results, it's often useful to work with a low number of points (coarse discretisation) and make plots that display your raw data points, i.e. not just directly draw lines between the points.

### My notation

$$\begin{aligned}x &\rightarrow x_i \\ u(x) &\rightarrow u(x_i) \equiv u_i \\ u(x \pm h) &\rightarrow u(x_i \pm h) \equiv u_{i \pm 1}\end{aligned}$$

- So far  $u_i$  is the exact  $u(x)$  at point  $x = x_i$
- Our numerical methods will find an *approximation to the exact*  $u_i$
- We will sometimes call this approximation  $v_i$ , to highlight that this approximation is not the same as the exact  $u_i$

### Basic relations

- $x_i = x_0 + ih$ , with  $i = 0, 1, 2, \dots, n$
- step size:  $h = x_1 - x_0 = \frac{x_2 - x_0}{2} = \dots = \frac{x_n - x_0}{n}$ .  
( $x_0 = x_{\min}, x_n = x_{\max}$ )
- Will sometimes use notation  $\Delta x$  for  $h$
- Remember:  $n$  **steps** corresponds to  $n + 1$  **points**
- Always make a sketch if you are unsure about the discretisation

## Numerical differentiation

See Chapter 3.1 in Morten's notes.

### Main results

#### First derivative:

$$\left. \frac{du}{dx} \right|_{x_i} = u'_i = \frac{u_{i+1} - u_i}{h} + \mathcal{O}(h), \quad (\text{two-point, forward difference})$$

$$\left. \frac{du}{dx} \right|_{x_i} = u'_i = \frac{u_i - u_{i-1}}{h} + \mathcal{O}(h), \quad (\text{two-point, backward difference})$$

$$\left. \frac{du}{dx} \right|_{x_i} = u'_i = \frac{u_{i+1} - u_{i-1}}{2h} + \mathcal{O}(h^2) \quad (\text{three-point})$$

#### Second derivative:

$$\left. \frac{d^2u}{dx^2} \right|_{x_i} = u''_i = \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} + \mathcal{O}(h^2)$$

### Derivation

- Starting point: Taylor expansion of  $u$  around a point  $x$

$$\begin{aligned} u(x+h) &= \sum_{n=0}^{\infty} \frac{1}{n!} u^{(n)}(x) h^n \\ &= u(x) + u'(x)h + \frac{1}{2}u''(x)h^2 + \frac{1}{6}u'''(x)h^3 + \mathcal{O}(h^4) \end{aligned}$$

An aside on notation:

- $u(x+h) = u(x) + u'(x)h + \mathcal{O}(h^2)$ , (exact)
- $u(x+h) \approx u(x) + u'(x)h$ , (approximation, with truncation error  $\mathcal{O}(h^2)$ )

- Can get expression for  $u'(x)$ :

$$u(x+h) = u(x) + u'(x)h + \mathcal{O}(h^2)$$

$$\Rightarrow u'(x) = \frac{u(x+h) - u(x) - \mathcal{O}(h^2)}{h}$$

$$u'(x) = \frac{u(x+h) - u(x)}{h} + \mathcal{O}(h), \quad (\text{note power of } h)$$

Discretise:

$$u(x) \rightarrow u_i$$

$$\Rightarrow u'_i = \frac{u_{i+1} - u_i}{h} + \mathcal{O}(h)$$

(Two-point, forward difference)

- Compare to definition of the first derivative:

$$u'(x) \equiv \lim_{h \rightarrow 0} \frac{u(x+h) - u(x)}{h}$$

- We could have used the points  $x$  and  $x - h$ , which would have given us

$$u'(x) = \frac{u(x) - u(x-h)}{h} + \mathcal{O}(h)$$

Discretise:

$$u(x) \rightarrow u_i$$

$$\Rightarrow u'_i = \frac{u_i - u_{i-1}}{h} + \mathcal{O}(h)$$

(Two-point, backward difference)

- Quick illustration of forward difference method:

- Example:  $u(x) = a_0 + a_1x + a_2x^2$



– Exact:  $u'(x) = a_1 + 2a_2x$

– Approximation:

$$\begin{aligned} u'(x) &\approx \frac{u(x+h) - u(x)}{h} \\ &= \frac{[a_0 + a_1(x+h) + a_2(x+h)^2] - [a_0 + a_1x + a_2x^2]}{h} \\ &= \frac{a_1h + a_2x^2 + 2a_2xh + a_2h^2 - a_2x^2}{h} \\ &= a_1 + 2a_2x + a_2h \end{aligned}$$

– Compare to the exact expression: our approximation is wrong by an  $\mathcal{O}(h)$  term, as expected

– This **truncation error** gets smaller when we take  $h \rightarrow 0$

– But doing this can lead to **roundoff errors** in the subtraction  $u(x+h) - u(x)$ , causing a **loss of precision**

– We will return to this topic later

• We can use more than two points to compute  $u'(x)$ :

– Starting point: Taylor expansions for  $u(x+h)$  and  $u(x-h)$ :

$$\begin{aligned} u(x+h) &= u(x) + u'(x)h + \frac{1}{2}u''(x)h^2 + \frac{1}{6}u'''(x)h^3 + \mathcal{O}(h^4) \\ u(x-h) &= u(x) - u'(x)h + \frac{1}{2}u''(x)h^2 - \frac{1}{6}u'''(x)h^3 + \mathcal{O}(h^4) \end{aligned}$$

– Subtract:

$$u(x+h) - u(x-h) = 2u'h + \frac{2}{6}u'''h^3 + \mathcal{O}(h^5) \quad (\text{note power } h^5)$$

– Rearrange:

$$u' = \frac{u(x+h) - u(x-h)}{2h} - \frac{1}{6}u'''h^2 - \mathcal{O}(h^4)$$

$$u'(x) = \frac{u(x+h) - u(x-h)}{2h} + \mathcal{O}(h^2)$$

Discretise:

$$u'_i = \frac{u_{i+1} - u_{i-1}}{2h} + \mathcal{O}(h^2)$$

(Three-point expression)

- The second derivative

- Add Taylor expansions for  $u(x+h)$  and  $u(x-h)$

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

- Rearrange to isolate  $u''(x)$

$$u''(x) = \frac{u(x+h) - 2u(x) + u(x-h)}{h^2} + \mathcal{O}(h^2)$$

Discretise:

$$u''_i = \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} + \mathcal{O}(h^2)$$

## In-lecture code discussion #2

- Hidden files on Unix systems
  - Hidden files have file names starting with a dot
  - Some relevant examples:
    - \* `.bashrc` and/or `.profile` in your home directory
    - \* `.gitignore` in your git repositories
  - Use the `-a` option to see hidden files in your file listings: `ls -a`
- Terminal-based text editors
  - Useful when you want to make quick file edits
  - Useful when you are logged into another system via a Unix terminal, e.g. if you are working on a supercomputer
  - Some people use the terminal-based editors as their main editors – can become very powerful and efficient tools
  - Two popular examples: `vim` and `nano`
- Short discussion of the `std::vector` class:  
[anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/containers](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/containers)
  - Note the use of `my_vector.at(10)` as a safe alternative to `my_vector[10]` for accessing vector elements.
- How (not) to use `using namespace` in C++ programs:  
[anderkve.github.io/FYS3150/book/introduction\\_to\\_cpp/source\\_files\\_and\\_header\\_files](https://anderkve.github.io/FYS3150/book/introduction_to_cpp/source_files_and_header_files)
- Integer vs floating-point division:
  - In Python, the statement `x = 7/10` will by default evaluate to `x = 0.7`
  - However, in C++ the statement `x = 7/10` will evaluate to `x = 0`
  - Since 7 and 10 are written as integers, C++ will do **integer division**
  - In integer division, it is correct that  $7/10 = 0$
  - If we instead write `7.` and `10.`, C++ will treat these as floating-point numbers and perform floating-point division
  - So `x = 7./10.` will give the result `x = 0.7`
  - (The combinations `x = 7./10` and `x = 7/10.` will also give `x = 0.7`)
  - *Question:* Given the variable assignment `double x = 7/10;`, what value will `x` get?
  - *Answer:* `x` will be set to `x = 0.0`, since the assignment is evaluated as `double x = 0;`

## Boundary value problems (BVPs)

- Our case in project 1:

$$-\frac{d^2u}{dx^2} = f(x)$$

- $u(x)$  is an *unknown* function  $\rightarrow$  what we want to find
- $f(x)$  is some *known* function
- $x \in [0, 1]$
- Boundary values:  $u(0) = 0$  and  $u(1) = 0$  (Dirichlet)

- Special case of:

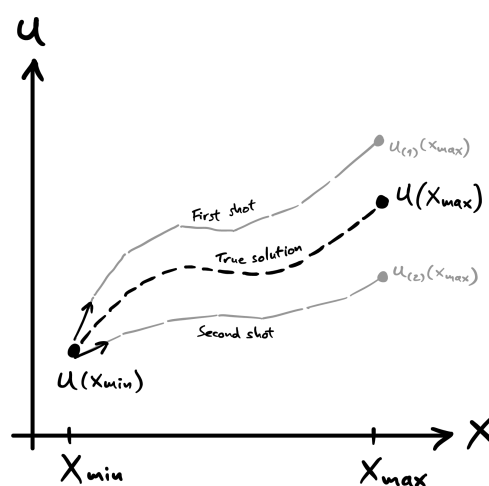
$$\alpha \frac{d^2u}{dx^2} + \beta \frac{du}{dx} + \gamma u(x) = f(x)$$

- *Ordinary* diff. eq., since there is only one independent variable ( $x$ )
- *Linear* diff. eq., since each term has maximum one power of  $u, u', u'', \dots$
- *Second order* diff. eq., since the highest-order derivative is  $u''$
- *Inhomogenous* diff. eq., when  $f(x) \neq 0$
- Many diff. eqs. in physics are linear
  - Then the sum of two solutions is a new, valid solution! (*superposition*)
  - Famous example: The Schrödinger eq. in quantum mechanics is linear  
 $\rightarrow$  superposition of quantum states!
- Many approaches to finding a solution
  - **Shooting methods** (described quickly below)
  - **Finite difference methods** (project 1, described below)
  - **Finite elements methods** (not covered)

## Quick description of shooting methods

- We want to solve a boundary value problem (BVP), where we start with known  $u(x_{\min})$  and  $u(x_{\max})$

- We'll do this by instead repeatedly solving an *initial value problem* (IVP), where we start with known  $u(x_{\min})$  and  $u'(x_{\min})$ :
  - Start from the known  $u(x_{\min})$
  - Guess a value for  $u'(x_{\min})$
  - Solve the corresponding IVP forward (“shoot”). (We will discuss IVPs later in the course.)
  - Repeat the previous two steps (in some clever way) until we find a solution  $u(x)$  that hits the known boundary condition at  $u(x_{\max})$
  - This solution  $u(x)$  is then a solution to our original BVP
- Things are easier when our diff. eq. is *linear*:
  - Guess a value for  $u'(x_{\min})$ , solve the IVP  
→ let's call this solution  $u_{(1)}(x)$
  - Guess another  $u'(x_{\min})$ , solve the IVP  
→ let's call this solution  $u_{(2)}(x)$
  - Since we have a linear diff. eq., a sum of solutions is a new solution:  
 $u_c(x) = cu_{(1)}(x) + (1 - c)u_{(2)}(x)$
  - Require that  $u_c(x_{\max})$  should equal the known  $u(x_{\max})$  (the second boundary condition)
  - Use this condition to determine a value for the free parameter  $c$   
→ this  $u_c(x)$  is then the solution  $u(x)$  to our BVP
- A drawback: Need to solve multiple IVPs to find the single solution to our BVP



**Figure 6:** Sketch of the shooting method

## Finite difference method

- Our problem: Find the function  $u(x)$  that solves this diff. eq.:

$$-\frac{d^2u}{dx^2} = f(x)$$

- We know  $u(0), u(1), f(x)$  and that  $x \in [0, 1]$
- Strategy:
  - **Step 1:** Express problem as a matrix eq.
  - **Step 2:** Solve the matrix eq.

