

Numerical Methods for Chemical Engineers

Study guide for 6E5X0, 2017-2018

Ivo Roghair, Martin van Sint Annaland

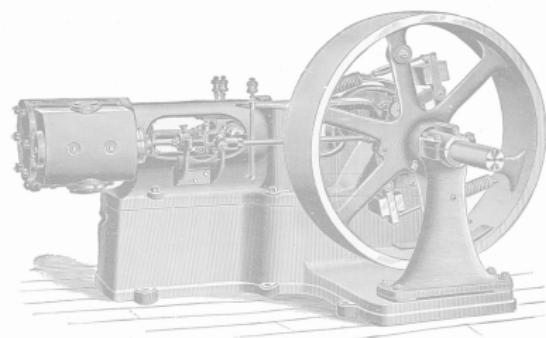
Chemical Process Intensification
Eindhoven University of Technology

13 November 2017

Numerical Methods

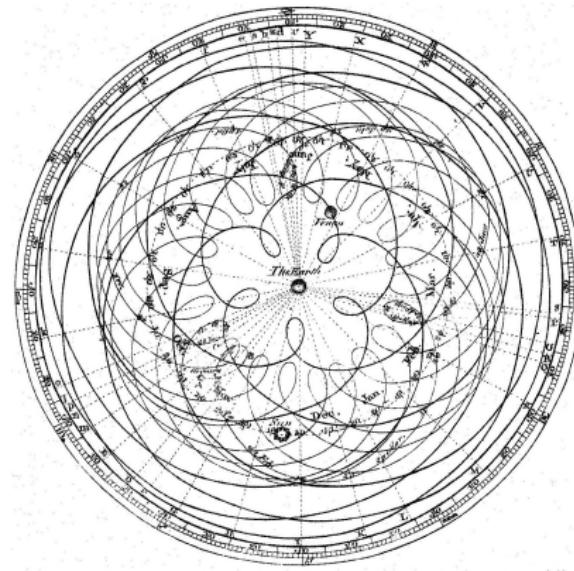
“Simulation and mathematical modeling will power the twenty-first century the way steam powered the nineteenth.”

— W.H. Press*



*Author of Numerical recipes, in “The Nature of Mathematical Modeling” by Neil Gershenfeld

Ptolemy and the Almagest



~150 AD. Development of numerical approximations to describe the motions of the heavenly bodies with accuracy matching reality sufficiently.

Numerical Methods

- Numerical analysis is concerned with obtaining approximate solutions to problems while maintaining reasonable bounds of error...
- ...because it is often impossible to obtain exact answers ...
- Numerical analysis makes use of algorithms to approximate solutions

Relevance

- Important to the world!
- E.g. in astronomy, construction, agriculture, architecture,
- And of course in Engineering!

...Chemical Engineering...

- Description of reactors and separators (dynamic and steady state)
- Computational fluid dynamics
- Thermodynamic equations of state
- Optimizing process performance
- Design and synthesis of processes
- Regression of data, e.g. isotherms, kinetics, ...

Course Schedule

Lecture	Date	Topic	Teacher
1	13/11/2017	Programming and algorithms (1)	IR
2	16/11/2017	Programming and algorithms (2)	IR
3	20/11/2017	Numerical errors	IR
4	23/11/2017	No lecture	
5	27/11/2017	Linear eqns: direct methods	IR
6	30/11/2017	Linear eqns: iterative methods	IR
7	04/12/2017	Non-linear equations	MSA
8	07/12/2017	Interpolation + integration	IR
9	11/12/2017	ODEs (1)	MSA
10	14/12/2017	ODEs (2) MSA	
11	18/12/2017	PDEs (1)	MSA
12	21/12/2017	No lecture	
13	08/01/2018	No lecture	
14	11/01/2018	Regression and optimization	IR

Course Objectives

- Gain experience with programming basics and algorithm design
- Acquire knowledge of and experience with different techniques for the numerical solution of systems of linear and non-linear algebraic and differential equations, as well as data analysis and optimization.
- Being able to solve various numerical problems using Matlab or Excel.

Prerequisites

The following subjects should give you enough hold-on to follow this course comfortably:

- Calculus A and B
- Linear Algebra
- Some basic MATLAB experience
 - We will shortly cover some aspects on MATLAB programming in the first lectures. Detailed documents and courses are provided on Canvas, for your own reference.
- Laptop with Matlab and Excel installed

Course Materials

- Lecture slides
- MATLAB scripts
- Additional articles
- There are some useful books:
 - Numerical recipes, W.H. Press et al.
 - Numerical methods for chemical engineering, K.J. Beers
 - Numerical methods for chemical engineers, A. Constantinides
 - Essential matlab-for engineers ,B.D. Hahn
 - Introduction to Numerical Methods and Matlab Programming for Engineers, T. Young and M.J. Mohlenkamp

Look on Canvas for the slides, exercises, scripts, assignments and additional documentation on MATLAB.

Assessment

5 assignments

- Each 20% of the final result
- Done in groups of 2 persons
 - Form groups via Canvas
 - Make sure that you have similar intentions!
- Short report (template provided, Canvas)

About the 5th assignment

- Short assignment + oral exam (in groups)
- Oral exam covers *all assignments and topics*
- Individual knowledge is assessed
- Grade needs to be at least a 5.0

Assignment grading

We will use rubrics to grade your reports. The following categories will be looked at.

- Use of numerical methods: e.g. built-in solvers vs. show implementation numerical methods
- Analysis of results: just the number is provided vs. high detail analysis and interpretation
- Programming skills: unstructured code, difficult to change vs. readable code with comments and UI
- Visualisation: unreadable graphs with no axes labels or legend vs. publication quality graphs, consistency between datasets

The first assignment will be graded via peer-review.

Assignment handout and deadlines

Hand-in your assignments via Canvas

- Deadlines are given on Canvas as well
- Deliver the report in PDF format
- Send along the scripts + necessities in a .zip

When delivering your final assignment, suggest a timeslot for the oral exam (e.g. via Canvas/assignment comment section or as Canvas message).

Contact information

Ivo Roghair

- E-mail: i.roghair@tue.nl
- Office: STW 0.37

Martin van Sint Annaland

- E-mail: m.v.sintannaland@tue.nl
- Office: STW 0.39

Milan Mihajlovic and Alessandro Battistella

For help with the exercises, Milan and Alessandro (STO 1.28) will help out during the lectures.

Some last remarks

- Tell us if something is not clear.
- We try to make the lectures interactive, working on examples and creating scripts as we go. It is advised that you work along with us to get the most out of this course!
- The exercises are meant to provide a jump start towards the assignments.
- We will always answer questions on the exercises. We may didactically answer questions on the assignments.
- During the lectures/tutorials we first and foremost work on the exercises. If they are done, you can work on the assignment if you want.

Some Acknowledgements



Some Real Acknowledgements

- To Roel Verstappen of Groningen University
- To Johan Hult of Cambridge University
- To Edwin Zondervan, now at Universität Bremen

Matlab and Programming 1

Programming basics and algorithms

Ivo Roghair, Martin van Sint Annaland

Chemical Process Intensification
Eindhoven University of Technology

Today's outline

① Introduction

② Variables

③ Creating algorithms

④ Functions

⑤ Conclusions

Programming

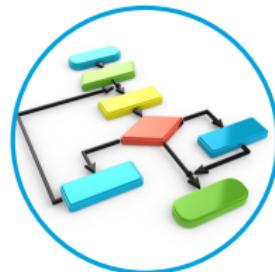
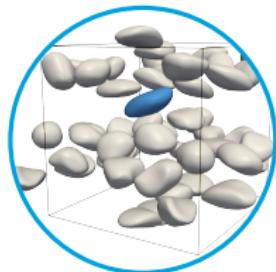
“Everybody in this country should learn to program a computer, because it teaches you to think..”

—Steve Jobs



Why?

- Scientific techniques depend in an increasing fashion upon computer programs and simulation methods
- Knowledge of programming allows you to automate routine tasks
- Ability to understand algorithms by inspection of the code
- Learn to think by dissecting a problem into smaller bits



Getting started

Start Matlab, and enter the following commands on the command line. Evaluate the output.

```
>> 2 + 3          % Some simple calculations
>> 2*3
>> 2*3^2         % Powers are done with ^
>> a = 2          % Storing values into the workspace
>> b = 3
>> c = (2*3)^2   % Parentheses set priority
>> 8/a-b
>> sin(a)         % Mathematical functions can be used
>> sin(0.5*pi)   % pi is an internal Matlab variable
>> 1/0            % Infinity is a thing ...
>> sqrt(-1)       % ... as are imaginary numbers
```

Introduction to programming

What is a program?

A program is a sequence of instructions that is written to perform a certain task on a computer.

- The computation might be something mathematical, such as solving a system of equations or finding the roots of a polynomial
- It can also be a symbolic computation, such as searching and replacing text in a document
- A program may even be used to compile another program
- A program consists of one or more *algorithms*

Getting started

- Use an *integrated development environment*
 - Matlab
 - MS Visual Studio/Code
 - Eclipse
 - Dev C++
 - IDLE, Canopy (express)
- Create a simple program:
 - Hello world
 - Find the roots of a parabola
 - Find the greatest common divisor of two numbers

Some often used programming languages

Python

- Many functionalities available
- Smooth learning curve
- Slow compared to compiled languages
- Many freely available editors

C / C++ / C#

- Many functionalities available
- Steeper learning curve
- Needs compilation, very fast (HPC)
- Freely available (gcc, MSVC)

Pascal

- Limited number of libraries available
- Steep learning curve
- Compiled language, may be fast
- Some free compilers (fpc)

Spreadsheet (Excel, Google Docs, ...)

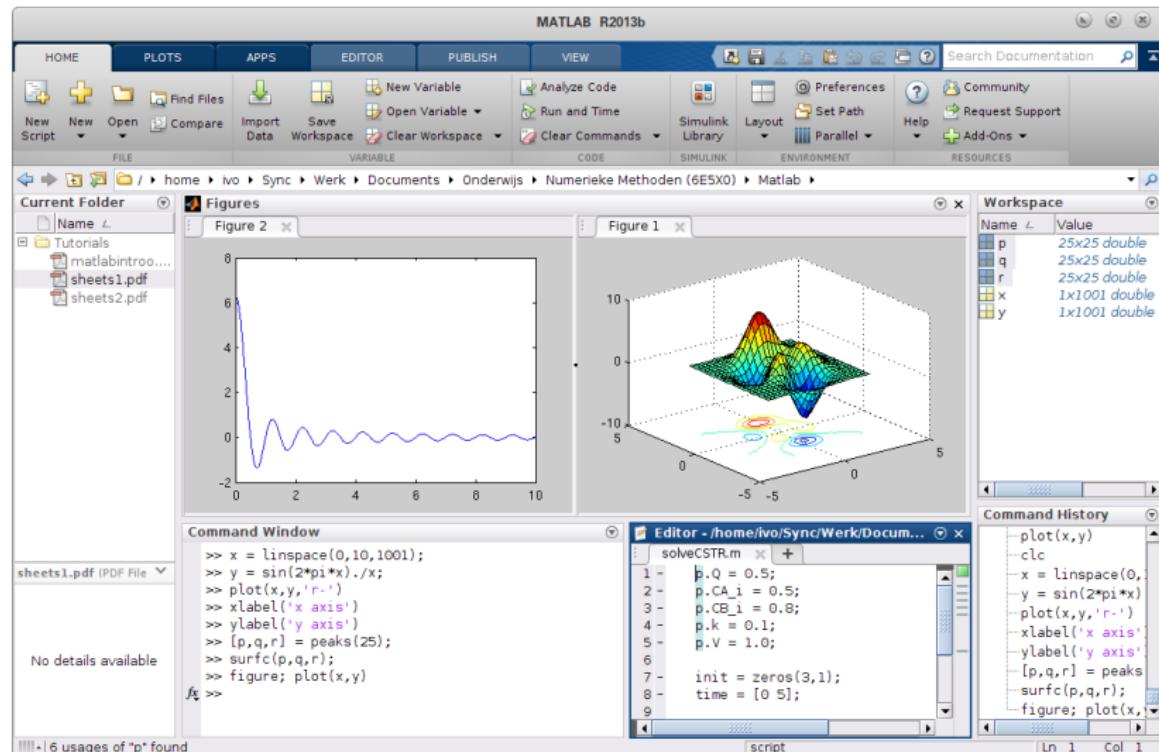
- High availability
- Low learning curve
- Very limited for larger problems, unbeatable for quick calculations
- Not always free

Matlab

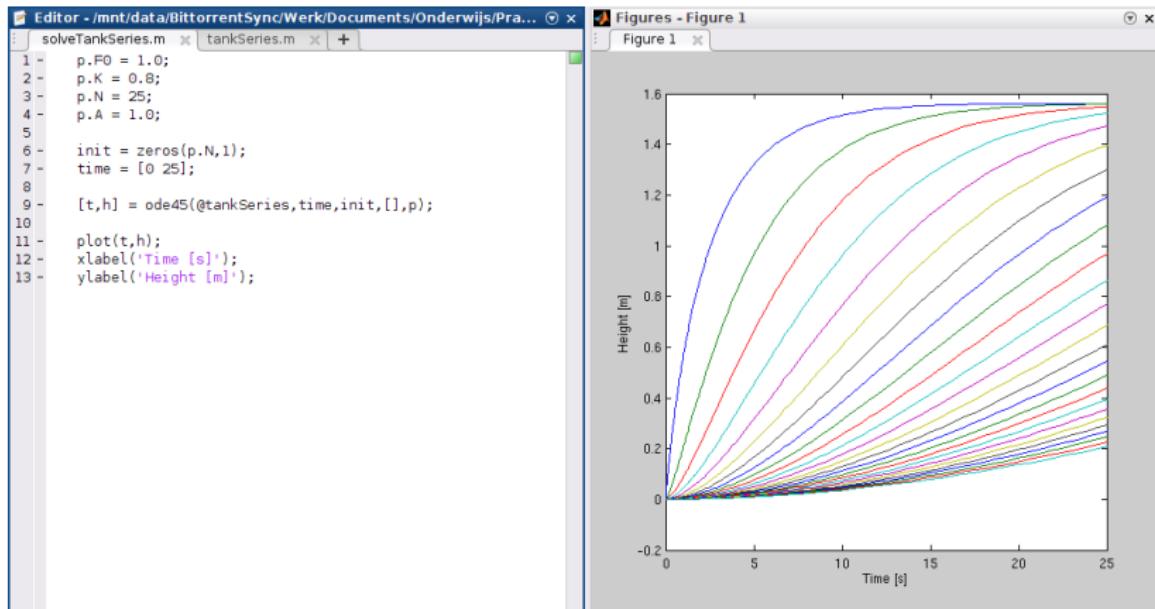
- Many functionalities built-in (80+ toolkits!)
- Slow compared to compiled languages

- Fairly smooth learning curve
- Needs a license (alternatives: SciLab, GNU Octave)

Versatility of Matlab



Versatility of Matlab: ODE solver

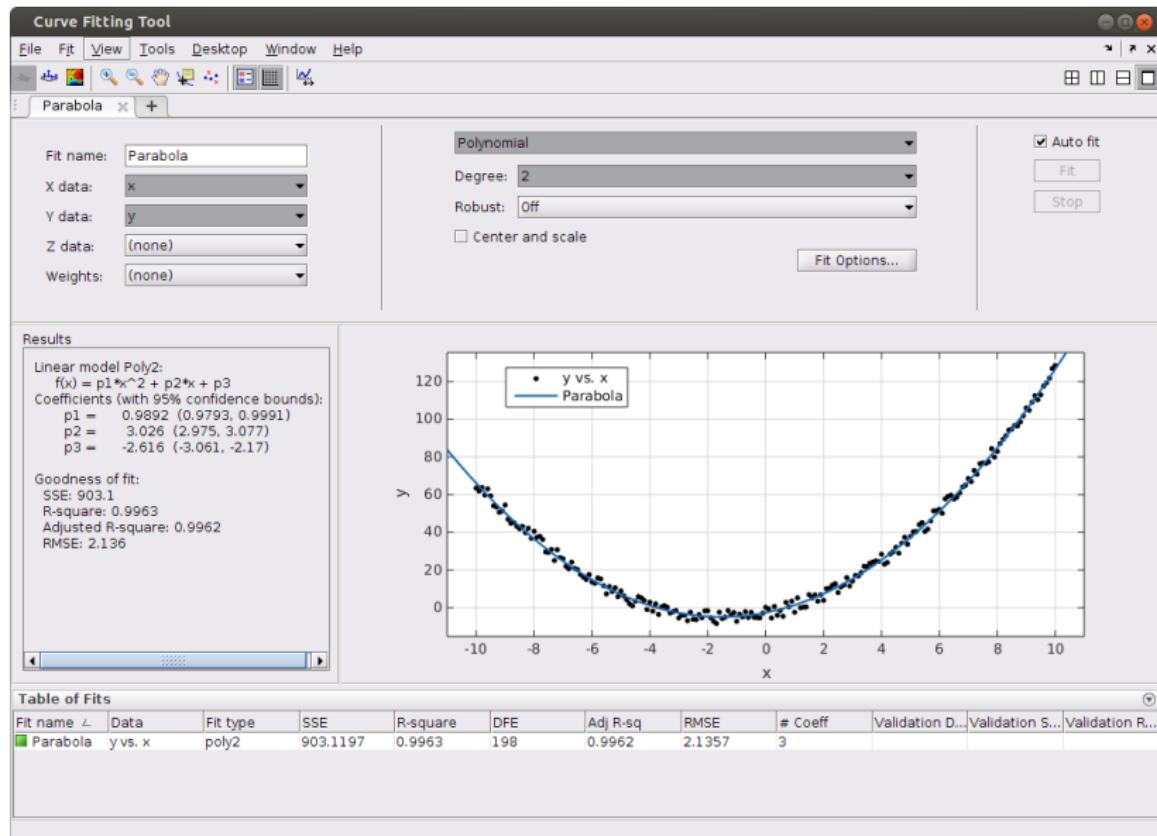


Versatility of Matlab: Image analysis

```
I = imread('bubbles.png');  
BW = rgb2gray(I);  
E = edge(BW, 'canny');  
F = imfill(E, 'holes');  
result = regionprops(F);
```



Versatility of Matlab: Curve fitting



Matlab help

- Matlab documentation: `doc` or `help` function
- Canvas page
- Introduction to Numerical Methods and Matlab Programming for Engineers. T. Young and M.J. Mohlenkamp (2015). GNU-licensed document, online
- Search the web!



Today's outline

① Introduction

② Variables

③ Creating algorithms

④ Functions

⑤ Conclusions

Terminology

Variable Piece of data stored in the computer memory, to be referenced and/or manipulated

Function Piece of code that performs a certain operation/sequence of operations on given input

Operators Mathematical operators (e.g. + - * or /), relational (e.g. < >or ==, and logical operators (&&, ||)

Script Piece of code that performs a certain sequence of operations without specified input/output

Expression A command that combines variables, functions, operators and/or values to produce a result.

Variables in Matlab

- Matlab stores variables in the *workspace*
- You should recognize the difference between the *identifier* of a variable (its name, e.g. `x`, `setpoint_p`), and the data that it actually stores (e.g. 0.5)
- Matlab also defines a number of variables by default, e.g. `eps`, `pi` or `i`.
- You can assign a variable by the = sign:

```
>> x = 4*3  
x =  
    12
```

- If you don't assign a variable, it will be stored in `ans`
- Clearing the workspace is done with `clear`.

Vectors in Matlab (1)

A row vector:

```
>> v = [0 1 2 3]
```

A column vector by separating elements with semi-colons:

```
>> u = [9; 10; 11; 12; 13; 14; 15]
```

Access (i.e. read) an entry in a vector:

```
>> u(2)
```

Manipulate the value of that entry:

```
>> u(2)=47
```

Get a slice of a vector:

```
>> u([2 3 4]) % With colon operator: u(2:4)
```

Transposing vectors:

```
>> w = v'
```

Vectors in Matlab (2)

Manual definition may be cumbersome. A colon (:) generates a list:

```
>> a = 1:10      % Default stride is 1  
>> x = -1:.1:1   % start:stride:stop specifies list
```

Or, when you prefer to set the *number of elements* instead of the step size:

```
>> y = linspace(0,10,11)  
>> p = logspace(2,6,5)
```

Manipulating multiple components:

```
>> y([1 4:7]) = 1
```

Or (by supplying a vector instead of a scalar):

```
>> y([1 4:7]) = 16:20 % equivalent to y([1 4 5 6 7]) =  
[16 17 18 19 20]
```

Practice

Given a vector

$$x = [2 \ 4 \ 6 \ 8 \ 10 \ 12 \ 14 \ 16 \ 18 \ 20 \ 30 \ 40 \ 50 \ 60 \ 70 \ 80]$$

- Find a way to define the vector without typing all individual elements
- Investigate the meaning of the following commands:

```
>> y = x(5:end)  
>> y(4)  
>> y(4) = []  
>> sum(x)  
>> mean(x)  
>> std(x)  
>> max(x)  
>> fliplr(x)  
>> diff(x)
```

Operations on vectors (1)

```
>> e = 1:5  
>> f = 2*e  
>> g = 4*f + 20  
>> h = e^2
```

... wait ... what's that?

```
Error using ^  
Inputs must be a scalar and a square matrix.  
To compute elementwise POWER, use POWER (.^) instead.
```

Matlab uses matrix operations by default, we should use a dot operator to make operations element-wise for *, / and ^.

```
>> e.^2
```

Operations on vectors (2)

To demonstrate the matrix product:

```
>> p = [1; 1; 1]
>> q = [1 2 3]
>> p*q    % which is not equal to q*p
```

All kinds of mathematical functions on vectors typically operate on elements:

```
>> x = linspace(0,2*pi,100);
>> s = sin(x)
>> e = exp(x)
```

Building blocks: Mathematics and number manipulation

Programming languages usually support the use of various mathematical functions (sometimes via a specialized library). Some examples of the most elementary functions in Matlab:

Command	Explanation
<code>cos(x)</code> , <code>sin(x)</code> , <code>tan(x)</code>	Cosine, sine or tangens of x
<code>mean(x)</code> , <code>std(x)</code>	Mean, st. deviation of vector x
<code>exp(x)</code>	Value of the exponential function e^x
<code>log10(x)</code> , <code>log(x)</code>	Base-10/Natural logarithm of x
<code>floor(x)</code>	Largest integer smaller than x
<code>ceil(x)</code>	Smallest integer that exceeds x
<code>abs(x)</code>	Absolute value of x
<code>size(x)</code>	Size of a vector x
<code>length(x)</code>	Number of elements in a vector x
<code>rem(x,y)</code>	Remainder of division of x by y

Printing results

You can prevent displaying the outcome of a command by adding a semi-colon at the end of a line:

```
>> c = linspace(0,10,11);  
>> length(c)  
>> c  
>> size(c)
```

Altering the display format can be done using the `format` command:

```
>> format compact % loose  
>> format long % short
```

Simple plotting

Make a plot of the following table

T (°C)	5	20	30	50	55
μ (Pa·s)	0.08	0.015	0.009	0.006	0.0055

```
>> x = [ 5 20 30 50 55 ]  
>> y = [ 0.08 0.015 0.009 0.006 0.0055]
```

```
>> plot(x,y)
```

```
>> plot(x,y, '*')
```

```
>> plot(x,y, 'r--')
```

```
>> plot(x,y, 'ko-', 'LineWidth', 2)
```

```
>> xlabel('Temperature [^\circ C]')  
>> ylabel('Viscosity [Pa s]')  
>> title('Experiment 1')
```

Practice

Create plots of the following functions in a single figure for $x \in \{0, 2\pi\}$:

$$y_1 = \cos x$$

$$y_2 = \arctan x$$

$$y_3 = \frac{\sin x}{x}$$

Strategies to draw multiple graphs in 1 figure:

```
>> plot(x,y1,x,y2,x,y3)
```

```
>> plot(x,y1)
>> hold on; % Maintain drawn plots in current figure
>> plot(x,y2)
>> plot(x,y3) % The 'hold-property' was already set
```

Matrices in Matlab

Matrix A is defined as:

$$A = \begin{bmatrix} 8 & 1 & 6 \\ 3 & 5 & 7 \\ 4 & 9 & 2 \end{bmatrix}$$

In Matlab:

```
>> A = [ 8 1 6; 3 5 7; 4 9 2 ]
```

Elements can be accessed/manipulated by the following syntax:

```
>> A(3,1) % Third row, first column, also A(3)
>> A(3,:) = [2 4 8] % Set entire third row
>> A(:,3) % Print third column
>> A(A>5) = 2 % Set elements by condition
```

There are a few functions that help creating matrices:

```
>> A = zeros(4) % A 4x4 matrix with zeros
>> A = ones(4,1) % A 4-element vector with ones
>> A = eye(3) % Identity matrix of 3x3
>> A = rand(3,4) % A 3x4 matrix with random numbers
```

Practice

- Find a *short* Matlab expression to create the following matrix:

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ 9 & 7 & 5 & 3 & 1 & -1 & -3 \\ 4 & 8 & 16 & 32 & 64 & 128 & 256 \end{bmatrix}$$

- Investigate the command `max(A)`. What does it give?
- How to obtain the maximum for each row?
- Use a vector multiplication to compute the following matrix:

$$A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \\ 4 & 4 & 4 & 4 \end{bmatrix}$$

Datatypes and variables

Matlab uses different types of variables:

Datatype	Example
string	'Wednesday'
integer	15
float	0.15
vector	[0.0; 0.1; 0.2]
matrix	[0.0 0.1 0.2; 0.3 0.4 0.5]
struct	sct.name = 'MyDataName' sct.number = 13
logical	0 (false) 1 (true)

About variables

- Matlab variables can change their type as the program proceeds (this is not common for other programming languages!):

```
>> s = 'This is a string'  
s =  
This is a string  
>> s = 10  
s =  
10
```

- Vectors and matrices are essentially *arrays* of another data type. A vector of *struct* is therefore possible.
- Variables are *local* to a function (more on this later).

Today's outline

① Introduction

② Variables

③ Creating algorithms

④ Functions

⑤ Conclusions

Building blocks: conditional statements

if-statement: Check whether a (set of) condition(s) is met.

```
num = floor (10 * rand + 1);
guess = input ('Your guess please : ');
if ( guess ~= num )
    disp (['Wrong, it was ',num2str(num),'. Kbye.']);
else
    disp ('Correct !');
end
```

Other relational operators

Combining conditional statements

==	is equal to
<=	is less than or equal to
>=	is greater than or equal to
<	is less than
>	is greater than

;&&	and
	or
xor	exclusive or

Building blocks: loops

for-loop: Performs a block of code a certain number of times.

```
>> p(1) = 1;
>> p(2) = 1;
>> for i = 2:10
p(i+1) = p(i)+p(i-1);
end
>> p
p =
    1     1     2     3     5     8    13    21    34    55    89
```

Building blocks: indeterminate repetition

while-loop: Performs and repeats a block of code until a certain condition.

```
num = floor (10* rand +1) ;
guess = input ('Your guess please : ');

while ( guess ~= num )
    guess = input ('That is wrong. Try again ... ');
end

if (isempty(guess))
    disp('No number supplied - exit');
else
    disp ('Correct!');
end
```

Example algorithm

Compute the factorial of N : $N! = N \cdot (N - 1) \cdot (N - 2) \cdots 2 \cdot 1$

How to deal with this?

Naive approach

```
Z = 1;  
Z = Z*2;  
Z = Z*3;  
Z = Z*4;  
... etc ...
```

For-loop

```
Z = 1;  
for i = 1:N  
    Z = Z*i;  
end
```

While-loop

```
Z = 1;  
i = 1;  
while (i<=N)  
    Z = Z*i;  
    i = i+1;  
end
```

Note: N must be set beforehand!

Note: Pay attention to the relational operators!

Building blocks: case selection

switch-statement: Selects and runs a block of code.

```
[dnum,dnam] = weekday(now);
switch dnum
    case {1,7}
        disp('Yay! It is weekend!');
    case 6
        disp('Hooray! It is Friday!');
    case {2,3,4,5}
        disp(['Today is ', dnam]);
    otherwise
        disp('Today is not a good day...');

end
```

Input and output

Many programs require some input to function correctly. A combination of the following is common:

- Input may be given in a parameters file (“hard-coded”)
- Input may be entered via the keyboard

```
>> a = input('Please enter the number ');
```

- Input may be read from a file, e.g.

```
>> data = getfield(importdata('myData.txt', ' ', 4), 'data');
>> numdata = xlsread('myExcelDataFile.xls');
```

- There are many more advanced functions, e.g. `fread`, `fgets`, ...

Input and output

Output of results to screen, storing arrays to a file or exporting a graphic are the most common ways of getting data out of Matlab:

- Results of each expression are automatically shown on screen as long as the line is not ended with a semi-colon;
- Output may be stored via the GUI:
 - Use the 'Export Setup' function
 - Save figure (use .fig, .eps or .png, not .jpg or .pcx)
 - Save variables (right click, save as)
- Save variables automatically (scripted):

```
>> savefile = 'test.mat';
>> p = rand(1,10);
>> q = ones(10);
>> save(savefile, 'p', 'q')
```

- More advanced functions can be found in e.g. `fwrite`, `fprintf`,
...

Today's outline

① Introduction

② Variables

③ Creating algorithms

④ Functions

⑤ Conclusions

Functions - general

A function in a programming language is a program fragment that performs a certain task. Creating functions keeps your code clean, re-usable and structured.

- You can use functions supplied by the programming language, and define functions yourself
- Functions take one or more input parameters (*arguments*), and *return* an output (*result*).
 - If functions do not return a result, it is called a procedure
- In Matlab, functions are defined as follows (2 output variables and 3 input arguments):

```
function [out1, out2] = myFunction(in1, in2, in3)
```

Functions - locality and arguments

- You are supplying arguments to a function because it does not have access to previously defined variables. This is called *locality*.
 - This does not include global variables - but they're evil!
 - Local variables created in a function are not accessible to other functions unless they are returned or supplied as an argument!

Exercise: write a function that takes 3 variables, and returns the average:

Approach 1

```
function res = avg1(a,b,c)
    mySum = a + b + c;
    res = mySum / 3;
end
```

Approach 2

```
function res = avg2(a,b,c)
    data = [a; b; c];
    res = mean(data);
end
```

Exercise: create a function

Compute $N! = N \cdot (N - 1) \cdot (N - 2) \cdots 2 \cdot 1$

Create a function of our while-loop approach with N the argument:

Original script

```
Z = 1;  
i = 1;  
while (i<=N)  
    Z = Z*i;  
    i = i+1;  
end
```

Function

```
function Z = fact_while(N)  
  
Z = 1;  
i = 1;  
while (i<=N)  
    Z = Z*i;  
    i = i+1;  
end  
  
end
```

Functions - checking input

The function we created computes the factorial correctly!

- When the supplied argument is positive and
- When the supplied argument is a natural number...



- In this case, we should check the user input to prevent an infinite loop:

```
if (fix(N) ~= N) | (N<0)
    disp 'Provide a positive
          integer number!'
    return;
end
```

- If no check can be done before a while-loop, you may want to stop after x loops

Functions - checking input

The whole factorial function, including comments:

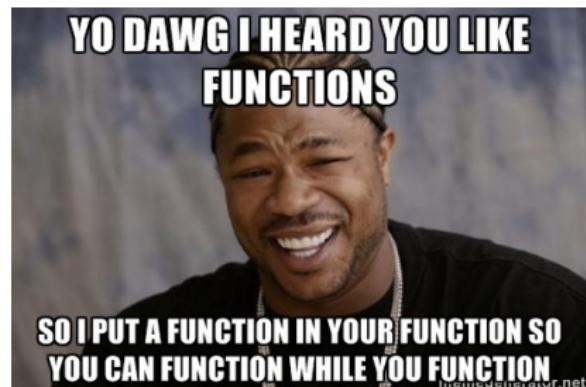
```
function Z = fact_while(N)
    % This function computes a factorial of input value N
    % Usage   : fact_while(N)
    % N       : value of which the factorial is computed
    % returns: factorial of N

    % Catch non-integer case
    if (fix(N) ~= N) | (N < 0)
        disp 'Provide a positive integer number!'
        return;
    end

    Z = 1;
    i = 1;
    while (i <= N)
        Z = Z*i;
        i = i+1;
    end
```

Recursion

- In order to understand recursion, one must first understand recursion
- A recursive function is called by itself (a function within a function)
 - This could lead to infinite calls;
 - A base case is required so that recursion is stopped;
 - Base case does not call itself, simply returns.



Recursion: example

```
function out = mystery(a,b)
if (b == 1)
    % Base case
    out = a;
else
    % Recursive function call
    out = a + mystery(a,b-1);
end
```

- What does this function do?
- Can you spot the error?
- How deep can you go? Which values of b don't work anymore?

Recursion: exercise

Create a function computing the factorial of N , based on recursion.

```
function res = fact_recursive(x)

% Catch non-integer case
if (fix(x) ~= x) | (x<0)
    disp 'You should provide a positive integer number
          only'
    return;
end

if (x > 1)
    res = x*fact_recursive(x-1);
else
    res = 1;
end

end
```

Today's outline

① Introduction

② Variables

③ Creating algorithms

④ Functions

⑤ Conclusions

In conclusion...

- Matlab: A versatile development environment, with excellent vector and matrix computations
- Programming basics: variables, operators and functions, locality of variables, recursive operations
- Next lecture: advanced practices (prepare, read instructions on Canvas)
- For now: exercises 1-4 (basics), 5+6 (advanced).

Matlab and Programming 2

Advanced programming techniques

Ivo Roghair, Martin van Sint Annaland

Chemical Process Intensification
Eindhoven University of Technology

Today's outline

- ① Coding in style
- ② Error management
- ③ Visualisation
- ④ Functions: the sequel
- ⑤ Excel
- ⑥ Algorithms
- ⑦ Conclusions

If anything sticks today, let it be this

Your code will not be understood by anyone

That includes future-you

Code organization

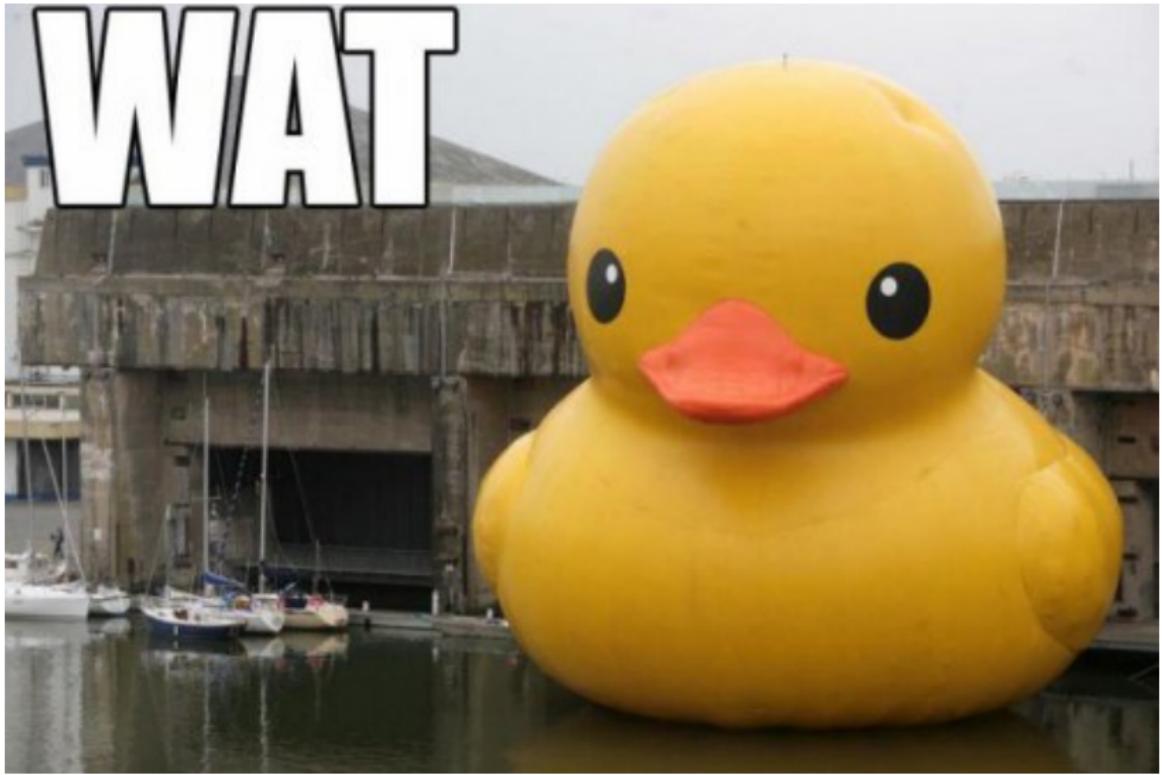
- Optimization of a code is time-consuming and complicated
- The more you optimize your code, the less readable it becomes
- But... You can write it in a such way that it will be flexible and easy to maintain
- Especially important in team work
- Any person has its own handwriting. Any programmer has its own coding style.