### Step 1: Express as matrix eq.

- Discretise equation:

$$-\left[\frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} + \mathcal{O}(h^2)\right] = f_i, \quad f_i \equiv f(x_i)$$

- Approximate (leave out the  $\mathcal{O}(h^2)$  terms) and change notation:  $v_i \approx u_i$
- Arrange terms:

$$-v_{i-1} + 2v_i - v_{i+1} = h^2 f_i$$

- Note: this is a collection of multiple equations, one for each value we can insert for  $i$
- New goal: Determine  $v_1, v_2, \dots, v_{n-1}$
- We know:  $v_0, v_n$  and all the  $f_i$
- Below we'll consider the special case with  $n_{\text{steps}} = 5$ 
  - $v_0, v_1, v_2, v_3, v_4, v_5$ : 6 points
  - $v_0$  and  $v_5$  are known boundary points
  - 4 unknowns:  $v_1, v_2, v_3, v_4$
  - $h = \frac{v_5 - v_0}{n_{\text{steps}}} = 0.2$  (very large, just for illustration)

- The boxed expression represents a set of four equations. Let's write them out in a suggestive manner...

$$\begin{array}{rclclclclcl}
 (i = 1) & -v_0 & +2v_1 & -v_2 & & & = & h^2 f_1 \\
 (i = 2) & & -v_1 & +2v_2 & -v_3 & & = & h^2 f_2 \\
 (i = 3) & & & -v_2 & +2v_3 & -v_4 & = & h^2 f_3 \\
 (i = 4) & & & & -v_3 & +2v_4 & -v_5 & = & h^2 f_4
 \end{array}$$

- $v_0$  and  $v_5$  are known – let's move them over to the right-hand side and define some simpler notation  $g_1, g_2, g_3, g_4$ :

$$\begin{array}{rclclclclcl}
 +2v_1 & -v_2 & & & & = & h^2 f_1 + v_0 & \equiv & g_1 \\
 -v_1 & +2v_2 & -v_3 & & & = & h^2 f_2 & \equiv & g_2 \\
 & -v_2 & +2v_3 & -v_4 & & = & h^2 f_3 & \equiv & g_3 \\
 & & -v_3 & +2v_4 & & = & h^2 f_4 + v_5 & \equiv & g_4
 \end{array}$$

- This can be written as

$$\begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{bmatrix} = \begin{bmatrix} g_1 \\ g_2 \\ g_3 \\ g_4 \end{bmatrix}$$

$\mathbf{A}\vec{v} = \vec{g}$

- $\mathbf{A}$  and  $\vec{g}$  are known, we want to solve for  $\vec{v}$
- Note that  $\mathbf{A}$  is a **tridiagonal matrix**.
- The *diagonal* has only 2's, while the *superdiagonal* and *subdiagonal* contain only  $-1$ 's.
- Note that the vector  $\vec{v} = [v_1, v_2, v_3, v_4]$  in this equation only contains the *unknown*  $v_i$ . The known values at the boundaries,  $v_0$  and  $v_5$ , are *not* included in  $\vec{v}$ .

**Step 2: Solve the matrix eq.**

- Overview of things we'll discuss:
  1. *Now:* Method for solving  $\mathbf{A}\vec{v} = \vec{g}$  when  $\mathbf{A}$  is a *general, tridiagonal* matrix
    - Gaussian elimination turns into the Thomas algorithm
  2. *A task for you in Project 1:* Method for solving  $\mathbf{A}\vec{v} = \vec{g}$  when  $\mathbf{A}$  is the *special, tridiagonal* matrix above (with only -1's and 2's along the diagonals)
  3. *Later in the course:* Methods for solving a general matrix equation  $\mathbf{A}\vec{x} = \vec{b}$ 
    - Gaussian elimination, LU decomposition, iterative methods

**Matrix equations: Gaussian elimination and the Thomas algorithm****Introduction**

- A matrix equation  $\mathbf{A}\vec{x} = \vec{b}$  ( $\mathbf{A}$  and  $\vec{b}$  known,  $\vec{x}$  unknown) represents a set of linear equations

$$\begin{array}{llllllllll}
 \text{(eq. 1)} & a_{11}x_1 & + & a_{12}x_2 & + & \dots & + & a_{1n}x_n & = & b_1 \\
 \text{(eq. 2)} & a_{21}x_1 & + & a_{22}x_2 & + & \dots & + & a_{2n}x_n & = & b_2 \\
 (\dots) & & & & & & & & & \\
 \text{(eq. } m) & a_{m1}x_1 & + & a_{m2}x_2 & + & \dots & + & a_{mn}x_n & = & b_m
 \end{array}$$

- $m$  equations, each with  $n$  terms — one for each unknown variable  $(x_1, \dots, x_n)$

$$\underset{(m \times n)}{\mathbf{A}} \underset{(n \times 1)}{\vec{x}} = \underset{(m \times 1)}{\vec{b}}$$

- We will focus on the case of a *square* matrix, i.e. when  $m = n$
- This means we have  $n$  equations and  $n$  unknowns
- If all our equations are *linearly independent*, i.e. when each equation represents information not contained in the other equations, we should be able to solve for all our  $n$  unknowns  $(x_1, \dots, x_n)$
- Some equivalent statements:
  - All the equations are linearly independent
  - $\mathbf{A}$  is *not* singular (all eigenvalues of  $\mathbf{A}$  are non-zero)
  - $\det \mathbf{A} \neq 0$



*Side note:* When we have *more* equations (constraints) than unknowns, there is generally *no exact solution*. But we can *fit* our unknowns such that all our equations are as close to solved as possible. This is the typical case in science: you have a model with a few free parameters and the model needs to match many observations (constraints) as closely as possible.

## Gaussian elimination, overview

- Start from general matrix equation

$$\begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix} \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

- **Step 1:** Forward substitution/elimination
  - Turn matrix into upper-triangular form

$$\begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ & \bullet & \bullet & \bullet \\ & & \bullet & \bullet \\ & & & \bullet \end{bmatrix} \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

- Can then read off solution for  $x_m$  (last row)

- **Step 2:** Back substitution/elimination
  - Use the now known  $x_m$  to find  $x_{m-1}$ , then use these to find  $x_{m-2}$ , and so on
  - End up with this

$$\begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{bmatrix} \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

## Thomas algorithm: Gaussian elimination on a tridiagonal matrix

- Let's go back to the notation of project 1:  $\mathbf{A}\vec{v} = \vec{g}$

- Let  $\mathbf{A}$  be a *general* tridiagonal matrix
  - For concreteness we look at the case with a  $4 \times 4$  matrix:

$$\begin{bmatrix} b_1 & c_1 & 0 & 0 \\ a_2 & b_2 & c_2 & 0 \\ 0 & a_3 & b_3 & c_3 \\ 0 & 0 & a_4 & b_4 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{bmatrix} = \begin{bmatrix} g_1 \\ g_2 \\ g_3 \\ g_4 \end{bmatrix}$$

- Note that we have used indices that correspond to the row numbers:
  - Subdiagonal:  $\vec{a} = [a_2, a_3, a_4]$
  - Diagonal:  $\vec{b} = [b_1, b_2, b_3, b_4]$
  - Superdiagonal:  $\vec{c} = [c_1, c_2, c_3]$

- Now let's do **step 1**, the *forward substitution*

- We start from the *augmented matrix*:

$$\begin{array}{lcl} R_1 : & b_1 & c_1 \quad 0 \quad 0 \quad \Bigg| \quad g_1 \\ R_2 : & a_2 & b_2 \quad c_2 \quad 0 \quad \Bigg| \quad g_2 \\ R_3 : & 0 & a_3 \quad b_3 \quad c_3 \quad \Bigg| \quad g_3 \\ R_4 : & 0 & 0 \quad a_4 \quad b_4 \quad \Bigg| \quad g_4 \end{array}$$

- To move towards an upper-triangular form, we want to set the  $a_2$  entry in  $R_2$  to 0
- Use a row operation with  $R_1$  to achieve this:  $R_2 \rightarrow R_2 - \frac{a_2}{b_1} R_1$
- This turns the  $a_2$  entry into  $a_2 - \frac{a_2}{b_1} b_1 = 0$

$$\begin{array}{lcl} R_1 : & b_1 & c_1 \quad 0 \quad 0 \quad \Bigg| \quad g_1 \\ R_2 : & 0 & (b_2 - \frac{a_2}{b_1} c_1) \quad c_2 \quad 0 \quad \Bigg| \quad (g_2 - \frac{a_2}{b_1} g_1) \\ R_3 : & 0 & a_3 \quad b_3 \quad c_3 \quad \Bigg| \quad g_3 \\ R_4 : & 0 & 0 \quad a_4 \quad b_4 \quad \Bigg| \quad g_4 \end{array}$$

- Introduce shorthand notation:

$$\begin{aligned} * \quad \tilde{b}_1 &= b_1 \\ * \quad \tilde{b}_2 &= b_2 - \frac{a_2}{b_1} c_1 \\ * \quad \tilde{g}_1 &= g_1 \\ * \quad \tilde{g}_2 &= g_2 - \frac{a_2}{b_1} g_1 \end{aligned}$$

- We then have

$$\begin{array}{lcl} R_1 : & \tilde{b}_1 & c_1 \quad 0 \quad 0 \quad \left| \quad \tilde{g}_1 \\ R_2 : & 0 & \tilde{b}_2 \quad c_2 \quad 0 \quad \left| \quad \tilde{g}_2 \\ R_3 : & 0 & a_3 \quad b_3 \quad c_3 \quad \left| \quad g_3 \\ R_4 : & 0 & 0 \quad a_4 \quad b_4 \quad \left| \quad g_4 \end{array}$$

- Now continue in the same way to turn the  $a_3$  entry to zero

- Row operation:  $R_3 \rightarrow R_3 - \frac{a_3}{\tilde{b}_2} R_2$

- Define notation:

$$\begin{aligned} * \quad \tilde{b}_3 &= b_3 - \frac{a_3}{\tilde{b}_2} c_2 \\ * \quad \tilde{g}_3 &= g_3 - \frac{a_3}{\tilde{b}_2} \tilde{g}_2 \end{aligned}$$

- We then get

$$\begin{array}{lcl} R_1 : & \tilde{b}_1 & c_1 \quad 0 \quad 0 \quad \left| \quad \tilde{g}_1 \\ R_2 : & 0 & \tilde{b}_2 \quad c_2 \quad 0 \quad \left| \quad \tilde{g}_2 \\ R_3 : & 0 & 0 \quad \tilde{b}_3 \quad c_3 \quad \left| \quad \tilde{g}_3 \\ R_4 : & 0 & 0 \quad a_4 \quad b_4 \quad \left| \quad g_4 \end{array}$$

- And once more, with feeling...

- Row operation:  $R_4 \rightarrow R_4 - \frac{a_4}{\tilde{b}_3} R_3$

- Define notation:

$$\begin{aligned} * \quad \tilde{b}_4 &= b_4 - \frac{a_4}{\tilde{b}_3} c_3 \\ * \quad \tilde{g}_4 &= g_4 - \frac{a_4}{\tilde{b}_3} \tilde{g}_3 \end{aligned}$$

$$\begin{array}{lcl} R_1 : & \tilde{b}_1 & c_1 \quad 0 \quad 0 \quad \left| \quad \tilde{g}_1 \\ R_2 : & 0 & \tilde{b}_2 \quad c_2 \quad 0 \quad \left| \quad \tilde{g}_2 \\ R_3 : & 0 & 0 \quad \tilde{b}_3 \quad c_3 \quad \left| \quad \tilde{g}_3 \\ R_4 : & 0 & 0 \quad 0 \quad \tilde{b}_4 \quad \left| \quad \tilde{g}_4 \end{array}$$

- The forward substitution is now done! Here's the summary:

**Forward substitution:**

$$\begin{aligned}\tilde{b}_1 &= b_1 \\ \tilde{b}_i &= b_i - \frac{a_i}{\tilde{b}_{i-1}} c_{i-1} \quad i = 2, 3, 4 \\ \tilde{g}_1 &= g_1 \\ \tilde{g}_i &= g_i - \frac{a_i}{\tilde{b}_{i-1}} \tilde{g}_{i-1} \quad i = 2, 3, 4\end{aligned}$$

The lecture on August 31 ended here.

- Now let's do **step 2**, the *back substitution*

- Starting point

$$\begin{array}{lcl} R_1 : & \tilde{b}_1 & c_1 \quad 0 \quad 0 \quad \left| \quad \tilde{g}_1 \\ R_2 : & 0 & \tilde{b}_2 \quad c_2 \quad 0 \quad \left| \quad \tilde{g}_2 \\ R_3 : & 0 & 0 \quad \tilde{b}_3 \quad c_3 \quad \left| \quad \tilde{g}_3 \\ R_4 : & 0 & 0 \quad 0 \quad \tilde{b}_4 \quad \left| \quad \tilde{g}_4 \end{array}$$

- We now want to get to an identity matrix form, starting from the bottom row
- Row operation:  $R_4 \rightarrow \frac{R_4}{\tilde{b}_4}$

$$\begin{array}{lcl} R_1 : & \tilde{b}_1 & c_1 \quad 0 \quad 0 \quad \left| \quad \tilde{g}_1 \\ R_2 : & 0 & \tilde{b}_2 \quad c_2 \quad 0 \quad \left| \quad \tilde{g}_2 \\ R_3 : & 0 & 0 \quad \tilde{b}_3 \quad c_3 \quad \left| \quad \tilde{g}_3 \\ R_4 : & 0 & 0 \quad 0 \quad 1 \quad \left| \quad \frac{\tilde{g}_4}{\tilde{b}_4} \rightarrow v_4 \end{array}$$

- We now have the solution for  $v_4$ :

$$v_4 = \frac{\tilde{g}_4}{\tilde{b}_4}$$

- Now we want to get  $R_3$  on the form (0,0,1,0)
- We can subtract  $c_3 R_4$  to get rid of the  $c_3$  entry in  $R_3$ , and then divide by  $\tilde{b}_3$  to set the third element to 1

- Row operation:  $R_4 \rightarrow \frac{R_3 - c_3 R_4}{\tilde{b}_3}$

$$\begin{array}{lcl} R_1 : & \tilde{b}_1 & c_1 \quad 0 \quad 0 \quad \left| \quad \tilde{g}_1 \\ R_2 : & 0 & \tilde{b}_2 \quad c_2 \quad 0 \quad \left| \quad \tilde{g}_2 \\ R_3 : & 0 & 0 \quad 1 \quad 0 \quad \left| \quad \frac{\tilde{g}_3 - c_3 v_4}{\tilde{b}_3} \rightarrow v_3 \\ R_4 : & 0 & 0 \quad 0 \quad 1 \quad \left| \quad v_4 \end{array}$$

- This gives us the solution for  $v_3$ :

$$v_3 = \frac{\tilde{g}_3 - c_3 v_4}{\tilde{b}_3}$$

- We can continue upwards like this to find all the remaining  $v_i$ . In summary:

**Back substitution:**

$$v_4 = \frac{\tilde{g}_4}{\tilde{b}_4}$$

$$v_i = \frac{\tilde{g}_i - c_i v_{i+1}}{\tilde{b}_i} \quad i = 3, 2, 1$$

- Let's summarise what we've done:
  - Given a general tridiagonal matrix  $\mathbf{A}$  and a vector  $\vec{g}$ , we have found the vector  $\vec{v}$  that solves the equation  $\mathbf{A}\vec{v} = \vec{g}$ .
  - We used Gaussian elimination, which has two steps:
    - \* *forward substitution*
    - \* *back substitution*
  - Because  $\mathbf{A}$  was tridiagonal, the Gaussian elimination procedure resulted in a fairly simple algorithm, which is known as the **Thomas algorithm** (Llewellyn Thomas, 1903–1992)

*Coding tip:* Note that we don't need to work with an entire matrix in memory here. To implement the Thomas algorithm above, we just need some arrays/vectors  $\vec{a}$ ,  $\vec{b}$ ,  $\vec{c}$ ,  $\vec{g}$ ,  $\vec{\tilde{b}}$ ,  $\vec{\tilde{g}}$  and  $\vec{v}$ .

## Back to our boundary value problem

- We now have the tools we need to use a **finite difference method** to solve a boundary value problem like

$$-\frac{d^2u}{dx^2} = f(x)$$

where  $f(x)$  is some known function, and we know  $u(x_{\min})$  and  $u(x_{\max})$

- Discretise the problem, using a discretised approximation for the second derivative
    - \* At this step we changed notation  $u_i \rightarrow v_i$
  - Formulate the resulting set of equations as a matrix equation  $\mathbf{A}\vec{v} = \vec{g}$ 
    - \* The second derivative in the diff. eq.  $\rightarrow$  the matrix  $\mathbf{A}$  will be tridiagonal, with a simple (-1,2,-1) form
  - Use the Thomas algorithm to solve the matrix equation
    - \* **However**, the Thomas algorithm is a method that can solve *any* tridiagonal matrix equation, but in the case of our BVP we are only interested in the case of a particularly simple, tridiagonal matrix. This means that we can simplify the Thomas algorithm for our usecase – something you will do in project 1.
- *A reasonable question:* Why are we doing all this? Why not rather find  $\mathbf{A}^{-1}$  and solve the equation as  $\vec{v} = \mathbf{A}^{-1}\vec{g}$ ?
    - Finding  $\mathbf{A}^{-1}$  numerically takes  $\mathcal{O}(n^3)$  operations for an  $n \times n$  matrix. This approach becomes useful if we need to solve *many* different equations ( $\mathbf{A}\vec{v}_1 = \vec{g}_1, \mathbf{A}\vec{v}_2 = \vec{g}_2, \dots$ ) that all involve the same matrix  $\mathbf{A}$ . But for solving a single equation  $\mathbf{A}\vec{v} = \vec{g}$ , other methods are quicker.

## Counting floating-point operations (FLOPs)

- Floating-point numbers, *floats*: (inexact) machine representation of the real numbers ( $\mathbb{R}$ )
- Floats are numbers where the decimal point can be placed anywhere (it can “float”) in a given string of digits, depending on which number we need to represent
  - Example: The digits 112358 can represent 11.2358 or 1123.58, depending on the placement of the decimal point
- Floating-point operations:  $\{+, -, \times, \div\}$  with floats
- Much slower than integer operations. (One FLOP consists of several integer operations.)
- Counting FLOPs is a way of estimating the *efficiency of an algorithm*
- Note: FLOPs** (Floating-point Operations) vs **FLOPS** (Floating-point Operations per Second). FLOPS is a measure of *computer performance*, which we will not discuss in this course.
  - So how long a given task will take on a given computer will depend both on the number of **FLOPs** required for the task *and* the number of **FLOPS** for the computer – and a bunch of other things...

### Examples

- Example 1:

$$y = ab + c, \quad 1 \text{ mult.}, 1 \text{ add.} \rightarrow 2 \text{ FLOPs}$$

- Example 2:

$$\begin{array}{ll} \text{for } i = 1, \dots, n : & n \text{ repetitions} \\ y_i = ay_{i-1} + i & 2 \text{ FLOPs} \\ & \rightarrow 2n \text{ FLOPs} \end{array}$$

- Example 3:

$$\begin{array}{ll} \text{for } i = 1, \dots, n : & n \text{ repetitions} \\ y_i = \frac{a}{b}y_{i-1} + i & 3 \text{ FLOPs} \\ & \rightarrow 3n \text{ FLOPs (silly!)} \end{array}$$

- A more efficient version of example 3:

|                         |   |
|-------------------------|---|
| $c = \frac{a}{b}$       | 1 FLOP  |
| for $i = 1, \dots, n$ : | $n$ repetitions   |
| $y_i = cy_{i-1} + i$    | 2 FLOPs   |
|                         | $\rightarrow (2n + 1) \text{ FLOPs} \approx 2n \text{ FLOPs}$ |

*Tip:* When code speed is important, avoid recomputing constants within a loop.



## Binary representation

- *In short*: How to represent numbers using only two different symbols
- Basic element: a **bit**
  - 1/0, on/off, true/false, yes/no, hole/not-hole (punched cards), red/blue, ...
  - The term *bit* is originally a contraction of *binary information digit*
- A bit doesn't have to be related to computers – it's a basic concept from information theory
  - A bit is the expected amount of *information* or *surprise* contained in the outcome of a 50/50 random draw. (The more surprising a result/message/signal is, the more information it contains – look up literature on *Shannon entropy* for more on this.)
- In principle, any physical system with two possible states can be used to represent the digits 0 and 1
- So we better use a numeral system that only needs two different digits to represent any number → the **binary system** or the **base 2 system**
- In base 10, we have ten different symbols (0–9) that can be used per position
- In base 2, we only have two different symbols per position
  - need to use more positions to express numbers
  - longer strings of symbols compared to the decimal system
- *Side note*: In computing and mathematics we also sometimes encounter the *hexadecimal* (base 16) system. In this case there are 16 different symbols (0–9 and A–F) per position.

## Integers

- Example: 137 in base 10 and base 2:  $(137)_{10} = (10001001)_2$
- Representation in the decimal (base 10) system:

$$\begin{aligned}(137)_{10} &= \frac{10^2}{1} + \frac{10^1}{3} + \frac{10^0}{7} = (1 \times 10^2) + (3 \times 10^1) + (7 \times 10^0) \\ &= 100 + 30 + 7 \\ &= 137\end{aligned}$$

- Representation in the binary (base 2) system:

$$\begin{aligned}
 (10001001)_2 &= \frac{2^7 \quad 2^6 \quad 2^5 \quad 2^4 \quad 2^3 \quad 2^2 \quad 2^1 \quad 2^0}{1 \quad 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 1} \\
 &= (1 \times 2^7) + (0 \times 2^6) + \dots + (1 \times 2^3) + \dots + (1 \times 2^0) \\
 &= 128 + 0 + 0 + 0 + 8 + 0 + 0 + 1 \\
 &= 137
 \end{aligned}$$

- How can we find the correct binary string of 0's and 1's for a given number?
  - The same way we (without thinking about it) identify the correct string of digits in the decimal system: by doing *integer division and keeping track of remainders*

|                         | Remainder | Position |
|-------------------------|-----------|----------|
| $137 \setminus 10 = 13$ | 7         | $10^0$   |
| $13 \setminus 10 = 3$   | 3         | $10^1$   |
| $3 \setminus 10 = 0$    | 1         | $10^2$   |

**Table 1:** Repeated integer division with 10 produces the base 10 representation of 137.

|                        | Remainder | Position |
|------------------------|-----------|----------|
| $137 \setminus 2 = 68$ | 1         | $2^0$    |
| $68 \setminus 2 = 34$  | 0         | $2^1$    |
| $34 \setminus 2 = 17$  | 0         | $2^2$    |
| $17 \setminus 2 = 8$   | 1         | $2^3$    |
| $8 \setminus 2 = 4$    | 0         | $2^4$    |
| $4 \setminus 2 = 2$    | 0         | $2^5$    |
| $2 \setminus 2 = 1$    | 0         | $2^6$    |
| $1 \setminus 2 = 0$    | 1         | $2^7$    |

**Table 2:** Repeated integer division with 2 produces the base 2 representation of 137.

- The more bits we have available, the longer the integer we can store
- If we are working with *signed* integers, we need one additional bit to represent the sign:  $(-1)^0$  or  $(-1)^1$

## Floating-point numbers

- How to represent the real numbers ( $\mathbb{R}$ ) in binary?

- Strategy: use *normalised, scientific notation in base 2*

- Example in decimal:

$$-9.90625 \times 10^0, \text{ or } -0.990625 \times 10^1$$

- The latter convention, where the first digit is always zero, is often used in computing
- General form:

$$\pm \left[ \text{number in } \left( \frac{1}{10}, 1 \right) \right] \times 10^{\text{integer exponent}}$$

- In binary (base 2):

$$\pm \left[ \text{number in } \left( \frac{1}{2}, 1 \right) \right] \times 2^{\text{integer exponent}}$$

- Terminology:

$$[\text{sign}][\text{mantissa}] \times 2^{\text{exponent}}$$

- Whether the mantissa should be a number within  $(\frac{1}{2}, 1)$  or within  $(1, 2)$  is a matter of convention
- Another common term for the mantissa is the **significand**
- We already know how to represent the integer exponent and the sign bit in binary
- Binary representation of the mantissa:

- Example:  $(0.5625)_{10}$

$$\begin{aligned} (0.1001)_2 &= \frac{2^0 \quad 2^{-1} \quad 2^{-2} \quad 2^{-3} \quad 2^{-4}}{0 \quad 1 \quad 0 \quad 0 \quad 1} \\ &= (0 \times 2^0) + (1 \times 2^{-1}) + (0 \times 2^{-2}) + (0 \times 2^{-3}) + (1 \times 2^{-4}) \\ &= 0 + 0.5 + 0 + 0 + 0.0625 \\ &= 0.5625 \end{aligned}$$

- **Single precision:** Using 32 bits (4 bytes) to represent a floating-point number

- Sign: 1 bit
- Exponent: 8 bits
- Mantissa: 23 bits

- Example: The number  $-3.25$

- In normalised, scientific notation in base 2:  $-0.8125 \times 2^2$

\* Sign: -1 (so the *sign bit* will be 1 since  $-1 = (-1)^1$ )

- \* Exponent: 2
- \* Mantissa: 0.8125

- In memory, something like this:

| sign bit | 8-bit exponent (2) | 23-bit mantissa (0.8125) |
|----------|--------------------|--------------------------|
| 1        | 00000010           | 1101000 ... 000          |

- On most computer systems, the type **float** in C++ will correspond to a 32-bit number
- **Double precision:** Using 64 bits (8 bytes) to represent a floating-point number
  - Sign: 1 bit
  - Exponent: 11 bits
  - Mantissa: 52 bits
- On most systems, the type **double** in C++ will correspond to a 64-bit number
- An 11-bit exponent gives an exponent range of  $(-1024, 1024)$ , since  $2^{11} = 2048$ 
  - Since  $2^{1024} \approx 10^{308}$ , the range of numbers that can be represented in double-precision is roughly  $(10^{-308}, 10^{308})$
  - You can test this quickly in your Python terminal, since a floating-point number in Python (the **float** type in Python) by default will be a 64-bit number on most systems:

```
>>> 2.**1023
8.98846567431158e+307
>>>
>>> 2.**1024
OverflowError: (34, 'Numerical result out of range')
```

- Finite number of bits → unavoidable problems with range and accuracy
- Limited number of bits for the **exponent:**
  - a limited **range** of  $\mathbb{R}$  can be represented
    - With 11 bits for the exponent, we get a range of  $\sim (10^{-308}, 10^{308})$
- Limited number of bits for the **mantissa:**
  - a limited **resolution/precision** in our representation of the continuous  $\mathbb{R}$ 
    - With 52 bits for the mantissa, we get a precision of around **15 digits** in the decimal system ( $\log_{10}(2^{52}) \approx 15.654$ )

**Hidden bit:** When using normalised, scientific notation in base 2, we know that the most significant digit, i.e. the first digit of the mantissa, will always be 0 (if the  $(\frac{1}{2}, 1)$ -convention is used), or always be 1 (if the  $(1, 2)$ -convention is used). So we don't need to explicitly store this bit in memory. This trick is referred to as the *hidden bit*, and it effectively increases the mantissa precision by one bit, e.g. from 52 bits to 53 bits for a double-precision number.