The coding style ≡ handwriting...

Interpret the following code

```
s=checksc();  
if(s==true)  
a=cb();  
b=cfrsp();  
if(a<5)  
if(b>5)  
a=gtbs();  
end  
if(a>b)  
ubx();  
end  
end  
else  
brn();  
gtbs();  
end
```

WAT



Let's change that a bit... Indentation

Shown here with 2 spaces of indentation, Matlab uses 4 by default!

```
s=checksc();  
if(s==true)  
a=cb();  
b=cfrsp();  
if(a<5)  
if(b>5)  
a=gtbs();  
end  
if(a>b)  
ubx();  
end  
end  
else  
brn();  
gtbs();  
end
```

```
s = checksc();  
if (s == true)  
    a = cb();  
    b = cfrsp();  
    if (a < 5)  
        if (b > 5)  
            a = gtbs();  
        end  
        if (a > b)  
            ubx();  
        end  
    end  
else  
    brn();  
    gtbs();  
end
```

Readable variables and function names

```
s = checksc();
if (s == true)
    a = cb();
    b = cfrsp();
    if (a < 5)
        if (b > 5)
            a = gtbs();
        end
        if (a > b)
            ubx();
        end
    end
else
    brn();
    gtbs();
end
```

```
IAmFree = checkSchedule();
if (IAmFree == true)
    books = countBooks();
    shelfSize =
        countFreeSpaceShelf();
    if (books < 5)
        if (shelfSize > 5)
            books = goToBookStore();
        end
        if (books > shelfSize)
            useBox();
        end
    end
else
    burnBooks();
    goToBookStore();
end
```

Get rid of obscure constants in the code

```
IAmFree = checkSchedule();
if (IAmFree == true)
    books = countBooks();
    shelfSize =
        countFreeSpaceShelf();
    if (books < 5)
        if (shelfSize > 5)
            books = goToBookStore();
        end
        if (books > shelfSize)
            useBox();
        end
    end
else
    burnBooks();
    goToBookStore();
end
```

```
IAmFree = checkSchedule();
if (IAmFree == true)
    books = countBooks();
    shelfSize =
        countFreeSpaceShelf();
    if (books < maxShelfSize)
        if (shelfSize >
            minBooksNeeded)
            books = goToBookStore();
        end
        if (books > shelfSize)
            useBox();
        end
    end
else
    burnBooks();
    goToBookStore();
end
```

That's more like it!

```
s=checksc();
if(s==true)
a=cb();
b=cfrrsp();
if(a<5)
if(b>5)
a=gtbs();
end
if(a>b)
ubx();
end
end
else
brn();
gtbs();
end
```

```
IAmFree = checkSchedule();
if (IAmFree == true)
books = countBooks();
shelfSize = countFreeSpaceShelf();
if (books < maxShelfSize)
if (shelfSize > minBooksNeeded)
books = goToBookStore();
end
if (books > shelfSize)
useBox();
end
end
else
burnBooks();
goToBookStore();
end
```

Writing readable code

Good code reads like a book.

- When it doesn't, make sure to use comments. In Matlab, everything following `% is a comment`
- Prevent “smart constructions” in the code
- Re-use working code (i.e. create functions for well-defined tasks).
- Documentation is also useful, but hard to maintain. (Matlab comes with a function that generates reports from comments)

How not to comment

- Useless:

```
% Start program
```

- Obvious:

```
if (a > 5)      % Check if a is greater than 5
    ...
end
else           % else add 1 to b
    b = b + 1;
end
```

- Too much about the life:

```
% Well... I do not know how to explain what is going on
% in the snippet below. I tried to code in the night
% with some booze and it worked then, but now I have a
% strong hangover and some parameters still need to be
% worked out...
```

Adding comments to our program

```
IAmFree = checkSchedule();
if (IAmFree == true)
% Count books and amount of free space on a shelf.
% If minimum number of books I need is less than a
% shelf capacity, go shopping and buy additional
% literature. If the amount of books after the
% shopping is too big, use boxes to store them.
books = countBooks();
shelfSize = countFreeSpaceShelf();

...
else
burnBooks();
goToBookStore();
end
```

If anything sticks today, let it be this

Your code will not be understood by anyone

That includes future-you

Use comments and code to document design and purpose (functionality), not mechanics (implementation).

Use consistent and sensible naming of functions and variables.

Today's outline

- ① Coding in style
- ② Error management
- ③ Visualisation
- ④ Functions: the sequel
- ⑤ Excel
- ⑥ Algorithms
- ⑦ Conclusions

Errors in computer programs

Computer programs often contain errors (bugs): buildings collapse, governments fall, kittens will die.



Errors in computer programs

The following symptoms can be distinguished:

- Unable to execute the program
- Program crashes, warnings or error messages
- Never-ending loops
- Wrong (unexpected) result

Three error categories:

Syntax errors You did not obey the language rules. These errors prevent running or compilation of the program.

Runtime errors Something goes wrong during the execution of the program resulting in an error message (problem with input, division by zero, loading of non-existent files, memory problems, etc.)

Semantic errors The program does not do what you expect, but does what have told it to do.

A convenient tool: the debugger

- No-one can write a 1000-line code without making errors
 - If you can, please come work for us
- One of the most important skills you will acquire is debugging.
- Although it can be frustrating, debugging is one of the most intellectually rich, challenging, and interesting parts of programming.
- In some ways, debugging is like detective work. You are confronted with clues, and you have to infer the processes and events that led to the results you see.

"When you have eliminated the impossible, whatever remains, however improbable, must be the truth."
— A. Conan Doyle, The Sign of Four

A convenient tool: the debugger

The debugger can help you to:

- Pause a program at a certain line: set a *breakpoint*
- Check the values of variables during the program
- Controlled execution of the program:
 - One line at a time
 - Run until a certain line
 - Run until a certain condition is met (conditional breakpoint)
 - Run until the current function exits
- Note: You may end up in the source code of Matlab functions!
- Check Canvas (Matlab Crash Course section) for a movie that demonstrates the debugger.

About testcases (validation)

- Testcases: run the program with parameters such that a known result is (should be) produced.
- Testcases: what happens when unforeseen input is encountered?
 - More or fewer arguments than anticipated? (Matlab uses `varargin` and `nargin` to create a varying number of input arguments, and to check the number of given input arguments)
 - Other data types than anticipated? How does the program handle this? Warnings, error messages (crash), NaN or worse (a continuing program)?
- For physical modeling, we typically look for analytical solutions
 - Sometimes somewhat stylized cases
 - Possible solutions include Fourier-series
 - Experimental data

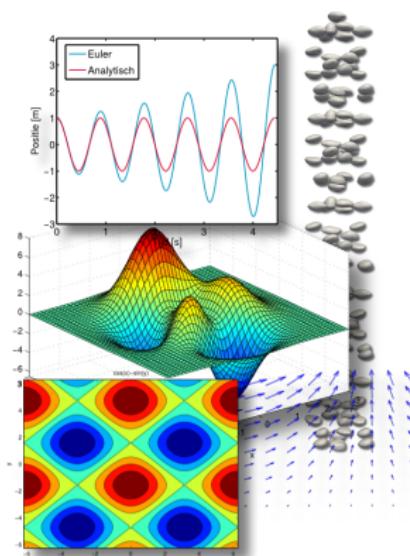
Today's outline

- ① Coding in style
- ② Error management
- ③ Visualisation
- ④ Functions: the sequel
- ⑤ Excel
- ⑥ Algorithms
- ⑦ Conclusions

Data visualisation

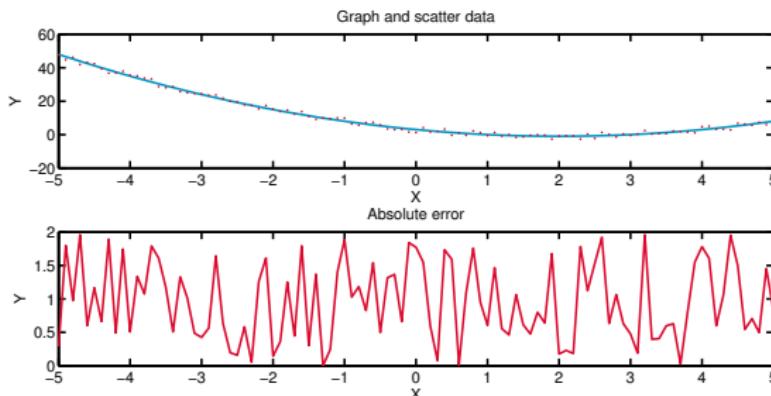
Modeling can lead to very large data sets, that require appropriate visualisation to convey your results.

- 1D, 2D, 3D visualisation
- Multiple variables at the same time (temperature, concentration, direction of flow)
- Use of colors, contour lines
- Use of stream lines or vector plots
- Animations



Plotting

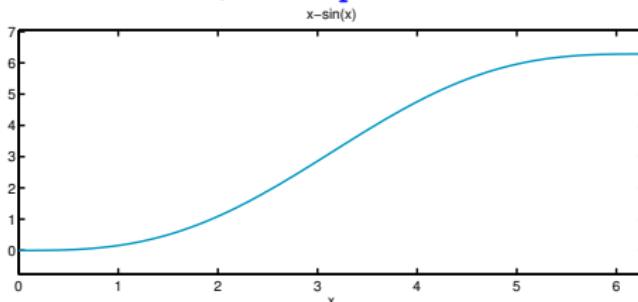
```
x = -5:0.1:5;
y = x.^2-4*x+3;
y2 = y + (2-4*rand(size(y)));
subplot(2,1,1); plot(x,y,'-',x,y2,'r.');
xlabel('X'); ylabel('Y'); title('Graph and Scatter');
subplot(2,1,2); plot(x,abs(y-y2),'r-');
xlabel('X'); ylabel('Y'); title('Absolute error');
```



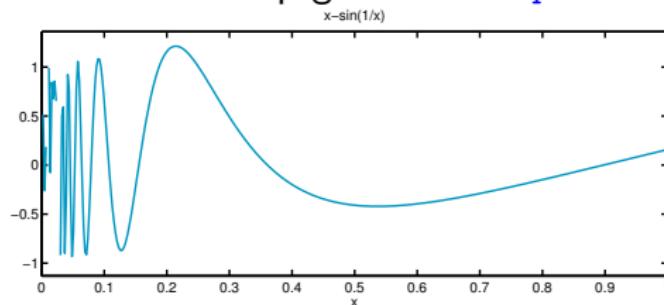
Plotting (2)

Easy plotting of functions can be done using the `ezplot` function:

`ezplot('x-sin(x)', [0 2*pi]):`

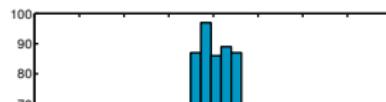
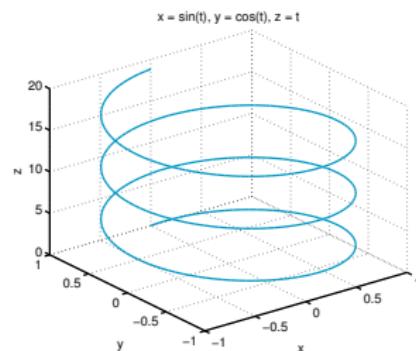
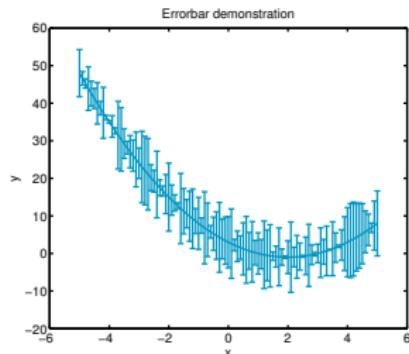


Be careful with steep gradients: `ezplot('x-sin(1/x)', [0 1])`



Other plotting tools

- Errorbars: `errorbar(x,y,err)`
- 3D-plots: `plot3(x,y,z)`
- Histograms: `histogram(x,20)`

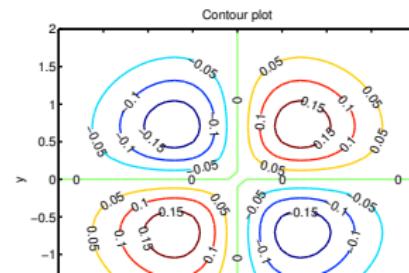
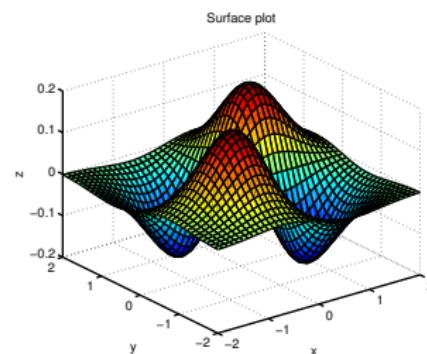


Multi-dimensional data

Matlab typically requires the definition of rectangular grid coordinates using `meshgrid`:

```
[x y] = meshgrid(-2:0.1:2,  
                    -2:0.1:2);  
z = x .* y .* exp(-x.^2 - y.^2);
```

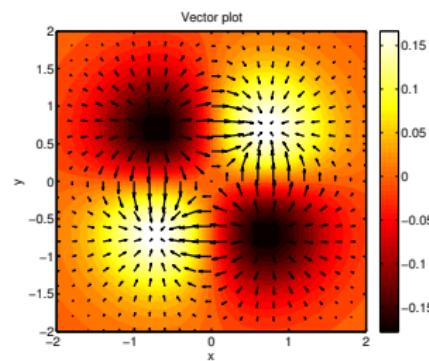
- Surface plot
- Contour plot
- Waterfall
- Ribbons



Vector data

The gradient operator, as expected, is used to obtain the gradient of a scalar field. Colors can be used in the background to simultaneously plot field data:

```
[x y] = meshgrid(-2:0.2:2,  
                   -2:0.2:2);  
z = x .* y .* exp(-x.^2 - y.^2)  
[dx dy] = gradient(z,8,8)  
  
% Background  
contourf(x,y,z,30,'LineColor','  
          none');  
colormap(hot); colorbar;  
  
axis tight; hold on;  
  
% Vectors  
quiver(x,y,dx,dy,'k');
```



Today's outline

- ① Coding in style
- ② Error management
- ③ Visualisation
- ④ Functions: the sequel
- ⑤ Excel
- ⑥ Algorithms
- ⑦ Conclusions

Functions: revisited

In MATLAB you can define your own functions to re-use certain functionalities. We now define the mathematical function $f = x^2 + e^x$:

```
function y = f(x)
y = x.^2 + exp(x);
```

Note:

- The first line of the file has to contain the `function` keyword
- The variables used are *local*. They will not be available in your Workspace
- The file needs to be saved with the same name as the function, i.e. “`f.m`”
- The semi-colon prevents that at each function evaluation output appears on the screen
- If `x` is an array, then `y` becomes an array of function values.

Anonymous functions

If you do not want to create a file, you can create an *anonymous function*:

```
>> f = @(x) (x.^2+exp(x))
```

- f: the name of the function
- @: the function handle
- x: the input argument
- x.^2+exp(x): the actual function

```
>> f(0:0.1:1)
```

Using function handles

A function handle points to a function. It behaves as a variable

```
>> myFunctionHandle = @exp  
>> myFunctionHandle(1)
```

Used a.o. for passing a function to another function, for instance for optimization functions.

$$f(x) = x^3 - x^2 - 3 \arctan x + 1$$

Matlab offers a function `fzero` that can find the roots of a function in a certain range:

```
>> f = @(x) x.^3 - x.^2 - 3*atan(x) + 1;  
>> fzero(f, [-2 2])  
>> ezplot(f)  
>> f(ans)
```

Practice function handles

Consider the function

$$f(x) = -x^2 - 3x + 3 + e^{x^2}$$

The built-in Matlab function `fminbnd` allows to find the minimum of a function in a certain range. Find the minimum of $f(x)$ on $-2 \leq x \leq 2$. Example usage:

```
x = fminbnd(fun,x1,x2)
```

Answer using an anonymous function:

```
>> f = @(x) -x.^2 - 3*x + 3 + exp(x.^2)
>> ezplot(f,[-2 2])
>> fminbnd(f,-2,2)
>> f(ans)
```

Various related functions are `fzero`, `feval`, `fsolve`, `fminsearch`. They will be discussed later in the course.

Today's outline

- ① Coding in style
- ② Error management
- ③ Visualisation
- ④ Functions: the sequel
- ⑤ Excel
- ⑥ Algorithms
- ⑦ Conclusions

Solver and goal-seek

Excel comes with a goal-seek and solver function. For Excel 2010:

- Install via Excel \Rightarrow File \Rightarrow Options \Rightarrow Add-Ins \Rightarrow Go (at the bottom) \Rightarrow Select solver add-in. You can now call the solver screen on the 'data' menu ('Oplosser' in Dutch)
- Select the goal-cell, and whether you want to minimize, maximize or set a certain value
- Enter the variable cells; Excel is going to change the values in these cells to get to the desired solution
- Specify the boundary conditions (e.g. to keep certain cells above zero)
- Click 'solve' (possibly after setting the advanced options).

Goal-seek: a simple example

Goal-Seek can be used to make the goal-cell to a specified value by changing another cell:

- Open Excel and type the following:

	A	B
1	x	3
2	f(x)	=-3*B1^2-5*B1+2
3		

- Go to Data ⇒ What-If Analysis ⇒ Goal Seek...
 - Set cell: B2
 - To value: 0
 - By changing cell: B1
- OK. You find a solution of 0.333 . . .

Solver: a simple example

The solver is used to change the value in a goal-cell, by changing the values in 1 or more other cells while keeping boundary conditions:

- Use the following sheet:

	A	B	C
1		x	$f(x)$
2	x_1	3	$=2*B2*B3-B3+2$
3	x_2	4	$=2*B3-4*B2-4$

- Go to Data ⇒ Solver
 - Goalfunction: C2 (value of: 0)
 - Add boundary condition: C3 = 0
 - By changing cells: \$B\$2:\$B\$3 (you can just select the cells)
- Solve. You will find $B2=0$ and $B3=2$.

Exercise

Use Excel functions to obtain the Antoine coefficients A , B and C for diethyl ether following the equation:

$$\ln P = A - \frac{B}{T + C}$$

P in kPa, T in °C. Experimental data is given (see Canvas for the xls):

P [mmHg]	T [K]
15.6	230.0
29.1	239.3
52.5	248.9
91.9	258.9
156.0	269.3
257.6	280.1
414.6	291.4
651.1	303.1
999.2	315.3
1501.0	328.0

- ① Dedicate three separate cells for A , B and C . Give an initial guess of 15, 2800, 200.
- ② Convert all values to proper units (hint: use e.g. `=CONVERT(A2, "mmHg", "Pa")`)
- ③ Compute P_{Antoine} (in Pa)
- ④ Compute $(P_{\text{exp}} - P_{\text{Antoine}})^2$, and sum this column
- ⑤ Start the solver, and aim for a value of 0 in the sum, by changing cells for A , B and C .

Today's outline

- ① Coding in style
- ② Error management
- ③ Visualisation
- ④ Functions: the sequel
- ⑤ Excel
- ⑥ Algorithms
- ⑦ Conclusions

Algorithm design

① Problem analysis

Contextual understanding of the nature of the problem to be solved

② Problem statement

Develop a detailed statement of the mathematical problem to be solved with the program

③ Processing scheme

Define the inputs and outputs of the program

④ Algorithm

A step-by-step procedure of all actions to be taken by the program (*pseudo-code*)

⑤ Program the algorithm

Convert the algorithm into a computer language, and debug until it runs

⑥ Evaluation

Test all of the options and conduct a validation study

Example: finding the roots of a parabola

We are writing a program that finds for us the roots of a parabola.
We use the form

$$y = ax^2 + bx + c$$

What is our program in pseudo-code?

① Input data (a , b and c)

② Identify special cases ($a = b = c = 0$, $a = 0$)

$a = b = c = 0$ Solution indeterminate

$a = 0$ Solution: $x = -\frac{c}{b}$

③ Find $D = b^2 - 4ac$

④ Decide, based on D :

$D < 0$ Display message: complex roots

$D = 0$ Display 1 root value

$D > 0$ Display 2 root values

Example: finding the roots of a parabola

```
function x = parabola(a,b,c)
% Catch exception cases
if (a==0)
    if(b==0)
        if(c==0)
            disp('Solution indeterminate'); return;
        end
        disp('There is no solution');
    end
    x = -c/b;
end

D = b^2 - 4*a*c;
if (D<0)
    disp('Complex roots'); return;
else if (D==0)
    x = -b/(2*a);
else if (D>0)
    x(1) = (-b + sqrt(D))/(2*a);
    x(2) = (-b - sqrt(D))/(2*a);
    x = sort(x);
end
end
end
```

Example: finding the roots of a parabola

```
>> roots([1 -4 -3])
ans =
    4.6458
   -0.6458
```

Advanced concepts

- Object oriented programming: classes and objects
- Memory management: some programming languages require you to allocate computer memory yourself (e.g. for arrays)
- External libraries: in many cases, someone already built the general functionality you are looking for
- Compiling and scripting (“interpreted”); compiling means converting a program to computer-language before execution. Interpreted languages do this on the fly.
- Profiling, optimization, parallelization: Checking where your program spends the most of its time, optimizing (or parallelizing) that part.

In conclusion...

- Algorithm design: define your problem, think ahead, make a scheme, sketch the interplay between variables and functions, then start programming
- Programming basics: variables, operators and functions, locality of variables, recursive operations
- Dealing with complex programs, verification of your algorithms, use of the debugger
- Visualisation: how to make 1D and 2D/3D plots, create a sensible and intuitive presentation of your data.
- Examples: a few practice cases

Errors in computer simulations

Ivo Roghair, Martin van Sint Annaland

Chemical Process Intensification,
Eindhoven University of Technology

Today's outline

① Introduction

② Roundoff and truncation errors

③ Break errors

④ Loss of digits

⑤ (Un)stable methods

⑥ Symbolic math

⑦ Summary

Example 1

Start your spreadsheet program (Excel, ...)

Enter:

Cell	Value
A1	0.1
A2	= (A1*10)-0.9
A3	= (A2*10)-0.9
A4	= (A3*10)-0.9

(repeat until A30)

What's happening?

Enterprise

Cell	Value
A1	2
A2	$=A1*10)-18$
A3	$=A2*10)-18$
A4	$=A3*10)-18$

(repeat until A30)

Errors in computer simulations

In this course we will outline different numerical errors that may appear in computer simulations, and how these errors can affect the simulation results.

- Errors in the mathematical model (physics)
- Errors in the program (implementation)
- Errors in the entered parameters
- Roundoff- and truncation errors
- Break errors

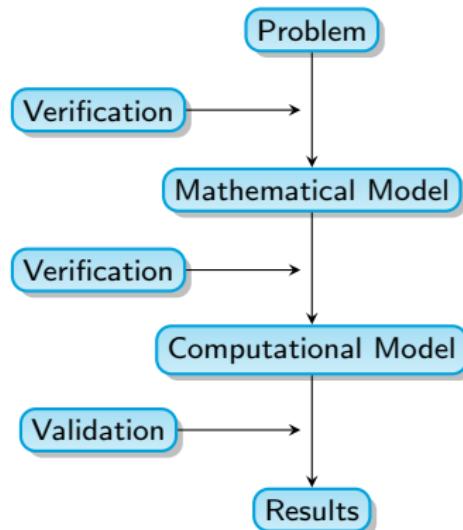
Verification and validation

Verification

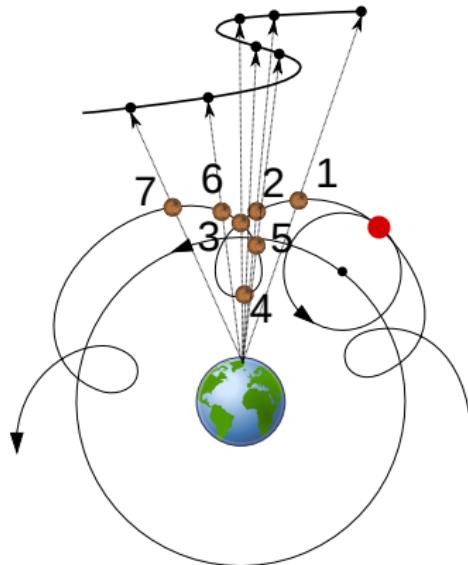
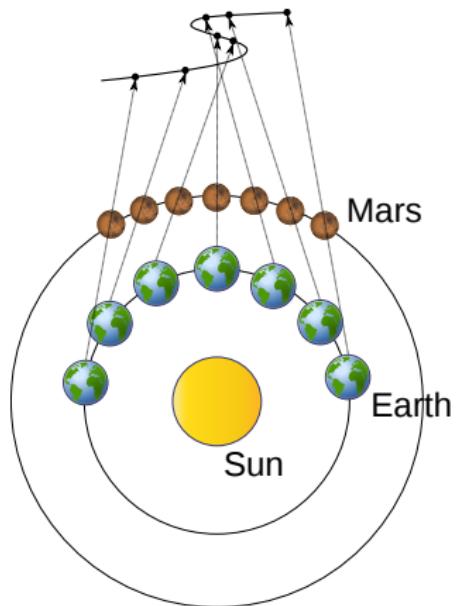
Verification is the process of mathematically and computationally assuring that the model computes what you have entered.

Validation

Validation is the process of determining the degree to which a model is an accurate representation of the real world from the perspective of the intended uses of the model



Verification of the physical model



- The perceived orbit of Mars from Earth shows a zig-zag (in contrast to the Sun, Mercury, Venus)
- Even though they were not 'right', Earth-centered models (Ptolemy) were still valid

Be aware of your uncertainties

Aleatory uncertainty

Uncertainty that arises due to inherent randomness of the system, features that are too complex to measure and take into account

Epistemic uncertainty

Uncertainty that arises due to lack of knowledge of the system, but could in principle be known

Errors in computer simulations

In this course we will outline different numerical errors that may appear in computer simulations, and how these errors can affect the simulation results.

- Errors in the mathematical model (physics)
- Errors in the program (implementation)
- Errors in the entered parameters
- Roundoff- and truncation errors
- Break errors

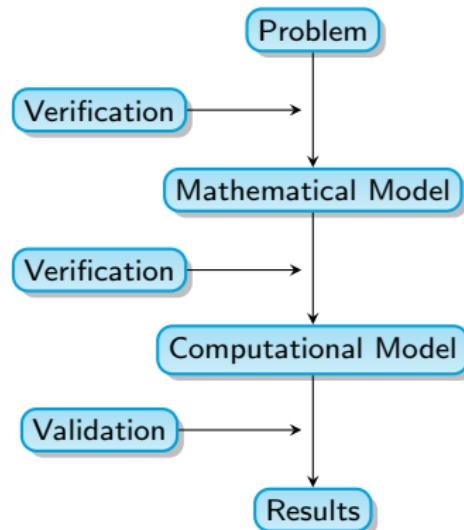
Verification and validation

Verification

Verification is the process of mathematically and computationally assuring that the model computes what you have entered.

Validation

Validation is the process of determining the degree to which a model is an accurate representation of the real world from the perspective of the intended uses of the model



Errors in computer simulations

In this course we will outline different numerical errors that may appear in computer simulations, and how these errors can affect the simulation results.

- Errors in the mathematical model (physics)
- Errors in the program (implementation)
- Errors in the entered parameters
- Roundoff- and truncation errors
- Break errors

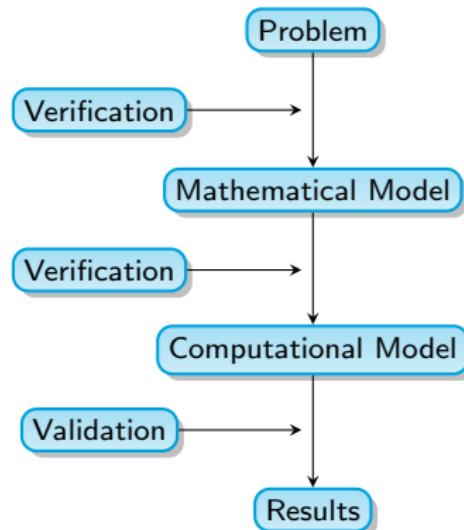
Verification and validation

Verification

Verification is the process of mathematically and computationally assuring that the model computes what you have entered.

Validation

Validation is the process of determining the degree to which a model is an accurate representation of the real world from the perspective of the intended uses of the model



Errors in computer simulations

In this course we will outline different numerical errors that may appear in computer simulations, and how these errors can affect the simulation results.

- Errors in the mathematical model (physics)
- Errors in the program (implementation)
- Errors in the entered parameters
- Roundoff- and truncation errors
- Break errors

Significant digits

A numerical result \tilde{x} is an approximation of the real value x .

- Absolute error

$$\delta = |\tilde{x} - x|, x \neq 0$$

- Relative error

$$\frac{\delta}{\tilde{x}} = \left| \frac{\tilde{x} - x}{\tilde{x}} \right|$$

- Error margin

$$\tilde{x} - \delta \leq x \leq \tilde{x} + \delta$$

$$x = \tilde{x} \pm \delta$$

Significant digits

- \tilde{x} has m significant digits if the absolute error in x is smaller or equal to 5 at the $(m + 1)$ -th position:

$$10^{q-1} \leq |\tilde{x}| \leq 10^q$$

$$|x - \tilde{x}| = 0.5 \times 10^{q-m}$$

- For example:

$$x = \frac{1}{3}, \tilde{x} = 0.333 \Rightarrow \delta = 0.00033333\dots$$

3 significant digits

Today's outline

① Introduction

② Roundoff and truncation errors

③ Break errors

④ Loss of digits

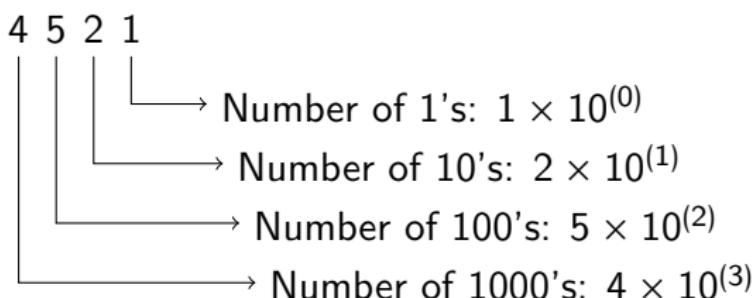
⑤ (Un)stable methods

⑥ Symbolic math

⑦ Summary

Representation of numbers

- Computers represent a number with a finite number of digits: each number is therefore an approximation due to roundoff and truncation errors.
- In the decimal system, a digit c at position n has a value of $c \times 10^{n-1}$



$$(4521)_{10} = 4 \times 10^3 + 5 \times 10^2 + 2 \times 10^1 + 1 \times 10^0$$

Representation of numbers

- You could use another basis, computers often use the basis 2:

$$\begin{aligned}(4521)_{10} &= 1 \times 2^{12} + 0 \times 2^{11} + 0 \times 2^{10} + 0 \times 2^9 + 1 \times 2^8 + \dots \\ &\quad \dots 1 \times 2^7 + 0 \times 2^6 + 1 \times 2^5 + 0 \times 2^4 + \dots \\ &\quad \dots 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 \\ &= (1000110101001)_2\end{aligned}$$

- In general:

$$(c_m \dots c_1 c_0)_q = c_0 q^0 + c_1 q^1 + \dots + c_m q^m, c \in \{0, 1, 2, \dots, q-1\}$$

Representation of numbers

- Numbers are stored in binary in the memory of a computer, in segments of a specific length (called a *word*).
- We distinguish multiple types of numbers:
 - Integers: $-301, -1, 0, 1, 96, 2293, \dots$
 - Floating points: $-301.01, 0.01, 3.14159265, 14498.2$
- A binary integer representation looks like the following bit sequence:

$$z = \sigma \left(c_0 2^0 + c_1 2^1 + \dots + c_{\lambda-1} 2^{\lambda-1} \right)$$

σ is the sign of z (+ or -), and λ is the length of the word

- Endianness: the order of bits stored by a computer

Excercise

- Convert the following decimal number to base-2: 214

$$214_{10} = 11010110_2$$

- Excel:
 - Decimal: =DEC2BIN(214)
 - Octal: =DEC2OCT(214)
 - Hexadecimal: =DEC2HEX(214)
- Matlab:
 - Decimal: dec2bin(214)
 - Other base: dec2base(214,<base>)

Arithmetic operations with binary numbers

Addition:

$$0 + 0 = 0$$

$$0 + 1 = 1$$

$$1 + 0 = 1$$

$$\textcolor{blue}{1 + 1 = 0}$$

(carry one)

$$\begin{array}{r}
 & 1 & 4 & 5 \\
 + & 2 & 3 \\
 \hline
 1 & 6 & 8
 \end{array}$$

$$\begin{array}{r}
 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\
 + & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 \\
 \hline
 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0
 \end{array}$$

Subtraction:

$$0 - 0 = 0$$

$$1 - 0 = 1$$

$$1 - 1 = 0$$

$$\textcolor{blue}{0 - 1 = 1}$$

(borrow one)

$$\begin{array}{r}
 & 1 & 4 & 5 \\
 - & 2 & 3 \\
 \hline
 1 & 2 & 2
 \end{array}$$

$$\begin{array}{r}
 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\
 - & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 \\
 \hline
 0 & 1 & 1 & 1 & 1 & 0 & 1 & 0
 \end{array}$$

- Multiplication and division are more expensive, and more elaborate

Excercise

Try the following commands in Matlab:

Command	Result
intmin	-2147483648
intmax	2147483647
i = int16(intmax)	i = 32767
whos i	int16 information
i = i + 100	i = 32767
realmax	1.7977e+308
f = 0.1	
whos f	double information
format long e	
realmax	1.797693134862316e+308
f	
fprintf("%0.16f",f)	0.1000000000000000
fprintf("%0.20f",f)	0.1000000000000000555

Representation of integer numbers

- In Matlab, integers of the type `int32` are represented by 32-bit words ($\lambda = 31$).
- The set of numbers that an `int32` z can represent is:

$$-2^{31} \leq z \leq 2^{31} - 1 \approx 2 \times 10^9$$

- If, during a calculation, an integer number becomes larger than $2^\lambda - 1$, the computer reports an [overflow](#)¹
- How can a computer identify an overflow?

¹Matlab does not perform actual integer overflows, it just stops at the maximum.

Representation of real (floating point) numbers

- Formally, a real number is represented by the following bit sequence

$$x = \sigma \left(2^{-1} + c_2 2^{-2} + \dots + c_m 2^{-m} \right) 2^{e-1023}$$

Here, σ is the sign of x and e is an integer value.