- *A silly example to illustrate the effect of limited range and precision:*

- Let's work in base 10
- Assume we only had memory for *one digit in the exponent* and *one digit in the mantissa*
- We could then only represent these numbers:

$$\dots, 1 \times 10^{-1}, 2 \times 10^{-1}, \dots, 1 \times 10^0, 2 \times 10^0, \dots, 1 \times 10^1, 2 \times 10^1, \dots$$

- We would have a range of  $\sim (10^{-5}, 10^5)$
- The only numbers we would be able to use would be

$$\dots, 0.1, 0.2, 0.3, \dots, 1, 2, 3, \dots, 10, 20, 30, \dots, 100, 200, 300, \dots$$

- So a number 17 would just end up as 10, and a computation like  $100 + 80$  would just give the result 100

## Errors

### Truncation errors

- Truncation errors are purely mathematical in origin
- Typical case: we cut of a series expansion at some point
- Example: leaving out the  $\mathcal{O}(h^2)$  terms in

$$u_i'' = \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} + \mathcal{O}(h^2)$$

- *Note:* here a smaller step size  $h$  will give a smaller truncation error

### Roundoff errors

- Numbers can only be stored with limited accuracy
- For *doubles*, our precision is ~15 digits
  - We will often refer to this as our *machine precision*
- So *almost* all numbers we store are **approximations** to the true number we intended to store
  - True number:  $a$
  - Floating-point representation of  $a$ :  $fl(a)$
  - Given a true  $a$ , your  $fl(a)$  will be in the range

$$a(1 - \delta_m) < fl(a) < a(1 + \delta_m)$$

where  $\delta_m$  is the machine precision (e.g.  $\delta_m \sim 10^{-15}$ )

- So given  $fl(a)$ , all you know is that the true number  $a$  is in the range

$$fl(a)(1 - \delta_m) < a < fl(a)(1 + \delta_m)$$



**Figure 7:** The continuous number line and the discretised floating-point representation

The lecture on September 1 ended here.

## Loss of numerical precision

- Also known as **loss of significance**
- Typical case: subtraction with similar numbers  
→ we lose the most significant digits, left with digits that are more affected by roundoff errors
- Example:

- True values:  $a = 1.0054321$ ,  $b = 1.0040001$
- Assume a machine precision of  $\delta_m \sim 10^{-4}$  (just for illustration)
- Approximate floating-point representations:  $fl(a) = 1.005$ ,  $fl(b) = 1.004$
- 4 significant digits
- Relative errors in the approximations:

$$\left| \frac{a - fl(a)}{a} \right| \approx 10^{-4}$$

$$\left| \frac{b - fl(b)}{b} \right| \approx 10^{-7}$$

- So  $fl(a)$  and  $fl(b)$  are clearly very reasonable approximations to  $a$  and  $b$ , given our assumed machine precision
- Now perform a subtraction:
- True value:  $a - b = 0.0014320$
- Approximate:  $fl(a) - fl(b) = 1.005 - 1.004 = 0.001$
- Now we only have 1 significant digit!
- Relative error:

$$\left| \frac{0.0014320 - 0.001}{0.0014320} \right| \approx 3 \times 10^{-1}$$

- Suddenly we have a **30% error**, even though our input numbers were reasonable representations of the true values

- Such **loss of precision** can easily happen in the middle of some long, complicate computation, and then all subsequent computations will end up with a large error.
- Another common term for this is **catastrophic cancellation**
- Note that when discussing errors, we are usually most interested in the *relative* error:
- Example from project 1:
  - Absolute error:  $\Delta = |v_i - u_i|$
  - Relative error:  $\epsilon = \left| \frac{v_i - u_i}{u_i} \right|$
  - It may be useful to e.g. study plots or tables of  $\log_{10}(\epsilon)$  vs  $\log_{10}(h)$
- Typical case for us:
  - If the step size is large: truncation error dominates
  - If the step size is tiny: roundoff errors lead to loss of precision → garbage results
  - So we expect that there is some optimal, intermediate step size that gives the smallest overall error

## An example error analysis

This topic includes a small in-lecture code discussion, using the `error_analysis` code example.

- Consider the function  $u(x) = e^{2x}$ 
  - (We choose this example function just because it's trivial to differentiate many times, and the 2 in the exponent ensures that all the derivatives are not exactly equal.)
- We will use our familiar expression

$$\frac{u_{i-1} - 2u_i + u_{i+1}}{h^2}$$

to implement a computer program that computes (an approximation to) the true second derivative  $u''_i$  at a given point  $x_i$

- *Our question:* How do we expect that the *relative error* of our code output will depend on our choice of step size  $h$ ?
- First of all, we know the exact answer for  $u''_i$ :

$$u''_i = 4e^{2x_i}$$



- In what follows we will first consider the **absolute error**, and then later the **relative error**

- Absolute error:

$$\Delta(h) \equiv |\text{approx.} - \text{true}| = \left| \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} - u_i'' \right|$$

- Relative error:

$$\epsilon(h) \equiv \left| \frac{\text{approx.} - \text{true}}{\text{true}} \right| = \left| \frac{\Delta(h)}{u_i''} \right|$$

- Now, as a first step towards answering our question, let's construct a simple *model* for the absolute error  $\Delta(h)$
- We will assume that  $\Delta(h)$  is the sum of two contributions, namely a truncation error  $\Delta_{\text{tr}}(h)$  and a roundoff error  $\Delta_{\text{ro}}(h)$ :

$$\Delta(h) = \Delta_{\text{tr}}(h) + \Delta_{\text{ro}}(h)$$

- Let's first look at the truncation error:

$$\begin{aligned} \Delta_{\text{tr}}(h) &= \left| \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} - u_i'' \right| \\ &= \left| \left( \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} \right) - \left( \frac{u_{i-1} - 2u_i + u_{i+1}}{h^2} + \mathcal{O}(u_i^{(4)} h^2) \right) \right| \\ &= \left| \mathcal{O}(u_i^{(4)} h^2) \right| \end{aligned}$$

- Note that here we have included in our big-O notation that the leading term is not only proportional to  $h^2$ , but also to the fourth derivative,  $u_i^{(4)}$ . (To see this, go back to our derivation of the discretised expression for the second derivative.)
- It is useful here to keep track of this dependence on the fourth derivative, since for other choices of the example function  $u(x)$  the different-order derivatives at  $x_i$  could have vastly different values – or indeed be zero, if our  $u(x)$  was a low-order polynomial.
- Now let's look at the roundoff error:
  - What we *want* to compute is

$$\frac{u_{i-1} - 2u_i + u_{i+1}}{h^2}$$

or to rewrite it slightly,

$$\frac{(u_{i+1} - u_i) - (u_i - u_{i-1}))}{h \times h}$$

- However, the computation we *actually* end up performing on the computer is something like this:

$$fl \left[ \frac{fl \left[ fl(u_{i+1}) - fl(u_i) \right] - fl \left( fl(u_i) - fl(u_{i-1}) \right)}{fl(h) \times fl(h)} \right]$$

- Consider the limit of small  $h$  and focus on the subtractions of near identical numbers

$$fl(u_{i+1}) - fl(u_i)$$

and similar for  $fl(u_i) - fl(u_{i-1})$

- Recall:

$$a(1 - \delta_m) < fl(a) < a(1 + \delta_m)$$

- We can then estimate an upper bound for the result of the subtraction

$$\begin{aligned} fl(u_{i+1}) - fl(u_i) &\leq u_{i+1}(1 + \delta_m) - u_i(1 - \delta_m) \\ &= (u_{i+1} - u_i) + (u_{i+1} + u_i)\delta_m \end{aligned}$$

- In the limit  $h \rightarrow 0$ , i.e. when  $u_{i+1} \rightarrow u_i$ , the first parenthesis vanishes, but the second parenthesis does not
- So we are left with

$$fl(u_{i+1}) - fl(u_i) \leq \mathcal{O}(u_i \delta_m)$$

- *What this means:* While we know that the *true* value of  $u_{i+1} - u_i$  goes to 0 when  $h \rightarrow 0$ , we have no guarantee that our actual computation  $fl(u_{i+1}) - fl(u_i)$  will go to exactly 0 in this limit.
- Assuming that this subtraction is the most “dangerous” part in our computation of

$$\frac{u_{i-1} - 2u_i + u_{i+1}}{h^2}$$

we can estimate the roundoff error  $\Delta_{\text{ro}}$  to be

$$\Delta_{\text{ro}}(h) = \mathcal{O}\left(\frac{u_i \delta_m}{h^2}\right)$$

- We can now put everything together in our simple model for the absolute error:

$$\begin{aligned}\Delta(h) &= |\Delta_{\text{tr}}(h) + \Delta_{\text{ro}}(h)| \\ &= \left| \mathcal{O}\left(u_i^{(4)} h^2\right) + \mathcal{O}\left(\frac{u_i \delta_m}{h^2}\right) \right|\end{aligned}$$

- Our model for the *relative* error  $\epsilon$  in our computation of  $u_i''$  then becomes

$$\begin{aligned}\epsilon(h) &= \left| \frac{\Delta(h)}{u_i''} \right| \\ &= \left| \mathcal{O}\left(\frac{u_i^{(4)}}{u_i''} h^2\right) + \mathcal{O}\left(\frac{u_i \delta_m}{u_i''} \frac{1}{h^2}\right) \right|\end{aligned}$$

- The first term grows when  $h$  *increases*, while the second term grows when  $h$  *decreases*
- For our choice of example function we also know that  $\mathcal{O}(u_i) \approx \mathcal{O}(u_i'') \approx \mathcal{O}(u_i^{(4)})$ , so the factors  $\frac{u_i^{(4)}}{u_i''}$  and  $\frac{u_i}{u_i''}$  won't suppress or enlarge the error terms much

- Let's now look at  $\log_{10} \epsilon(h)$

- Collecting the stuff that doesn't depend on  $h$  in two constants  $C_1$  and  $C_2$ , we can write

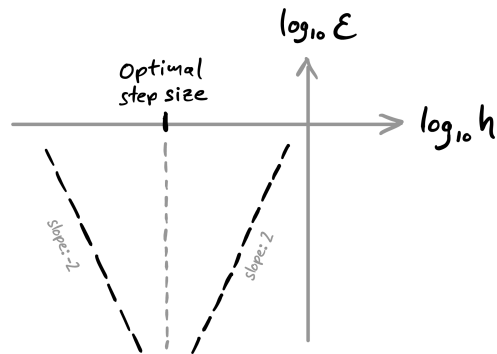
$$\log_{10} \epsilon(h) = \log_{10} \left| C_1 h^2 + C_2 h^{-2} \right|$$

- Look at the behaviour in the limits  $h \rightarrow \infty$  (first term dominates) and  $h \rightarrow 0$  (second term dominates)

$$\log_{10} \epsilon(h) \approx \begin{cases} -2 \log_{10} h + \log_{10} C_2 & \text{for } h \rightarrow 0, \text{ i.e. } \log_{10} h \rightarrow -\infty \\ 2 \log_{10} h + \log_{10} C_1 & \text{for } h \rightarrow \infty, \text{ i.e. } \log_{10} h \rightarrow \infty \end{cases}$$

- Note that these are the equations for two straight lines (slopes 2 and  $-2$ ) in a plot of  $\log_{10} \epsilon(h)$  vs  $\log_{10} h$
- So we see that our model for the error suggests the qualitative behaviour we expected, namely that there should be some *optimal, intermediate choice for the step size* that gives the smallest overall error

- We also see that we get a quantitative prediction for how quickly the error will grow when we move far away from the optimal step size choice



**Figure 8:** Sketch of a  $\log_{10} \epsilon(h)$  vs  $\log_{10} h$  plot, as suggested by our simple error model

## Recap: solving matrix equations

- We have discussed how to solve matrix equations  $\mathbf{A}\vec{x} = \vec{b}$ 
  - (We used the notation  $\mathbf{A}\vec{v} = \vec{g}$  in project 1)
- **Gaussian elimination:**
  - Can be used to solve  $\mathbf{A}\vec{x} = \vec{b}$  for a *general* (dense)  $\mathbf{A}$
  - In that case it requires  $\mathcal{O}(n^3)$  FLOPs, or more accurately  $\mathcal{O}\left(\frac{2}{3}n^3\right)$
  - We have only looked at the special case for a *tridiagonal*  $\mathbf{A}$  (more efficient)
- Next up: **LU decomposition**
- Later: **Iterative methods**

## Classification of methods for solving matrix equations

### Direct methods

- Examples:
  - Gaussian elimination
  - LU decomposition
- In theory, these methods give the *exact* answer in a *finite number of steps*
- In practice, these methods can suffer from numerical instabilities
- They typically work with the entire matrix at once
  - keeps the full matrix stored in memory

### Indirect methods

- Examples:
  - Jacobi's iterative method
  - Gauss-Seidel
  - Relaxation methods
- *Iterate* closer and closer to the exact answer, but *will never get there exactly*
- Can often work without keeping the full matrix in memory
- Are often less susceptible to roundoff errors

## Lower-upper (LU) decomposition

- Also commonly known as **lower-upper (LU) factorisation**
- We will introduce LU decomposition as an approach for solving  $\mathbf{A}\vec{x} = \vec{b}$
- Actually a starting point for several different matrix tasks, as we will see
- Our plan:
  1. What is LU decomposition?
  2. What is it good for? (And what's the difficulty?)
  3. An algorithm for LU decomposition

### 1. What is LU decomposition?

- We will only consider *square* matrices
- A matrix  $\mathbf{A}$  is said to admit an **LU decomposition** if it can be written as a product of a lower-triangular matrix ( $\mathbf{L}$ ) and an upper-triangular matrix ( $\mathbf{U}$ )

$$\mathbf{A} = \mathbf{L}\mathbf{U}$$

$$\begin{bmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{bmatrix} = \begin{bmatrix} \bullet & & \\ \bullet & \bullet & \\ \bullet & \bullet & \bullet \end{bmatrix} \begin{bmatrix} \bullet & \bullet & \bullet \\ & \bullet & \bullet \\ & & \bullet \end{bmatrix}$$

- Consider  $3 \times 3$  example:
  - $\mathbf{A}$  contains 9 elements  $a_{ij}$
  - $\mathbf{L}$  contains 6 elements  $l_{ij}$  and  $\mathbf{U}$  contains 6 elements  $u_{ij}$
  - So the relation  $\mathbf{A} = \mathbf{L}\mathbf{U}$  implies 9 equations (one for each known element  $a_{ij}$ ) involving 12 unknowns (the  $l_{ij}$ 's and  $u_{ij}$ 's)
  - This is an *underdetermined* (underconstrained) set of equations (infinitely many solutions)
  - We can choose 3 elements to get a unique solution
  - Common to set the diagonal elements of  $\mathbf{L}$  to 1

$$\mathbf{L} = \begin{bmatrix} 1 & & \\ \bullet & 1 & \\ \bullet & \bullet & 1 \end{bmatrix}$$

- If  $\mathbf{A} = \mathbf{LU}$  and all the diagonal elements of  $\mathbf{L}$  are 1's, the matrix  $\mathbf{A}$  can also be factorised in the form  $\mathbf{A} = \mathbf{LDU}'$ , where
  - $\mathbf{D}$  is a diagonal matrix with the diagonal elements  $u_{ii}$  of the original  $\mathbf{U}$  matrix
  - $\mathbf{U}'$  is the matrix generated by taking  $\mathbf{U}$  and multiplying each row  $i$  with  $\frac{1}{u_{ii}}$ , so that both  $\mathbf{L}$  and  $\mathbf{U}'$  have only 1's along their diagonals

$$\mathbf{A} = \mathbf{LDU}'$$

$$\begin{bmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{bmatrix} = \begin{bmatrix} 1 & & \\ \bullet & 1 & \\ \bullet & \bullet & 1 \end{bmatrix} \begin{bmatrix} \bullet & & \\ & \bullet & \\ & & \bullet \end{bmatrix} \begin{bmatrix} 1 & \bullet & \bullet \\ & 1 & \bullet \\ & & 1 \end{bmatrix}$$

\* This is (unsurprisingly) called **LDU decomposition** or **LDU factorisation**

- Computational complexity:
  - It takes  $\mathcal{O}\left(\frac{2}{3}n^3\right)$  operations to determine  $\mathbf{L}$  and  $\mathbf{U}$  for a given  $\mathbf{A}$  ( $n \times n$ )
  - So the computational complexity of performing the decomposition  $\mathbf{A} = \mathbf{LU}$  is the same as that of solving  $\mathbf{A}\vec{x} = \vec{b}$  with Gaussian elimination
  - LU decomposition can be seen as the matrix representation of Gaussian elimination
- Existence:
  - If a square matrix  $\mathbf{A}$ 
    - \* is *non-singular* (invertible), and
    - \* *all its leading principal minors* are non-zero (see note below)

then it admits an LU (or LDU) decomposition
  - If a square matrix  $\mathbf{A}$ 
    - \* is *singular* (not invertible),
    - \* has rank  $k$ , and
    - \* *the first  $k$*  of the leading principal minors are non-zero

then it admits an LU (or LDU) decomposition
  - (For more details about this, see textbooks on linear algebra)

*Side note: The leading principal minors of  $\mathbf{A}$  are the determinants of the square submatrices you get from  $\mathbf{A}$  if you start in the upper left-hand corner and grow the submatrix by one row and column at a time*

- Example: The leading principal minors of the matrix

$$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$$

are the determinants

$$\det \begin{bmatrix} a \end{bmatrix}, \det \begin{bmatrix} a & b \\ d & e \end{bmatrix} \text{ and } \det \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$$

## 2. What is it good for?

- Assume we *have performed* the LU decomposition
- We can now
  - solve **matrix equations**,  $\mathbf{A}\vec{x} = \vec{b}$ , at  $\mathcal{O}(n^2)$  cost
  - easily compute **the determinant**,  $\det \mathbf{A}$ , at  $\mathcal{O}(n)$  cost
  - find **the inverse**,  $\mathbf{A}^{-1}$ , at  $\mathcal{O}(n^3)$  cost
    - \* Finding  $\mathbf{A}^{-1}$  would have cost  $\mathcal{O}(n^4)$  if we did it by treating  $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$  as  $n$  equations of the form  $\mathbf{A}\vec{x} = \vec{b}$  and solved each one with Gaussian elimination

### Solving matrix equations after LU decomposition

- We have  $\mathbf{A} = \mathbf{LU}$
- Want to solve  $\mathbf{A}\vec{x} = \vec{b}$  for  $\vec{x}$
- We will solve  $\mathbf{A}\vec{x} = \mathbf{LU}\vec{x} = \vec{b}$  in two steps
  - First we define some notation:  $\vec{w} \equiv \mathbf{U}\vec{x}$ .
    - \* Since  $\vec{x}$  is unknown,  $\vec{w}$  is unknown.
    - \* We can now write  $\mathbf{LU}\vec{x} = \mathbf{L}\vec{w} = \vec{b}$
  - 1. Solve  $\mathbf{L}\vec{w} = \vec{b}$  for  $\vec{w}$
  - 2. Solve  $\mathbf{U}\vec{x} = \vec{w}$  for  $\vec{x}$
- **Step 1:** Solve  $\mathbf{L}\vec{w} = \vec{b}$  for  $\vec{w}$



- Consider example with  $4 \times 4$  matrices

$$\begin{bmatrix} l_{11} & 0 & 0 & 0 \\ l_{21} & l_{22} & 0 & 0 \\ l_{31} & l_{32} & l_{33} & 0 \\ l_{41} & l_{42} & l_{43} & l_{44} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}$$

- Solve by forward substitution:

- \* From  $l_{11}w_1 = b_1$  we immediately get the solution for  $w_1$ :

$$w_1 = \frac{1}{l_{11}}b_1$$

- \* From  $l_{21}w_1 + l_{22}w_2 = b_2$ , and using the now known  $w_1$ , we get

$$w_2 = \frac{1}{l_{22}}[b_2 - l_{21}w_1]$$

- \* Continuing the same way, we get

$$w_3 = \frac{1}{l_{33}}[b_3 - l_{31}w_1 - l_{32}w_2]$$

$$w_4 = \frac{1}{l_{44}}[b_4 - l_{41}w_1 - l_{42}w_2 - l_{43}w_3]$$

- \* In general, when  $\mathbf{A}$  is  $n \times n$ :

$$w_i = \frac{1}{l_{ii}} \left[ b_i - \sum_{j=1}^{i-1} l_{ij}w_j \right]$$

- \* Counting FLOPs (here we assume  $l_{ii} = 1$ ):

$$\sum_{i=1}^n (2i - 1) = n^2$$

(which is less than the  $\mathcal{O}(n^3)$  cost of doing the LU decomposition in the first place)

- Now we know  $\vec{w}$  and can so step 2

- **Step 2:** Solve  $\mathbf{U}\vec{x} = \vec{w}$  for  $\vec{x}$

$$\begin{bmatrix} u_{11} & u_{12} & u_{13} & u_{14} \\ 0 & u_{22} & u_{23} & u_{24} \\ 0 & 0 & u_{33} & u_{34} \\ 0 & 0 & 0 & u_{44} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

– Solve for  $\vec{x}$  by back substitution

\* From  $u_{44}x_4 = w_4$  we immediately get the solution for  $x_4$ :

$$x_4 = \frac{1}{u_{44}}w_4$$

\* From  $u_{33}x_3 + u_{34}x_4 = w_3$  we then get

$$x_3 = \frac{1}{u_{33}}[w_3 - u_{34}w_4]$$

\* And so on...

$$x_2 = \frac{1}{u_{22}}[w_2 - u_{23}x_3 - u_{24}x_4]$$

$$x_1 = \frac{1}{u_{11}}[w_1 - u_{12}x_2 - u_{13}x_3 - u_{14}x_4]$$

\* In general, when  $\mathbf{A}$  is  $n \times n$ :

$$x_n = \frac{1}{u_{nn}}w_n$$

$$x_i = \frac{1}{u_{ii}}\left[w_i - \sum_{j=i+1}^n u_{ij}x_j\right], \quad i = (n-1), (n-2), \dots, 1$$

\* This also takes  $\mathcal{O}(n^2)$  FLOPs, so the combined task of forward + back substitution to find  $\vec{x}$  has an  $\mathcal{O}(n^2)$  cost.

• In summary:

– If we already have  $\mathbf{A} = \mathbf{LU}$ , we can solve  $\mathbf{A}\vec{x} = \vec{b}$  at a total  $\mathcal{O}(n^2)$  cost as follows:

1. From  $\mathbf{L}\vec{w} = \vec{b}$ , find  $\vec{w}$  by forward substitution
2. From  $\mathbf{U}\vec{x} = \vec{w}$ , find  $\vec{x}$  by back substitution

**A difficulty**

- We need to store the full matrix  $\mathbf{A}$  ( $n \times n$ ) in memory for the LU decomposition
  - That's  $n^2$  floating-point numbers
  - At double precision (64 bits = 8 bytes per number), this requires  $n^2 \times 8$  bytes of memory
  - Example:
    - \* Assume  $n = 10^4$
    - \* We then need  $8 \times 10^8$  bytes  $\approx 10^9$  bytes = 1 GB of memory
    - \* So we can quite quickly run out of memory
  - Also, since the decomposition is an  $\mathcal{O}(n^3)$  operation, it will be slow when  $n$  is large

**Finding the determinant after LU decomposition**

- Once we have the decomposition  $\mathbf{A} = \mathbf{L}\mathbf{U}$ , computing the determinant of  $\mathbf{A}$  is trivial:

$$\begin{aligned}\det(\mathbf{A}) &= \det(\mathbf{L}\mathbf{U}) \\ &= \det(\mathbf{L}) \det(\mathbf{U}) \\ &= (1)(u_{11}u_{22} \dots u_{nn})\end{aligned}$$

where we have assumed that  $\mathbf{L}$  is on the standard form with 1's on the diagonal

- So in summary:

$$\det(\mathbf{A}) = \prod_{i=1}^n u_{ii}$$

or equivalently

$$\log(\det(\mathbf{A})) = \sum_{i=1}^n \log u_{ii}$$

- For large matrices it is often useful to work numerically with the logarithm of the determinant, rather than the determinant itself

### Finding the inverse after LU decomposition

- Once we have  $\mathbf{A} = \mathbf{LU}$ , we can find  $\mathbf{A}^{-1}$  at  $\mathcal{O}(n^3)$
- We know

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{I} = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{bmatrix} = \mathbf{A}\mathbf{A}^{-1}$$