- A floating point number hence contains sections that contain the sign, the exponent and the mantissa

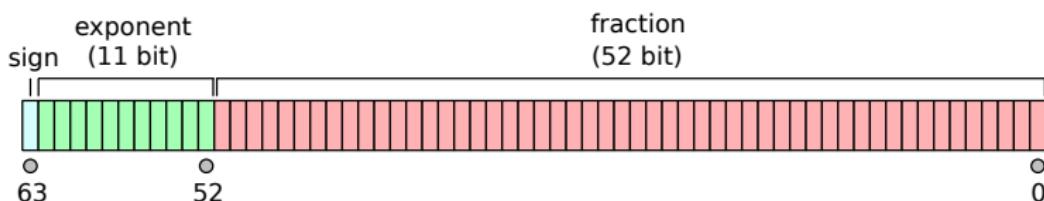


Image: Wikimedia Commons CC by-SA

Representation of real (floating point) numbers

- Example: $\lambda = 3$, $m = 2$, $x = \frac{2}{3}$

$$x = \pm \left(2^{-1} + c_2 2^{-2} \right) 2^e$$

- $c_0 \in \{0, 1\}$
- $e = \pm a_0 2^0$
- $a_0 \in \{0, 1\}$
- Truncation: $f_l(x) = 2^{-1} = 0.5$
- Round off: $f_u(x) = 2^{-1} + 2^{-2} = 0.75$

Today's outline

① Introduction

② Roundoff and truncation errors

③ Break errors

④ Loss of digits

⑤ (Un)stable methods

⑥ Symbolic math

⑦ Summary

Trigonometric, Logarithmic, and Exponential computations

- Processors can do logic and arithmetic instructions
- Trigonometric, logarithmic and exponential calculations are “higher-level” functions:
 \exp , \sin , \cos , \tan , \sec , \arcsin , \arccos , \arctan , \log , \ln , ...
- Such functions can be performed using these “low level” instructions, for instance using a Taylor series:

$$\sin(x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{2n+1} = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

Trigonometric, Logarithmic, and Exponential computations

- These operations involve many multiplications and additions, and are therefore *expensive*
- Computations can only take finite time, for infinite series, calculations are interrupted at N

$$\sin(x) = \sum_{n=0}^N \frac{(-1)^n}{(2n+1)!} x^{2n+1} = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots + \frac{(-1)^N}{(2N+1)!} x^{2N+1}$$

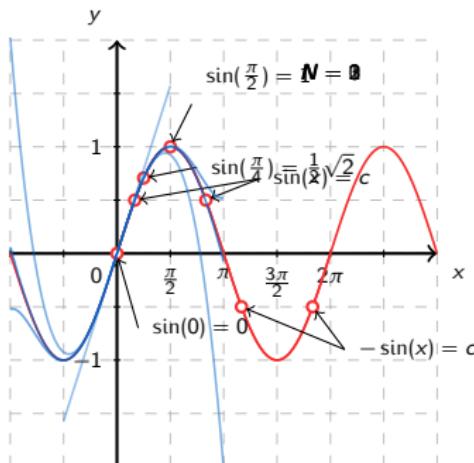
$$e^x = \sum_{n=0}^N \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots + \frac{x^N}{N!}$$

- This results in a *break error*

Algorithm for sine-computation

A computer may use a clever algorithm to limit the number of operations required to perform a higher-level function. A (fictional!) example for the computation of $\sin(x)$:

- ① Use periodicity so that
 $0 \leq x \leq 2\pi$
- ② Use symmetry ($0 \leq x \leq \frac{\pi}{2}$)
- ③ Use lookup tables for known values
- ④ Perform taylor expansion



Today's outline

① Introduction

② Roundoff and truncation errors

③ Break errors

④ Loss of digits

⑤ (Un)stable methods

⑥ Symbolic math

⑦ Summary

Loss of digits

- During operations such as $+$, $-$, \times , \div , an error can add up
- Consider the summation of x and y

$$\tilde{x} - \delta \leq x \leq \tilde{x} + \delta \quad \text{and} \quad \tilde{y} - \varepsilon \leq y \leq \tilde{y} + \varepsilon$$

$$(\tilde{x} + \tilde{y}) - (\delta + \varepsilon) \leq x + y \leq (\tilde{x} + \tilde{y}) + (\delta + \varepsilon)$$

Loss of digits: Example 1

$$\left. \begin{array}{l} x = \pi, \tilde{x} = 3.1416 \\ y = 22/7, \tilde{y} = 3.1429 \end{array} \right\} \Rightarrow \left. \begin{array}{l} \delta = \tilde{x} - x = 7.35 \times 10^{-6} \\ \varepsilon = \tilde{y} - y = 4.29 \times 10^{-5} \end{array} \right\}$$

$$x + y = \tilde{x} + \tilde{y} \pm (\delta + \varepsilon) \approx 6.2845 - 5.025 \times 10^{-5}$$

$$x - y = \tilde{x} - \tilde{y} \pm (\delta + \varepsilon) \approx -0.0013 + 3.55 \times 10^{-5}$$

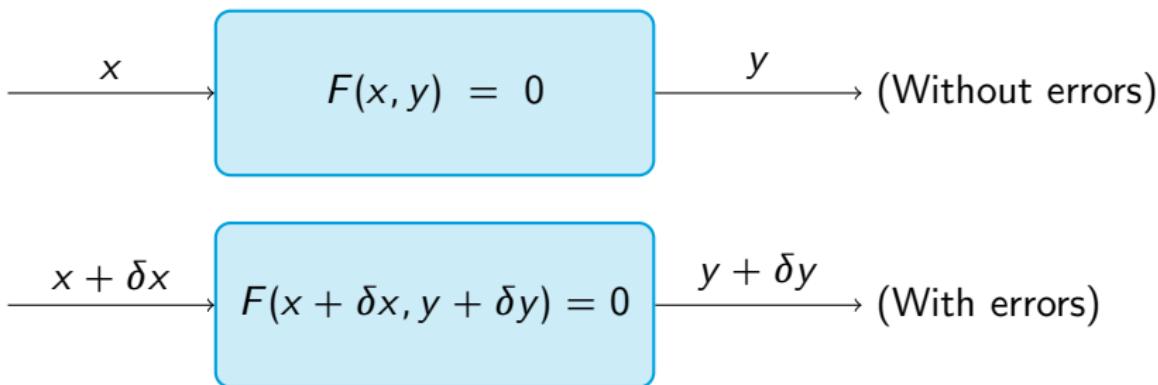
- The absolute error is small ($\approx 10^{-5}$), but the relative error is much bigger (0.028).
- Adding up the errors results in a loss of significant digits!

Loss of digits: Example 2

- Calculate e^{-5}
 - Use the Taylor series
 - Calculate the first 26 terms ($N = 26$)
- Now repeat the calculation, but use for each calculation only 4 digits. What do you find?
Use: `str2double(sprintf('%.4g', term))`
- Without errors you would find: $e^{-5} = 0.006738$
- If you only use 4 digits in the calculations, you'll find 0.00998

Badly (ill) conditioned problems

We consider a system $F(x, y)$ that computes a solution from input data. The input data may have errors:



$$y(x + \delta x) - y(x) \approx y'(x)\delta x$$

$$C = \max_{\delta x} \left(\left| \frac{\delta y / y}{\delta x / x} \right| \right)$$

Propagated error on the basis of
Taylor expansion

Condition criterion, $C < 10$ error
development small

Badly (ill) conditioned problems: Example

Solve the following linear system in Matlab using double and single precision:

$$A = \begin{bmatrix} 1.2969 & 0.8648 \\ 0.2161 & 0.1441 \end{bmatrix}, \quad x = \begin{bmatrix} 0.8642 \\ 0.1440 \end{bmatrix}, \quad y = \begin{bmatrix} 2.0 \\ 2.0 \end{bmatrix}$$

Double precision

```
>> clear;clc;format long e;
>> A = [[1.2969 0.8648];
   [0.2161 0.1441]];
>> x = [0.8642; 0.1440];
>> y = A\x
y =
    2.00000002400302e+00
   -2.00000003599621e+00
```

Single precision

```
>> clear;clc;format long e;
>> A = single(
    [[1.2969 0.8648];
    [0.2161 0.1441]] );
>> x = single(
    [0.8642; 0.1440] );
>> y = A\x
y =
    1.3331791e+00
   -1.0000000e+00
```

Badly (ill) conditioned problems: Example

- Matlab already warned us about the bad condition number:

Warning: Matrix is `close` to singular or badly scaled.

Results may be inaccurate. `RCOND = 1.148983e-08`.

- The `RCOND` is the reciprocal condition number
- A small error in x results in a big error in y . This is called an ill conditioned problem.

Today's outline

① Introduction

② Roundoff and truncation errors

③ Break errors

④ Loss of digits

⑤ (Un)stable methods

⑥ Symbolic math

⑦ Summary

(Un)stable methods

- The condition criterion does not tell you anything about the quality of a numerical solution method!
- It is very well possible that a certain solution method is more sensitive for one problem than another
- If the method propagates the error, we call it an *unstable method*. Let's look at an example.

The Golden mean

- Let's evaluate the following recurrent relationship:

$$y_{n+1} = y_{n-1} - y_n$$

$$y_0 = 1, \quad y_1 = \frac{2}{1 + \sqrt{5}}$$

- You can prove (by substitution) that:

$$y_n = x^{-n}, \quad n = 0, 1, 2, \dots, \quad x = \frac{1 + \sqrt{5}}{2}$$

The Golden mean

Recurrent version

```
% initialise  
y(1) = 1;  
y(2) = 2 / (1 + sqrt(5));  
  
% Perform recurrent  
    approach  
for n = 2:39  
    y(n+1) = y(n-1)-y(n);  
end
```

Powerlaw version

```
% initialise  
x = (1 + sqrt(5))/2;  
y2(1) = x^0; % n = 1  
  
% Perform powerlaw approach  
for n = 0:39  
    y2(n+1) = x^-n  
end
```

The Golden mean

n	Recurrent	Powerlaw
1	1.0000	1.0000
1	0.6180	0.6180
2	0.3820	0.3820
3	0.2361	0.2361
...
37	$3.080 \cdot 10^{-8}$	$2.995 \cdot 10^{-8}$
38	$1.714 \cdot 10^{-8}$	$1.851 \cdot 10^{-8}$
39	$1.366 \cdot 10^{-8}$	$1.144 \cdot 10^{-8}$
40	$3.485 \cdot 10^{-8}$	$7.071 \cdot 10^{-8}$

- The recurrent approach enlarges errors from earlier calculations!

Example 1: Explanation

Recall example 1, where the errors blew up our computation of 0.1, whereas they did not for 2. Why did we see these results?

- The number 0.1 is not exactly represented in binary
 - A tiny error can accumulate up to catastrophic proportions!
- The number 2 does have an exact binary representation

Example 2

Start your calculation program of choice (Excel, Matlab, ...)

Calculate the result of y :

$$y = e^\pi - \pi = 19.999099979 \neq 20$$

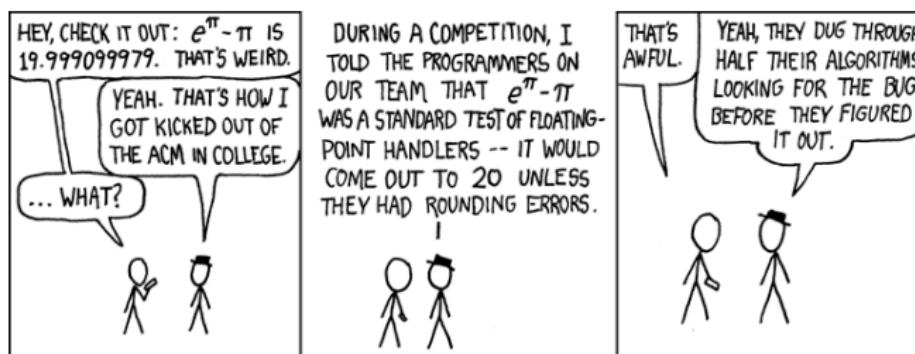


Image: [xkcd](#)

Today's outline

① Introduction

② Roundoff and truncation errors

③ Break errors

④ Loss of digits

⑤ (Un)stable methods

⑥ Symbolic math

⑦ Summary

Symbolic math packages

Definition

The use of computers to manipulate mathematical equations and expressions in symbolic form, as opposed to manipulating the numerical quantities represented by those symbols.

- Symbolic integration or differentiation, substitution of one expression into another
- Simplification of an expression, change of subject etc.
- Packages and toolboxes:

Symbolic math packages

Mathematica Well known software package, license available via [TU/e](#)

Maple Well known, license available via [TU/e](#)

Wolfram|Alpha Web-based interface by Mathematica developer.
Less powerful in mathematical respect, but more accessible and has a broad application range (unit conversion, semantic commands).

Sage Open-source alternative to Maple, Mathematica, Magma, and MATLAB.

Matlab Symbolic math toolbox

Symbolic math: simplify

$$f(x) = (x - 1)(x + 1)(x^2 + 1) + 1$$

```
>> syms x
>> f = (x - 1)*(x + 1)*(x^2 + 1) + 1
f =
(x^2 + 1)*(x - 1)*(x + 1) + 1
>> f2 = simplify(f)
f2 =
x^4
```

Symbolic math: integration and differentiation

$$f(x) = \frac{1}{x^3 + 1}$$

```
>> syms x
>> f = 1/(x^3+1);
>> my_f_int = int(f)

my_f_int = log(x + 1)/3 - log((x - 1/2)^2 + 3/4)/6 +
(3^(1/2)*atan((2*3^(1/2)*(x - 1/2))/3))/3

>> my_f_diff = diff(my_f_int)

my_f_diff = 1/(3*(x + 1)) + 2/(3*((4*(x - 1/2)^2)/3 +
1)) - (2*x - 1)/(6*((x - 1/2)^2 + 3/4))

>> simplify(my_f_diff)

ans = 1/(x^3 + 1)
```

Symbolic math: exercises

Exercise 1

Simplify the following expression:

$$f(x) = \frac{2 \tan x}{(1 + \tan^2 x)} = \sin 2x$$

```
>> simplify(2*tan(x)/(1 + tan(x)^2))
```

Exercise 2

Calculate the *value* of p :

$$p = \int_0^{10} \frac{e^x - e^{-x}}{\sinh x} dx$$

```
>> f = ((exp(x)- exp(-x))/sinh(x));  
>> p = int(f, 0, 10)  
p = 20
```

Symbolic math: root finding

A root finding method searches for the values where a function reaches zero. We will cover the numerical methods later, here we show how to use root finding with symbolic math in Matlab.

Symbolic math function

$$f(x) = \frac{3}{x^2 + 3x} - 2$$

```
>> syms x
>> f = 3 / (x^2 + 3*x) - 2;
>> solve(f)
ans =
15^(1/2)/2 - 3/2
- 15^(1/2)/2 - 3/2
```

Function as a string

$$f(x) = x^2 - 4x + 2$$

```
>> solve('x^2 - 4*x + 2')
ans =
2^(1/2) + 2
2 - 2^(1/2)
```

Symbolic math toolbox: variable precision arithmetic

Variable precision can be used to specify the number of significant digits.

```
>> p = vpa(1/3,16)
p = 0.3333333333333333
>> p = vpa(1/3,4)
p = 0.3333
>> a = vpa(0.1, 30)
a = 0.1
>> b = vpa(0.1, 5);
b = 0.1
>> a-b

ans = 0.000000000000056843418860808014869689938467514
```

Today's outline

① Introduction

② Roundoff and truncation errors

③ Break errors

④ Loss of digits

⑤ (Un)stable methods

⑥ Symbolic math

⑦ Summary

Summary

- Numerical errors may arise due to truncation, roundoff and break errors, which may seriously affect the accuracy of your solution
- Errors may propagate and accumulate, leading to smaller accuracy
- Ill-conditioned problems and unstable methods have to be identified so that proper measures can be taken
- Symbolic math computations may be performed to solve certain equations algebraically, bypassing numerical errors, but this is not always possible.

Linear equations

Linear algebra basics

Ivo Roghair, Martin van Sint Annaland

Chemical Process Intensification,
Eindhoven University of Technology

Today's outline

① Introduction

② Matrix inversion

③ Solving a linear system

④ Towards larger systems

⑤ Summary

Today's outline

① Introduction

② Matrix inversion

③ Solving a linear system

④ Towards larger systems

⑤ Summary

Overview

Goals

- Different ways of looking at a system of linear equations
- Determination of the inverse, determinant and the rank of a matrix
- The existence of a solution to a set of linear equations

Different views of linear systems

- Separate equations:

$$x + y + z = 4$$

$$2x + y + 3z = 7$$

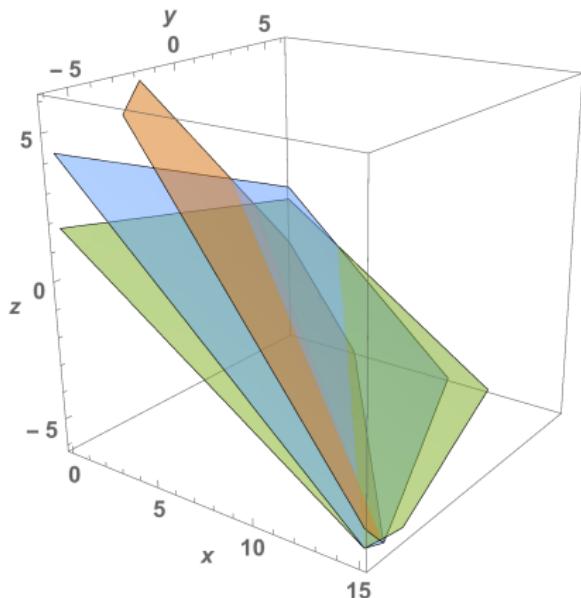
$$3x + y + 6z = 5$$

- Matrix mapping $Mx = b$:

$$\begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix}$$

- Linear combination:

$$x \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + y \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + z \begin{bmatrix} 3 \\ 6 \\ 6 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix}$$



Today's outline

① Introduction

② Matrix inversion

③ Solving a linear system

④ Towards larger systems

⑤ Summary

Inverse of a matrix

- The inverse M^{-1} is defined such that:

$$MM^{-1} = I \quad \text{and} \quad M^{-1}M = I$$

- Use the inverse to solve a set of linear equations:

$$Mx = b$$

$$M^{-1}Mx = M^{-1}b$$

$$Ix = M^{-1}b$$

$$x = M^{-1}b$$

How to calculate the inverse?

- The inverse of an $N \times N$ matrix can be calculated using the co-factors of each element of the matrix:

$$M^{-1} = \frac{1}{\det |M|} \begin{bmatrix} C_{11} & C_{12} & C_{13} \\ C_{21} & C_{22} & C_{23} \\ C_{31} & C_{32} & C_{33} \end{bmatrix}^T$$

- $\det |M|$ is the *determinant* of matrix M .
- C_{ij} is the *co-factor* of the ij^{th} element in M .

Computing the co-factors

Consider the following example matrix: $M = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix}$

Computing the co-factors

Consider the following example matrix: $M = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix}$

A co-factor (e.g. C_{11}) is the determinant of the elements left over when you cover up the row and column of the element in question, multiplied by ± 1 , depending on the position.

$$\begin{bmatrix} 1 & \times & \times \\ \times & 1 & 3 \\ \times & 1 & 6 \end{bmatrix}$$

Computing the co-factors

Consider the following example matrix:

$$M = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix}$$

A co-factor (e.g. C_{11}) is the determinant of the elements left over when you cover up the row and column of the element in question, multiplied by ± 1 , depending on the position.

$$\begin{bmatrix} 1 & \times & \times \\ \times & 1 & 3 \\ \times & 1 & 6 \end{bmatrix} \quad \begin{bmatrix} + & - & + \\ - & + & - \\ + & - & + \end{bmatrix}$$

Computing the co-factors

Consider the following example matrix:

$$M = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix}$$

A co-factor (e.g. C_{11}) is the determinant of the elements left over when you cover up the row and column of the element in question, multiplied by ± 1 , depending on the position.

$$\begin{bmatrix} 1 & \times & \times \\ \times & 1 & 3 \\ \times & 1 & 6 \end{bmatrix}$$

$$\begin{bmatrix} + & - & + \\ - & + & - \\ + & - & + \end{bmatrix}$$

$$\begin{aligned} C_{11} &= +1 \cdot \det \begin{vmatrix} 1 & 3 \\ 1 & 6 \end{vmatrix} \\ &= 6 \times 1 - 3 \times 1 = 3 \end{aligned}$$

Computing the co-factors

Back to our example:

$$M^{-1} = \frac{1}{\det |M|} \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix}^{-1} = \frac{1}{\det |M|} \begin{bmatrix} 3 & -3 & -1 \\ -5 & 3 & 2 \\ 2 & -1 & -1 \end{bmatrix}^T$$

Computing the co-factors

Back to our example:

$$M^{-1} = \frac{1}{\det |M|} \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix}^{-1} = \frac{1}{\det |M|} \begin{bmatrix} 3 & -3 & -1 \\ -5 & 3 & 2 \\ 2 & -1 & -1 \end{bmatrix}^T$$

- The determinant is very important
- If $\det |M| = 0$, the inverse does not exist (singular matrix)

Calculating the determinant

Compute the determinant by multiplication of each element on a row (or column) by its cofactor and adding the results:

$$\det \begin{vmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{vmatrix} = +\det \begin{vmatrix} 1 & 3 \\ 1 & 6 \end{vmatrix} - \det \begin{vmatrix} 2 & 3 \\ 3 & 6 \end{vmatrix} + \det \begin{vmatrix} 2 & 1 \\ 3 & 1 \end{vmatrix} = -1$$

Calculating the determinant

Compute the determinant by multiplication of each element on a row (or column) by its cofactor and adding the results:

$$\det \begin{vmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{vmatrix} = +\det \begin{vmatrix} 1 & 3 \\ 1 & 6 \end{vmatrix} - \det \begin{vmatrix} 2 & 3 \\ 3 & 6 \end{vmatrix} + \det \begin{vmatrix} 2 & 1 \\ 3 & 1 \end{vmatrix} = -1$$

$$\det \begin{vmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{vmatrix} = +\det \begin{vmatrix} 2 & 1 \\ 3 & 1 \end{vmatrix} - 3\det \begin{vmatrix} 1 & 1 \\ 3 & 1 \end{vmatrix} + 6\det \begin{vmatrix} 1 & 1 \\ 2 & 1 \end{vmatrix} = -1$$

Today's outline

① Introduction

② Matrix inversion

③ Solving a linear system

④ Towards larger systems

⑤ Summary

Solving a linear system

- Our example:

$$\begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix}$$

Solving a linear system

- Our example:

$$\begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix}$$

- The solution is:

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = M^{-1}b = \frac{1}{-1} \begin{bmatrix} 3 & -5 & 2 \\ -3 & 3 & -1 \\ -1 & 2 & -1 \end{bmatrix} \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix} = \frac{1}{-1} \begin{bmatrix} -13 \\ 4 \\ 5 \end{bmatrix} = \begin{bmatrix} 13 \\ -4 \\ -5 \end{bmatrix}$$

Solving a linear system

- Our example:

$$\begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix}$$

- The solution is:

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = M^{-1}b = \frac{1}{-1} \begin{bmatrix} 3 & -5 & 2 \\ -3 & 3 & -1 \\ -1 & 2 & -1 \end{bmatrix} \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix} = \frac{1}{-1} \begin{bmatrix} -13 \\ 4 \\ 5 \end{bmatrix} = \begin{bmatrix} 13 \\ -4 \\ -5 \end{bmatrix}$$

- The inverse exists, because $\det |M| = -1$.

Solving a linear system in Matlab using the inverse

- Create the matrix:

```
>> A = [1 1 1; 2 1 3; 3 1 6];
```

Solving a linear system in Matlab using the inverse

- Create the matrix:

```
>> A = [1 1 1; 2 1 3; 3 1 6];
```

- Create solution vector:

```
>> b = [4; 7; 5];
```

Solving a linear system in Matlab using the inverse

- Create the matrix:

```
>> A = [1 1 1; 2 1 3; 3 1 6];
```

- Create solution vector:

```
>> b = [4; 7; 5];
```

- Get the matrix inverse:

```
>> Ainv = inv(A);
```

Solving a linear system in Matlab using the inverse

- Create the matrix:

```
>> A = [1 1 1; 2 1 3; 3 1 6];
```

- Create solution vector:

```
>> b = [4; 7; 5];
```

- Get the matrix inverse:

```
>> Ainv = inv(A);
```

- Compute the solution:

```
>> x = Ainv * b
```

Solving a linear system in Matlab using the inverse

- Create the matrix:

```
>> A = [1 1 1; 2 1 3; 3 1 6];
```

- Create solution vector:

```
>> b = [4; 7; 5];
```

- Get the matrix inverse:

```
>> Ainv = inv(A);
```

- Compute the solution:

```
>> x = Ainv * b
```

- Matlab's internal direct solver:

```
>> x = A\b
```

- These are black boxes! We are going over some methods later!

Exercise: performance of inverse computation

Create a script that generates matrices with random elements of various sizes $N \times N$. Compute the inverse of each matrix, and use `tic` and `toc` to see the computing time for each inversion. Plot the time as a function of the matrix size N .

Exercise: performance of inverse computation

Create a script that generates matrices with random elements of various sizes $N \times N$. Compute the inverse of each matrix, and use `tic` and `toc` to see the computing time for each inversion. Plot the time as a function of the matrix size N .

```
% Generate random matrices of various sizes 's'.
% Invert the matrices and store the time required
% for the inversion. Plot the times vs 's'
s = [10:10:90 100:100:1000 2000:1000:5000 10000]
for n = 1:length(s)
```

Exercise: performance of inverse computation

Create a script that generates matrices with random elements of various sizes $N \times N$. Compute the inverse of each matrix, and use `tic` and `toc` to see the computing time for each inversion. Plot the time as a function of the matrix size N .

```
% Generate random matrices of various sizes 's'.
% Invert the matrices and store the time required
% for the inversion. Plot the times vs 's'
s = [10:10:90 100:100:1000 2000:1000:5000 10000]
for n = 1:length(s)
    s(n)
    A = rand(s(n));
```

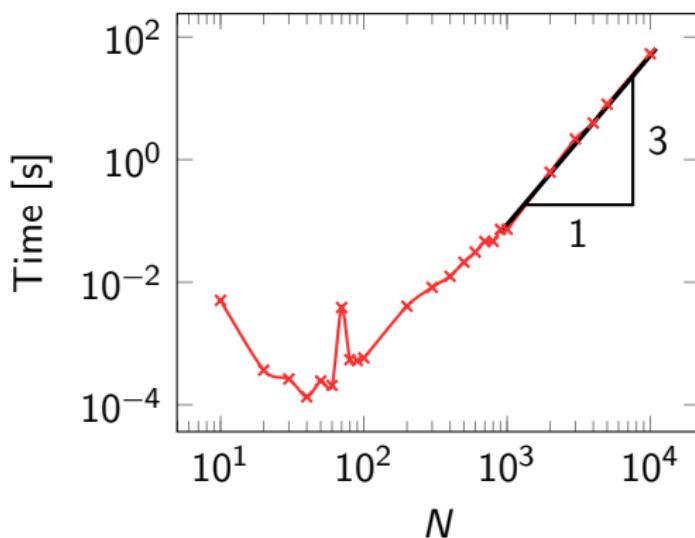
Exercise: performance of inverse computation

Create a script that generates matrices with random elements of various sizes $N \times N$. Compute the inverse of each matrix, and use `tic` and `toc` to see the computing time for each inversion. Plot the time as a function of the matrix size N .

```
% Generate random matrices of various sizes 's'.
% Invert the matrices and store the time required
% for the inversion. Plot the times vs 's'
s = [10:10:90 100:100:1000 2000:1000:5000 10000]
for n = 1:length(s)
    s(n)
    A = rand(s(n));
    tic;
    Ainv = inv(A);
    t_inv(n) = toc;
end
loglog(s,t_inv)
xlabel('N')
ylabel('Time [s]')
```

Exercise: sample results

Each computer produces slightly different results because of background tasks, different matrices, etc. This is especially noticeable for small systems.



The time increases by 3 orders of magnitude, for every magnitude in N . A matrix inversion scales with $\mathcal{O}(N^3)$!

Solving a linear system in Excel using the inverse

$$Ax = b \quad \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix}$$

- Create matrix A in 3×3 cells
- Create right hand side vector b in 3 vertical cells

¹In Dutch Excel: INVERSEMAT

²In Dutch Excel: PRODUCTMAT. The semicolon may be a comma.

Solving a linear system in Excel using the inverse

$$Ax = b \quad \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix}$$

- Create matrix A in 3×3 cells
- Create right hand side vector b in 3 vertical cells
- Compute the inverse I :
 - Select an empty area of 3×3 cells
 - Type `=MINVERSE(B2:D4)`¹
 - Close with `Ctrl+Shift+Enter`

¹In Dutch Excel: INVERSEMAT

²In Dutch Excel: PRODUCTMAT. The semicolon may be a comma.

Solving a linear system in Excel using the inverse

$$Ax = b \quad \begin{bmatrix} 1 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 5 \end{bmatrix}$$

- Create matrix A in 3×3 cells
- Create right hand side vector b in 3 vertical cells
- Compute the inverse I :
 - Select an empty area of 3×3 cells
 - Type `=MINVERSE(B2:D4)`¹
 - Close with `Ctrl+Shift+Enter`
- Solution:
 - Select 3 vertical cells
 - Type `=MMULT(H2:J4 ; B6:B8)`²
 - Close with `Ctrl+Shift+Enter`

¹In Dutch Excel: INVERSEMAT

²In Dutch Excel: PRODUCTMAT. The semicolon may be a comma.

Today's outline

① Introduction

② Matrix inversion

③ Solving a linear system

④ Towards larger systems

⑤ Summary

Towards larger systems

Computation of determinants and inverses of large matrices in this way is too difficult (slow), so we need other methods to solve large linear systems!

Towards larger systems

- Determinant of upper triangular matrix:

$$\det |M_{\text{tri}}| = \prod_{i=1}^n a_{ii} \quad M = \begin{bmatrix} 5 & 3 & 2 \\ 0 & 9 & 1 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \det |M| = 5 \times 9 \times 1 = 45$$

- Matrix multiplication:

$$\det |AM| = \det |A| \times \det |M|$$

- When A is an identity matrix ($\det |A| = 1$):

$$\det |AM| = \det |A| \times \det |M| = 1 \times \det |M|$$

- With rules like this, we can use row-operations so that we can compute the determinant more cheaply.

Solutions of linear systems

Rank of a matrix: the number of linearly independent columns (columns that can not be expressed as a linear combination of the other columns) of a matrix.

$$M = \begin{bmatrix} 5 & 3 & 2 \\ 0 & 9 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

$$M = \begin{bmatrix} 1 & 2 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

- 3 independent columns
- In Matlab:

```
>> rank(M)
```

- col 2 = 2 × col 1
- col 4 = col 3 - col 1
- 2 independent columns:
rank = 2

Solutions of linear systems

The solution of a system of linear equations may or may not exist, and it may or may not be unique. Existence of solutions can be determined by comparing the rank of the Matrix M with the rank of the augmented matrix M_a :

```
>> rank(A)  
>> rank([A b])
```

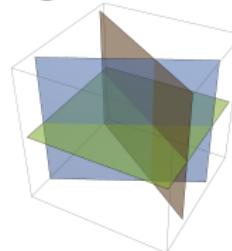
Our system: $Mx = b$

$$M = \begin{bmatrix} M_{11} & M_{12} & M_{13} \\ M_{21} & M_{22} & M_{23} \\ M_{31} & M_{32} & M_{33} \end{bmatrix}, b = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \Rightarrow M_a = \begin{bmatrix} M_{11} & M_{12} & M_{13} & b_1 \\ M_{21} & M_{22} & M_{23} & b_2 \\ M_{31} & M_{32} & M_{33} & b_3 \end{bmatrix}$$

Existence of solutions for linear systems

For a matrix M of size $n \times n$, and augmented matrix M_a :

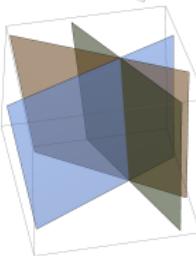
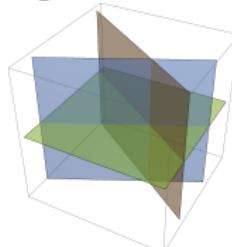
- $\text{Rank}(M) = n$:
Unique solution



Existence of solutions for linear systems

For a matrix M of size $n \times n$, and augmented matrix M_a :

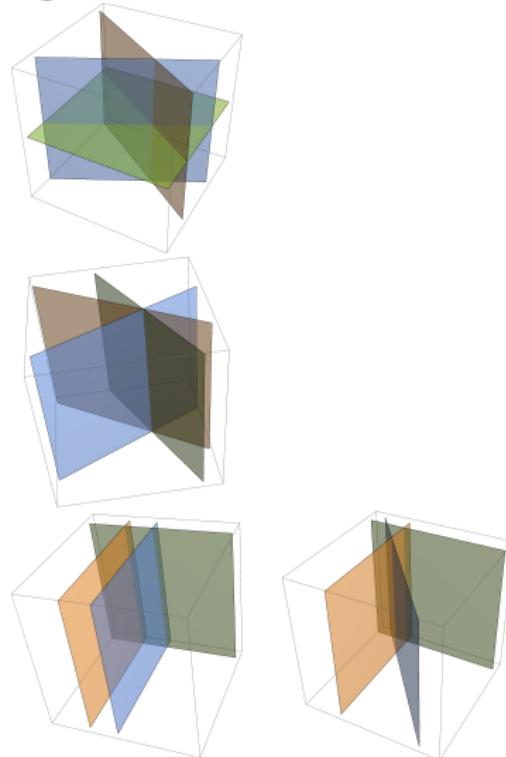
- $\text{Rank}(M) = n$:
Unique solution
- $\text{Rank}(M) = \text{Rank}(M_a) < n$:
Infinite number of solutions



Existence of solutions for linear systems

For a matrix M of size $n \times n$, and augmented matrix M_a :

- $\text{Rank}(M) = n$:
Unique solution
- $\text{Rank}(M) = \text{Rank}(M_a) < n$:
Infinite number of solutions
- $\text{Rank}(M) < n$,
 $\text{Rank}(M) < \text{Rank}(M_a)$:
No solutions



Two examples

$$M = \begin{bmatrix} 1 & 1 & 2 \\ 0 & 3 & 1 \\ 0 & 0 & 2 \end{bmatrix} \quad b = \begin{bmatrix} 17 \\ 11 \\ 4 \end{bmatrix} \Rightarrow M_a = \begin{bmatrix} 1 & 1 & 2 & 17 \\ 0 & 3 & 1 & 11 \\ 0 & 0 & 2 & 4 \end{bmatrix}$$

$\text{rank}(M) = 3 = n \Rightarrow \text{Unique solution}$

Two examples

$$M = \begin{bmatrix} 1 & 1 & 2 \\ 0 & 3 & 1 \\ 0 & 0 & 2 \end{bmatrix} \quad b = \begin{bmatrix} 17 \\ 11 \\ 4 \end{bmatrix} \Rightarrow M_a = \begin{bmatrix} 1 & 1 & 2 & 17 \\ 0 & 3 & 1 & 11 \\ 0 & 0 & 2 & 4 \end{bmatrix}$$

$\text{rank}(M) = 3 = n \Rightarrow \text{Unique solution}$

$$M = \begin{bmatrix} 1 & 1 & 2 \\ 0 & 3 & 1 \\ 0 & 0 & 0 \end{bmatrix} \quad b = \begin{bmatrix} 17 \\ 11 \\ 0 \end{bmatrix} \Rightarrow M_a = \begin{bmatrix} 1 & 1 & 2 & 17 \\ 0 & 3 & 1 & 11 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$\text{rank}(M) = \text{rank}(M_a) = 2 < n \Rightarrow \text{Infinite number of solutions}$

Today's outline

① Introduction

② Matrix inversion

③ Solving a linear system

④ Towards larger systems

⑤ Summary

Summary

- Linear equations can be written as matrices
- Using the inverse, the solution can be determined
 - Inverse via cofactors
 - Inverse and solution in Matlab
 - Inverse and solution in Excel
- Inversion scales with N^3
- A solution depends on the rank of a matrix

Linear equation solvers

Direct methods (elimination methods)

Ivo Roghair, Martin van Sint Annaland

Chemical Process Intensification,
Eindhoven University of Technology

Today's outline

① Introduction

② Gauss elimination

③ Partial Pivoting

④ LU decomposition

⑤ Summary

Today's outline

① Introduction

② Gauss elimination

③ Partial Pivoting

④ LU decomposition

⑤ Summary

Overview

Goals

Today we are going to write a program, which can solve a set of linear equations

- The first method is called Gaussian elimination
- We will encounter some problems with Gaussian elimination
- Then LU decomposition will be introduced

Today's outline

① Introduction

② Gauss elimination

③ Partial Pivoting

④ LU decomposition

⑤ Summary

Define the linear system

Consider the system:

$$Ax = b$$

In general:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

Desired solution:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b'_1 \\ b'_2 \\ b'_3 \end{bmatrix}$$

Using row operations

- Use row operations to simplify the system. Eliminate element A_{21} by subtracting $A_{21}/A_{11} = d_{21}$ times row 1 from row 2.
- In this case, Row 1 is the pivot row, and A_{11} is the pivot element.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ A_{21} & A_{22} & A_{23} & b_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right]$$

Using row operations

Eliminate element A_{21} using $d_{21} = A_{21}/A_{11}$.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ A_{21} & A_{22} & A_{23} & b_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right]$$

Using row operations

Eliminate element A_{21} using $d_{21} = A_{21}/A_{11}$.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ A_{21} & A_{22} & A_{23} & b_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right]$$

- $d_{21} \rightarrow A_{21}/A_{11}$
- $A_{21} \rightarrow A_{21} - A_{11}d_{21}$
- $A_{22} \rightarrow A_{22} - A_{12}d_{21}$
- $A_{23} \rightarrow A_{23} - A_{13}d_{21}$
- $b_2 \rightarrow b_2 - b_1 d_{21}$

Using row operations

Eliminate element A_{21} using $d_{21} = A_{21}/A_{11}$.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ A_{21} & A_{22} & A_{23} & b_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right]$$

- $d_{21} \rightarrow A_{21}/A_{11}$
- $A_{21} \rightarrow A_{21} - A_{11}d_{21}$
- $A_{22} \rightarrow A_{22} - A_{12}d_{21}$
- $A_{23} \rightarrow A_{23} - A_{13}d_{21}$
- $b_2 \rightarrow b_2 - b_1 d_{21}$

```
d21 = A(2,1)/A(1,1);  
A(2,1) = A(2,1) - A(1,1)*d21;  
A(2,2) = A(2,2) - A(1,2)*d21;  
A(2,3) = A(2,3) - A(1,3)*d21;  
b(2) = b(2) - b(1)*d21;
```

Using row operations

Eliminate element A_{31} using $d_{31} = A_{31}/A_{11}$.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ 0 & A'_{32} & A'_{33} & b'_3 \end{array} \right]$$

Using row operations

Eliminate element A_{31} using $d_{31} = A_{31}/A_{11}$.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ A_{31} & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ 0 & A'_{32} & A'_{33} & b'_3 \end{array} \right]$$

- $d_{31} \rightarrow A_{31}/A_{11}$
- $A_{31} \rightarrow A_{31} - A_{11}d_{31}$
- $A_{32} \rightarrow A_{32} - A_{12}d_{31}$
- $A_{33} \rightarrow A_{33} - A_{13}d_{31}$
- $b_3 \rightarrow b_3 - b_1d_{31}$

```
d31 = A(3,1)/A(1,1);  
A(3,1) = A(3,1) - A(1,1)*d31;  
A(3,2) = A(3,2) - A(1,2)*d31;  
A(3,3) = A(3,3) - A(1,3)*d31;  
b(3) = b(3) - b(1)*d31;
```

Using row operations

Eliminate element A_{32} using $d_{32} = A_{32}/A'_{22}$. Note that now the second row has become the pivot row.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ 0 & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ 0 & 0 & A''_{33} & b''_3 \end{array} \right]$$

Using row operations

Eliminate element A_{32} using $d_{32} = A_{32}/A'_{22}$. Note that now the second row has become the pivot row.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ 0 & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ 0 & 0 & A''_{33} & b''_3 \end{array} \right]$$

- $d_{32} \rightarrow A_{32}/A'_{22}$
- $A_{32} \rightarrow A_{32} - A'_{22}d_{32}$
- $A_{33} \rightarrow A_{33} - A'_{23}d_{32}$
- $b_3 \rightarrow b_3 - b'_2 d_{32}$

```
d32 = A(3,2)/A(2,2);  
A(3,2) = A(3,1) - A(2,2)*d32;  
A(3,3) = A(3,2) - A(2,3)*d32;  
b(3) = b(3) - b(2)*d32;
```

Using row operations

Eliminate element A_{32} using $d_{32} = A_{32}/A'_{22}$. Note that now the second row has become the pivot row.

$$\left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ 0 & A_{32} & A_{33} & b_3 \end{array} \right] \longrightarrow \left[\begin{array}{ccc|c} A_{11} & A_{12} & A_{13} & b_1 \\ 0 & A'_{22} & A'_{23} & b'_2 \\ 0 & 0 & A''_{33} & b''_3 \end{array} \right]$$

- $d_{32} \rightarrow A_{32}/A'_{22}$
- $A_{32} \rightarrow A_{32} - A'_{22}d_{32}$
- $A_{33} \rightarrow A_{33} - A'_{23}d_{32}$
- $b_3 \rightarrow b_3 - b'_2 d_{32}$

```
d32 = A(3,2)/A(2,2);  
A(3,2) = A(3,1) - A(2,2)*d32;  
A(3,3) = A(3,2) - A(2,3)*d32;  
b(3) = b(3) - b(2)*d32;
```

The matrix is now a triangular matrix — the solution can be obtained by back-substitution.

Backsubstitution

The system now reads:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & 0 & A''_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b'_2 \\ b''_3 \end{bmatrix}$$

Backsubstitution

The system now reads:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & 0 & A''_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b'_2 \\ b''_3 \end{bmatrix}$$

Start at the last row N , and work upward until row 1.

$$x_3 = b''_3 / A''_{33}$$

$$x_2 = (b'_2 - A'_{23}x_3) / A'_{22}$$

$$x_1 = (b_1 - A_{12}x_2 - A_{13}x_3) / A_{11}$$

Backsubstitution

The system now reads:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & 0 & A''_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b'_2 \\ b''_3 \end{bmatrix}$$

Start at the last row N , and work upward until row 1.

$$x_3 = b''_3 / A''_{33}$$

$$x_2 = (b'_2 - A'_{23}x_3) / A'_{22}$$

$$x_1 = (b_1 - A_{12}x_2 - A_{13}x_3) / A_{11}$$

$$\begin{aligned} x(3) &= b(3) / A(3,3) \\ x(2) &= (b(2) - A(2,3)*x(3)) / A(2,2) \\ x(1) &= (b(1) - A(1,2)*x(2) - A(1,3)*x(3)) / A(1,1) \end{aligned}$$

In general:

$$x_N = \frac{b_N}{A_{NN}} \quad x_i = \frac{b_i - \sum_{j=i+1}^N A_{ij}x_j}{A_{ii}}$$

Writing the program

- Create a function that provides the framework: take matrix A and vector b as an input, and return the solution x as output:

```
function [x,A,b] = GaussianEliminate(A,b)
```

- We will use *for-loops* instead of typing out each command line.
- Useful Matlab shortcuts:
 - $A(1,:) = [A_{11}, A_{12}, A_{13}]$
 - $A(:,2) = [A_{12}, A_{22}, A_{32}]^T$
 - $A(1,2:end) = [A_{12}, A_{13}]$
- A row operation could look like:

```
A(i,:) = A(i,:) - d*A(1,:)
```

The program: elimination

```
function [x,A,b] = GaussianEliminate(A,b)

% Perform elimination to obtain an upper triangular
matrix
N = length(b);
for column=1:(N-1) % Select pivot
    for row=(column+1):N % Loop over subsequent rows (
        below pivot)
        d=A(row,column)/A(column,column);
        A(row,:)=A(row,:)-d*A(column,:);
        b(row)= b(row)-d*b(column);
    end
end
```

The program: Backsubstitution

```
% Assign b to x
x=b;

% Perform backsubstitution
for row=N:-1:1
    x(row) = b(row);
    for i =(row+1):N
        x(row)=x(row)-A(row,i)*x(i);
    end
    x(row)=x(row)/A(row,row);
end
```

$$x_N = \frac{b_N}{A_{NN}} \quad x_i = \frac{b_i - \sum_{j=i+1}^N A_{ij}x_j}{A_{ii}}$$

Exercise: Gaussian Elimination

- The function we just made can be found on Canvas
- Use `help GaussianEliminate` to find out how it works
- Solve the following system of equations:

$$\begin{bmatrix} 9 & 9 & 5 & 2 \\ 6 & 7 & 1 & 3 \\ 6 & 4 & 3 & 5 \\ 2 & 6 & 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 7 \\ 4 \\ 10 \\ 1 \end{bmatrix}$$

- Compare your solution with `A\b`

Today's outline

① Introduction

② Gauss elimination

③ Partial Pivoting

④ LU decomposition

⑤ Summary

Partial pivoting

- Now try to run the algorithm with the following system:

$$\begin{bmatrix} 0 & 2 & 1 \\ 3 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 4 \\ 3 \\ 10 \end{bmatrix}$$

Partial pivoting

- Now try to run the algorithm with the following system:

$$\begin{bmatrix} 0 & 2 & 1 \\ 3 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 4 \\ 3 \\ 10 \end{bmatrix}$$

- It does not work! Division by zero, due to $A_{11} = 0$.
- Solution: Swap rows to move largest element to the diagonal.

Partial pivoting: implementing row swaps

- Find maximum element row below pivot in current column

```
[dummy, index] = max(abs(A(column:end, column)));
Index = index+column-1;
```

Partial pivoting: implementing row swaps

- Find maximum element row below pivot in current column
- Store current row

```
[dummy, index] = max(abs(A(column:end, column)));
Index = index+column-1;
```

```
temp = A(column,:);
```

Partial pivoting: implementing row swaps

- Find maximum element row below pivot in current column
- Store current row
- Swap pivot row and desired row in A

```
[dummy, index] = max(abs(A(column:end, column)));
Index = index+column-1;
```

```
temp = A(column,:);
```

```
A(column,:) = A(index,:);
A(index,:) = temp;
```

Partial pivoting: implementing row swaps

- Find maximum element row below pivot in current column
- Store current row
- Swap pivot row and desired row in A
- Do the same for b: store and swap

```
[dummy, index] = max(abs(A(column:end, column)));
Index = index+column-1;
```

```
temp = A(column,:);
```

```
A(column,:) = A(index,:);
A(index,:) = temp;
```

```
temp = b(column);
b(column) = b(index);
b(index) = temp;
```

Improve the program by using re-usable functions

```
function [x] = GaussianEliminate(A,b)
% GaussianEliminate(A,b): solves x in Ax=b
N = length(b);
for c=1:(N-1)
    [dummy , index]=max(abs(A(c:end,c)));
    index=index+c-1;
    A = SWAP(A,c,index); % Created swap function
    b = SWAP(b,c,index);
    for row=(column+1):N
        d=A(row,column)/A(column,column);
        A(row,:)=A(row,:)-d*A(column,:);
        b(row)= b(row)-d*b(column);
    end
end
x = backwardSub(A,b); % Created BS function
return
```

This function is also provided (named GaussianEliminate_v2 and GaussianEliminate_v3 on Canvas).

Alternatives to this program

- MATLAB can compute the solution to $Ax=b$ with its own solvers (more efficient) $A\b$
- Too many loops. Loops make MATLAB slow.
- There are fundamental problems with Gaussian elimination

Alternatives to this program

- MATLAB can compute the solution to $Ax=b$ with its own solvers (more efficient) $A\b$
- Too many loops. Loops make MATLAB slow.
- There are fundamental problems with Gaussian elimination
 - You can add a counter to the algorithm to see how many subtraction and multiplication operations it performs for a given size of matrix A .
 - The number of operations to perform Gaussian elimination is $\mathcal{O}(2N^3)$ (where N is the number of equations)
 - Exercise: verify this for our script

Alternatives to this program

- MATLAB can compute the solution to $Ax=b$ with its own solvers (more efficient) $A\b$
- Too many loops. Loops make MATLAB slow.
- There are fundamental problems with Gaussian elimination
 - You can add a counter to the algorithm to see how many subtraction and multiplication operations it performs for a given size of matrix A .
 - The number of operations to perform Gaussian elimination is $\mathcal{O}(2N^3)$ (where N is the number of equations)
 - Exercise: verify this for our script
 - LU decomposition takes $\mathcal{O}(2N^3/3)$ flops, 3 times less!
 - Forward and backward substitution each take $\mathcal{O}(N^2)$ flops (both cases)

Today's outline

- ① Introduction
- ② Gauss elimination
- ③ Partial Pivoting
- ④ LU decomposition
- ⑤ Summary

LU Decomposition

Suppose we want to solve the previous set of equations, but with several right hand sides:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} \begin{bmatrix} \vdots & \vdots & \vdots \\ x_1 & x_2 & x_3 \\ \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots \\ b_1 & b_2 & b_3 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

LU Decomposition

Suppose we want to solve the previous set of equations, but with several right hand sides:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} \begin{bmatrix} \vdots & \vdots & \vdots \\ x_1 & x_2 & x_3 \\ \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots \\ b_1 & b_2 & b_3 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

Factor the matrix A into two matrices, L and U, such that $A = LU$:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ \times & 1 & 0 \\ \times & \times & 1 \end{bmatrix} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \end{bmatrix}$$

Now we can solve for each right hand side, using only a forward followed by a backward substitution!

Substitutions

- Define a lower and upper matrix L and U so that $A = LU$
- Therefore $LUX = b$
- Define a new vector $y = UX$ so that $Ly = b$
- Solve for y , use L and forward substitution
- Then we have y , solve for x , use $UX = y$
- Solve for x , use U and backward substitution
- But how to get L and U?

Decomposing the matrix (1)

When we eliminate the element A_{21} we can keep multiplying by a matrix that undoes this row operations, so that the product remains equal to A .

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{21} & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A_{22} - d_{21}A_{12} & A_{23} - d_{21}A_{13} \\ A_{31} & A_{32} & A_{33} \end{bmatrix}$$

Decomposing the matrix (2)

When we eliminate the element A_{31} we can keep multiplying by a matrix that undoes this row operations, so that the product remains equal to A .

$$A = \begin{bmatrix} 1 & 0 & 0 \\ d_{21} & 1 & 0 \\ d_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} = A_{22} - d_{21}A_{12} & A'_{23} = A_{23} - d_{21}A_{13} \\ 0 & A'_{32} = A_{32} - d_{31}A_{12} & A'_{33} = A_{33} - d_{31}A_{11} \end{bmatrix}$$

Decomposing the matrix (3)

When we eliminate the element A_{32} we can keep multiplying by a matrix that undoes this row operations, so that the product remains equal to A .

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{21} & 1 & 0 \\ d_{31} & d_{32} & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & A'_{32} & A''_{33} = A'_{33} - d_{32}A'_{23} \end{bmatrix}$$

Decomposing the matrix (3)

When we eliminate the element A_{32} we can keep multiplying by a matrix that undoes this row operations, so that the product remains equal to A .

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{21} & 1 & 0 \\ d_{31} & d_{32} & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & A'_{32} & A''_{33} = A'_{33} - d_{32}A'_{23} \end{bmatrix}$$

This finishes the LU decomposition!

Pivoting during decomposition

Suppose we have arrived at the situation below, where $A'_{32} > A'_{22}$:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{21} & 1 & 0 \\ d_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & A'_{32} & A'_{33} \end{bmatrix}$$

Pivoting during decomposition

Suppose we have arrived at the situation below, where $A'_{32} > A'_{22}$:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{21} & 1 & 0 \\ d_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & A'_{32} & A'_{33} \end{bmatrix}$$

Exchange rows 2 and 3 to get the largest value on the main diagonal. Use a permutation matrix to store the swapped rows:

Pivoting during decomposition

Suppose we have arrived at the situation below, where $A'_{32} > A'_{22}$:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{21} & 1 & 0 \\ d_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & A'_{32} & A'_{33} \end{bmatrix}$$

Exchange rows 2 and 3 to get the largest value on the main diagonal. Use a permutation matrix to store the swapped rows:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{31} & 0 & 1 \\ d_{21} & 1 & 0 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{32} & A'_{33} \\ 0 & A'_{22} & A'_{23} \end{bmatrix}$$

Pivoting during decomposition

Suppose we have arrived at the situation below, where $A'_{32} > A'_{22}$:

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{21} & 1 & 0 \\ d_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{22} & A'_{23} \\ 0 & A'_{32} & A'_{33} \end{bmatrix}$$

Exchange rows 2 and 3 to get the largest value on the main diagonal. Use a permutation matrix to store the swapped rows:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ d_{31} & 0 & 1 \\ d_{21} & 1 & 0 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A'_{32} & A'_{33} \\ 0 & A'_{22} & A'_{23} \end{bmatrix}$$

Multiplying with a permutation matrix will swap the rows of a matrix. The permutation matrix is just an identity matrix, whose rows have been interchanged.

Recipe for LU decomposition

When decomposing matrix A into $A = LU$, it may be beneficial to swap rows to get the largest values on the diagonal of U (pivoting). A permutation matrix P is used to store row swapping such that:

$$PA = LU$$

- Write down a permutation matrix and the linear system
- Promote the largest value in the column diagonal
- Eliminate all elements below diagonal
- Move on to the next column and move largest elements to diagonal
- Eliminate elements below diagonal
- Repeat 5 and 6
- Write down L,U and P

LU decomposition example (1)

Write down a permutation matrix, which starts as the identity matrix, and the linear system:

$$PA = LU$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 1 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 1 \\ 1 & 2 & 0 \end{bmatrix}$$

LU decomposition example (1)

Write down a permutation matrix, which starts as the identity matrix, and the linear system:

$$PA = LU$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 1 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 1 \\ 1 & 2 & 0 \end{bmatrix}$$

Promote the largest value into the diagonal of column 1 — swap row 1 and 2:

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 1 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 2 & 0 \end{bmatrix}$$

LU decomposition example (2)

Eliminate all elements below the diagonal — row 2 already contains a zero in column 1, row 3 = row 3 - 0.5 row 1. Record the multiplier 0.5 in L :

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 1 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0.5 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 1.5 & -0.5 \end{bmatrix}$$

LU decomposition example (2)

Eliminate all elements below the diagonal — row 2 already contains a zero in column 1, row 3 = row 3 - 0.5 row 1. Record the multiplier 0.5 in L :

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 1 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0.5 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 1.5 & -0.5 \end{bmatrix}$$

Elimination of column 1 is done. Step to the next column, and move the largest value below/on the diagonal to the diagonal (swap rows 2 and 3). Adjust P and lower triangle of L accordingly:

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 0 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1.5 & -0.5 \\ 0 & 1 & 1 \end{bmatrix}$$

LU decomposition example (3)

Eliminate all elements below the diagonal —

row 3 = row 3 - $\frac{2}{3}$ row 2. Record the multiplier $\frac{2}{3}$ in L:

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 0 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0 & \frac{2}{3} & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1.5 & -0.5 \\ 0 & 0 & \frac{4}{3} \end{bmatrix}$$

LU decomposition example (3)

Eliminate all elements below the diagonal —

row 3 = row 3 - $\frac{2}{3}$ row 2. Record the multiplier $\frac{2}{3}$ in L:

$$\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 0 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0 & \frac{2}{3} & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1.5 & -0.5 \\ 0 & 0 & \frac{4}{3} \end{bmatrix}$$

We have obtained the matrices from $PA = LU$:

$$P = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \quad L = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0 & \frac{2}{3} & 1 \end{bmatrix} \quad U = \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1.5 & -0.5 \\ 0 & 0 & \frac{4}{3} \end{bmatrix}$$

Proceed with solving for x .

Substitutions

$$Ax = b \Rightarrow PAx = Pb \equiv d$$

$$PA = LU \Rightarrow LUx = d$$

- Define a new vector $y = Ux$
 - $Ly = b \Rightarrow Ly = d$
 - Solve for y , forward substitution:

$$y_1 = \frac{d_1}{L_{11}}$$

$$y_i = \frac{d_i - \sum_{j=1}^{i-1} L_{ij}y_j}{L_{ii}}$$

- Then solve $Ux = y$:
 - Solve for x , backward substitution:

$$x_N = \frac{y_N}{U_{NN}}$$

$$x_i = \frac{y_i - \sum_{j=i+1}^{N-1} U_{ij}x_j}{U_{ii}}$$

How to use the solver in Matlab

```
A = rand(5,5); % Get random matrix
[L, U, P] = lu(A); % Get L, U and P
b = rand(5,1); % Random b vector
d = P*b; % Permute b vector
y = forwardSub(L,d); % Can also do y=L\d
x = backwardSub(U,y); % Can also do x=U\y
rnorm = norm(A*x - b); % Residual

% Compare results to internal Matlab solver
x = A\b
```

How to use the solver in Matlab

```
A = rand(5,5); % Get random matrix
[L, U, P] = lu(A); % Get L, U and P
b = rand(5,1); % Random b vector
d = P*b; % Permute b vector
y = forwardSub(L,d); % Can also do y=L\d
x = backwardSub(U,y); % Can also do x=U\y
rnorm = norm(A*x - b); % Residual

% Compare results to internal Matlab solver
x = A\b
```

- Use this as a basis to create a function that takes A and b , and returns x .
- Use the function to check the performance for various matrix sizes and inspect the performance.

Today's outline

① Introduction

② Gauss elimination

③ Partial Pivoting

④ LU decomposition

⑤ Summary

Summary

- This lecture covered direct methods using elimination techniques.
- Gaussian elimination can be slow ($\mathcal{O}(N^3)$)
- Back substitution is often faster ($\mathcal{O}(N^2)$)
- LU decomposition means that we don't have to do Gaussian elimination every time (saves time and effort), but the matrix has to stay the same.
- Matlab has build in routines for solving linear equations (backslash operator \) and LU decomposition ([lu](#)).
- Advanced techniques such as (preconditioned) conjugate gradient or biconjugate gradient solvers are also available.
- Next part covers iterative approaches

Linear equation solvers

Iterative methods

Ivo Roghair, Martin van Sint Annaland

Chemical Process Intensification,
Eindhoven University of Technology

Today's outline

① Introduction

② Sparse matrices

③ Laplace's equation

④ Creating a sparse system

⑤ Iterative methods

⑥ Summary

Sparse matrices

- In many engineering cases, we deal with sparse matrices (as opposed to dense matrices)
 - A matrix is sparse when it mostly consists of zeros
 - Linear systems where equations depend on a limited number of variables (e.g. spatial discretization)
 - Storing zeros is not very efficient:

```
>> A = eye(10000);  
>> whos A  
>> S = sparse(A);  
>> whos S
```

- Can you think of a way to achieve this?
 - Sparse matrix formats: Yale, CRS, CCS

Sparse matrix storage format

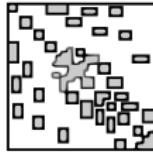
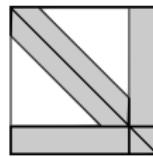
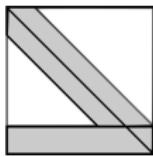
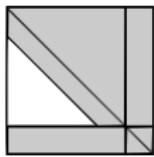
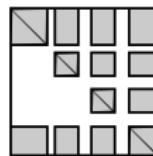
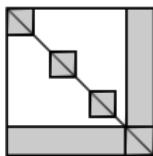
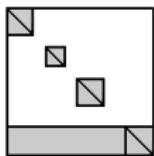
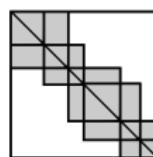
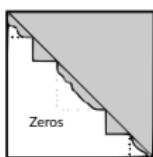
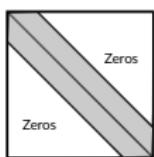
- Example: Yale storage format, storing 3 vectors:

- $A = [5 \ 8 \ 3 \ 6]$
- $IA = [0 \ 0 \ 2 \ 3 \ 4]$
- $JA = [0 \ 1 \ 2 \ 1]$

$$A = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 5 & 8 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 6 & 0 & 0 \end{bmatrix}$$

- A stores the non-zero values
- IA stores the index in A of the first non-zero in row i
- JA stores the column index
- Note: zero-based indices are used here!

Sparse matrix layout examples



Today's outline

- ① Introduction
- ② Sparse matrices
- ③ Laplace's equation
- ④ Creating a sparse system
- ⑤ Iterative methods
- ⑥ Summary

Laplace's equation

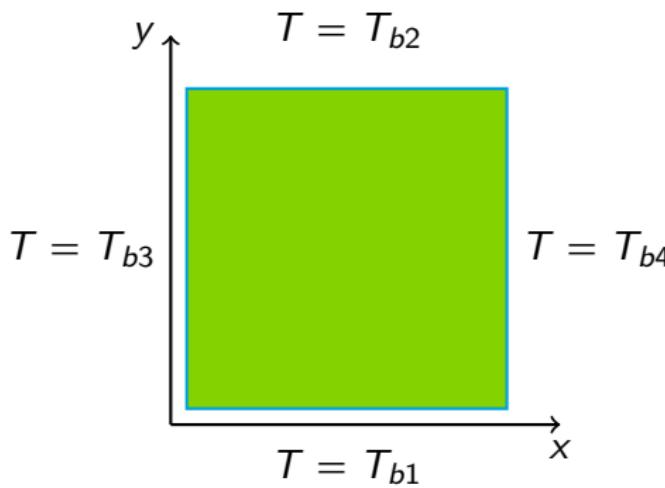
$$\frac{\partial T}{\partial t} = \alpha \nabla^2 T$$

T = Temperature

α = Thermal diffusivity

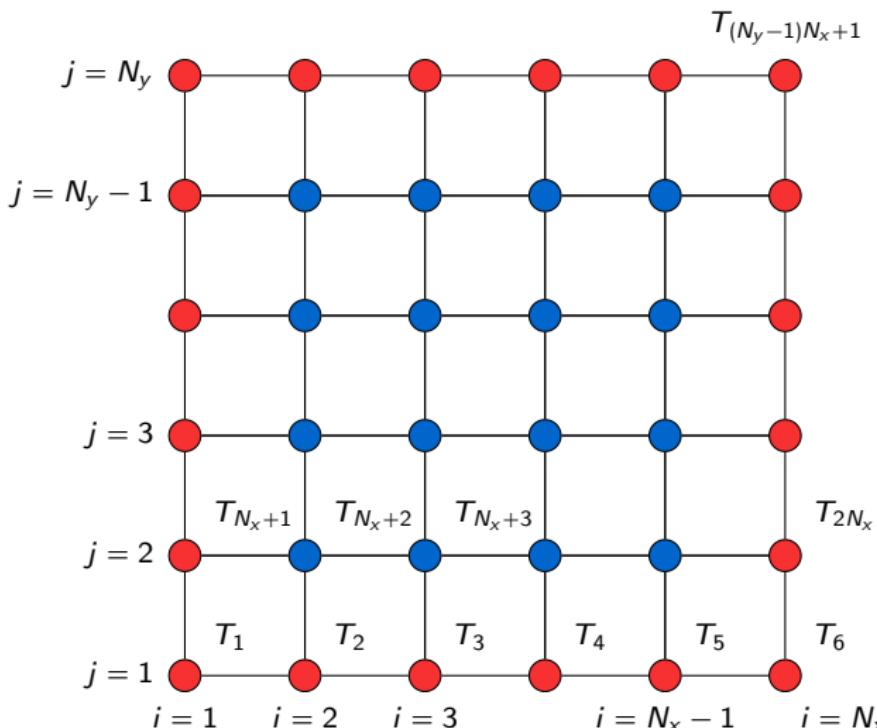
In steady state:

$$\nabla^2 T = 0$$



$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = 0$$

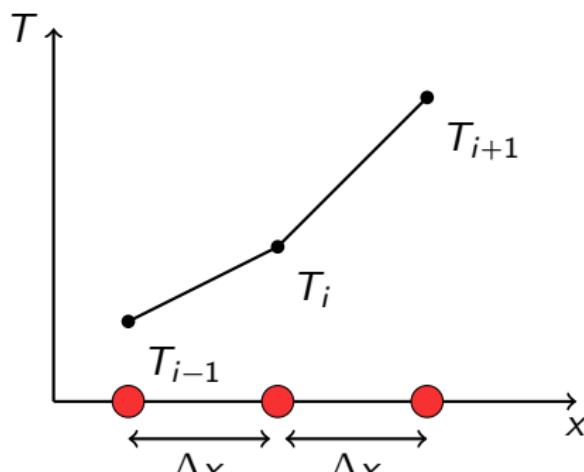
Discretization of Laplace's equation (I)



- Define a grid of points in x and y
- Index of the grid points using 2D coordinates i and j
- Set up the equations using a 1D index system:
 $T_{i,j} = T_{i+N_x(j-1)}$

Discretization of Laplace's equation (II)

Estimate the second-order differentials: assume a piece-wise linear profile in the temperature:



$$\begin{aligned}\frac{\partial^2 T}{\partial x^2} &\approx \frac{\frac{\partial T}{\partial x}\Big|_{i+\frac{1}{2}} - \frac{\partial T}{\partial x}\Big|_{i-\frac{1}{2}}}{\Delta x} \\ &\approx \frac{\frac{(T_{i+1,j} - T_{i,j})}{\Delta x} - \frac{(T_{i,j} - T_{i-1,j})}{\Delta x}}{\Delta x} \\ &= \frac{T_{i+1,j} - 2T_{i,j} + T_{i-1,j}}{(\Delta x)^2}\end{aligned}$$

Discretization of Laplace's equation (III)

The y -direction is derived analogously, so that the 2D Laplace's equation is discretized as:

$$\frac{T_{i+1,j} - 2T_{i,j} + T_{i-1,j}}{(\Delta x)^2} + \frac{T_{i,j+1} - 2T_{i,j} + T_{i,j-1}}{(\Delta y)^2} = 0$$

Use a single index counter $k = i + N_x(j - 1)$, so that the equation becomes:

$$\frac{T_{k+1} - 2T_k + T_{k-1}}{(\Delta x)^2} + \frac{T_{k+N_x} - 2T_k + T_{k-N_x}}{(\Delta y)^2} = 0$$

For an equal spaced grid $\Delta x = \Delta y = 1$:

$$T_{k-N_x} + T_{k-1} - 4T_k + T_{k+1} + T_{k+N_x} = 0$$

$$\Rightarrow AT = b$$

Today's outline

- ① Introduction
 - ② Sparse matrices
 - ③ Laplace's equation
 - ④ Creating a sparse system
 - ⑤ Iterative methods
 - ⑥ Summary

Creating the linear system

$$T_{k-N_x} + T_{k-1} - 4T_k + T_{k+1} + T_{k+N_x} = 0$$

Create a *banded* matrix A : the main diagonal k contains -4, whereas the bands at $k - 1$, $k + 1$, $k - N_x$ and $k + N_x$ contain a 1. Boundary cells just contain a 1 on the main diagonal so that the temperature is equal to T_b (e.g. $T_1 = 1T_b$).

$$\left[\begin{array}{ccccccccc} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots \\ \cdots & 1 & \cdots & 1 & -4 & 1 & \cdots & 1 & \ddots & 0 \\ 0 & \cdots & 1 & \cdots & 1 & -4 & 1 & \cdots & 1 & \vdots \\ \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right] \left[\begin{array}{c} T_1 \\ T_2 \\ \vdots \\ T_k \\ T_{k+1} \\ \vdots \\ T_{(N_y-1)N_x} \\ T_{(N_y-1)N_x+1} \end{array} \right] = \left[\begin{array}{c} T_b \\ T_b \\ \vdots \\ 0 \\ 0 \\ \vdots \\ T_b \\ T_b \end{array} \right]$$

Creating the linear system

$$T_{k-N_x} + T_{k-1} - 4T_k + T_{k+1} + T_{k+N_x} = 0$$

Create a *banded* matrix A in Matlab, by setting the coefficients for the internal cells:

```
Nx=5; %number of points along x direction
Ny=5; %number of points in the y direction
Nc=Nx*Ny; % Total number of points

e = ones(Nc,1);
A = spdiags([e,e,-4*e,e,e],[-Nx,-1,0,1,Nx],Nc,Nc);
b = zeros(Nc,1);
```

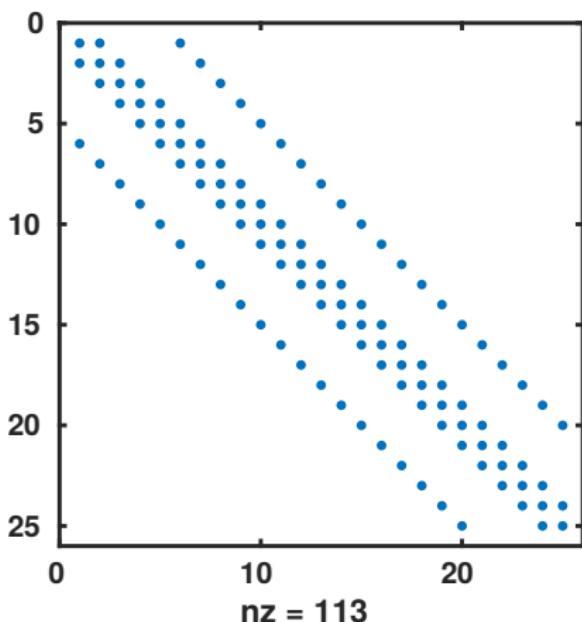
The function `spdiags` uses the following arguments:

- The coefficients that have to be put on the diagonals arranged as columns in a matrix
- The position of the bands with respect to the main diagonal
- Size of the resulting matrix (in our case square $N_x N_y \times N_x N_y$)

Matrix sparsity

- Let's check the matrix layout:

```
>> spy(A)
```
- This command shows the non-zero values of a matrix
- Apart from the main diagonal, there are offset bands!



About boundary conditions

- For the nodes on the boundary, we have a simple equation:

$$T_{k,\text{boundary}} = \text{Some fixed value}$$

- However, we have set all nodes to be a function of their neighbors...
- Find the boundary node indices using $k = i + Nx(j - 1)$
 - i = 1, j = 1:Ny
 - i = Nx, j = 1:Ny
 - j = 1, i = 1:Nx
 - j = Ny, i = 1:Nx
- Reset the row in A to zeros, set $A_{kk} = 1$
- Set value in rhs: $b_k = T_{k,\text{boundary}}$
- Boundary conditions are often more elaborate to implement!
See `setBoundaryConditions.m`.

Partial implementation of the boundary conditions

See `setBoundaryConditions.m`.

```
function [A,b] = setBoundaryConditions(A,b,Tb,Nx,Ny)

% Set boundary conditions over x-direction
for i=1:Nx
    j = 1;
    ind = i + Nx * (j-1);
    A(ind,:) = 0; % Reset matrix for boudary cells
    A(ind,ind) = 1; % Add a 1 on the diagonal
    b(ind) = Tb(1);
    j = Ny;
    ind = i + Nx * (j-1);
    A(ind,:) = 0; % Reset matrix for boudary cells
    A(ind,ind) = 1; % Add a 1 on the diagonal
    b(ind) = Tb(2);
end

%% Repeat for y-direction
```

How applying boundary conditions affects the linear system

```
function [A,b] = setBoundaryConditions(A,b,Tb,Nx,Ny)
```

- Make sure that matrix A and right hand side vector b are in your workspace, as well as N_x and N_y
- Create a vector that holds the temperature at each boundary:

```
>> T = [10 20 30 40];
```

- Call the function, store A and b in new variables:

```
>> [A2,b2] = setBoundaryConditions(A,b,T,Nx,Ny);
```

- Check the new structure of the matrix and the right hand side:

```
>> subplot(1,2,1); spy(A2);
>> subplot(1,2,2); spy(b2);
```

A full program, including solver

The program and auxiliary functions are on Canvas
(`solveLaplaceEq.m`)

```
function [x,y,T,A] = solveLaplaceEq(Nx,Ny)
% Solves the steady-state Laplace equation

Tb = [10 20 30 40]; % Fixed boundary temperatures

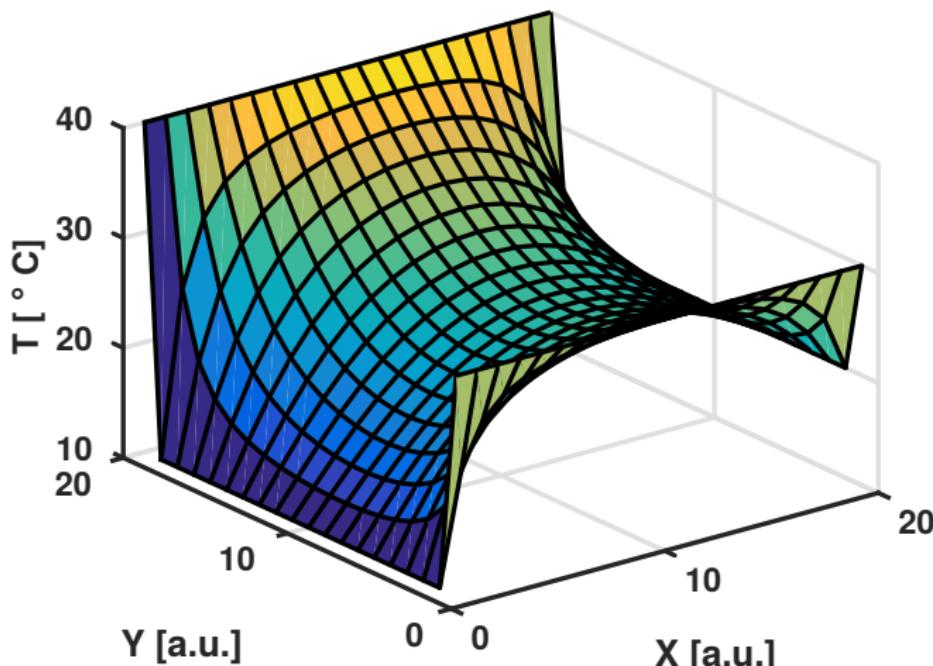
% Fill sparse matrix with [1 1 -4 1 1]
e = ones(Nx*Ny,1);
A = spdiags([e,e,-4*e,e,e],[-Nx,-1,0,1,Nx],Nx*Ny,Nx*Ny);
b = zeros(Nx*Ny,1);

[A,b] = setBoundaryConditions(A,b,Tb,Nx,Ny);

T = A\b; % Solve matrix
Tc = reshape(T,[Nx,Ny]); % Reshape x-vec to mat Nx,Ny
[xc yc] = meshgrid(1:Nx,1:Ny); % Get position arrays
surf(xc,yc,Tc); % Surface plot
```

Sample results

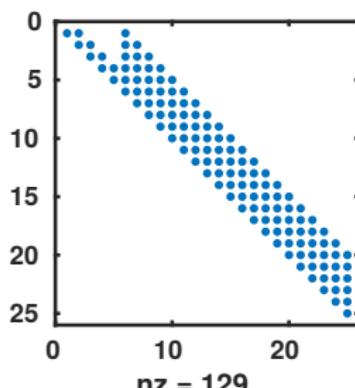
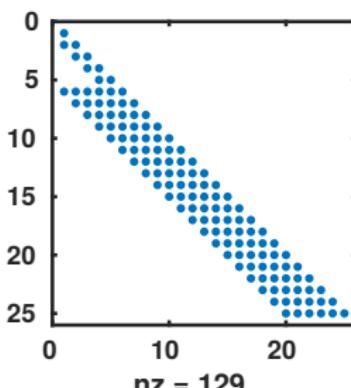
Solved for a 20×20 system with $T_b = [10 \ 20 \ 30 \ 40]$.