- Write  $\mathbf{A}^{-1}$  as column vectors

$$\mathbf{A}^{-1} = \left[ \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} \dots \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} \right] = [\vec{\alpha}_1 \quad \vec{\alpha}_2 \quad \vec{\alpha}_3 \quad \vec{\alpha}_4]$$

where we have used a notation where e.g.  $\vec{\alpha}_1$  contains the first column of  $\mathbf{A}^{-1}$

$$\vec{\alpha}_1 \equiv \begin{bmatrix} (\mathbf{A}^{-1})_{11} \\ (\mathbf{A}^{-1})_{21} \\ (\mathbf{A}^{-1})_{31} \\ (\mathbf{A}^{-1})_{41} \end{bmatrix}$$

- Using  $\mathbf{A} = \mathbf{LU}$  and our notation for  $\mathbf{A}^{-1}$  we then have

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{LU} [\vec{\alpha}_1 \quad \vec{\alpha}_2 \quad \vec{\alpha}_3 \quad \vec{\alpha}_4] = \mathbf{I} = \begin{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} & \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \end{bmatrix}$$

- We can read this as four matrix equations,  $(\mathbf{LU})\vec{\alpha}_i = \hat{e}_i$ , where  $\hat{e}_i$  are the unit vectors and  $\vec{\alpha}_i$  are the unknown vectors we want to determine:

$$(\mathbf{LU})\vec{\alpha}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \dots \quad (\mathbf{LU})\vec{\alpha}_4 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

- Note that each matrix equation is on the familiar form  $(\mathbf{LU})\vec{x} = \vec{b}$ , which we have already seen how to solve
- This was for a  $4 \times 4$  example. In the general case, with  $\mathbf{A}$  being  $n \times n$ :
  - If we have  $\mathbf{A} = \mathbf{LU}$ , then finding  $\mathbf{A}^{-1}$  requires solving  $n$  matrix equations on the form  $(\mathbf{LU})\vec{x} = \vec{b}$
  - Solving each equation has an  $\mathcal{O}(n^2)$  cost
  - Thus, we can find  $\mathbf{A}^{-1}$  in  $\mathcal{O}(n^3)$  FLOPs
  - Finding the decomposition  $\mathbf{A} = \mathbf{LU}$  in the first place also had an  $\mathcal{O}(n^3)$  cost
  - So the **total cost of LU decomposition + finding  $\mathbf{A}^{-1}$  is  $\mathcal{O}(n^3)$**
  - (Recall that finding  $\mathbf{A}^{-1}$  would have required  $\mathcal{O}(n^4)$  FLOPs if we had naively solved each of the  $n$  equations  $\mathbf{A}\vec{\alpha}_i = \hat{e}_i$  “from scratch” using Gaussian elimination)

Finding the inverse of a large matrix can be the main bottleneck in many applications, e.g. in various methods in machine learning. In such cases matrix decomposition techniques like LU decomposition, or the related *Cholesky decomposition*  $\mathbf{A} = \mathbf{LL}^*$ , typically play a key role.

### 3. An algorithm for LU decomposition

- How to determine the elements of  $\mathbf{L}$  and  $\mathbf{U}$  such that  $\mathbf{A} = \mathbf{LU}$ ?
- Consider a  $4 \times 4$  case:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & \bullet & \bullet & \bullet \\ a_{31} & \bullet & \bullet & \bullet \\ a_{41} & \bullet & \bullet & \bullet \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ l_{21} & 1 & 0 & 0 \\ l_{31} & l_{32} & 1 & 0 \\ l_{41} & l_{42} & l_{43} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} & u_{14} \\ 0 & u_{22} & u_{23} & u_{24} \\ 0 & 0 & u_{33} & u_{34} \\ 0 & 0 & 0 & u_{44} \end{bmatrix}$$

- Let's look at the first column of  $\mathbf{A}$ , i.e. the elements  $a_{i1}$ :
  - From  $a_{11}$  we get:

$$a_{11} = \begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} u_{11} \\ 0 \\ 0 \\ 0 \end{bmatrix} = u_{11}$$

$$u_{11} = a_{11}$$

- From  $a_{21}$  we get:

$$a_{21} = \begin{bmatrix} l_{21} & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} u_{11} \\ 0 \\ 0 \\ 0 \end{bmatrix} = l_{21}u_{11}$$

$$l_{21} = \frac{a_{21}}{u_{11}} \quad (\text{where } u_{11} \text{ is now known})$$

- Similarly, from the equations  $a_{31} = l_{31}u_{11}$  and  $a_{41} = l_{41}u_{11}$  we get

$$l_{31} = \frac{a_{31}}{u_{11}}$$

$$l_{41} = \frac{a_{41}}{u_{11}}$$

- Now for the second column **A** (elements  $a_{i2}$ ):

- From  $a_{12}$ :

$$a_{12} = \begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} u_{12} \\ u_{22} \\ 0 \\ 0 \end{bmatrix} = u_{12}$$

$$u_{12} = a_{12}$$

- From  $a_{22}$ :

$$a_{22} = \begin{bmatrix} l_{21} & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} u_{12} \\ u_{22} \\ 0 \\ 0 \end{bmatrix} = l_{21}u_{12} + u_{22} \quad (\text{where } u_{22} \text{ is the unknown})$$

$$u_{22} = a_{22} - l_{21}u_{12}$$

– From  $a_{32}$ :

$$a_{32} = \begin{bmatrix} l_{31} & l_{32} & 1 & 0 \end{bmatrix} \begin{bmatrix} u_{12} \\ u_{22} \\ 0 \\ 0 \end{bmatrix} = l_{31}u_{12} + l_{32}u_{22} \quad (\text{where } l_{32} \text{ is the unknown})$$

$$l_{32} = \frac{a_{32} - l_{31}u_{12}}{u_{22}}$$

– Similarly, from  $a_{42} = l_{42}u_{22} + l_{41}u_{12}$ :

$$l_{42} = \frac{a_{42} - l_{41}u_{12}}{u_{22}}$$

- We can continue like this for the third and fourth columns
- The general pattern that appears gives us the following recipe:

$$l_{ij} = \frac{1}{u_{jj}} \left[ a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \right] \quad \text{with } i > j$$

$$u_{ij} = a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj} \quad \text{with } i \leq j$$

## Pivoting

- For numerical stability we need to avoid  $u_{11} \approx 0$
- Use a permutation matrix to interchange rows
- Instead of  $\mathbf{A} = \mathbf{LU}$  we will then have

$$\mathbf{A} = \mathbf{PLU}$$

or equivalently

$$\mathbf{P}^T \mathbf{A} = \mathbf{L} \mathbf{U}$$

- Here  $\mathbf{P}$  is called the *pivot matrix*
- It satisfies the relation  $\mathbf{P}^T = \mathbf{P}^{-1}$ , or equivalently,  $\mathbf{P} \mathbf{P}^T = \mathbf{I}$

The lecture on September 9 ended here.



## Topics in project 2

- Scaling equations
- Physics case: “the buckling beam”
  - Will lead to a two-point boundary value problem of the form

$$\begin{aligned}-\frac{d^2u}{dx^2} &= -\lambda u(x) \\ x &\in [0, 1] \\ u(0) &= 0 \\ u(1) &= 0\end{aligned}$$

- Eigenvalue problems, connection to two-point boundary value problem
- Jacobi’s rotation method for eigenvalue problems
- Coding:
  - Unit testing
  - Using the Armadillo library

## Scaling equations

- Also known **nondimensionalisation** or **using natural units** (but don't confuse this with the term *natural units* as used in e.g. particle physics)
- We can only represent a limited range of numbers on a computer
- We therefore want to avoid too large/small numbers in our codes.
  - First, simple approach: choose some sensible units
  - Better approach: scale away units, i.e. work with **dimensionless** variables
  - This is also useful for debugging, since it's then easier to notice if a result is surprisingly small/large

### Example 1: A simulation of the solar system

- Silly choice of units: kg, m
  - In these units we would get e.g.  
 $m_{\text{sun}} \approx 1.989 \times 10^{30} \text{ kg}$   
 $r_{\text{earth-sun}} \approx 1.496 \times 10^{11} \text{ m}$
- Sensible choice of units:  $M_{\text{sun}}$ , au
  - Then we would have  
 $m_{\text{sun}} = 1 M_{\text{sun}}$   
 $r_{\text{earth-sun}} = 1 \text{ au}$

### Example 2: Exponential decay of radioactive nuclei

- Consider the differential equation

$$\frac{dN(t)}{dt} = -\lambda N(t)$$

or, written on a standard form

$$\frac{dN(t)}{dt} + \lambda N(t) = 0$$

- Initial value:  $N(t = 0) = N_0$

- Units:

- $N$ : number of nuclei,  $[N] = 1$
- $t$ : time,  $[t] = \text{s}$
- $\lambda$ : decay constant,  $[\lambda] = \text{s}^{-1}$

For this simple example we know that the analytical solution is  $N(t) = N_0 e^{-\lambda t}$

- Define a scaled, dimensionless time variable (the **independent** variable):

$$\hat{t} \equiv \lambda t$$

- We can then rewrite the time derivative as

$$\frac{d}{dt} = \frac{d\hat{t}}{dt} \frac{d}{d\hat{t}} = \lambda \frac{d}{d\hat{t}}$$

- If we insert this in our original differential equation (and remember that  $\lambda$  is just a constant), we get

$$\lambda \frac{dN}{d\hat{t}} + \lambda N = 0$$

$$\frac{dN}{d\hat{t}} + N = 0$$

- Note that once we have changed variable from  $t$  to  $\hat{t}$ , we should think of  $N$  as new function, compared to the original  $N(t)$ .

- Technically we should have indicated this with some new notation, e.g.  $N_{\hat{t}}(\hat{t})$
- This new function  $N_{\hat{t}}$  is defined by the requirement  $N_{\hat{t}}(\hat{t}) = N(t)$
- So our differential equation above should more properly be written out as something like this:

$$\frac{dN_{\hat{t}}(\hat{t})}{d\hat{t}} + N_{\hat{t}}(\hat{t}) = 0$$

- However, this is often not written out so explicitly
- It is common to just let the argument in  $N(\hat{t})$  tell us that we should read this as  $N_{\hat{t}}(\hat{t})$ 
  - \* This is similar to how we often use the simple notation  $p(x)$  and  $p(y)$  for two different

probability distributions, when we technically should have written something like  $p_x(x)$  and  $p_y(y)$ , or  $f(x)$  and  $g(y)$

Since  $\hat{t}$  is dimensionless and defined using the typical time scale  $(1/\lambda)$  of our problem, we know

- that a step size  $h \ll 1$  along the  $\hat{t}$  axis will actually correspond to a *small* step size for this problem
- that solving our differential equation for a time span  $\hat{t} \in [0, \text{a few}]$  will be a reasonable starting point for studying the problem (rather than starting with a time span of say  $\hat{t} \in [0, 0.1]$  or  $\hat{t} \in [0, 100]$ )

- In the differential equation above, both  $N$  and  $\hat{t}$  are dimensionless numbers
- But the typical value for  $N$  may still impractical, e.g. some very large number
- So we may also consider scaling the **dependent** variable ( $N$ )
- Sensible choice here: use the initial value  $N_0$  to set a scale for  $N$

$$\hat{N} \equiv \frac{N}{N_0}$$

- The initial value in terms of our new variable  $\hat{N}$  (let's call the initial value  $\hat{N}_0$ ) is simply

$$\hat{N}_0 = \frac{N(t=0)}{N_0} = \frac{N_0}{N_0} = 1$$

- If we insert  $N = N_0 \hat{N}$  in our differential equation, we get

$$N_0 \frac{d\hat{N}}{d\hat{t}} + N_0 \hat{N} = 0$$

or simply

$$\frac{d\hat{N}}{d\hat{t}} + \hat{N} = 0$$

In this final form of the differential equation, all quantities are dimensionless, and both the *independent* variable  $\hat{t}$  and the *dependent* variable  $\hat{N}$  have typical values of  $\mathcal{O}(1)$ .