LU decomposition of a sparse matrix

- With LU decomposition we produce matrices that are less sparse than the original matrix.
- Sparse storage often required, and also numerical techniques that fully utilizes this!

```
>> [L,U,P] = lu(A)
>> subplot(1,2,1)
>> spy(L)
>> subplot(1,2,2)
>> spy(U)
```



LU decomposition

- LU decomposition and Gaussian elimination on a matrix like A requires more memory (with 3D problems, the offset in the diagonal would even be bigger!)
- In general extra memory allocation will not be a problem for MATLAB
- MATLAB is clever, in that sense that it attempts to reorder equations, to move elements closer to the diagonal)

Alternatives for elimination methods

- Use iterative methods when systems are large and sparse.
- Often such systems are encountered when we want to solve PDEs of higher dimensions

Today's outline

- ① Introduction
- ② Sparse matrices
- ③ Laplace's equation
- ④ Creating a sparse system
- ⑤ Iterative methods
- ⑥ Summary

Examples of iterative methods

- Jacobi method
 - Gauss-Seidel method
 - Successive over relaxation
-
- `bicg` — Bi-conjugate gradient method
 - `pcg` — preconditioned conjugate gradient method
 - `gmres` — generalized minimum residuals method
 - `bicgstab` — Bi-conjugate gradient method

The Jacobi method

- In our example we derived the following equation:

$$T_{k-N_x} + T_{k-1} - 4T_k + T_{k+1} + T_{k+N_x} = 0$$

- Rearranging gives:

$$T_k = \frac{T_{k-N_x} + T_{k-1} + T_{k+1} + T_{k+N_x}}{4}$$

- In the Jacobi scheme the iteration proceeds as follows:
 - Start with an initial guess for the values of T at each node
 - Compute updated values and store a new vector:

$$T_k^{\text{new}} = \frac{T_{k-N_x}^{\text{old}} + T_{k-1}^{\text{old}} + T_{k+1}^{\text{old}} + T_{k+N_x}^{\text{old}}}{4}$$

- Do this for all nodes
- Repeat the procedure until converged

Jacobi method for Laplace's equation

See `laplace_jacobi.m` (from Canvas)

```
% Grid size
nx = 40; ny = 40;
% The temperature field + boundaries at old and new times
T = zeros(nx,ny);
T(1,:) = 40; % Left
T(nx,:) = 60; % Right
T(:,1) = 20; % Bottom
T(:,ny) = 30; % Top
Tnew = T;
% For plotting
[x y] = meshgrid(1:nx, 1:ny);
for iter = 1:1000
    for i = 2:nx-1
        for j = 2:ny-1
            Tnew(i,j) = (T(i-1,j)+T(i+1,j)+T(i,j-1)+T(i,j+1))/4.0;
        end
    end
    surf(x,y,Tnew);
    title(['Iteration: ', num2str(iter)]);
    drawnow
    T = Tnew; % Update T
end
```

About the straightforward implementation

- The method as implemented works fine for a simple Laplace equation
- For generic systems of linear equations, the implementation cannot be used.

We will now introduce the Jacobi method so it can be used for generic systems of linear equations.

The Jacobi method with matrices

We can split our (banded) matrix A into a diagonal matrix D and a remainder R :

$$A = D + R$$

$$\begin{bmatrix} \times & \times & & & & \times & & & \\ \times & \times & \times & & & \times & & & \\ & \times & \times & \times & & & & & \\ & & \times & \times & \times & & & & \\ & & & \times & \times & \times & & & \\ & & & & \times & \times & \times & & \\ & & & & & \times & \times & \times & \\ & & & & & & \times & \times & \times \\ & & & & & & & \times & \times \\ & & & & & & & & \times \end{bmatrix} = \begin{bmatrix} \times & & & & & & & & \\ & \times & & & & & & & \\ & & \times & & & & & & \\ & & & \times & & & & & \\ & & & & \times & & & & \\ & & & & & \times & & & \\ & & & & & & \times & & \\ & & & & & & & \times & \\ & & & & & & & & \times \end{bmatrix} + \begin{bmatrix} \times & & & & & & & & \\ & \times & & & & & & & \\ & & \times & & & & & & \\ & & & \times & & & & & \\ & & & & \times & & & & \\ & & & & & \times & & & \\ & & & & & & \times & & \\ & & & & & & & \times & \\ & & & & & & & & \times \end{bmatrix}$$

Jacobi method: solving a system

- We can solve $AT = b$, now written generally as $Ax = b$, by:

$$Ax = b$$

$$(D + R)x = b$$

$$Dx = b - Rx$$

$$Dx^{\text{new}} = b - Rx^{\text{old}}$$

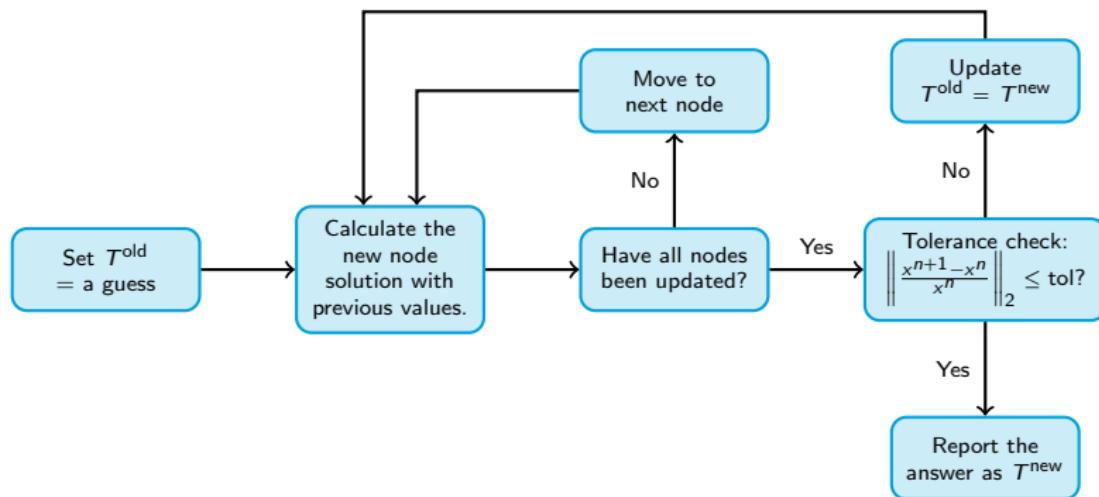
$$x^{\text{new}} = D^{-1}(b - Rx^{\text{old}})$$

- Using the n and $n + 1$ notation for old and new time steps, we find in general:

$$x^{n+1} = D^{-1}(b - Rx^n)$$

$$x_i^{n+1} = \frac{1}{A_{ii}} \left(b_i - \sum_{j \neq i} A_{ij} x_j^n \right)$$

Diagram of the Jacobi method



The core of the solver

The full file is on Canvas, solveJacobi.m.

```
1 while ( xDiff > tol && it_jac < 1000 )
2     x_old = x;
3     for i=1:N
4         s = 0;
5         for j = 1:N
6             if (j ~= i)
7                 s = s+A(i,j)*x_old(j);
8             end
9         end
10        x(i) = (b(i)-s)/A(i,i);
11    end
12    it_jac = it_jac+1;
13    xDiff = norm((x-x_old)./x,2);
14 end
15 it_jac
```

Try to call it from the solveLaplaceEq.m file, instead of using \.

A few details on this algorithm

- The while loop holds two aspects
 - A convergence criterion (`norm((x-x_old)./x,2)> tol)`). Some considerations are:
 - L_1 -norm (sum)
 - L_2 -norm (Euclidian distance)
 - L_∞ -norm (max)
 - Protection against infinite loops (no convergence)
- Reset the sum for each row, before summing for the new unknown node
- Start vector x is not shown in the example, but should be there!
- It can have huge impact on performance!
- The for-loops also have a large performance penalty!

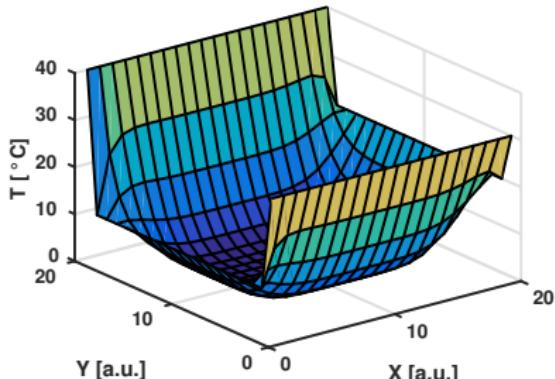
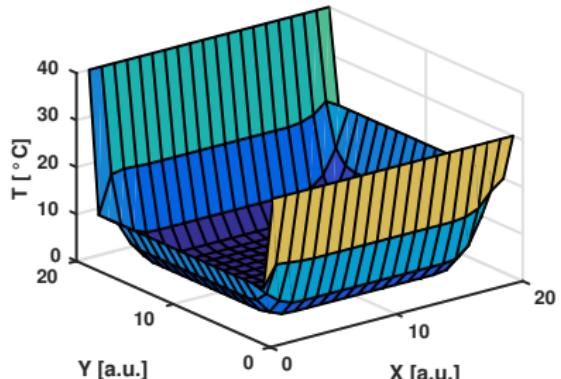
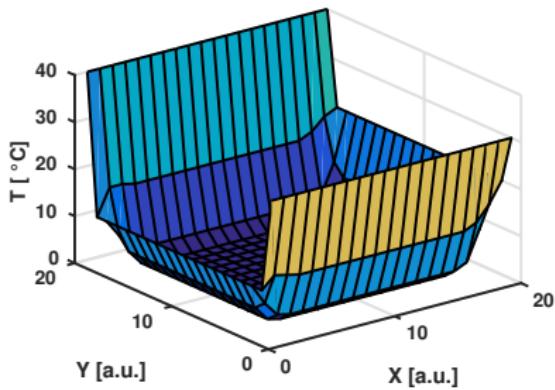
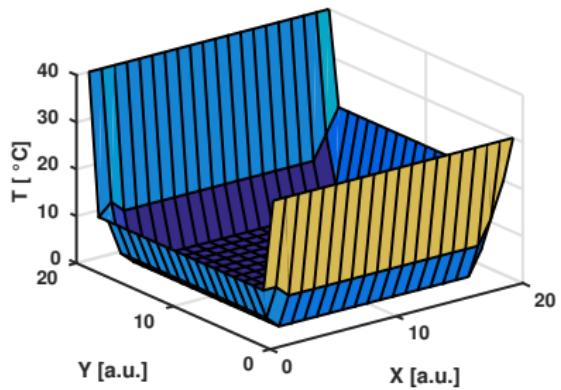
The solver using array indices

Make a copy of the Jacobian solver, and replace the for-loop by a vector-operation:

```
% While not converged or max_it not reached
while ( xDiff > tol && it_jac < 1000 )
    x_old = x;
    for i=1:N
        % Sum off-diagonal*x_old
        offDiagonalIndex = [1:(i-1) (i+1):N];
        Aij_Xj = A(i,offDiagonalIndex)*x_old(offDiagonalIndex);

        % Compute new x value
        x(i) = (b(i)-Aij_Xj)/A(i,i);
    end
    it_jac = it_jac+1;
    xDiff = norm((x-x_old)./x,2);
end
```

Iterations 1, 2, 3 and 10



Gauss-Seidel method

The Gauss-Seidel method is quite similar to Jacobi method

- The only difference is that the new estimate x^{new} is returned to the solution x^{old} as soon as it is completed
- For following nodes, the updated solution is used immediately
- Our straightforward script (from the Jacobi method) is therefore changed easily:
 - Do not create a T_{new} array (save memory!)
 - Do not store the solution in T_{new} , but simply in T
 - Do not perform the update step $T=T_{\text{new}}$
 - See `laplace_gaussseidel.m` for the algorithm.
- The straightforward script works well for the current Laplace equation, but we define the generic Gauss-Seidel algorithm on the following slides.

Gauss-Seidel method

- Define a lower and strictly upper triangular matrix, such that $A = L + U$
- Now we can solve $AT = b$ by:

$$(L + U)T = b$$

$$LT = b - UT$$

$$LT^{\text{new}} = b - UT^{\text{old}}$$

$$T^{\text{new}} = L^{-1}(b - UT^{\text{old}})$$

- Using the n and $n + 1$ notation for old and new time steps, we find in for the general Gauss-Seidel method:

$$x^{n+1} = L^{-1}(b - Ux^n)$$

$$x_i^{n+1} = \frac{1}{A_{ii}} \left(b_i - \sum_{j < i} A_{ij}x_j^{n+1} - \sum_{j > i} A_{ij}x_j^n \right)$$

Today's outline

- ① Introduction
- ② Sparse matrices
- ③ Laplace's equation
- ④ Creating a sparse system
- ⑤ Iterative methods
- ⑥ Summary

Summary

- Partial differential equations can be written as sparse systems of linear equations
- Sparse systems can be handled with a direct method like Gaussian elimination
- If you have systems of more than 1 dimension, a direct method still can be used, if there are no memory issues, otherwise an iterative method may be attractive.
- The Jacobi method was introduced. Many other methods are based on the Jacobi method (SOR method, for example)

Numerical methods for Chemical Engineers:

Non-linear equations

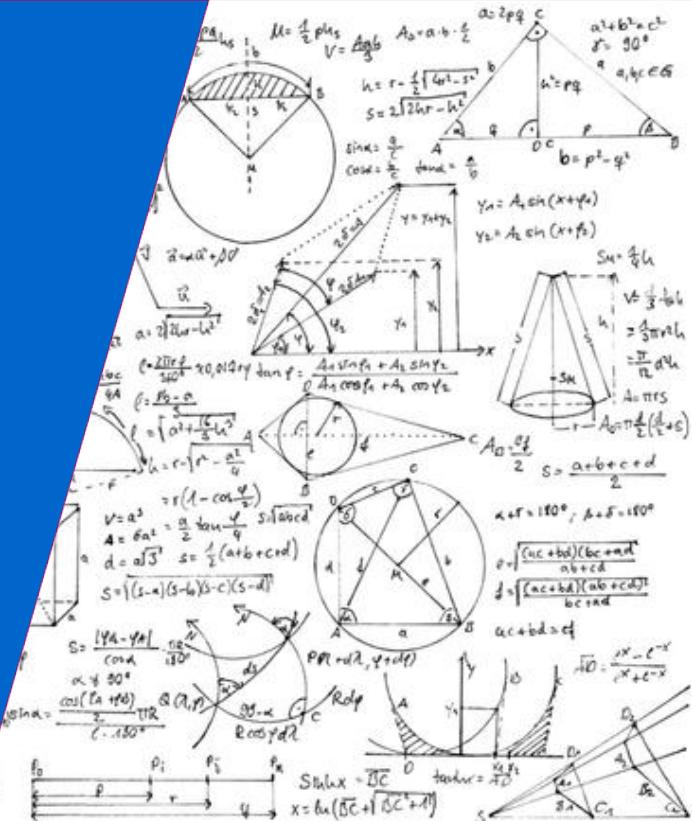
Prof.dr.ir. Martin van Sint Annaland
Dr.ir. Ivo Roghair

Chemical process Intensification

TU/e

Technische Universiteit
Eindhoven
University of Technology

Where innovation starts



Content

- **How to solve:**

$f(x) = 0$ for arbitrary functions f

“Root finding”

(i.e. move all terms to the left)

- **One dimensional case:** $f(x) = 0$
“Bracket or ‘trap’ a root between bracketing values,
then hunt it down like a rabbit.”
- **Multi-dimensional case:** $f(x) = 0$
 - N equations in N unknowns:
You can only *hope* to find a solution.
It may have no (real) solution, or more than one solution!
 - Much more difficult!!
“You never know whether a root is near,
unless you have found it”

Outline

- **One-dimensional case:**

- Bisection method
- Secant and false position method
- Brent's method
- Newton-Raphson method

- **Multi-dimensional case:**

- Newton-Raphson method
- Broyden's method

Do not use routines
as black boxes without
understanding them!!!

- Introduction to underlying ideas and algorithms
- Exercises in how to program the methods in Excel and MATLAB.

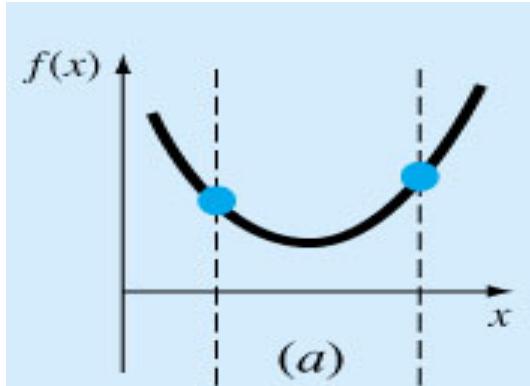
General idea

- **Root finding proceeds by iteration:**
 - Start with a good initial guess (crucially important!!)
 - Use an algorithm to improve the solution until some predetermined convergence criterion is satisfied
 - **Pitfalls:**
 - Convergence to the wrong root...
 - Fails to converge because there is no root...
 - Fails to converge because your initial estimate was not close enough...
- It never hurts to inspect your function graphically
- Pay attention to carefully select initial guesses

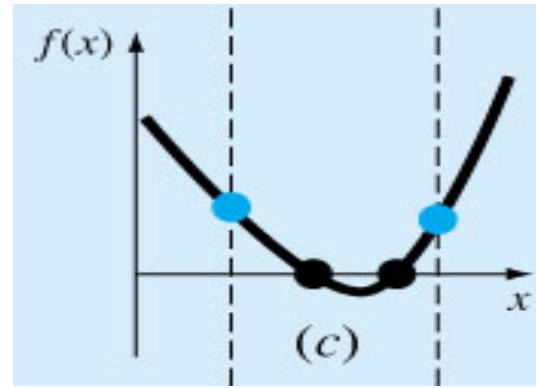
Hamming's motto:
the purpose of computing
is insight, not numbers!!

General idea

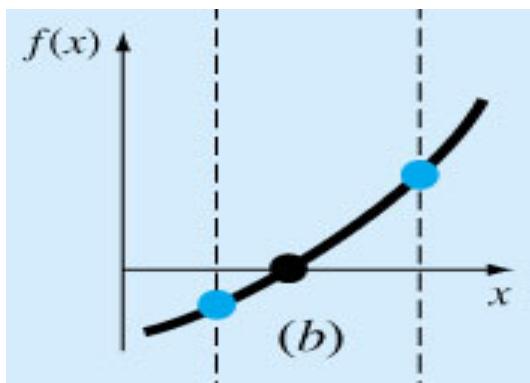
- Examples of pitfalls of root finding...



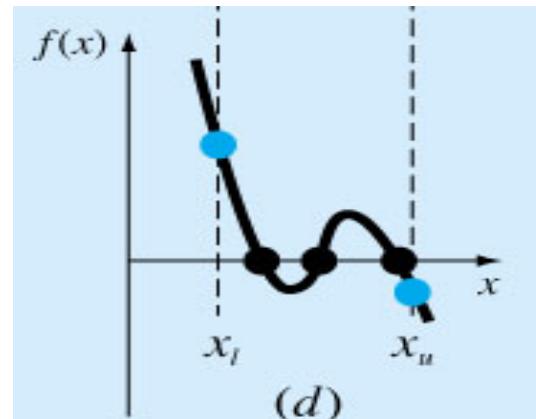
No answer (no root)



Oops!! (two roots!!)



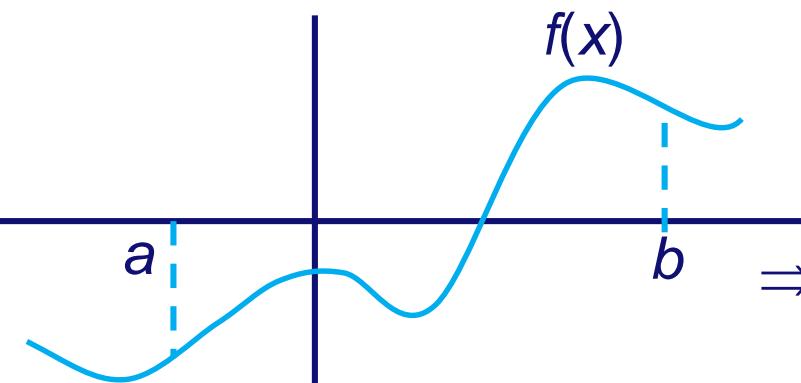
Nice case (one root)



Three roots (might work for a while!)

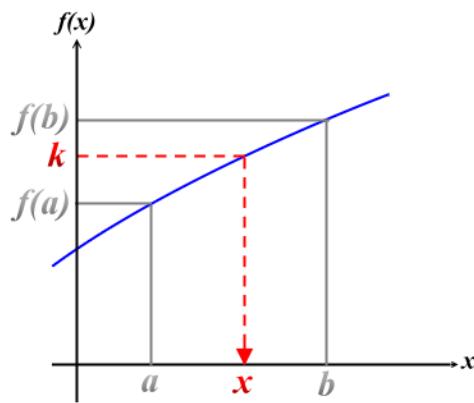
Bracketing

Bracketing a root = knowing that the function changes sign in an identified interval



A root is bracketed in the interval (a,b) , if $f(a)$ and $f(b)$ have opposite signs

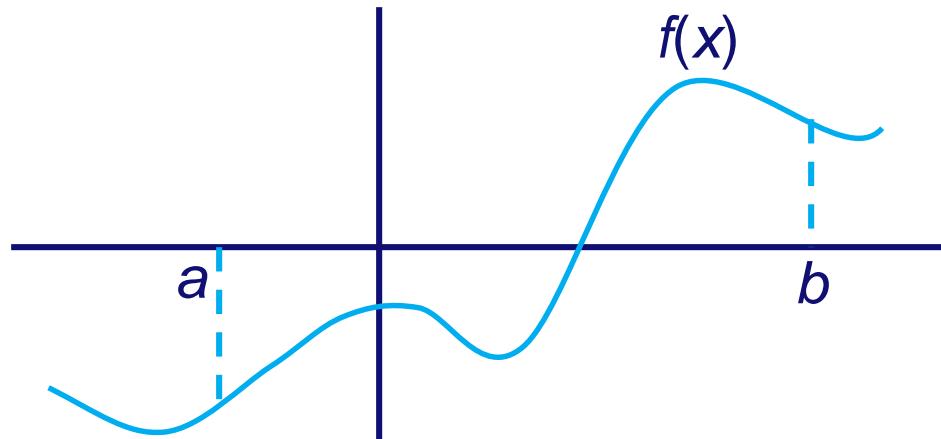
⇒ At least one root must lie in this interval, if the function is continuous



Intermediate Value Theorem
If $f(x)$ is continuous on $[a,b]$ and k is a constant that lies between $f(a)$ and $f(b)$, then there is a value $x \in [a,b]$ such that $f(x) = k$

Bracketing

Bracketing a root = knowing that the function changes sign in an identified interval



- **General best advise:**
 - Always bracket a root before trying to convergence...
 - Never allow your iteration to method to get outside the best bracketing bounds...

Bracketing

Exercise 1:

- Write a function in MATLAB to bracket a function given an initial guessed range x_1 and x_2 .
(via expansion of the interval)
- Write a program to find out how many roots exist (at minimum) in the interval x_1 and x_2 .

Of course these functions can then be combined to create a function that returns bracketing intervals for different roots.

Passing functions in Matlab

- In MATLAB function names can be passed as arguments to functions, this is called a **function handle** (other programs would call this pointer).
- For example: to solve $f(x) = x^2 - 4x + 2 = 0$ numerically, we can write a function that returns the value of f :

```
function f = MyFunc(x)
f = x^2 - 4*x + 2;           (Note: case sensitive!!)
return
```

- The function handle can be used as an alias
 - >> f = @MyFunc; a = 4; b = f(a)
- We can then call a solving routine (e.g. fzero):

```
>> ans = fzero(@MyFunc,5)
>> fzero(@(x) x^2-4*x+2,5)
```

Bracketing

- **Exercise 1: Function to bracket a function**

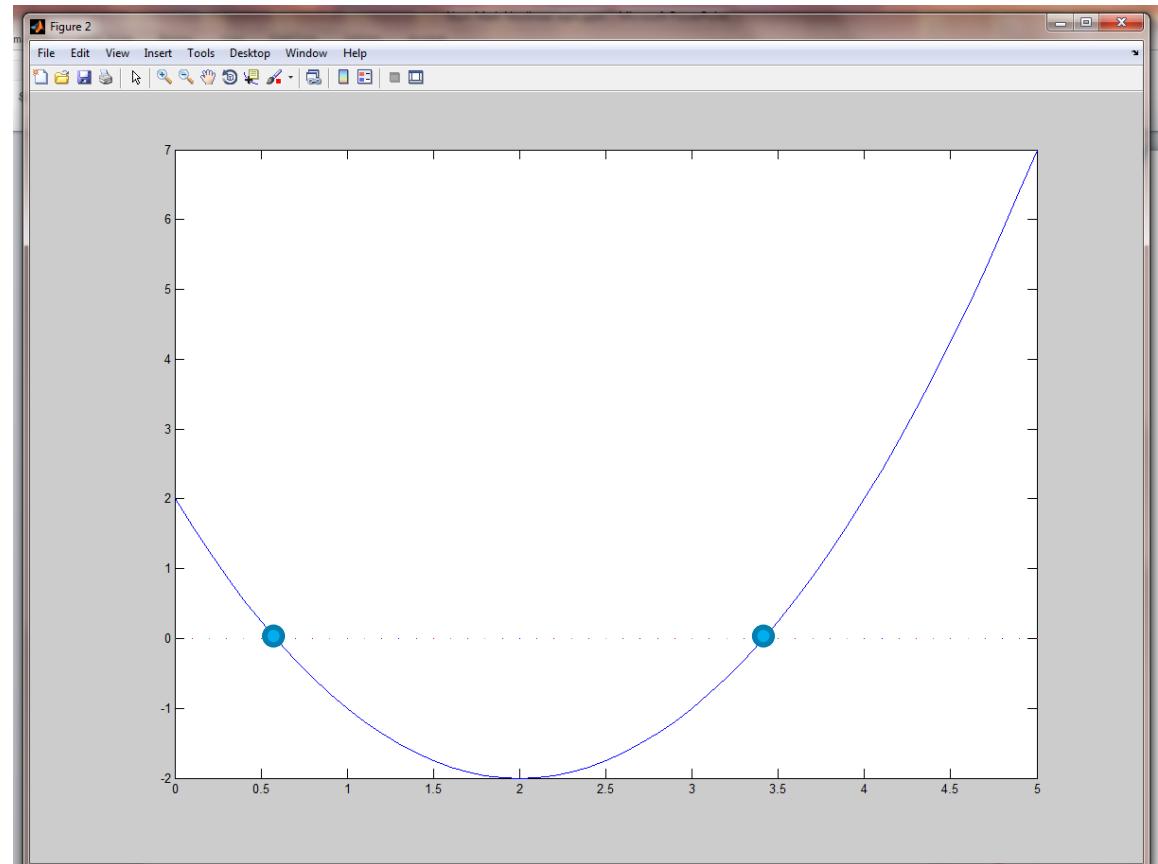
If possible, first make a graph: for example via

```
>> x=0:0.1:5;  
>> y=x.^2-4*x+2;  
>> figure;  
>> plot(x,y,x,0);
```

Makes immediately clear that there are two roots.

$$x_1 = 2 - \sqrt{2} \approx 0.59$$

$$x_2 = 2 + \sqrt{2} \approx 3.41$$



Bracketing

• Exercise 1: Function to bracket a function

```
1  [-] function found = brac(func, x1, x2)
2  -    ntry = 50;
3  -    factor = 1.6;
4
5  -    found = false;
6  -    if (x1~==x2)
7  -        f1 = func(x1);
8  -        f2 = func(x2);
9  -        for i = 1:ntry
10 -            if (f1*f2<0)
11 -                found = true
12 -                break;
13 -            end;
14 -            if (abs(f1)<abs(f2))
15 -                x1 = x1 + factor*(x1-x2);
16 -                f1 = func(x1);
17 -            else
18 -                x2 = x2 + factor*(x2-x1);
19 -                f2 = func(x2);
20 -            end;
21 -        end;
22 -    else
23 -        disp('Bad initial range!');
24 -    end;
25
26 -    if found
27 -        disp(sprintf('The bracketing interval = [%f, %f]\n', [x1,x2]));
28 -    else
29 -        disp('No bracketing interval found!');
30 -    end;
31 -    return
```

a function to expand the interval (x_1, x_2) maximally $2^{50} \sim 10^{15}$, until a root is found

returns true when root is found and false otherwise

displays results

Bracketing

- **Exercise 1: Function to bracket a function**

```
1  function nroot = brak(func, x1, x2, n)
2  nroot = 0;
3  dx = (x2 - x1)/n;
4  x = x1;
5  fp = func(x1);
6  for i = 0:n
7      x = x + dx;
8      fc = func(x);
9      if (fc*fp<=0)
10         nroot = nroot + 1;
11         xb1(nroot) = x - dx;
12         xb2(nroot) = x;
13     end;
14     fp = fc;
15   end;
16   if n>0
17     for i = 1:nroot
18       disp(sprintf('Root %d in bracketing interval [%f, %f]', [i,xb1(i),xb2(i)]));
19     end
20   else
21     disp('No roots found!');
22   end;
23
24   return;
```

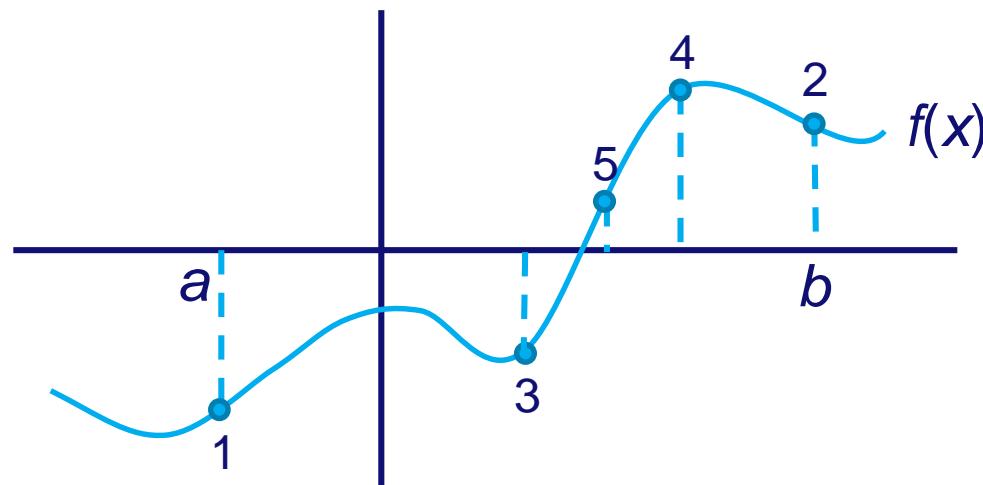
a function to subdivide the interval (x_1, x_2) in n parts and examines whether there is at least one root

Returns the left and right boundaries of the intervals of the roots in xb1, xb2

Bisection method

- **Bisection algorithm:**

- Over some interval it is known that the function will pass through zero, because the function changes sign
- Evaluate function value at the interval's midpoint and examine its sign
- Use the midpoint to replace whichever limit has the same sign



It cannot fail,
but relatively
slow convergence!

Bisection

Exercise 2:

- Write a function in Excel to find a root of a function using the bisection method
 - Assume that an initial bracketing interval (x_1, x_2) is provided
 - Also the required tolerance is specified (which tolerance?)
 - Also output the required number of iterations
- Do the same in MATLAB

Bisection method

- Exercise 2: Bisection method in Excel

it	x1	x2	f1	f2	xmid	fmid	interval size
0	-2	2	14	-2	0	2	4
1	0	2	2	-2	1	-1	2
2	0	1	-1	-1	0.5	-0.25	1
3	0.5	1	0.25	-1	0.75	-0.4375	0.5
4	0.5	1	0.25	-1	0.75	-0.4375	0.25
5	0.5	1	0.25	-1	0.75	-0.4375	0.125
6	0.5	1	0.066406	-1	0.5	-0.578125	0.0625
7	0.5625	0.59375	0.066406	-0.02246	0.5625	-0.03125	0.03125
8	0.578125	0.59375	0.021729	-0.02246	0.578125	-0.02246	5625
9	0.578125	0.585938	0.021729	-0.00043	0.578125	-0.00043	7813
10	0.582031	0.585938	0.010635	-0.00043	0.582031	-0.00043	3906
11	0.583984	0.585938	0.0051	-0.00043	0.583984	-0.00043	1953
12	0.584961	0.585938	0.002336	-0.00043	0.584961	-0.00043	0.000977
13	0.585449	0.585938	0.000954	-0.00043	0.585449	0.000954	0.000488
14	0.585693	0.585938	0.000263	-0.00043	0.585693	0.000263	0.000244
15	0.585693	0.585815	0.000263	-8.2E-05	0.585693	-8.2E-05	0.000122
16	0.585754	0.585815	9.06E-05	-8.2E-05	0.585754	9.06E-05	6.1E-05
17	0.585785	0.585815	4.31E-06	-8.2E-05	0.585785	4.31E-06	3.05E-05
18	0.585785	0.5858	4.31E-06	-3.9E-05	0.585785	-3.9E-05	1.53E-05
19	0.585785	0.585793	4.31E-06	-1.7E-05	0.585785	-1.7E-05	7.63E-06
20	0.585785	0.585789	4.31E-06	-6.5E-06	0.585785	-6.5E-06	3.81E-06
21	0.585785	0.585787	4.31E-06	-1.1E-06	0.585785	-1.1E-06	1.91E-06
22	0.585786	0.585787	1.62E-06	-1.1E-06	0.585786	1.62E-06	9.54E-07
23	0.585786	0.585787	2.69E-07	-1.1E-06	0.585786	2.69E-07	4.77E-07
24	0.585786	0.585787	2.69E-07	-4.1E-07	0.585786	-4.1E-07	2.38E-07
25	0.585786	0.585786	2.69E-07	-6.8E-08	0.585786	1E-07	1.19E-07
26	0.585786	0.585786	1E-07	-6.8E-08	0.585786	1.58E-08	5.96E-08

=IF(f1*fmid<0;x1;xmid)

=IF(f2*fmid<0;x2;xmid)

$$\begin{aligned}x_{\text{mid}} &= 0.5 * (x_1 + x_2) \\f_{\text{mid}} &= f(x_{\text{mid}})\end{aligned}$$

Bisection method

- Exercise 2: Bisection method in MATLAB

```
1 function [p] = bisection(f, x1, x2, tol_step, tol_func)
2     f1 = f(x1);
3     f2 = f(x2);
4     fp = f2;
5     if (f1*f2>0)
6         error('Root must be bracketed!');
7     else
8         it = 0;
9         while ((abs(fp)>tol_func) && (abs(x2 - x1)>tol_step))
10            it = it + 1;
11            p = 0.5*(x1 + x2);
12            fp = f(p);
13            if (f1*fp<0)
14                x2 = p;
15                f2 = fp;
16            else
17                x1 = p;
18                f1 = fp;
19            end
20        end
21        disp(sprintf('Root found in %d iterations at x = %e\n (function value = %e)', [it,p,fp]));
22    end
23 end
```

```
>> bisection(@(x) x^2-4*x+2,0,2,1e-7,1e-7);
```

Note1: We have used a criterion for the function value and the step size!

Note2: usually while loop needs protection for maximum number of iterations
(but here bisection is sure to convergence...)

Root found in 24 iterations required.
Can we do better?

Bisection method

- **Required number of iterations?**

- After each iteration the interval bounds containing the root decrease by a factor of 2:

$$\epsilon_{n+1} = \frac{1}{2} \epsilon_n \quad \Rightarrow \quad n = \log_2 \frac{\epsilon_0}{tol}$$

ϵ_0 = initial bracketing interval
 tol = desired tolerance

i.e. after 50 iterations the interval is decreased by factor $2^{50} = 10^{15}!$
(Mind machine accuracy when setting tolerance!)

- Order of convergence = 1

$$\epsilon_{n+1} = K(\epsilon_n)^m$$

$m = 1$: linear convergence
 $m = 2$: quadratic convergence

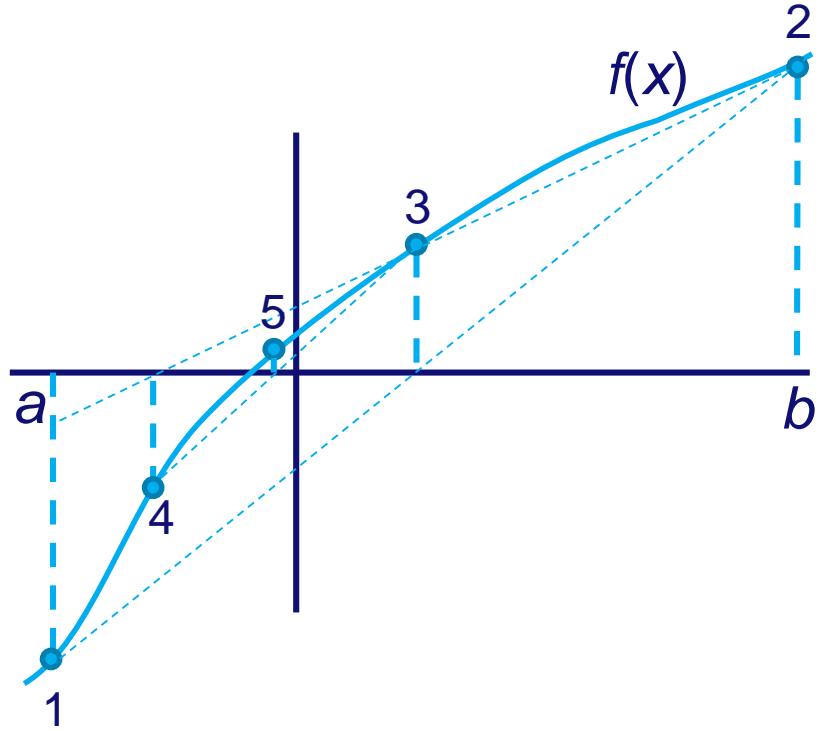
- Must succeed:
 - More than root \Rightarrow bisection will find one of them
 - No root, but singularity \Rightarrow bisection will find singularity

Secant and False position method

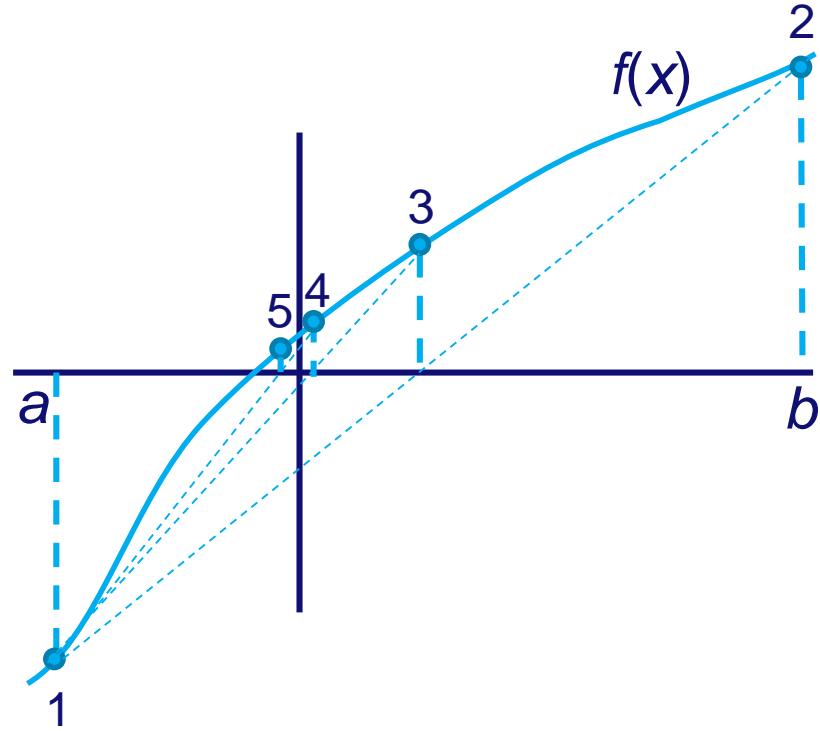
- **Secant/False position (= Regula Falsi) method**
 - Faster convergence (provided sufficiently smooth behaviour)
 - Difference with bisection method in choice of next point:
 - Bisection: mid-point of interval
 - Secant/False position: point where the approximating line crosses the axis
 - One of the boundary points is discarded in favor of the latest estimate of
 - Secant: retains the most recent of the prior estimates
 - False position: retains prior estimate with opposite sign, so that the points continue to bracket the root

Secant and False position method

Secant method



False position method



Secant: slightly faster convergence: $\lim_{n \rightarrow \infty} |\epsilon_{n+1}| = K|\epsilon_n|^{1.618}$

False position: guaranteed convergence

Secant and False position method

Exercise 3:

- Write a function in Excel and MATLAB to find a root of a function using the Secant and the False position methods
 - Assume that an initial bracketing interval (x_1, x_2) is provided
 - Also the required tolerance is specified
 - Also output the required number of iterations
 - Compare the bisection, false position and secant methods

Secant and False position method

Exercise 3:

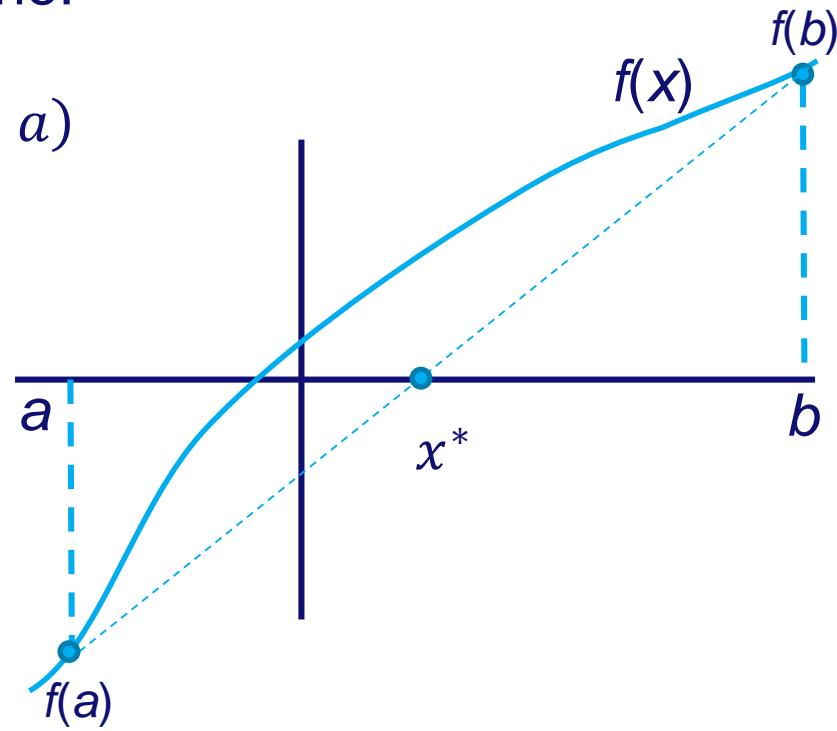
- **Determination of the abscissa of the approximating line:**
 - Determine the approximating line:

$$f(x) \approx f(a) + \frac{f(b) - f(a)}{b - a} (x - a)$$

- Determine abscissa:

$$f(x^*) = 0$$

$$\Rightarrow x^* = a - \frac{f(a)(b - a)}{f(b) - f(a)}$$
$$= \frac{af(b) - bf(a)}{f(b) - f(a)}$$



Secant and False position method

- Exercise 3: False position method in Excel

it	x1	x2	f1	f2	x absc	f absc	interval size
0	-2	2	14	-2	1.5	-1.75	4
1	-2	1.5	14	-1.75	1.111111	-1.20988	0.388889
2	-2	1.111111	14	-1.20988	0.863636	-0.70868	0.247475
3	-2	0.863636	14	-0.70868	0.725664	-0.37607	0.137973
4	-2	0.725664	14	-0.37607	0.654362	-0.18926	0.071301
5	-2	0.654362	14	-0.18926	0.618958	-0.09272	0.035404
6	-2	0.618958	14	-0.09272	0.601727	-0.04483	0.017231
7	-2	0.601727	14	-0.04483	0.593422	-0.02154	0.008305
8	-2	0.593422	14	-0.02154	0.589438	-0.01032	0.003984
9	-2	0.589438	14	-0.01032	0.587532	-0.00493	0.001907
10	-2	0.587532	14	-0.00493	0.586662	-0.00236	0.000911
11	-2	0.586662	14	-0.00236	0.586185	-0.00113	0.000436
12	-2	0.586185	14	-0.00113	0.586977	-0.00054	0.000208

=IF(f1*fabs<0;
x1;xabsc)

=IF(f2*fabs<0;
x2;xabsc)

$x \text{ absc} = x_1 - f_1 * (x_2 - x_1) / (f_2 - f_1)$
 $f \text{ absc} = f(x \text{ absc})$

Secant and False position method

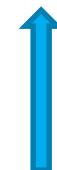
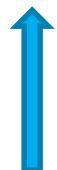
- Exercise 3: Secant method in Excel

it	x1	x2	f1	f2	x absc	f absc	interval size
0	-2	2	14	-2	1.5	-1.75	4
1	-2	1.5	14	-1.75	1.111111	-1.20988	3.111111
2	1.111111	1.5	-1.20988	-1.75	0.24	1.0976	0.388889
3	0.24	1.111111	1.0976	-1.20988	0.654362	-0.18926	0.871111
4	0.24	0.654362	1.0976	-0.18926	0.593422	-0.02154	0.414362
5	0.593422	0.654362	-0.02154	-0.18926	0.585596	0.000538	0.060941
6	0.585596	0.593422	0.000538	-0.02154	0.585787	-1.5E-06	0.007826
7	0.585596	0.585787	0.000538	-1.5E-06	0.585786	-9.8E-11	0.000191
8	0.585786	0.585787	-9.8E-11	-1.5E-06	0.585786	0	5.15E-07
9	0.585786	0.585786	0	-9.8E-11	0.585786	0	3.46E-11



$$= \min(x_{it-2}, x_{it-1})$$

$$= \max(x_{it-2}, x_{it-1})$$



$$x \text{ absc} = x_1 - f_1 * (x_2 - x_1) / (f_2 - f_1)$$
$$f \text{ absc} = f(x \text{ absc})$$

Secant and False position method

- Exercise 3: False position method in MATLAB

```
1 function [p] = falseposition(f, x1, x2, tol_step, tol_func)
2 -     f1 = f(x1);
3 -     f2 = f(x2);
4 -     fp = f2;
5 -     if (f1*f2>0)
6 -         error('Root must be bracketed!');
7 -     else
8 -         it = 1;
9 -         while ((abs(fp)>tol_func) && (abs(x2 - x1)>tol_step))
10 -             it = it + 1;
11 -             p = (x1*f2 - x2*f1) / (f2 - f1); ← The only difference with bisection!
12 -             fp = f(p);
13 -             if (f1*fp<0)
14 -                 x2 = p;
15 -                 f2 = fp;
16 -             else
17 -                 x1 = p;
18 -                 f1 = fp;
19 -             end
20 -         end
21 -         disp(sprintf('Root found in %d iterations at x = %e\n (function value = %e)', [it,p,fp]));
22 -     end
23 - end
```

>> falseposition(@(x) x^2-4*x+2,0,2,1e-7,1e-7);

Root found in 12 iterations!
(Bisection needed 24 iterations)

Secant and False position method

• Exercise 3: Secant method in MATLAB

```
1 function [p] = secant(f, x1, x2, tol_step, tol_func)
2 -     f1 = f(x1);
3 -     f2 = f(x2);
4 -     fp = f2;
5 -     if (f1*f2>0)
6 -         error('Root must be bracketed!');
7 -     else
8 -         it = 1;
9 -         while ((abs(fp)>tol_func) && (abs(x2 - x1)>tol_step))
10 -             it = it + 1;
11 -             p = (x1*f2 - x2*f1)/(f2 - f1);
12 -             fp = f(p);
13 -             x1 = x2;
14 -             f1 = f2;
15 -             x2 = p;
16 -             f2 = fp;
17 -         end
18 -         disp(sprintf('Root found in %d iterations at x = %e\n (function value = %e)', [it,p,fp]));
19 -     end
20 - end
```



The only difference with
False position method!

```
>> secant(@(x) x^2-4*x+2,0,2,1e-7,1e-7);
```

Secant method: 8 iterations
False position: 12 iterations
Bisection: 24 iterations

Secant and False position method

- **Comparison of methods**

$$f(x) = x^2 - 4x + 2 = 0$$

tol_eps, tol_func = 1e-15, and $(x_1, x_2) = (0,2)$

Method	Nr. iterations
Bisection	51
False position	22
Secant	9

Compare with:

```
>> fzero(@(x) x^2-4*x+2,2,optimset('TolX',1e-15,'Display','iter'))
```

Note the initial bracketing steps in fzero!

Brent's method

- **Superlinear convergence + sureness of bisection**
 - Keep track of superlinear convergence, and if not, intersperse with bisection steps (assures at least linear convergence)
 - Brent's method (is implemented in MATLAB's fzero): root-bracketing + bisection + inverse quadratic interpolation
 - Inverse quadratic interpolation: uses 3 prior points to fit an inverse quadratic function (i.e. $x(y)$) with contingency plans, if root falls outside brackets:
$$x = b + P/Q \quad R = f(b)/f(c)$$
$$P = S[T(R - T)(c - b) - (1 - R)(b - a)] \quad S = f(b)/f(a)$$
$$Q = (T - 1)(R - 1)(S - 1) \quad T = f(a)/f(c)$$
$$b = \text{current best estimate}$$
$$P/Q = \text{ought to be a 'small' correction}$$
 - When P/Q does not land within the bounds or when bounds are not collapsing fast enough \Rightarrow take bisection step

Brent's method

```
1  function [root] = brent(f, x1, x2, tol)
2  ITMAX = 100;
3  EPS = 3e-8;
4  a = x1; b = x2; c = x2;
5  fa = f(a);
6  fb = f(b);
7  fc = fb;
8  if (fa*fb>0)
9    error('Root must be bracketed!');
10 else
11   for iter=1:ITMAX
12     if (fb*fc>0)
13       c = a; fc = fa;      % Rename a, b, c and
14       d = b - a; e = d;    % adjust bounding interval d
15     end;
16     if (abs(fc)<abs(fb))
17       a = b; fa = fb;
18       b = c; fb = fc;
19       c = a; fc = fa;
20     end;
21     tol1 = 2.0*EPS*abs(b) + 0.5*tol; % Convergence check.
22     xm = 0.5*(c - b);
23     if ((abs(xm)<=tol1) || (fb == 0)) |
24       root = b;
25       disp(sprintf('\nRoot found in %d iterations at x = %e (f(x) = %e)', [iter,b,fb]));
26       break;
27     end;
28     if ((abs(e)>=tol1) && (abs(fa)>abs(fb)))
29       % Attempt inverse quadratic interpolation.
30       s = fb/fa;
31       if (a==c)
32         p = 2.0*xm*s;
33         q = 1.0 - s;
34       else
35         q = fa/fc;
36         r = fb/fc;
37         p = s*(2.0*xm*q*(q - r) - (b - a)*(r - 1.0));
38         q = (q - 1.0)*(r - 1.0)*(s - 1.0);
39       end;
```

Brent's method

```
40 -         if (p>0.0)
41 -             q = -q; % Check whether in bounds.
42 -         end;
43 -         p = abs(p);
44 -         min1 = 3.0*xm*q - abs(tol1*q);
45 -         min2 = abs(e*q);
46 -         if (2.0*p<min(min1,min2))
47 -             e = d; % Accept interpolation.
48 -             d = p/q;
49 -         else
50 -             d = xm; % Interpolation failed, use bisection.
51 -             e = d;
52 -         end;
53 -     else
54 -         d = xm; % Bounds decreasing too slowly, use bisection.
55 -         e = d;
56 -     end;
57 -     a = b; % Move last best guess to a.
58 -     fa = fb;
59 -     if (abs(d)>tol1) % Evaluate new trial root.
60 -         b = b + d;
61 -     else
62 -         if (xm<0)
63 -             b = b - tol1;
64 -         else
65 -             b = b + tol1;
66 -         end;
67 -     end;
68 -     fb = f(b);
69 -     if (d == xm)
70 -         disp(sprintf('Iteration: %d => x = %e, f(x) = %e (bisection)', [iter,b,fb]));
71 -     else
72 -         disp(sprintf('Iteration: %d => x = %e, f(x) = %e (inverse quadratic interpolation)', [iter,b,fb]));
73 -     end;
74 - end;
75 - if (iter==ITMAX)
76 -     disp('Maximum number of iterations exceeded in brent!');
77 - end;
78 - end;
79 - end
```

Newton-Raphson method

- Requires the evaluation of the function $f(x)$ and the derivative $f'(x)$ at arbitrary points

- Algorithm:

- Extend tangent line at current point x_i till it crosses zero
 - Set next guess x_{i+1} to the abscissa of that zero crossing

$$f(x + \delta) \approx f(x) + f'(x)\delta + \frac{1}{2}f''\delta^2 + \dots \quad (\text{Taylor series at } x)$$

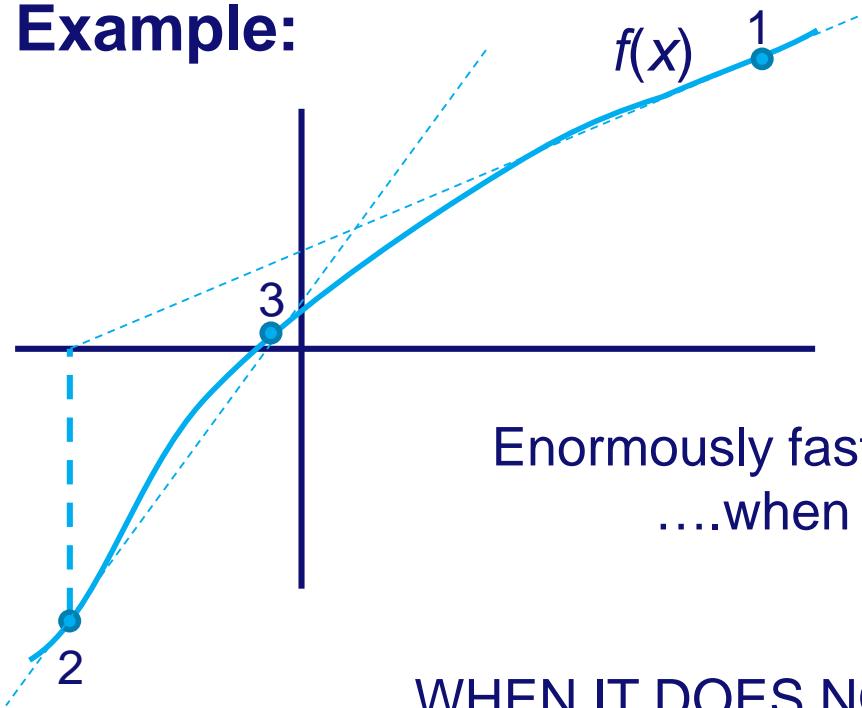
For small enough values of δ and for well-behaved functions, the non-linear terms become unimportant

$$\Rightarrow \delta = -\frac{f(x)}{f'(x)}$$

- Can be extended to higher dimensions
- Requires an initial guess sufficiently close to the root! (otherwise even failure!!)

Newton-Raphson method

- Example:

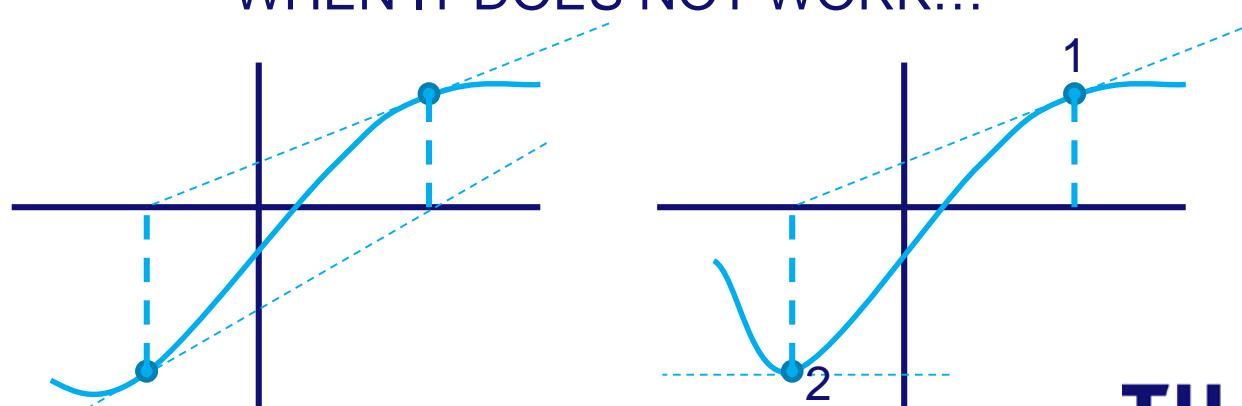


Newton-Raphson method:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Enormously fast convergence,
....when it works

WHEN IT DOES NOT WORK...



Newton-Raphson method

- **Basic algorithm:**

Given initial x , required tolerance $\varepsilon > 0$

Repeat

1. Compute $f(x)$ and $f'(x)$.
2. If $|f(x)| \leq \varepsilon$, return x
3. $x := x - f(x)/f'(x)$

until maximum number of iterations is exceeded

Newton-Raphson method

- Why is Newton-Raphson so powerful?
⇒ High rate of convergence

Newton-Raphson method:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Subtracting the solution x^* :

$$x_{n+1} - x^* = x_n - x^* - \frac{f(x_n)}{f'(x_n)}$$

Defining the error $\epsilon_n = x_n - x^*$: $\epsilon_{n+1} = \epsilon_n - \frac{f(x_n)}{f'(x_n)}$

$$\epsilon_{n+1} = \epsilon_n - \frac{f(x^*) + f'(x^*)\epsilon_n + \frac{1}{2}f''(x^*)\epsilon_n^2 + \dots}{f'(x^*) + \dots}$$

$$\epsilon_{n+1} = \epsilon_n - \epsilon_n - \frac{1}{2} \frac{f''(x^*)}{f'(x^*)} \epsilon_n^2 \Rightarrow$$

$\epsilon_{n+1} \sim K \epsilon_n^2$
Quadratic convergence!!

Newton-Raphson method

- Order of convergence

$$\lim_{n \rightarrow \infty} \frac{|\epsilon_{n+1}|}{|\epsilon_n|^m} = K \quad m = \text{order of convergence}$$
$$K = \text{asymptotic error constant}$$

$$\epsilon_n = x_n - x^* \quad \text{with } x^* \text{ the solution}$$

When the solution is not known a priori: $\epsilon_{n+1} \approx x_{n+1} - x_n$

$$\frac{|\epsilon_{n+1}|}{|\epsilon_n|} = \frac{K|\epsilon_n|^m}{K|\epsilon_{n-1}|^m} \Rightarrow \frac{|\epsilon_{n+1}|}{|\epsilon_n|} = \left(\frac{|\epsilon_n|}{|\epsilon_{n-1}|} \right)^m$$

$$\Rightarrow \ln \left(\frac{|\epsilon_{n+1}|}{|\epsilon_n|} \right) = m \ln \left(\frac{|\epsilon_n|}{|\epsilon_{n-1}|} \right)$$

$$m = \frac{\ln \left(\frac{|\epsilon_{n+1}|}{|\epsilon_n|} \right)}{\ln \left(\frac{|\epsilon_n|}{|\epsilon_{n-1}|} \right)}$$

for $n \rightarrow \infty$

Newton-Raphson method

Exercise 4:

- Write a function in MATLAB to find a root of a function using the Newton-Raphson method
 - Assume that an initial guess x_0 is provided
 - Also the required tolerance is given
 - Output the results for every iteration
 - Verify that at every iteration the number of significant digits double, and compute the order of convergence

Newton-Raphson method

Exercise 4: Newton-Raphson in MATLAB

```
1 function [p] = newton1D(func, grad, x, tol_x, tol_f)
2 %
3 % ITMAX = 100;
4 % error = 2*tol_f;
5 % it = 0;
6 % f = func(x);
7 % while (((error>tol_f) || (dx>tol_x)) && (it<ITMAX))
8 %     it = it + 1;
9 %     g = grad(x);
10 %     dx = -f/g;
11 %     x = x + dx;
12 %     f = func(x);
13 %     error = abs(f);
14 % end;
15 % if it<=ITMAX
16 %     disp(sprintf('Root found in %d iterations at x = %e\n (function value = %e)', [it,x,f]));
17 % else
18 %     disp(sprintf('No root found after %d iterations!', [it]));
19 % end
```

```
>> newton1D(@(x) x^2-4*x+2, @(x) 2*x-4, 1, 1e-12, 1e-12)
```

Convergence in 6 iterations.

Why does it not work with an initial guess of $x_0 = 2???$

Newton-Raphson method

- **Modifications to the basic algorithm**

- If the first derivative $f'(x)$ is not known or cumbersome to compute/program, we can use the local num. approximation:

$$f'(x) \approx \frac{f(x + dx) - f(x)}{dx} \quad (dx \sim 10^{-8})$$

dx should be small (otherwise the method reduces to first order)

But not too small (otherwise you will be wiped out by roundoff!)

- Unless you know that the initial guess is close to the solution, the Newton-Raphson method should be combined with:
 - a bracketing method, to reject the solution if it wanders outside of the bounds;
 - Reduced Newton step method (= relaxation) for more robustness. Don't take the entire step if the error does not decrease (enough)
 - More sophisticated step size control: Local line searches and backtracking using cubic interpolation (for global convergence)

Content

- **How to solve:**

$f(x) = 0$ for arbitrary functions f

“Root finding”

(i.e. move all terms to the left)

- **One dimensional case:** $f(x) = 0$
“Bracket or ‘trap’ a root between bracketing values,
then hunt it down like a rabbit.”
- **Multi-dimensional case:** $f(x) = 0$
 - N equations in N unknowns:
You can only *hope* to find a solution.
It may have no (real) solution, or more than one solution!
 - Much more difficult!!
“You never know whether a root is near,
unless you have found it”

Newton-Raphson method

- **Extensions to multi-dimensional case:**

Let's first consider the two-dimensional case:

$$f(x, y) = 0$$

$$g(x, y) = 0$$

Multi-variate Taylor series expansion:

$$f(x + \delta x, y + \delta y) \approx f(x, y) + \frac{\partial f}{\partial x} \delta x + \frac{\partial f}{\partial y} \delta y + O(\delta x^2, \delta y^2) = 0$$

$$g(x + \delta x, y + \delta y) \approx g(x, y) + \frac{\partial g}{\partial x} \delta x + \frac{\partial g}{\partial y} \delta y + O(\delta x^2, \delta y^2) = 0$$

Neglecting higher order terms:

$$\frac{\partial f}{\partial x} \delta x + \frac{\partial f}{\partial y} \delta y = -f(x, y)$$

$$\frac{\partial g}{\partial x} \delta x + \frac{\partial g}{\partial y} \delta y = -g(x, y)$$



Two linear equations in the
two unknowns δx and δy .

Newton-Raphson method

- **Extensions to multi-dimensional case:**

Newton-Raphson method:

$$\frac{\partial f}{\partial x} \delta x + \frac{\partial f}{\partial y} \delta y = -f(x, y)$$

$$\frac{\partial g}{\partial x} \delta x + \frac{\partial g}{\partial y} \delta y = -g(x, y)$$

Or in matrix notation:

$$\begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \\ \frac{\partial g}{\partial x} & \frac{\partial g}{\partial y} \end{bmatrix} \cdot \begin{bmatrix} \delta x \\ \delta y \end{bmatrix} = - \begin{bmatrix} f(x, y) \\ g(x, y) \end{bmatrix}$$



Jacobian matrix

Solution via Cramer's rule:

$$\delta x = \frac{\begin{vmatrix} -f & \frac{\partial f}{\partial y} \\ -g & \frac{\partial g}{\partial y} \end{vmatrix}}{\begin{vmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \\ \frac{\partial g}{\partial x} & \frac{\partial g}{\partial y} \end{vmatrix}} = \frac{-f \frac{\partial g}{\partial y} + g \frac{\partial f}{\partial y}}{\frac{\partial f}{\partial x} \frac{\partial g}{\partial y} - \frac{\partial f}{\partial y} \frac{\partial g}{\partial x}}$$

$$\delta y = \frac{\begin{vmatrix} \frac{\partial f}{\partial x} & -f \\ \frac{\partial g}{\partial x} & -g \end{vmatrix}}{\begin{vmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \\ \frac{\partial g}{\partial x} & \frac{\partial g}{\partial y} \end{vmatrix}} = \frac{-g \frac{\partial f}{\partial x} + f \frac{\partial g}{\partial x}}{\frac{\partial f}{\partial x} \frac{\partial g}{\partial y} - \frac{\partial f}{\partial y} \frac{\partial g}{\partial x}}$$

Newton-Raphson method

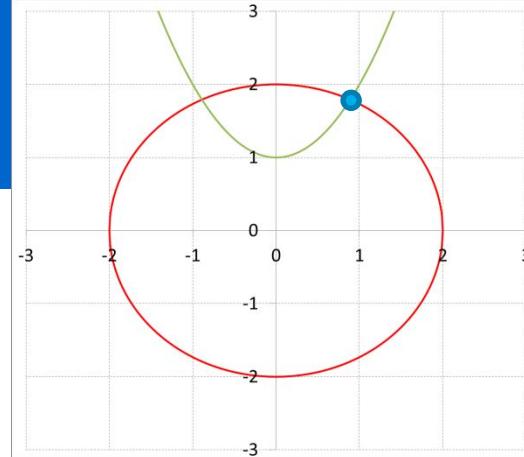
- **Extensions to multi-dimensional case:**

Example: intersection of circle with parabola:

$$\begin{aligned}x^2 + y^2 = 4 \\y = x^2 + 1\end{aligned}$$

In matrix form:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mathbf{f} = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \begin{bmatrix} x_1^2 + x_2^2 - 4 \\ x_1^2 - x_2 + 1 \end{bmatrix} \quad \mathbf{J} = \begin{bmatrix} 2x_1 & 2x_2 \\ 2x_1 & -1 \end{bmatrix}$$

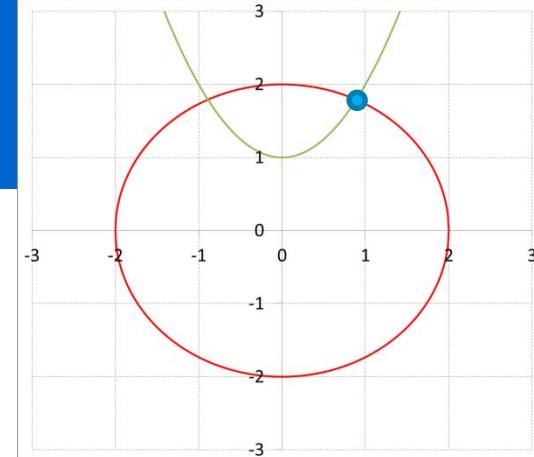


	$\mathbf{x}^{(i)}$	$\mathbf{f}^{(i)}$	$\mathbf{J}^{(i)}$	$\delta\mathbf{x}^{(i)}$
$i = 1:$	$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 2 & 4 \\ 2 & -1 \end{bmatrix}$	$\begin{bmatrix} -0.1 \\ -0.2 \end{bmatrix}$
$i = 2:$	$\begin{bmatrix} 0.9 \\ 1.8 \end{bmatrix}$	$\begin{bmatrix} 0.05 \\ 0.01 \end{bmatrix}$	$\begin{bmatrix} 1.8 & 3.6 \\ 1.8 & -1 \end{bmatrix}$	$\begin{bmatrix} -0.01039 \\ -0.0087 \end{bmatrix}$
$i = 3:$	$\begin{bmatrix} 0.889614 \\ 1.791304 \end{bmatrix}$	$\begin{bmatrix} 0.000183 \\ 0.0000108 \end{bmatrix}$	$\begin{bmatrix} 1.7792 & 3.5826 \\ 1.7792 & -1 \end{bmatrix}$	$\begin{bmatrix} -6.99 \cdot 10^{-5} \\ -1.65 \cdot 10^{-5} \end{bmatrix}$
$i = 4:$	$\begin{bmatrix} 0.8895436 \\ 1.7912878 \end{bmatrix}$	$\begin{bmatrix} 5.16 \cdot 10^{-9} \\ 4.89 \cdot 10^{-9} \end{bmatrix}$	$\begin{bmatrix} 1.779087 & 3.582576 \\ 1.779087 & -1 \end{bmatrix}$	$\begin{bmatrix} -2.78 \cdot 10^{-9} \\ -5.94 \cdot 10^{-11} \end{bmatrix}$

Newton-Raphson method

- **Extensions to multi-dimensional case:**

Example: intersection of circle with parabola:



Check order of convergence:

it	x1	x2	eps1	eps2	m1	m2
1	1.0000000000000000	2.0000000000000000				
2	0.9000000000000000	1.8000000000000000	0.1000000000000000	0.2000000000000000		
3	0.8896135265700480	1.7913043478260900	0.0103864734299518	0.0086956521739132	1.983532	2.948192
4	0.8895436203043770	1.7912878475373300	0.0000699062656710	0.0000165002887549	2.094992	2.32082
5	0.8895436175241320	1.7912878474779200	0.0000000027802448	0.0000000000594120	2.058946	2.138235

Quadratic convergence!
= doubling number of significant digits every iteration

$$\epsilon_{n+1} \approx x_{n+1} - x_n$$

$$m = \frac{\ln\left(\frac{|\epsilon_{n+1}|}{|\epsilon_n|}\right)}{\ln\left(\frac{|\epsilon_n|}{|\epsilon_{n-1}|}\right)}$$

Newton-Raphson method

- **Extensions to multi-dimensional case:**

Generalization to the N -dimensional case:

$$f_i(x_1, x_2, \dots, x_N) = 0 \quad \text{for } i = 1, 2, \dots, N$$

Define: $\mathbf{x} = [x_1, x_2, \dots, x_N]$ and $\mathbf{f} = [f_1, f_2, \dots, f_N] \Rightarrow \mathbf{f}(\mathbf{x}) = \mathbf{0}$

Multi-variate Taylor series expansion:

$$f_i(\mathbf{x} + \delta\mathbf{x}) = f_i(\mathbf{x}) + \sum_{j=1}^N \frac{\partial f_i}{\partial x_j} \delta x_j + O(\delta\mathbf{x}^2)$$

Define Jacobian matrix:
$$J_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$\Rightarrow \mathbf{f}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{f}(\mathbf{x}) + \mathbf{J} \cdot \delta\mathbf{x} + O(\delta\mathbf{x}^2)$$

Neglect higher order terms:

$$\mathbf{J} \cdot \delta\mathbf{x} = -\mathbf{f}(\mathbf{x})$$
$$\mathbf{x}_{new} = \mathbf{x}_{old} + \delta\mathbf{x}$$

Newton-Raphson method

Multi-variate Newton-Raphson in MATLAB

```
1 function [f] = MyFunc(x)
2 -    f(1) = x(1)^2 + x(2)^2 - 4;
3 -    f(2) = x(1)^2 - x(2) + 1;
4 -    f = f';
5 - end
```

```
1 function [jac] = MyJac(x)
2 -    jac(1,1) = 2*x(1);
3 -    jac(1,2) = 2*x(2);
4 -    jac(2,1) = 2*x(1);
5 -    jac(2,2) = -1;
6 - end
```

```
1 function [p] = newton(func, jac, x, tol_x, tol_f)
2 -    ITMAX = 100;
3 -    error = 2*tol_f;
4 -    it = 0;
5 -    f = feval(func,x);
6 -    while (((error>tol_f) || (max(abs(dx))>tol_x)) && (it<ITMAX))
7 -        it = it + 1;
8 -        j = feval(jac,x);
9 -        dx = j\b(-f);
10 -       x = x + dx;
11 -       f = func(x);
12 -       error = max(abs(f));
13 -       disp(sprintf('iteration %d: x[1] = %e, x[2] = %e with f[1] = %e, f[2] = %e', [it,x(1),x(2),f(1),f(2)]));
14 -    end;
15 -    if it<=ITMAX
16 -        disp(sprintf('\nRoot found in %d iterations at x[1] = %e, x[2] = %e with f[1] = %e; f[2] = %e\n', [it,x(1),x(2),f(1),f(2)]));
17 -    else
18 -        disp(sprintf('\nNo root found after %d iterations!\n', [it]));
19 -    end;
20 - end
```

Solve $\mathbf{A}^{-1} \cdot \mathbf{b}$ simply with “ $\mathbf{A}\backslash\mathbf{b}$ ”
This is the strength of MATLAB!