The analytical solution for our scaled differential equation is simply

$$\hat{N}(\hat{t}) = \hat{N}_0 e^{-\hat{t}} = e^{-\hat{t}} \quad (\text{recall that } \hat{N}_0 = 1)$$

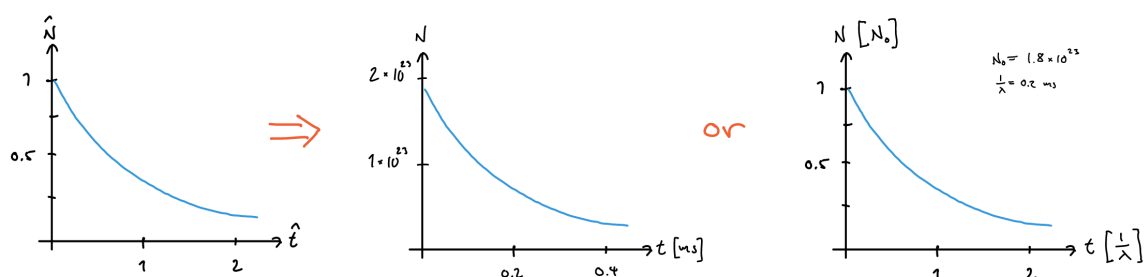
Using our definitions  $\hat{N} = \frac{N}{N_0}$  and  $\hat{t} = \lambda t$ , we can get the solution in terms of our original variables  $N$  and  $t$ :

$$\frac{N}{N_0} = e^{-\lambda t}$$

$$N(t) = N_0 e^{-\lambda t}$$

### Presenting results

- When presenting our results, we should either *use the original, dimensionful variables, or specify the natural units used*



**Figure 9:** Our raw numerical solution will be like the left-hand plot, but the result we present should be like the middle or right-hand plot

## Physics case for project 2: the buckling beam



**Figure 10:** The buckling beam

- Description of the setup:
  - A beam of some material is extended between two points  $x = 0$  and  $x = L$
  - The function  $u(x)$  denotes the shape of the beam (displacement from the  $x$  axis)
  - A force  $F$  is applied at the endpoint  $x = L$ , directed *into* the beam (i.e. along the  $x$  axis)
  - The material properties of the beam are collected in a constant  $\gamma$
  - The beam can not be displaced from the  $x$  axis at the end points:
    - \* Boundary conditions:  $u(0) = 0$  and  $u(L) = 0$
  - We'll consider a so-called *pin endpoints* case:
    - \* The beam shape  $u(x)$  can have non-zero derivative  $u'(x)$  at  $x = 0$  and  $x = L$
- If the force  $F$  is large enough, the configuration is *unstable*
  - Any tiny perturbation will cause the beam to *buckle* (bend) into some shape
- **Our question:** What beam shapes (i.e. what function forms  $u(x)$ ) can arise?
- We are considering this as a *static* problem – there is no time dependence here
  - We are looking for what static beam shapes are theoretically allowed under the conditions described above

- Differential equation:

$$\gamma \frac{d^2 u}{dx^2} = -F u(x)$$

- Do the following steps:
  - Scale the equation, to use dimensionless position  $\hat{x} \equiv \frac{x}{L}$
  - Define new notation:  $\lambda_c \equiv \frac{FL^2}{\gamma}$
  - Discretise with  $n$  steps (so we get  $n + 1$  points *including endpoints*, which means we have  $n - 1$  *interior* points)
- As you will see when you work on the project, after the above steps we end up with the following **eigenvalue problem**:

$$\mathbf{A} \vec{v} = \lambda \vec{v}$$

- $\mathbf{A}$  has size  $N \times N = (n - 1) \times (n - 1)$
- $\lambda = \lambda(h)$ , which should go to  $\lambda_c$  in the limit  $h \rightarrow 0$
- The elements  $v_i$  of  $\vec{v}$  are approximations to the exact  $u_i = u(\hat{x}_i)$
- $\vec{v}$  contains the *interior* points:

$$\vec{v} = \begin{bmatrix} v_1 \\ v_2 \\ \dots \\ v_N \end{bmatrix} = \begin{bmatrix} v_1 \\ v_2 \\ \dots \\ v_{n-1} \end{bmatrix}$$

- Our complete, approximate solution to the BVP:  $\vec{v}^* = [v_0, v_1, v_2, \dots, v_{n-1}, v_n]$
- The solutions we find for  $\mathbf{A} \vec{v} = \lambda \vec{v}$  are eigenvector-eigenvalue pairs,  $(\vec{v}^{(i)}, \lambda^{(i)})$
- In the continuous limit ( $h \rightarrow 0$ , i.e.  $n \rightarrow \infty$ ) these correspond to eigenvalues  $\lambda_c^{(i)}$  and *eigenfunctions*  $u^{(i)}(\hat{x})$
- Like in project 1 we have that  $\mathbf{A}$  is a tridiagonal matrix (due to the second derivative), but this time we keep the factor  $1/h^2$  from the second derivative as part of  $\mathbf{A}$

- $4 \times 4$  example:

$$\mathbf{A} = \begin{bmatrix} \frac{2}{h^2} & -\frac{1}{h^2} & 0 & 0 \\ -\frac{1}{h^2} & \frac{2}{h^2} & -\frac{1}{h^2} & 0 \\ 0 & -\frac{1}{h^2} & \frac{2}{h^2} & -\frac{1}{h^2} \\ 0 & 0 & -\frac{1}{h^2} & \frac{2}{h^2} \end{bmatrix} = \frac{1}{h^2} \begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix}$$

- Keeping  $\frac{1}{h^2}$  on the left-hand side makes it easier to see how our approach would work also for a more general case, where we had a more complicated differential equation:

$$[\text{some general operator}] u(x) = \lambda u(x)$$

$$\left[ \frac{d^2 u}{dx^2} + \dots \right] u(x) = \lambda u(x)$$

## Eigenvalue problems

TODO

## Jacobi rotation method

TODO

## In-lecture code discussion #3

TODO

- Topic: Debugging

## Iterative methods for solving matrix equations

- We now return to the topic of how to solve matrix equations of the general form  $\mathbf{A}\vec{x} = \vec{b}$



## Direct vs iterative methods

- We have previously looked at **direct methods** for solving  $\mathbf{A}\vec{x} = \vec{b}$ :
  - Gaussian elimination
    - \* In the case of a tridiagonal matrix  $\mathbf{A}$  this became the Thomas algorithm
  - LU decomposition
    - \* A starting point for many matrix tasks, including finding the inverse  $\mathbf{A}^{-1}$
    - \* When we have  $\mathbf{A}^{-1}$ , we get  $\vec{x}$  from matrix-vector multiplication:  $\vec{x} = \mathbf{A}^{-1}\vec{b}$
- The direct methods in principle give **the exact answer** in a **finite number of steps**
- But due e.g. to their reliance on exact equalities/cancellations in the mathematics, these methods can be susceptible to numerical problems
- Alternative class of methods: **iterative methods**
- Iterative methods iterate closer and closer to the true solution  $\vec{x}$ , but will generally never get there exactly
- Need some convergence criteria for deciding when to stop the iterations
- Compared to direct methods, iterative methods are often faster and with a smaller memory footprint
  - Particularly important for very large matrices
- Added bonus: the simplest iterative methods are often very easy to implement in code
  - As we will see, when an iterative method is expressed in matrix form it can look very complicated
  - But often it boils down to some very simple set of equations for the individual components  $x_i$  of the solution vector, and it is usually these component-level equations we would implement in our code

At this point you may be thinking: *Why didn't we use an iterative method to solve the matrix equation in project 1?* The answer is that we definitely could have done that – coding-wise it would have been very easy! But since Gaussian elimination is such an important stepping stone for many more advanced algorithms, it's an important algorithm to get some hands-on experience with, and that's why we focused on that (in the form of the Thomas algorithm) in project 1.

- In the following we will look at three examples of iterative methods:

- **The Jacobi method** (not to be confused with with Jacobi’s rotation method for eigenvalue problems  $\mathbf{A}\vec{x} = \lambda\vec{x}$ )
  - **Gauss-Seidel**
  - **Successive over-relaxation**
- But first we need to discuss how to check for convergence when using an iterative approach

### Checking convergence for an iterative method

- An iterative method has the conceptual form of a **while** loop in a program: we keep doing some steps over and over again until we have reached some stopping criterion
- Below are two approaches to how that stopping criterion can be formulated

#### Alternative 1: monitor the relative change in our estimate for $\vec{x}$

- We let  $\vec{x}^{(m)}$  denote our estimate for the true  $\vec{x}$  after  $m$  steps of our iterative method
- Let  $\epsilon$  denote the *relative change in the vector norm* from one iteration step ( $m$ ) to the next ( $m + 1$ ):

$$\epsilon = \left| \frac{|\vec{x}^{(m+1)}|_l - |\vec{x}^{(m)}|_l}{|\vec{x}^{(m)}|_l} \right|$$

- Here the subscript  $l$  indicates that we are using some particular  $l$ -norm to measure the length of a vector:

$$|\vec{x}|_l \equiv \left[ \sum_{i=1}^N |x_i|^l \right]^{\frac{1}{l}}$$

where  $N$  is the number of elements  $x_i$  in the vector  $\vec{x}$

- Some common choices for  $l$ :

$$l = 1 : \quad |\vec{x}|_1 = |x_1| + |x_2| + \dots + |x_N| \quad (\text{sum of absolute values of the vector elements})$$

$$l = 2 : \quad |\vec{x}|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_N^2} \quad (\text{the familiar Euclidian vector length})$$

$$l = \infty : \quad |\vec{x}|_\infty = \max_i |x_i| \quad (\text{the largest vector element in absolute value})$$

To see why  $|\vec{x}|_\infty$  ends up picking out the largest vector element, just check what happens in the general expression for  $|\vec{x}|_l$  when you set  $l$  to some very large value, say  $l = 1000$ . In that case, the

sum  $|x_1|^{1000} + |x_2|^{1000} + \dots + |x_N|^{1000}$  will be highly dominated by the term containing the largest vector element. For example, consider the vector  $\vec{x} = [0.9, 1.2, 1.19, 0.7]$ . Then we would have that  $0.9^{1000} + 1.2^{1000} + 1.19^{1000} + 0.7^{1000} \approx 1.2^{1000}$ . So the 1000-norm of this vector  $\vec{x}$  would be  $|\vec{x}|_{1000} = [0.9^{1000} + 1.2^{1000} + 1.19^{1000} + 0.7^{1000}]^{\frac{1}{1000}} \approx [1.2^{1000}]^{\frac{1}{1000}} = 1.2 = \max_i |x_i|$ . As we increase  $l$  towards  $l = \infty$  the approximation in the previous expression becomes more and more precise.

- When we have chosen a particular  $l$ -norm, we can decide on some small threshold value such that when  $\epsilon$  falls below this value, we stop the iterations
- In other words, we stop when  $\vec{x}^{(m+1)}$  is sufficiently similar to our previous  $\vec{x}^{(m)}$ , so that we don't think we will gain much from running more iterations
- *Pro:* Computing the  $l$ -norm of the new vector  $\vec{x}^{(m+1)}$  is a quick operation
- *Con:* We are here just checking whether or not our estimate for  $\vec{x}$  has (almost) stopped changing – we are *not* measuring how close our estimate  $\vec{x}^{(m)}$  is to actually solving the equation  $\mathbf{A}\vec{x} = \vec{b}$

### Alternative 2: monitor the residual

- Recall that a true solution  $\vec{x}$  should satisfy  $\mathbf{A}\vec{x} = \vec{b}$
- Let **the residual**  $\vec{r}$  be the vector difference between  $\vec{b}$  and  $\mathbf{A}\vec{x}^{(m)}$  (our current estimate for  $\mathbf{A}\vec{x}$ )

$$\vec{r} = \mathbf{A}\vec{x}^{(m)} - \vec{b}$$

- At every iteration we can then check the ratio of the lengths of  $\vec{r}$  and  $\vec{b}$ :

$$\frac{|\vec{r}|_l}{|\vec{b}|_l}$$

- We stop the iterations when this ratio falls below a threshold value that we have decided
- *Pro:* Compared to our previous approach to checking for convergence, this approach more directly monitors how far our estimate  $\vec{x}^{(m)}$  is from satisfying the equation  $\mathbf{A}\vec{x}^{(m)} = \vec{b}$
- *Con:* This approach is computationally more expensive than the previous approach, since it requires a matrix-vector multiplication.

### The Jacobi method

Despite the similar name, do not confuse this iterative method for solving  $\mathbf{A}\vec{x} = \vec{b}$  with *Jacobi's rotation method*, which is a method for solving an eigenvalue-eigenvector problem  $\mathbf{A}\vec{x} = \lambda\vec{x}$ .