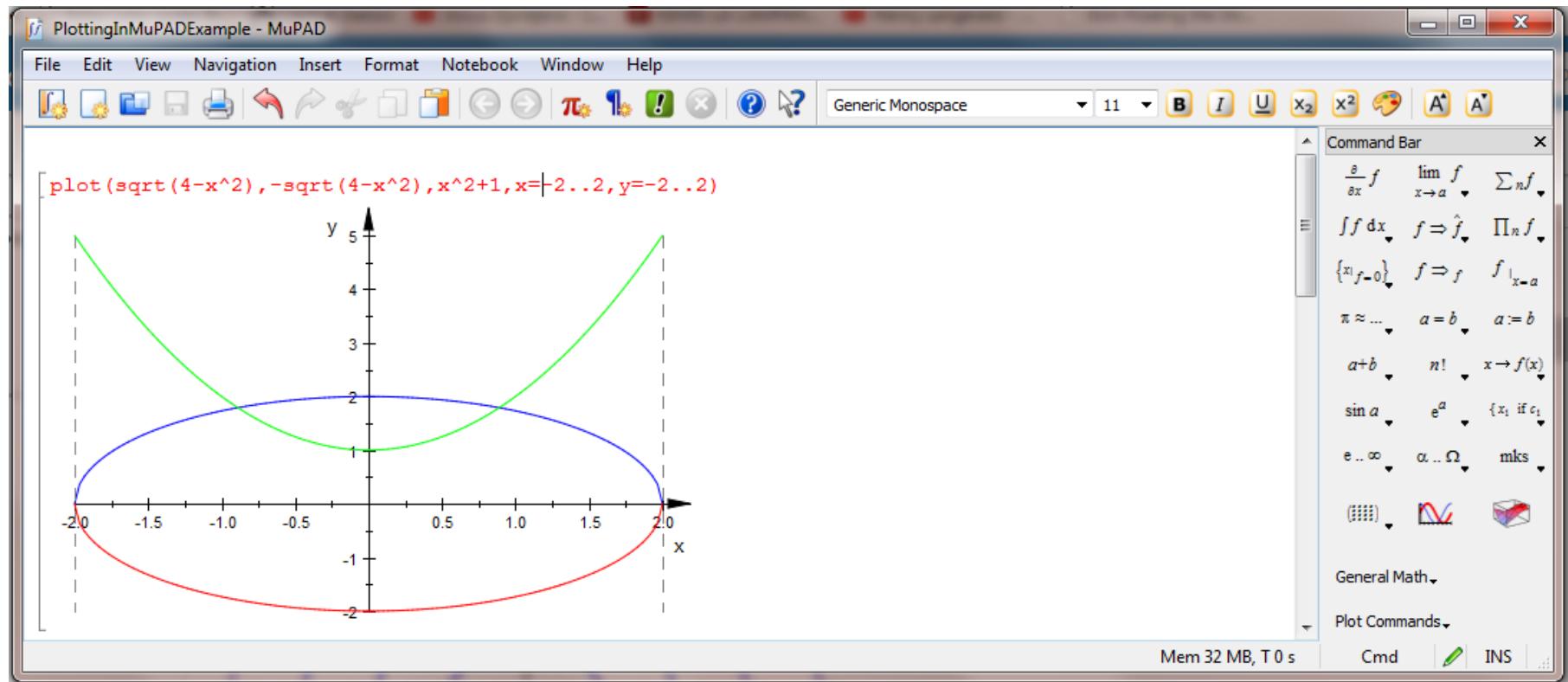
```
>> newton(@MyFunc,@MyJac,[1;2],1e-12,1e-12)
```

Newton-Raphson method

Multi-variate Newton-Raphson in MATLAB

Plotting the functions:

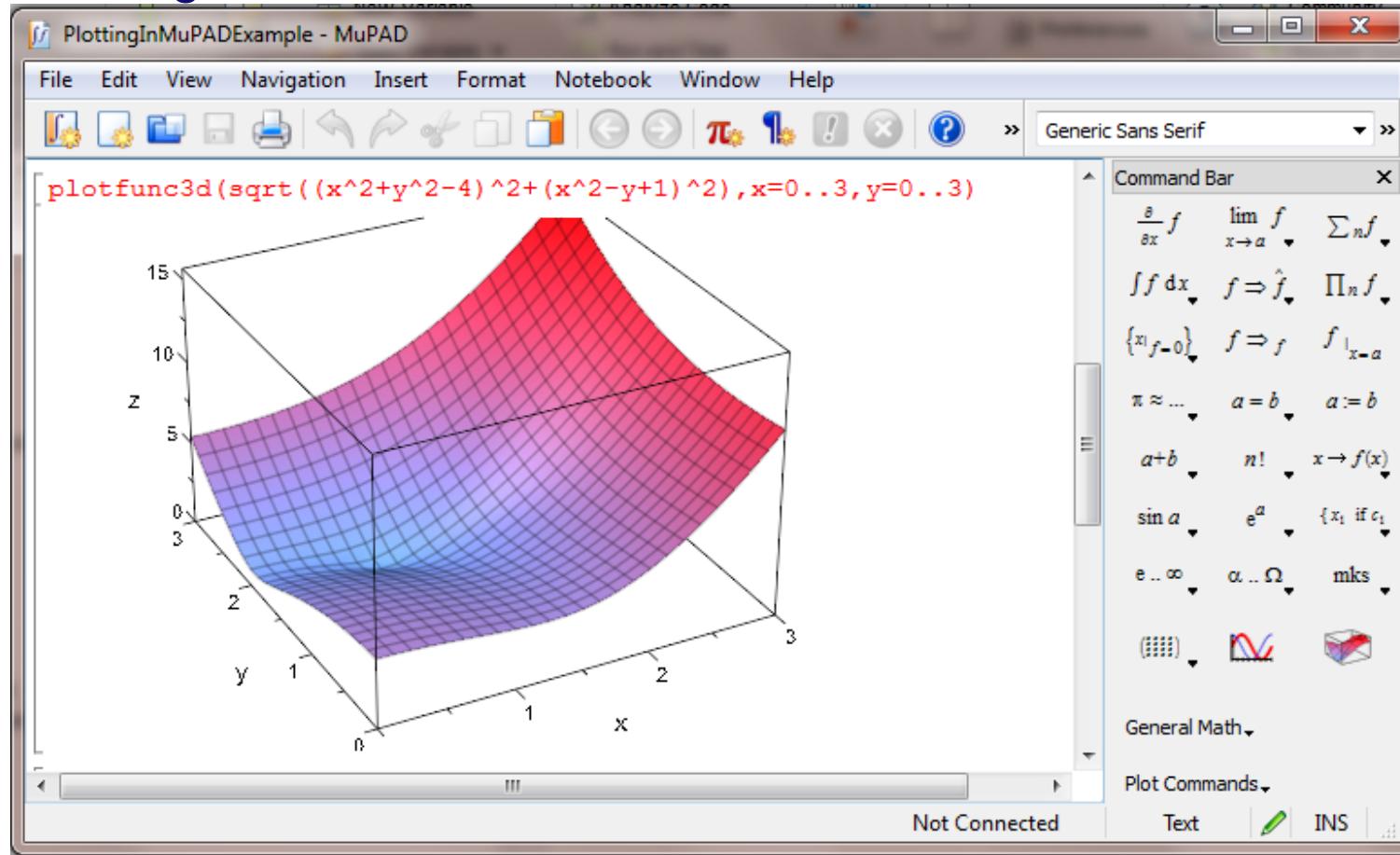
```
>> mphandle = mupad
```



Newton-Raphson method

Multi-variate Newton-Raphson in MATLAB

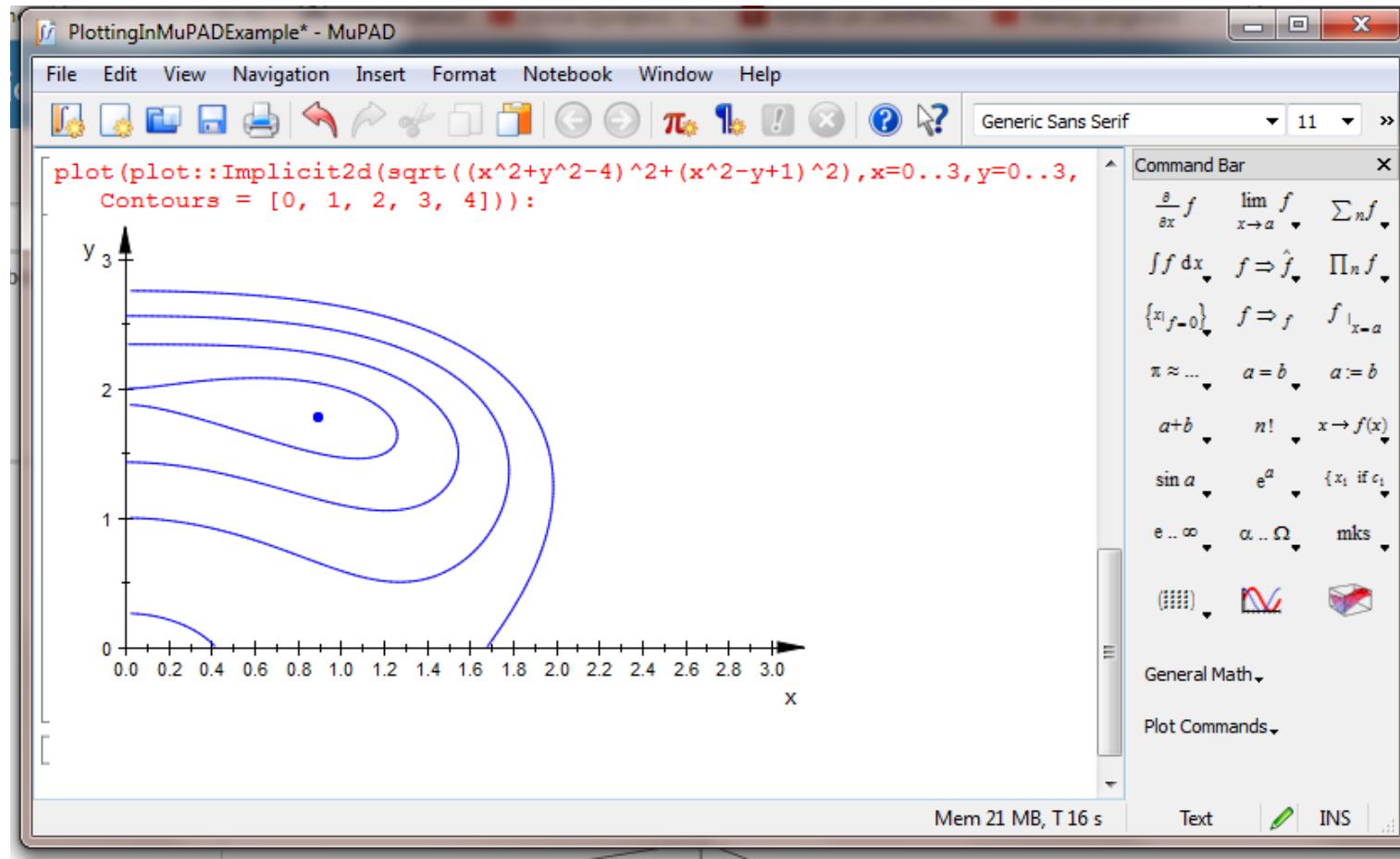
Plotting the norm of the function:



Newton-Raphson method

Multi-variate Newton-Raphson in MATLAB

Plotting contours of the norm of the function:



Broyden's method

- **Multi-dimensional secant method ('quasi-Newton'):**

Disadvantage of the Newton-Raphson method:

It requires the Jacobian matrix

- In many problems no analytical jacobian available
- If the function evaluation is expansive, the numerical approximation using finite differences can be prohibitive!

⇒ use cheap approximation of the Jacobian!

(= secant, or 'quasi-Newton' method)

Newton-Raphson:

$$\mathbf{J}^n \cdot \delta \mathbf{x}^n = -\mathbf{f}^n(\mathbf{x}^n)$$

$$\mathbf{x}^{n+1} = \mathbf{x}^n + \delta \mathbf{x}^n$$

Secant method:

$$\mathbf{B}^n \cdot \delta \mathbf{x}^n = -\mathbf{f}^n(\mathbf{x}^n)$$

$$\mathbf{x}^{n+1} = \mathbf{x}^n + \delta \mathbf{x}^n$$

\mathbf{B}^n = approximation
of the Jacobian

Broyden's method

- **Multi-dimensional secant method ('quasi-Newton'):**

Secant equation (generalisation of 1D case):

$$\mathbf{B}^{n+1} \cdot \delta \mathbf{x}^n = \delta \mathbf{f}^n \quad \delta \mathbf{x}^n = \mathbf{x}^{n+1} - \mathbf{x}^n \quad \delta \mathbf{f}^n = \mathbf{f}^{n+1} - \mathbf{f}^n$$

Underdetermined (i.e. not unique: n equations with n^2 unknowns)
⇒ we need another condition to pin down \mathbf{B}^{n+1}

Broyden's method: determine \mathbf{B}^{n+1} by making the least change to \mathbf{B}^n that is consistent with the secant condition

Updating formula:

$$\mathbf{B}^{n+1} = \mathbf{B}^n + \frac{(\delta \mathbf{f}^n - \mathbf{B}^n \cdot \delta \mathbf{x}^n)}{\delta \mathbf{x}^n \cdot \delta \mathbf{x}^n} \otimes \delta \mathbf{x}^n$$

(Note: sometimes \mathbf{B}^{-1} is updated directly)

Broyden's method

- Multi-dimensional secant method ('quasi-Newton'):

Background of Broyden's method:

Secant equation: $\mathbf{B}^{n+1} \cdot \delta \mathbf{x}^n = \delta \mathbf{f}^n$

Broyden's method: Since there is no update on derivative info, why would \mathbf{B}^n change in a direction \mathbf{w} orthogonal to $\delta \mathbf{x}^n$

$$\Rightarrow (\delta \mathbf{x}^n)^T \cdot \mathbf{w} = 0$$

$$\left. \begin{array}{l} \mathbf{B}^{n+1} \cdot \mathbf{w} = \mathbf{B}^n \cdot \mathbf{w} \\ \mathbf{B}^{n+1} \cdot \delta \mathbf{x}^n = \delta \mathbf{f}^n \end{array} \right\} \Rightarrow \boxed{\mathbf{B}^{n+1} = \mathbf{B}^n + \frac{(\delta \mathbf{f}^n - \mathbf{B}^n \cdot \delta \mathbf{x}^n)}{\delta \mathbf{x}^n \cdot \delta \mathbf{x}^n} \otimes \delta \mathbf{x}^n}$$

Initialize \mathbf{B}^0 with identity matrix (or with finite difference approx.)

Broyden's method

- Same example as before but now with Broyden's method

```
function [p] = broyden(func, x, tol_x, tol_f)
    ITMAX = 100;
    error = 2*tol_f;
    it = 0;
    f = feval(func,x);
    b = eye(2); % create identity matrix
    while (((error>tol_f) || (max(abs(dx))>tol_x)) && (it<ITMAX))
        it = it + 1;
        dx = b\(-f);
        x = x + dx;
        f0 = f;
        f = func(x);
        df = f - f0;
        b = b + ((df - b*dx)*dx.') / (dx.*dx); % Broyden's updating
        error = max(abs(f));
        disp(sprintf('iteration %d: x[1] = %e, x[2] = %e with f[1] = %e, f[2] = %e', [it,x(1),x(2),f(1),f(2)]));
    end;
    if it<=ITMAX
        disp(sprintf('\nRoot found in %d iterations at x[1] = %e, x[2] = %e with f[1] = %e; f[2] = %e\n', [it,x(1),x(2),f(1),f(2)]));
    else
        disp(sprintf('\nNo root found after %d iterations!\n', [it]));
    end;
end
```

```
>> broyden(@MyFunc,[1;2],1e-12,1e-12)
```

Requires 12 iterations (compare with Newton: 5 iterations)

But much fewer function evaluations per iteration!

Slower convergence with
Broyden's method should be
offset by improved efficiency
of each iteration!

Conclusions

- **Recommendations for root finding:**
 - **One-dimensional cases:**
 - If it is not easy/cheap to compute the function's derivative
⇒ use Brent's algorithm
 - If derivative information is available
⇒ use Newton-Raphson's method + bookkeeping on bounds provided you can supply a good enough initial guess!!
 - There are specialized routines for (multiple) root finding of polynomials (but not covered in this course)
 - **Multi-dimensional cases:**
 - Use Newton-Raphson method, but make sure that you provide an initial guess close enough to achieve convergence
 - In case derivative information is expensive
⇒ use Broyden's method (but slower convergence!)

Numerical interpolation and integration

Ivo Roghair, Martin van Sint Annaland

Chemical Process Intensification,
Eindhoven University of Technology

Part I

Numerical interpolation

Today's outline

① Introduction

② Piecewise constant

③ Linear

④ Polynomial

⑤ Splines

Interpolation problem

Definition

Given a set of points x_k , $k = 0, \dots, n$, $x_i \neq x_j$ with associated function values f_k , $k = 0, \dots, n$, or simply: $\{x_k, f_k\}_{k=0}^n$. The interpolation problem is defined as: find a polynomial p_n such that this interpolates the values of f_k on the points x_k :

$$p_n(x_k) = f_k, \quad k = 0, \dots, n$$

Theorem

The interpolation problem for $\{x_k, f_k\}_{k=0}^n$ has a unique solution when $x_i \neq x_j$ for $i \neq j$. Note that we cannot allow multiple function values f_k for the same value of x_k .

What is interpolation?

Interpolation means constructing additional data points within the range of, and using, a discrete set of known data points.

It is typically performed on a uniformly spread data set, but this is not strictly necessary for all methods

Is interpolation the same as curve fitting?

NO

- Curve-fitting requires additionally some way of computing the error between function (curve) and data
 - Curve-fitting does not strictly enforce the function to match the data exactly
 - Curve-fitting may be done on multiple datapoints at one position
 - Curve-fitting is much more expensive to do, requires optimisation

Why do chemical engineers need interpolation?

- Comparison of two data sets which are given at different positions
 - An experimental data set may have been recorded at a constant rate, but the numerical solution is computed at irregular intervals
- Reconstruction of field values distant of computing nodes
 - A CFD simulation on a regular grid containing structures that are not grid-conformant requires interpolation to the structures
- Calculation of a physical property at a condition between those of a lookup table
 - The viscosity of a substance may have been measured at 20°C and 30°C, but not at the desired 28.5°C

General

Several important numerical interpolation methods are discussed today:

- Piecewise constant interpolation
- Linear interpolation
 - Bilinear interpolation
- Polynomial interpolation (Newton's method)
- Spline interpolation

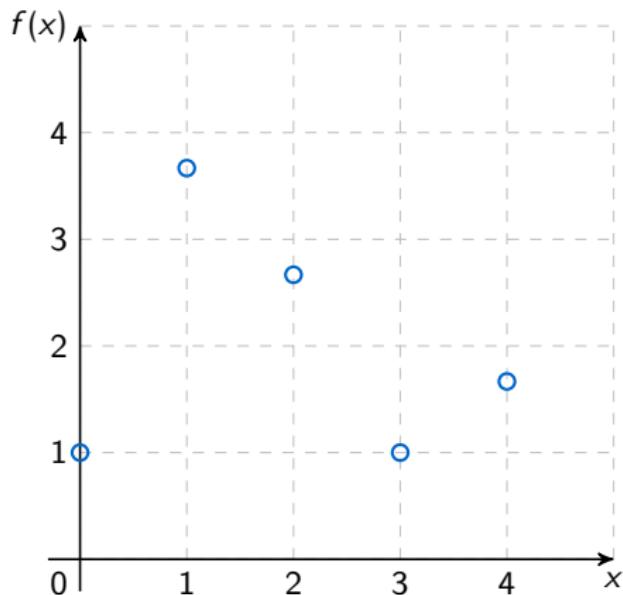
Today's data set

Download the datafile
`interpolation-dataset.mat`,
which contains multiple data sets.

We start with x_1 and y_1 :

x_k	f_k
0	1.00
1	$\frac{11}{3} = 3.67$
2	$\frac{8}{3} = 2.67$
3	1.00
4	$\frac{5}{3} = 1.67$
5	$\frac{23}{3} = 7.67$

Data set $f_n(x_n)$ represented by ○ at discrete intervals $x_n \in \{0, 5\}$



Piecewise constant interpolation

- Nearest-neighbor interpolation in the continuous range $x \in [0, 5]$
- How to treat the point halfway (e.g. at $x = 2.5$)?

$$x \in [0, 0.5] \rightarrow f(x) = f(0)$$

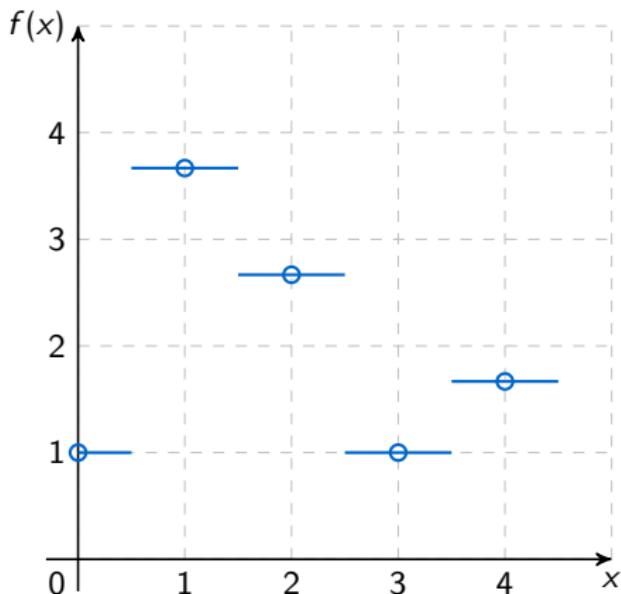
$$x \in]0.5, 1.5] \rightarrow f(x) = f(1)$$

$$x \in]1.5, 2.5] \rightarrow f(x) = f(2)$$

$$x \in]2.5, 3.5] \rightarrow f(x) = f(3)$$

$$x \in]3.5, 4.5] \rightarrow f(x) = f(4)$$

Data set $f_n(x_n)$ represented by \circ at discrete intervals $x_n \in \{0, 5\}$



- Not often used for simple problems, but e.g. for 2D (Voronoi)

Linear interpolation

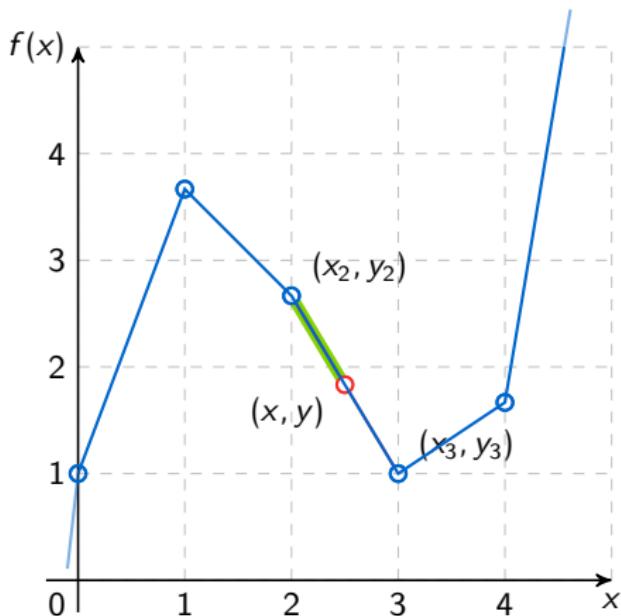
Data set $f_n(x_n)$ represented by \circ at discrete intervals $x_n \in \{0, 5\}$

- Linear interpolation to (x, y) between 2 data points (x_2, y_2) and (x_3, y_3) :

$$\frac{y - y_2}{x - x_2} = \frac{y_3 - y_2}{x_3 - x_2}$$

- Reordered, and more formally:

$$y = y_n + (y_{n+1} - y_n) \frac{x - x_n}{x_{n+1} - x_n}$$

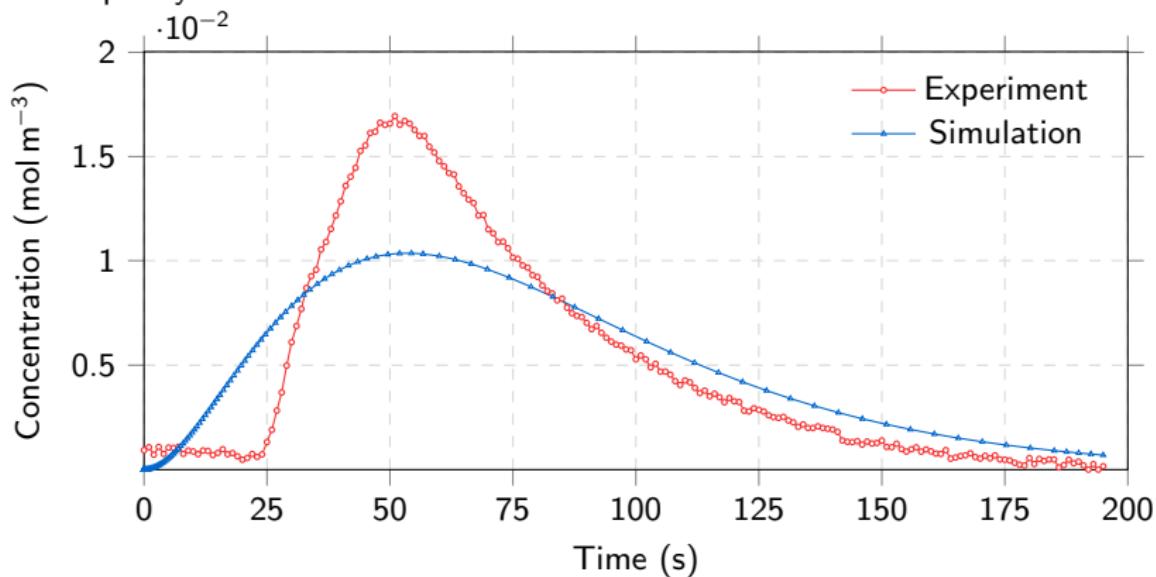


Linear interpolation

- While linear interpolation is fast, and relatively easy to program, it is not very accurate
- At the nodes, the derivatives are discontinuous i.e. not differentiable
- Error is proportional to the square of the distance between nodes

Example: Linear interpolation in Matlab

Consider the data set in `sim_exp_dataset.mat`, containing a normalized concentration and time vector for an experiment and a simulation. The simulation was performed with adaptive node distance to save computation time, thus the concentration is not known at the same times. We are not able to compare yet.



Example: Linear interpolation in Matlab

Consider the data set in `sim_exp_dataset.mat`, containing a normalized concentration and time vector for an experiment and a simulation. The simulation was performed with adaptive node distance to save computation time, thus the concentration is not known at the same times. We are not able to compare yet.

```
% Linear interpolation
c_sim_new = interp1(t_sim,c_sim,t_exp,'linear');
diff = abs(c_exp-c_sim_new);
% Plot the solution
subplot(2,1,1);
plot(t_exp,c_exp,'b-x',t_exp,c_sim_new,'r-o');
subplot(2,1,2);
stem(t_exp,diff);
% Compute the L2-norm
norm(diff)
```

Bi-linear interpolation

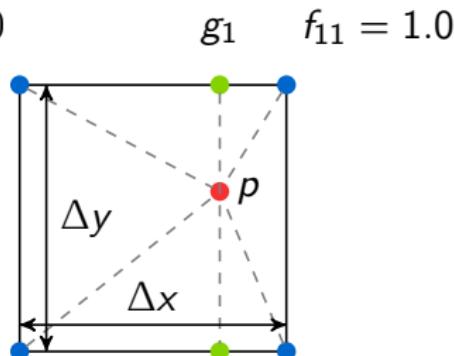
When a 2D field of some quantity is known, we can interpolate the solution to an arbitrary position in the 2D domain $p(x, y)$ using 4 field values f_{00} , f_{10} , f_{01} and f_{11} .

$$\begin{aligned}g_1 &= f_{01} \frac{x_1 - x}{x_1 - x_0} + f_{11} \frac{x - x_0}{x_1 - x_0} \\&= f_{01} \frac{x_1 - x}{\Delta x} + f_{11} \frac{x - x_0}{\Delta x}\end{aligned}$$

$$g_2 = f_{00} \frac{x_1 - x}{\Delta x} + f_{10} \frac{x - x_0}{\Delta x}$$

$$p = g_2 \frac{y_1 - y}{\Delta y} + g_1 \frac{y - y_0}{\Delta y}$$

$$f_{01} = 8.0$$

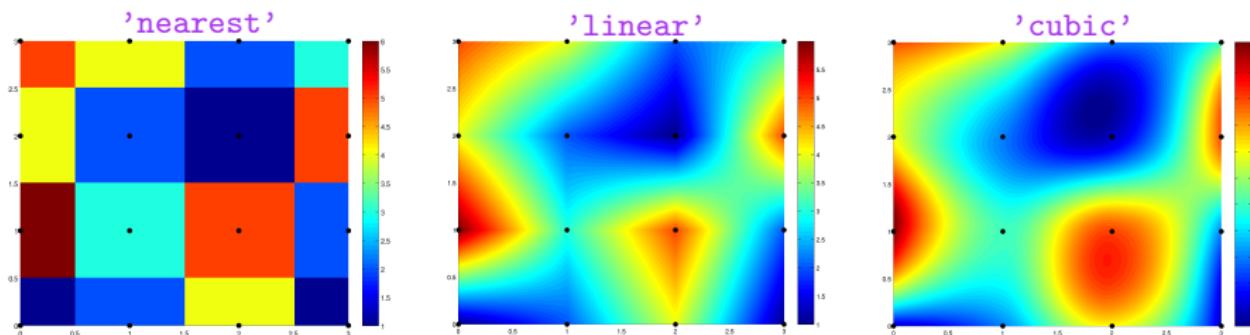


$$f_{00} = 4.0$$

- The order of interpolation (x or y direction first) does not matter; the results are equal

Higher-dimensional field interpolation in Matlab

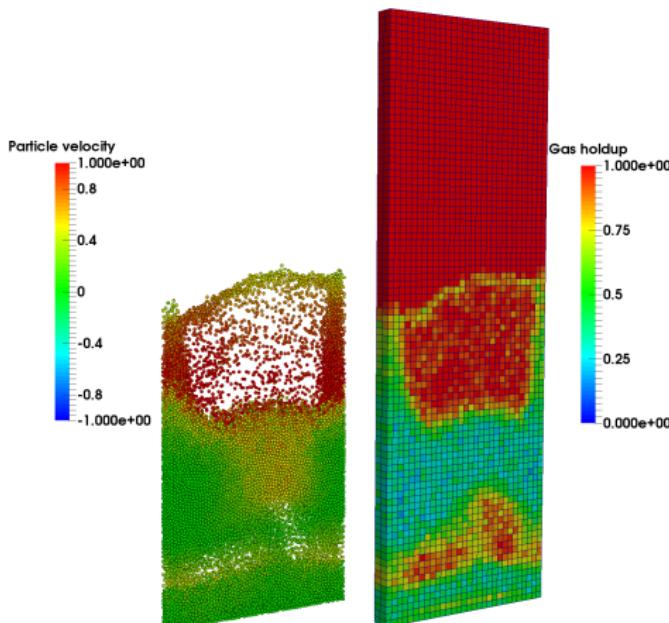
2D or higher-dimensional fields of data can be interpolated in Matlab using the `interp2`, `interp3` or even `interpn` functions, the method can be adjusted:



- Similar to 1D linear interpolation, the derivatives are discontinuous on the grid nodes
- Also consider tri-linear interpolation (for 3D fields), or bicubic interpolation (2D, but third order)

A practical example

Field interpolation is used in e.g. CFD simulations, e.g. a fluidized bed simulation using a *discrete particle model*, where particles are found in between the grid nodes used for velocity computation.



Polynomial interpolation

The examples that we have seen, are simplified forms of *Newton polynomials*. We can interpolate our data with a polynomial of degree n :

$$p_n(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_2 x^2 + a_1 x + a_0$$

Polynomial interpolation via Vandermonde matrix

Consider the data points $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$, the Vandermonde matrix V , coefficient vector a and function value vector y :

$$V_{m,n} = \begin{pmatrix} x_1^0 & x_1^1 & x_1^2 & \cdots & x_1^{n-1} \\ x_2^0 & x_2^1 & x_2^2 & \cdots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_m^0 & x_m^1 & x_m^2 & \cdots & x_m^{n-1} \end{pmatrix} \quad a = \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_{n-1} \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

The coefficients of a polynomial through the data points can be obtained by solving the linear system $Va = y$.

```
>> x = [0 1 2];
>> y = [1.0000; 3.6667;
          2.6667];
>> V = vander(x);
>> a = V\y
a =
    -1.8333
     4.5000
    1.0000
```

So we found the equation:

$$p_2(x) = -1.8333x^2 + 4.5x - 1$$

These Vandermonde-systems are often *ill-conditioned*, so we need another, more stable, method!

Construction of Newton polynomials

Formally, the polynomials $p_n(x)$ are described using prefactors $f[x_0, \dots, x_k]$ and polynomial terms $w_m(x)$:

$$p_n(x) = \sum_{k=0}^n f[x_0, \dots, x_k] w_k(x)$$

The polynomial terms are computed via:

$$w_0(x) = 1, \quad w_1(x) = (x - x_0), \quad w_2(x) = (x - x_0) \cdot (x - x_1),$$

$$w_m(x) = (x - x_0) \cdot (x - x_1) \cdots (x - x_{m-1}) = w_{m-1} \cdot (x - x_{m-1})$$

$$w_m(x) = \prod_{j=0}^{m-1} (x - x_j), \quad m = 0, \dots, n$$

The prefactors are *forward divided differences*, which can be computed as:

$$f[x_{r-k}, \dots, x_r] \equiv \frac{f[x_{r-k+1}, \dots, x_r] - f[x_{r-k}, \dots, x_{r-1}]}{x_r - x_{r-k}}$$

Construction of Newton polynomials: example

Sample data

x_k	f_k
0	1.00
1	$\frac{11}{3} = 3.67$
2	$\frac{8}{3} = 2.67$

$$p_n(x) = \sum_{k=0}^n f[x_0, \dots, x_k] w_k(x)$$

$$f[x_{r-k}, \dots, x_r] \equiv \frac{f[x_{r-k+1}, \dots, x_r] - f[x_{r-k}, \dots, x_{r-1}]}{x_r - x_{r-k}}$$

$$w_m(x) = \prod_{j=0}^{m-1} (x - x_j)$$

x_k	f_k
x_0	$f[x_0] = f_0$
x_1	$f[x_1] = f_1$
x_2	$f[x_2] = f_2$
	$f[x_0, x_1] = \frac{f_1 - f_0}{x_1 - x_0}$
	$f[x_1, x_2] = \frac{f_2 - f_1}{x_2 - x_1}$
	$f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0}$

x_k	f_k
0	1
1	3.67
2	2.67

$$\frac{\frac{11}{3} - 1}{1 - 0} = \frac{8}{3}$$

$$\frac{\frac{8}{3} - \frac{11}{3}}{2 - 1} = \frac{-1}{1} = -1$$

$$\frac{(-1) - \frac{8}{3}}{2 - 0} = -\frac{11}{6}$$

Construction of Newton polynomials: example

Sample data

x_k	f_k
0	1.00
1	$\frac{11}{3} = 3.67$
2	$\frac{8}{3} = 2.67$

$$p_n(x) = \sum_{k=0}^n f[x_0, \dots, x_k] w_k(x)$$

$$f[x_{x-k}, \dots, x_r] = \frac{f[x_{r-k+1}, \dots, x_r] - f[x_{r-k}, \dots, x_{r-1}]}{x_r - x_{r-k}}$$

$$w_m(x) = \prod_{j=0}^{m-1} (x - x_j)$$

x_k	f_k
0	1
1	3.67
2	2.67

$\frac{\frac{11}{3} - 1}{1 - 0} = \frac{8}{3}$
 $\frac{\frac{8}{3} - \frac{11}{3}}{2 - 1} = \frac{-1}{1} = -1 \quad \frac{(-1) - \frac{8}{3}}{2 - 0} = -\frac{11}{6}$

$$\begin{aligned}
 p_2(x) &= 1 \cdot w_m(0) + \frac{8}{3} \cdot w_m(1) + \left(-\frac{11}{6} \right) \cdot w_m(2) \\
 &= 1 \cdot 1 + \frac{8}{3} \cdot (x - 0) + \left(-\frac{11}{6} \right) \cdot (x - 0)(x - 1) = -\frac{11}{6}x^2 + 4\frac{1}{2}x + 1
 \end{aligned}$$

Construction of Newton polynomials: example

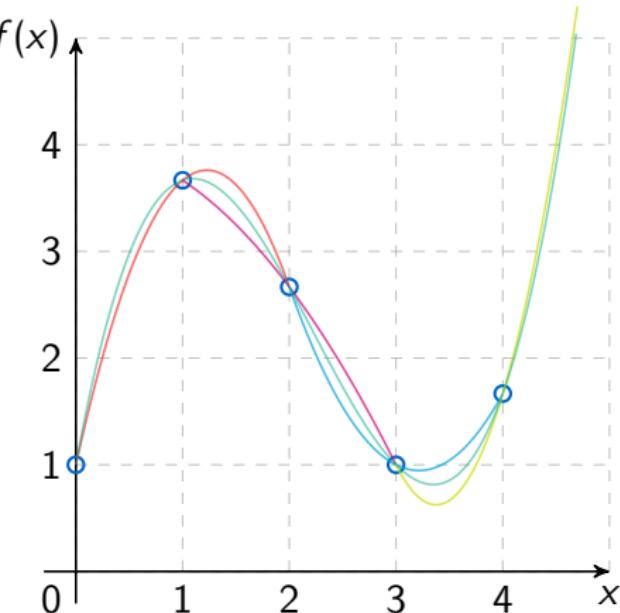
For each three points, a new polynomial interpolant can be derived:

$$p_2(x) = -\frac{11}{6}x^2 + 4\frac{1}{2}x + 1$$

$$p_2(x) = 4 - \frac{x^2}{3}$$

$$p_2(x) = \frac{7x^2}{6} - 7\frac{1}{2}x + 13$$

$$p_2(x) = \frac{8}{3}x^2 - 18x + 31$$

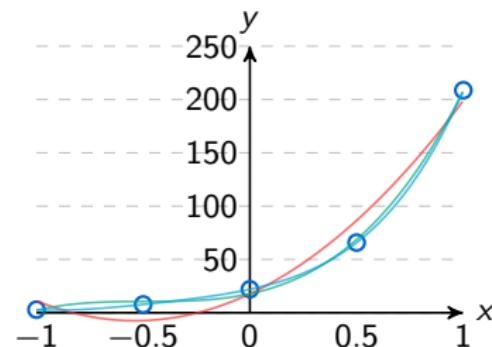


$$f(x) = \frac{x^3}{2} - \frac{10x^2}{3} + \frac{11x}{2} + 1$$

Polynomial fitting in Matlab: example

Develop the $p_2(x)$, $p_3(x)$ and $p_4(x)$ from the following data set (example data x_2 and y_2):

x_k	y_k
-1.0	2.8677
-0.5	7.7530
0.0	22.0000
0.5	65.7863
1.0	208.6744



We use the built-in `polyfit(x,y,n)` and `polyval(p,x)` functions:

```
x_cont = linspace(-1,1,1001);
p2 = polyfit(x2,y2,2);
p3 = polyfit(x2,y2,3);
p4 = polyfit(x2,y2,4);

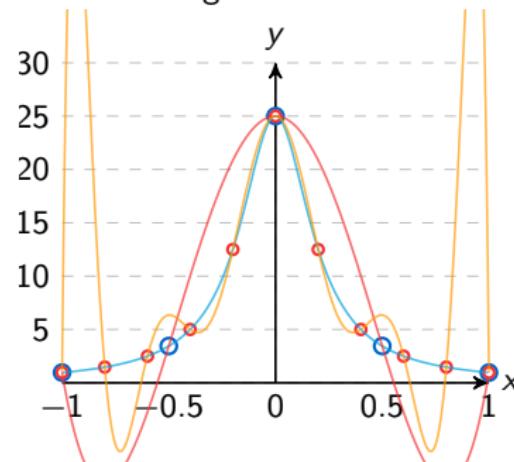
y_cont2 = polyval(p2,x_cont);
y_cont3 = polyval(p3,x_cont);
y_cont4 = polyval(p4,x_cont);
plot(x2,y2,'o',x_cont,y_cont2,x_cont,y_cont3,
x_cont,y_cont4)
```

Exercise

Develop the $p_4(x)$ and $p_{10}(x)$ interpolants from the following data sets:

$$f(x) = \frac{1}{x^2 + \frac{1}{25}} \quad x \in [-1, 1]$$

```
x3a = linspace(-1, 1, 5);
x3b = linspace(-1, 1, 11);
y3a = 1 ./ (x3a.^2 + (1/25));
y3b = 1 ./ (x3b.^2 + (1/25));
```



```
x_cont = linspace(-1, 1, 1001);
p4 = polyfit(x3a, y3a, 4);
p10 = polyfit(x3b, y3b, 10);
y_cont4 = polyval(p4, x_cont);
y_cont10 = polyval(p10, x_cont);
ezplot('1./(x.^2+(1/25))', [-1 1]); hold on;
plot(x3a, y3a, 'o', x3b, y3b, 'x', x_cont, y_cont4, x_cont,
y_cont10);
```

Final thoughts on polynomial interpolation

- An polynomial interpolant of order n requires $n + 1$ data points
 - More data points: interpolant does *not always* cross the points
 - Fewer data points: interpolant is not unique
- Higher-degree polynomials at equidistant points may cause strong oscillatory behaviour (Runge's phenomenon)
 - Mitigation of the problem on Chebyshev (i.e. non uniform grid)...
 - ... or by performing piecewise interpolation (next topic)
- Matlab functions `polyfit(x,y,n)` and `polyval(p,x_new)` were demonstrated.