- The starting point is to rewrite our matrix  $\mathbf{A}$  as the sum  $\mathbf{L} + \mathbf{D} + \mathbf{U}$ , where  $\mathbf{L}$  is a strictly lower-triangular matrix,  $\mathbf{D}$  is a diagonal matrix and  $\mathbf{U}$  is a strictly upper-triangular matrix
- Here is a  $3 \times 3$  example:

$$\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U}$$

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} & & \\ a_{21} & & \\ a_{31} & a_{32} & \end{bmatrix} \begin{bmatrix} a_{11} & & \\ & a_{22} & \\ & & a_{33} \end{bmatrix} \begin{bmatrix} & a_{12} & a_{13} \\ & & a_{23} \\ & & \end{bmatrix}$$

- Note that this simple decomposition is *not* the same as *LU decomposition*, which was the much more complicated task of writing  $\mathbf{A}$  as a *product*  $\mathbf{A} = \mathbf{L}\mathbf{U}$  (or alternatively as  $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{U}$ )
- We can now rewrite our matrix problem as follows

$$\begin{aligned} \mathbf{A}\vec{x} &= \vec{b} \\ (\mathbf{L} + \mathbf{D} + \mathbf{U})\vec{x} &= \vec{b} \\ \mathbf{D}\vec{x} &= -(\mathbf{L} + \mathbf{U})\vec{x} + \vec{b} \\ \vec{x} &= -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\vec{x} + \mathbf{D}^{-1}\vec{b} \end{aligned}$$

- This last expression has the form  $\vec{x} = [\text{some matrix}]\vec{x} + [\text{stuff independent of } \vec{x}]$
- We will take this as inspiration and suggest the following iterative recipe for how to change some current estimate  $\vec{x}^{(m)}$  to a new (and hopefully improved) estimate  $\vec{x}^{(m+1)}$ :

$$\boxed{\vec{x}^{(m+1)} = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\vec{x}^{(m)} + \mathbf{D}^{-1}\vec{b}}$$

- Note that since  $\mathbf{D}$  is diagonal, we already know that the inverse  $\mathbf{D}^{-1}$  is simply the diagonal matrix with the reciprocal matrix elements

$$\mathbf{D}^{-1} = \text{diag}\left(\frac{1}{a_{11}}, \frac{1}{a_{22}}, \dots, \frac{1}{a_{NN}}\right)$$

- The matrix  $\mathbf{T} = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})$ , which operates on  $\vec{x}^{(m)}$  on the right-hand side of our iteration rule, is called the **iteration matrix** or the **update matrix**
- Our iterative matrix recipe above might look complicated, but it turns into some very simple update rules when expressed in terms of the individual vector components  $x_i^{(m+1)}$

- To see this, let's write out the matrix multiplication above for a  $4 \times 4$  example:

$$x_1^{(m+1)} = \frac{1}{a_{11}} \left[ b_1 - a_{12}x_2^{(m)} - a_{13}x_3^{(m)} - a_{14}x_4^{(m)} \right]$$

$$x_2^{(m+1)} = \frac{1}{a_{22}} \left[ b_2 - a_{21}x_1^{(m)} - a_{23}x_3^{(m)} - a_{24}x_4^{(m)} \right]$$

$$x_3^{(m+1)} = \frac{1}{a_{33}} \left[ b_3 - a_{31}x_1^{(m)} - a_{32}x_2^{(m)} - a_{34}x_4^{(m)} \right]$$

$$x_4^{(m+1)} = \frac{1}{a_{44}} \left[ b_4 - a_{41}x_1^{(m)} - a_{42}x_2^{(m)} - a_{43}x_3^{(m)} \right]$$

- So for a general case with a length- $N$  vector, the iterative recipe on component form is simply

$$x_i^{(m+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{\substack{j=1 \\ j \neq i}}^N a_{ij}x_j^{(m)} \right]$$

- To start our iterative method, we need to make an **initial guess**,  $\vec{x}^{(0)}$
- In general, the closer our initial guess  $\vec{x}^{(0)}$  is to the true  $\vec{x}$ , the fewer iterations we need to perform
- Starting from *any* initial guess  $\vec{x}^{(0)}$ , the method will converge towards the true solution if the largest eigenvalue (in absolute value) of the update matrix  $\mathbf{T}$  is less than 1
  - This is an example of a general result regarding convergence of iterative methods:
    - \* The **spectral radius**  $\rho(\mathbf{M})$  of some matrix  $\mathbf{M}$  is defined as  $\rho(\mathbf{M}) = \max_i |\lambda_i|$
    - \* If the spectral radius of the update matrix  $\mathbf{T}$  satisfies  $\rho(\mathbf{T}) < 1$ , the iterative method will converge starting from any initial guess
- A useful special case to remember is that the requirement  $\rho(\mathbf{T}) < 1$  is always satisfied if the original matrix  $\mathbf{A}$  is **diagonally dominant**
  - The matrix  $\mathbf{A}$  is diagonally dominant if
    - \*  $|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|$  for every row  $i$  in  $\mathbf{A}$ , and
    - \*  $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$  for at least one of the rows in  $\mathbf{A}$

## Gauss-Seidel

- The Gauss-Seidel method builds on the Jacobi method, but it immediately starts using each computed estimate  $x_i^{(m+1)}$  in subsequent computations during the same iteration
- Let's see this in action for an example where  $\mathbf{A}$  is  $4 \times 4$ :

$$x_1^{(m+1)} = \frac{1}{a_{11}} \left[ b_1 - a_{12}x_2^{(m)} - a_{13}x_3^{(m)} - a_{14}x_4^{(m)} \right]$$

$$x_2^{(m+1)} = \frac{1}{a_{22}} \left[ b_2 - a_{21}x_1^{(m+1)} - a_{23}x_3^{(m)} - a_{24}x_4^{(m)} \right]$$

$$x_3^{(m+1)} = \frac{1}{a_{33}} \left[ b_3 - a_{31}x_1^{(m+1)} - a_{32}x_2^{(m+1)} - a_{34}x_4^{(m)} \right]$$

$$x_4^{(m+1)} = \frac{1}{a_{44}} \left[ b_4 - a_{41}x_1^{(m+1)} - a_{42}x_2^{(m+1)} - a_{43}x_3^{(m+1)} \right]$$

- So we are effectively doing a form of forward substitution
- For the general case where  $\mathbf{A}$  is  $N \times N$ , the update rule on component form is

$$x_i^{(m+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(m+1)} - \sum_{j=i+1}^N a_{ij}x_j^{(m)} \right]$$

- Note that the elements  $x_i^{(m+1)}$  must be computed *sequentially*, in the order  $i = 1, 2, \dots, N$
- In contrast, for the Jacobi method we could have computed all the  $x_i^{(m+1)}$  in parallel, since these computations were independent of each other
- On matrix form, the Gauss-Seidel method corresponds to starting from the matrix equation

$$\begin{aligned} \mathbf{A}\vec{x} &= \vec{b} \\ (\mathbf{L} + \mathbf{D} + \mathbf{U})\vec{x} &= \vec{b} \\ (\mathbf{L} + \mathbf{D})\vec{x} &= -\mathbf{U}\vec{x} + \vec{b} \\ \vec{x} &= -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{U}\vec{x} + (\mathbf{L} + \mathbf{D})^{-1}\vec{b} \end{aligned}$$

and from this suggest the update rule

$$\vec{x}^{(m+1)} = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{U}\vec{x}^{(m)} + (\mathbf{L} + \mathbf{D})^{-1}\vec{b}$$

- The Gauss-Seidel method typically converges faster than the Jacobi method

### Successive over-relaxation

- Successive over-relaxation (SOR) is a modified version of the Gauss-Seidel method
- SOR can achieve faster convergence than Gauss-Seidel, but it comes with a free *weight parameter*  $\omega$ , and the optimal choice for  $\omega$  is unknown in most cases
- SOR can be seen as the following schematic generalisation of the Gauss-Seidel method:

$$[\text{new}] = [\text{old}] + [\text{weight}] * [\text{the change } [\text{new}]-[\text{old}] \text{ from Gauss-Seidel}]$$

- If we choose the weight  $\omega = 1$  we simply get back the Gauss-Seidel method
- For  $\omega > 2$  the method will *not* converge
- For  $1 < \omega \leq 2$  the method *will* converge, but the optimal choice is problem-specific (and generally unknown)
  - Note that SOR will also converge for  $0 < \omega < 1$ , but such *under-relaxation* is typically not useful for improving convergence compared to e.g. Gauss-Seidel
- On component form, the update rule for SOR is

$$x_i^{(m+1)} = x_i^{(m)} + \frac{\omega}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(m+1)} - \sum_{j=i+1}^N a_{ij} x_j^{(m)} - a_{ii} x_i^{(m)} \right]$$

- As with the Gauss-Seidel method, the vector components  $x_i^{(m+1)}$  must be computed *sequentially*, in the order  $i = 1, 2, \dots, N$
- The update rule on matrix form is given by the following (rather confusing-looking) expression:

$$\vec{x}^{(m+1)} = [\omega \mathbf{L} + \mathbf{D}]^{-1} \left[ -[\omega \mathbf{U} + (\omega - 1)\mathbf{D}] \vec{x}^{(m)} + \omega \vec{b} \right]$$

- The best choice for  $\omega$  would be the choice that gives the smallest spectral radius for the update matrix  $\mathbf{T} = [\omega \mathbf{L} + \mathbf{D}]^{-1} [-[\omega \mathbf{U} + (\omega - 1)\mathbf{D}]]$

## In-lecture code discussion #4

TODO

- Topic: Classes in C++

## How to write a scientific report

TODO

## Grading system for reports

TODO

## Topics in project 3

- Methods for solving *initial value problems* (ODEs)
  - (In previous projects we have focused on *boundary value problems*)
- Main algorithm: Runge-Kutta, 4th order
- Coding:
  - Object-oriented programming, classes in C++
  - More use of the Armadillo library
- Writing a proper scientific report
- Physics case: simulating a *Penning trap*

## Physics of project 3: the Penning trap

TODO



## Code design for simulations

TODO

## Initial value problems

- Previously we have worked on solving *boundary value problems*, i.e. ODEs where we know the solution at the boundaries of our domain
  - We did this in both projects 1 and 2
  - After discretising our ODE using a finite difference scheme, we ended up with some matrix problem
- Now we will turn to *initial value problems*, i.e. ODEs where we know how the solution starts, and we need to evolve it forwards
  - This will be our focus in project 3
  - For initial value problems the independent variable is often time – but it doesn't have to be!
- Later in the course (project 5) we will look at *partial differential equations* (PDEs)
- We will discuss the following methods for solving initial value problems:
  - Forward Euler
  - Predictor-Corrector
  - Runge-Kutta
  - Leapfrog
  - Verlet, including *velocity Verlet*
- But before we look at the different methods, let's recap some basics about ODEs

## Classification

- Since different methods work well for different types of differential equations, it's useful to remind ourselves of some basic ODE classification

- First order ODE:

$$\frac{dy}{dt} = f(t, y)$$

- This equation is *first order*, since the highest-order derivative of  $y$  is  $\frac{dy}{dt}$
- In an initial value problem we would typically know  $y(t_0)$

- Second order ODE:

$$\frac{d^2y}{dt^2} = f\left(t, y, \frac{dy}{dt}\right)$$

- This equation is *second order*, since the highest-order derivative of  $y$  is  $\frac{d^2y}{dt^2}$
- In an initial value problem we would typically know both  $y(t_0)$  and  $y'(t_0)$
- A typical example from physics would be Newton's second law:

$$\frac{d^2x}{dt^2} = \frac{1}{m}F\left(t, x, \frac{dx}{dt}\right)$$

- Linear versus non-linear ODEs:

- Example of a *linear*, first-order ODE:

$$\frac{dy}{dt} = g^3(t)y(t)$$

- \* Linear, since the highest power of  $y$  in any term is 1

- Example of a *non-linear*, first-order ODE:

$$\frac{dy}{dt} = g^3(t)y(t) - h(t)y^2(t)$$

- \* Non-linear, since the dependence on  $y$  is more complicated than a simple linear dependence (note the  $y^2$  term)
- In project 3, the Coulomb interaction between the charged particles in the Penning trap give rise to a set of non-linear differential equations
- Non-linear problems typically require a numerical approach. (Often a closed-form analytical solution does not exist.)

## From a second-order equation to coupled first-order equations

- If we have to solve an  $m^{\text{th}}$ -order ODE, we can always rewrite it as a set of  $m$  first-order ODEs
- The resulting first-order equations will typically be *coupled*
- Let's take Newton's second law in one dimension as example:

$$\frac{d^2x}{dt^2} = \frac{1}{m} F\left(t, x, \frac{dx}{dt}\right)$$

- This is a second-order differential equation for  $x(t)$
- We can turn this into two first-order differential equations as follows:
  - First we *define* a new variable with a rather suggestive name:  $v \equiv \frac{dx}{dt}$
  - We can now write  $\frac{d^2x}{dt^2}$  as  $\frac{dv}{dt}$
  - We are then left with the following two first-order differential equations:

$$\frac{dx}{dt} = v(t) \quad (\text{from the definition of } v)$$

$$\frac{dv}{dt} = \frac{1}{m} F(t, x, v) \quad (\text{from the original diff. eq.})$$

- These equations are coupled, since the unknown  $v(t)$  appears in the diff. eq. for  $x(t)$ , and the unknown  $x(t)$  appears in the diff. eq. for  $v(t)$
- In project 3 you will use this approach to obtain a set of coupled, first-order equations that your numerical programs will solve

## Forward Euler

## Predictor-Corrector

## Runge-Kutta

## Leapfrog

## Verlet

## Algorithm classification

- Consistency, order of global error, one-step vs multi-step, stability

**Stability**

- Includes a small in-lecture code discussion: code example: [IVP\\_comparison](#)

**In-lecture code discussion #4**

TODO

- Topic: static variables

**Intro to probability**

TODO

**Properties of probabilities and probability density functions**

TODO

**My notation****Properties****Some important one-dimensional probability distributions****Probability density functions of many variables****Expectation values**

TODO

**Moments****Summarising probability distributions with a single number****Introduction to Monte Carlo methods**

TODO

**Physics of Project 4: the Ising model**

TODO

**Markov chains**

TODO

**Markov chain Monte Carlo (MCMC)**

TODO