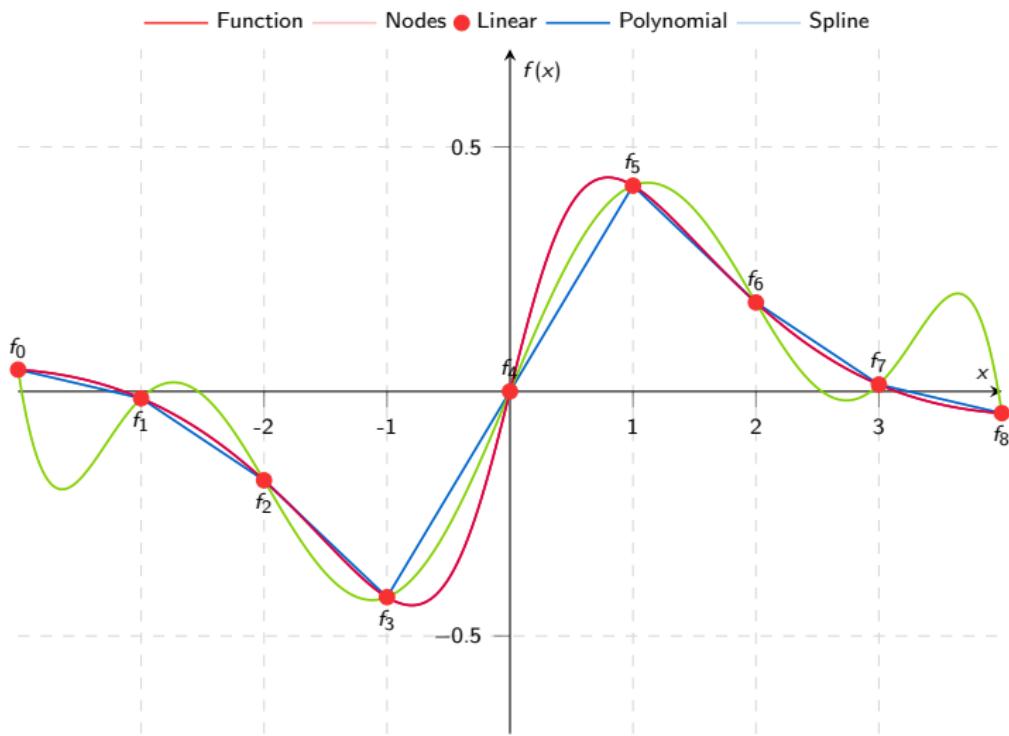
Spline interpolation

A spline is a numerical function that represents a **smooth, higher order, piecewise polynomial** interpolants of a data set.

- Smooth: the interpolant is continuous in the first and second derivatives
- Higher order: The most common type of splines uses third-order polynomials (cubic splines)
- Piecewise polynomial: The interpolant is constructed between each two consecutive tabulated points

Splines: comparison to other interpolation techniques

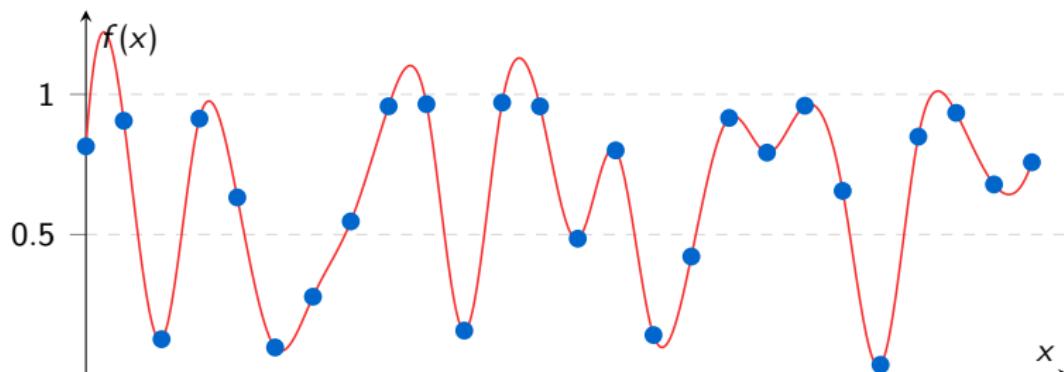
Interpolation of $f(x) = \frac{\sin x}{1 + x^2}$



Spline interpolation in Matlab

We can generate a random data set, and interpolate using `interp1`:

```
% Generate random data set
x=0:25;
y = rand(size(x));
% Interpolant on a fine mesh
xc = linspace(0,25,1001);
yc = interp1(x,y,xc,'spline');
plot(x,y,'o',xc,yc,'-r')
```



Part II

Numerical integration

Today's outline

⑥ Introduction

⑦ Riemann integrals

⑧ Trapezoid rule

⑨ Simpson's rule

⑩ Conclusion

What is numerical integration?

To determine the integral $I(x)$ of an integrand $f(x)$, which can be used to compute the area underneath the integrand between $x = a$ and $x = b$.

$$I(x) = \int_a^b f(x) dx$$

Today we will outline different numerical integration methods.

- Riemann integrals
- Trapezoidal rule
- Simpson's rule

Why do chemical engineers need integration?

- Obtaining the cumulative particle size distribution from a particle size distribution
- The concentration outflow over time may be integrated to yield the residence time distribution
- Integration of a varying product outflow yields the total product outflow
- Quantitative analysis of mixture components via e.g. GC/MS
- Not all functions have an explicit antiderivative, e.g. $\int e^{x^2} dx$ or $\int \frac{1}{\ln x} dx$

Today's outline

⑥ Introduction

⑦ Riemann integrals

⑧ Trapezoid rule

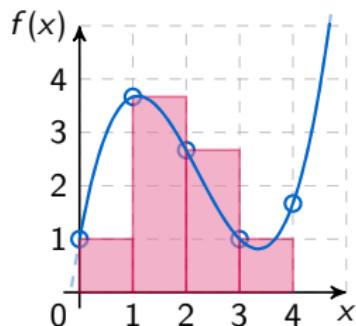
⑨ Simpson's rule

⑩ Conclusion

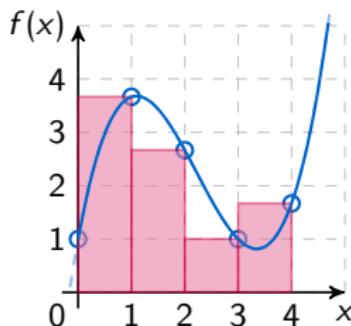
Riemann integrals

Basic idea: Subdivide the interval $[a, b]$ into n subintervals of equal length $\Delta x = \frac{b-a}{n}$ and use the sum of area to approximate the integral.

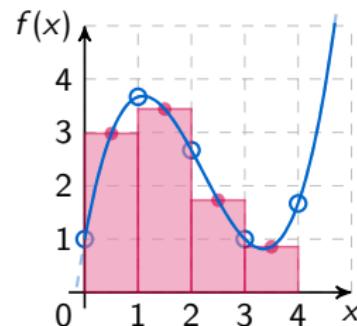
Left endpoint rule



Right endpoint rule



Midpoint rule



$$L_n = \sum_{i=1}^n f(x_{i-1})\Delta x_i$$

$$R_n = \sum_{i=1}^n f(x_i)\Delta x_i$$

$$M_n = \sum_{i=1}^n f(\bar{x}_i)\Delta x_i$$

$$\text{with } \bar{x}_i = \frac{x_{i-1}+x_i}{2}$$

Errors in Riemann integrals

We define the exact integral as $I = \int_a^b f(x)dx$, and L_n , R_n and M_n represent the left, right and midpoint rule approximations of I based on n intervals.

Writing $f_{\max}^{(k)}$ for the maximum value of the k -th derivative, the upper-bounds of the errors by Riemann integrals are:

- $|I - L_n| \leq \frac{f_{\max}^{(1)}(b-a)^2}{2n}$
- $|I - R_n| \leq \frac{f_{\max}^{(1)}(b-a)^2}{2n}$
- $|I - M_n| \leq \frac{f_{\max}^{(2)}(b-a)^3}{24n^2}$

Note that while $|I - L_n|$ and $|I - R_n|$ give the same *upper-bounds* of the error, this does not mean the same error. Rather, the error is of opposite sign!

Today's outline

⑥ Introduction

⑦ Riemann integrals

⑧ Trapezoid rule

⑨ Simpson's rule

⑩ Conclusion

Trapezoid rule

Since the sign of the approximation error of the left and right endpoint rules is opposite, we can take the average of these approximations:

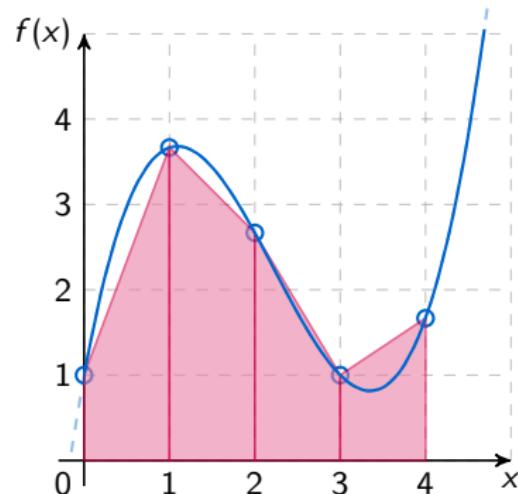
$$T_n = \frac{L_n + R_n}{2}$$

The total area is obtained by geometric reconstruction of trapezoids:

$$T_n = \sum_{i=1}^n \frac{f(x_{i+1}) + f(x_i)}{2} \Delta x_i$$

Note that this can be rewritten for equidistant intervals:

$$T_n = \frac{b-a}{2n} \left(f(x_0) + 2f(x_1) + \dots + 2f(x_{n-1}) + f(x_n) \right)$$



Error in trapezoid integration

The trapezoid rule result over n intervals T_n approximates the exact integral $I = \int_a^b f(x)dx$. The upper-bounds of the error is given as:

$$|I - T_n| \leq \frac{f_{\max}^{(2)}(b-a)^3}{12n^2}$$

Recall that the midpoint rule approximates with an upper-bound error of

$$|I - M_n| \leq \frac{f_{\max}^{(2)}(b-a)^3}{24n^2}$$

The midpoint rule approximation has lower error bounds than the trapezoid rule. A linear function is, however, better approximated by the trapezoid rule.

Today's outline

⑥ Introduction

⑦ Riemann integrals

⑧ Trapezoid rule

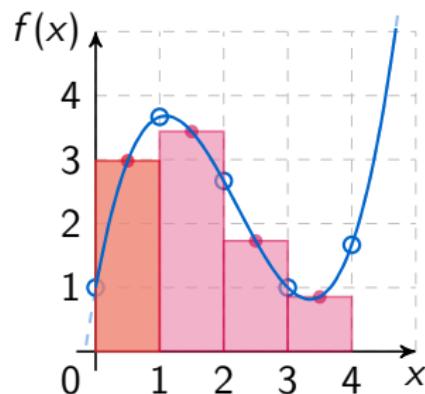
⑨ Simpson's rule

⑩ Conclusion

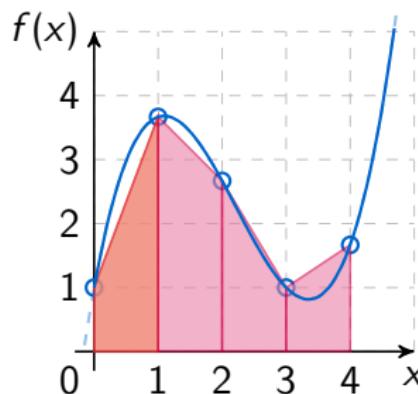
Towards higher-order integration

Compare how the midpoint and trapezoid functions behave on convex and concave parts of a graph.

Midpoint rule



Trapezoid rule



In convex parts (bending down), the midpoint rule tends to overestimate the integral (trapezoid underestimates).

In concave parts (bending up), the midpoint rule tends to underestimate the integral (trapezoid overestimates).

Towards higher-order integration

The errors of the midpoint rule and trapezoid rule behave in a similar way, but have opposite signs.

- Midpoint: $|I - M_n| \leq \frac{f_{\max}^{(2)}(b-a)^3}{24n^2}$
- Trapezoid: $|I - T_n| \leq \frac{f_{\max}^{(2)}(b-a)^3}{12n^2}$

For a quadratic function, the errors relate as:

$$|I - M_n| = \frac{1}{2}|I - T_n|$$

Taking the weighted average of these two yields the Simpson's rule:

$$S_{2n} = \frac{2}{3}M_n + \frac{1}{3}T_n$$

The $2n$ means we have $2n$ subintervals: the n trapezoid intervals are subdivided by the midpoint rule.

Simpson's rule

Consider the interval $i \in [x_0, x_2]$, subdivided in three equidistant interpolation points: x_0, x_1, x_2 .

- Midpoint: $M_i = f\left(\frac{x_0 + x_2}{2}\right)2\Delta x = f(x_1)2\Delta x$
- Trapezoid: $T_i = \frac{f(x_0) + f(x_2)}{2}2\Delta x$
- Simpson: $S_i = \frac{2}{3}M_i + \frac{1}{3}T_i$

Note that M_i and T_i were computed on interval $x_2 - x_0 = 2\Delta x$.

Now we have:

$$\begin{aligned} S_i &= \frac{2}{3}[f(x_1)2\Delta x] + \frac{1}{3}\left[\frac{f(x_0) + f(x_2)}{2}2\Delta x\right] \\ &= \frac{4\Delta x}{3}f(x_1) + \frac{\Delta x}{3}f(x_0) + f(x_2) = \frac{\Delta x}{3}(f(x_0) + 4f(x_1) + f(x_2)) \end{aligned}$$

Simpson's rule

We write $f(x_k) = f_k$. The integral of an interval $i \in [x_0, x_2]$ is approximated as:

$$S_i = \frac{\Delta x}{3} (f_0 + 4f_1 + f_2)$$

The next interval, S_j with $j \in [x_2, x_4]$ with midpoint $x_3 = \frac{x_2+x_4}{2}$ is approximated as:

$$S_j = \frac{\Delta x}{3} (f_2 + 4f_3 + f_4)$$

If we sum these two intervals we obtain:

$$\begin{aligned} I \approx S_i + S_j &= \left[\frac{\Delta x}{3} (f_0 + 4f_1 + f_2) \right] + \left[\frac{\Delta x}{3} (f_2 + 4f_3 + f_4) \right] \\ &= \frac{\Delta x}{3} (f_0 + 4f_1 + 2f_2 + 4f_3 + f_4) \end{aligned}$$

Simpson's rule

In general, Simpson's rule can be written as:

$$\int_a^b f(x)dx \approx \sum_{\substack{k=2 \\ k \text{ even}}}^n \frac{\Delta x}{3} (f_{k-2} + 4f_{k-1} + f_k)$$
$$= \frac{\Delta x}{3} (f_0 + 4f_1 + 2f_2 + 4f_3 + 2f_4 + \dots + 2f_{n-2} + 4f_{n-1} + f_n)$$

The error is given by:

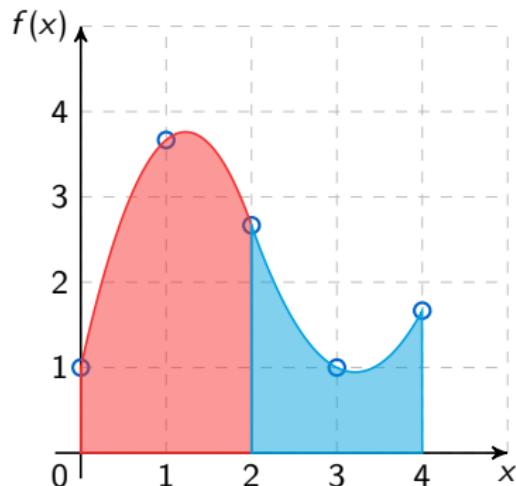
$$|I - S_n| \leq \frac{f_{\max}^{(4)}(b-a)^5}{180n^4}$$

if integrand f is differentiable on $[a, b]$.

Simpson's rule: example

Recall our example data, described by $f(x) = \frac{x^3}{2} - \frac{10x^2}{3} + \frac{11x}{2} + 1$
 $I = \int_0^4 \frac{x^3}{2} - \frac{10x^2}{3} + \frac{11x}{2} + 1 = \frac{80}{9} \approx 8.888\dots$

- Interpolating x_0, x_1 and x_2 :
 $p_{2a}(x) = -\frac{11}{6}x^2 + 4\frac{1}{2}x + 1$
 $\int_0^2 p_{2a} = \frac{55}{9} \approx 6.1111$
- Interpolating x_2, x_3 and x_4 :
 $p_{2b}(x) = \frac{7x^2}{6} - 7\frac{1}{2}x + 13$
 $\int_2^4 p_{2b} = \frac{25}{9} \approx 2.777\dots$
- Adding the separate integrals:
 $\int_0^2 p_{2a} + \int_2^4 p_{2b} = \frac{80}{9}$



Using Simpson's rule: $I \approx \frac{\Delta x}{3} (f_0 + 4f_1 + 2f_2 + 4f_3 + f_4) =$
 $\frac{1}{3} (1 + 4 \cdot 3.6667 + 2 \cdot 2.6667 + 4 \cdot 1.0000 + 1.6667) = 8.88888 = \frac{80}{9}$

Simpson's method is of fourth order, and it gives exact approximations of third order polynomials!

Integration in Matlab

Integration can be done numerically in Matlab.

- `trapz(x,y)` uses the trapezoid rule to integrate the data. Make sure you use the `x` variable if your data is not spaced with $\Delta x = 1$. Can handle non-equidistant data.
- Integration of functions can be done using the `integral(fun,xmin,xmax)` function:

```
fun = @(x) exp(-x.^2);
I = integral(fun,0,10)
I =
    0.886226925452758
```

Today's outline

⑥ Introduction

⑦ Riemann integrals

⑧ Trapezoid rule

⑨ Simpson's rule

⑩ Conclusion

What hasn't been discussed?

This course is by no means complete, and further reading is possible.

- Legendre polynomials: Another way of performing the polynomial interpolation
- Gaussian quadrature: A third-order integration method that requires only two base points (in contrast to the third order Simpson's method, which requires three points)
- Adaptive techniques: Parts of a function that are relatively steady (no wild oscillations) and differentiable can be integrated with much larger step sizes than other parts of the function.
- Simpson's 3/8-rule: Yet another integration technique, requiring an additional data point

Summary

- Interpolation is used to obtain data between existing data points
 - (Bi-)Linear, polynomial and spline interpolation methods
 - Construction of Newton polynomials
 - Oscillations of high-order polynomials
- Several techniques for numerical integration were discussed:
 - Riemann sums, trapezoid rule, Simpson's rule
 - Upper-bound errors were given for each technique

Ordinary differential equations

Martin van Sint Annaland, Ivo Roghair

m.v.sintannaland@tue.nl

Chemical Process Intensification,
Eindhoven University of Technology

Introduction
oooooooo

Euler's method
oooooooooooo

Rates of convergence
ooooooo

Runge-Kutta methods
oooooooooooo

Step size control
ooooo

Solving ODEs in Matlab
ooooooo

Part I

Explicit and Implicit methods

Today's outline

① Introduction

② Euler's method

Forward Euler

③ Rates of convergence

④ Runge-Kutta methods

RK2 methods

RK4 method

⑤ Step size control

⑥ Solving ODEs in Matlab

Overview

Ordinary differential equations

An equation containing a function of one independent variable and its derivatives, in contrast to a *partial differential equation*, which contains derivatives with respect to more independent variables.

Main question

How to solve

$$\frac{dy}{dx} = f(y(x), x) \quad \text{with} \quad y(x=0) = y_0$$

accurately and efficiently?

What is an ODE?

- Algebraic equation:

$$f(y(x), x) = 0 \quad \text{e.g. } -\ln(K_{eq}) = (1 - \zeta)$$

- First order ODE:

$$f \left(\frac{dy}{dx}(x), y(x), x \right) = 0 \quad \text{e.g. } \frac{dc}{dt} = -kc^n$$

- Second order ODE:

$$f \left(\frac{d^2y}{dx^2}(x), \frac{dy}{dx}(x), y(x), x \right) = 0 \quad \text{e.g.} \quad \mathcal{D} \frac{d^2c}{dx^2} = -\frac{kc}{1+Kc}$$

About second order ODEs

Very often a second order ODE can be rewritten into a system of first order ODEs (whether it is handy depends on the boundary conditions!)

More general

Consider the second order ODE:

$$\frac{d^2y}{dx^2} + q(x)\frac{dy}{dx} = r(x)$$

Now define and solve using z as a new variable:

$$\frac{dy}{dx} = z(x)$$

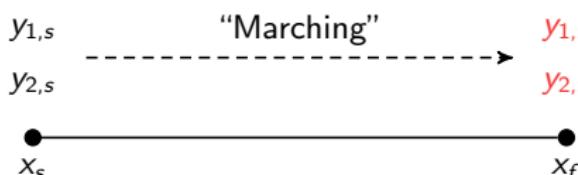
$$\frac{dz}{dx} = r(x) - q(x)z(x)$$

Importance of boundary conditions

The nature of boundary conditions determines the appropriate numerical method. Classification into 2 main categories:

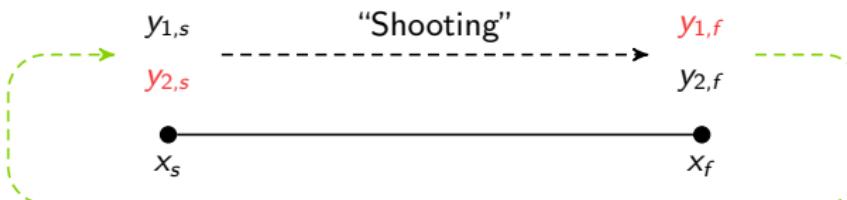
- *Initial value problems (IVP)*

We know the values of all y_i at some starting position x_s , and it is desired to find the values of y_i at some final point x_f .



- *Boundary value problems (BVP)*

Boundary conditions are specified at more than one x . Typically, some of the BC are specified at x_s and the remainder at x_f .



Overview

Initial value problems:

- Explicit methods
 - First order: forward Euler
 - Second order: improved Euler (RK2)
 - Fourth order: Runge-Kutta 4 (RK4)
 - Step size control
- Implicit methods
 - First order: backward Euler
 - Second order: midpoint rule

Boundary value problems

- Shooting method

Today's outline

① Introduction

② Euler's method

Forward Euler

③ Rates of convergence

④ Runge-Kutta methods

RK2 methods

RK4 method

⑤ Step size control

⑥ Solving ODEs in Matlab

Euler's method

Consider the following single initial value problem:

$$\frac{dc}{dt} = f(c(t), t) \quad \text{with} \quad c(t=0) = c_0 \quad (\text{initial value problem})$$

Easiest solution algorithm: Euler's method, derived here via Taylor series expansion:

$$c(t_0 + \Delta t) \approx c(t_0) + \left. \frac{dc}{dt} \right|_{t_0} \Delta t + \frac{1}{2} \left. \frac{d^2 c}{dt^2} \right|_{t_0} (\Delta t)^2 + \mathcal{O}(\Delta t^3)$$

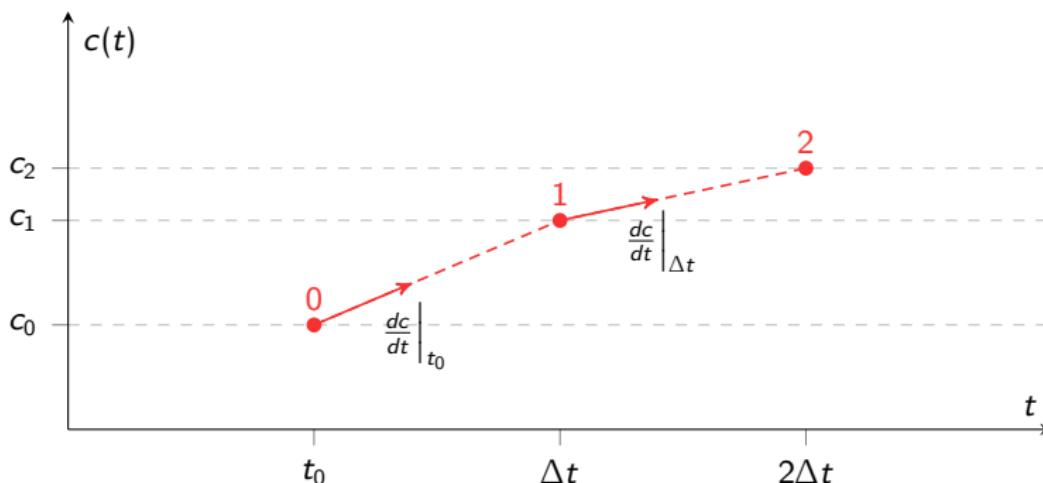
Neglect terms with higher order than two: $\left. \frac{dc}{dt} \right|_{t_0} = \frac{c(t_0 + \Delta t) - c(t_0)}{\Delta t}$

Substitution:

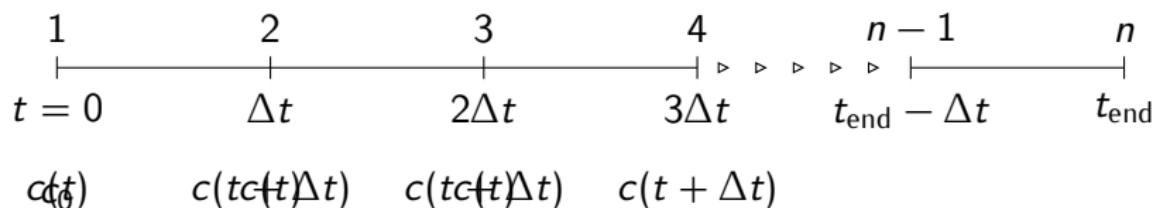
$$\frac{c(t_0 + \Delta t) - c(t_0)}{\Delta t} = f(c_0, t_0) \Rightarrow c(t_0 + \Delta t) = c(t_0) + \Delta t f(c_0, t_0)$$

Euler's method: graphical example

$$\frac{c(t_0 + \Delta t) - c(t_0)}{\Delta t} = f(c_0, t_0) \Rightarrow c(t_0 + \Delta t) = c(t_0) + \Delta t f(c_0, t_0)$$



Euler's method - solution method



Start with $t = t_0$, $c = c_0$, then calculate at discrete points in time:
 $c(t_1 = t_0 + \Delta t) = c(t_0) + \Delta t f(c_0, t_0)$.

Pseudo-code Euler's method: $\frac{dy}{dx} = f(x, y)$ and $y(x_0) = y_0$.

- ① Initialize variables, functions; set $h = \frac{x_1 - x_0}{N}$
- ② Set $x = x_0$, $y = y_0$
- ③ While $x < x_{\text{end}}$ do

$$x_{i+1} = x_i + h; \quad y_{i+1} = y_i + h f(x_i, y_i)$$

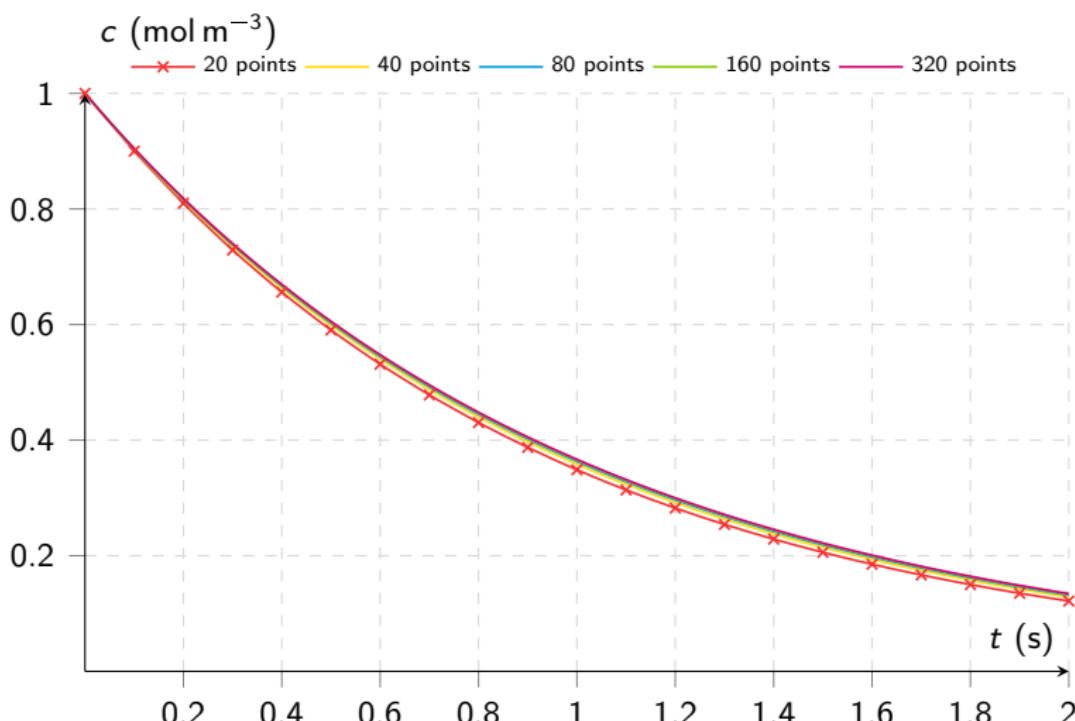
Euler's method - example

First order reaction in a batch reactor:

$$\frac{dc}{dt} = -kc \quad \text{with} \quad c(t=0) = 1 \text{ mol m}^{-3}, \quad k = 1 \text{ s}^{-1}, \quad t_{\text{end}} = 2 \text{ s}$$

Time [s]	Concentration [mol m ⁻³]
$t_0 = 0$	$c_0 = 1.00$
$t_1 = t_0 + \Delta t$ $= 0 + 0.1 = 0.1$	$c_1 = c_0 + \Delta t \cdot (-kc_0)$ $= 1 + 0.1 \cdot (-1 \cdot 1) = 0.9$
$t_2 = t_1 + \Delta t$ $= 0.1 + 0.1 = 0.2$	$c_2 = c_1 + \Delta t \cdot (-kc_1)$ $= 0.9 + 0.1 \cdot (-1 \cdot 0.9) = 0.81$
$t_3 = t_2 + \Delta t$ $= 0.2 + 0.1 = 0.3$	$c_3 = c_2 + \Delta t \cdot (-kc_2)$ $= 0.81 + 0.1 \cdot (-1 \cdot 0.81) = 0.729$
...	...
$t_{i+1} = t_i + \Delta t$	$c_{i+1} = c_i + \Delta t \cdot (-kc_i)$
...	...
$t_{20} = 2.0$	$c_{20} = c_{19} + \Delta t \cdot (-kc_{19}) = 0.121577$

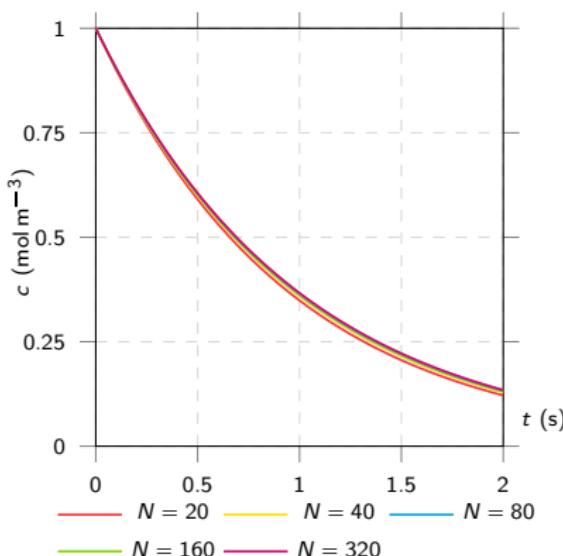
Euler's method - example



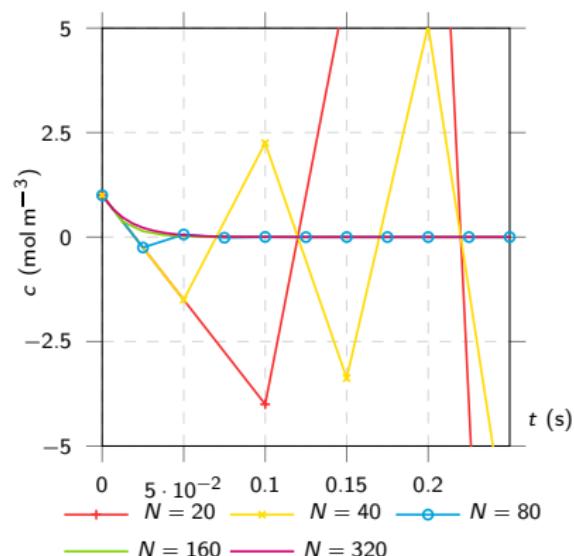
Problems with Euler's method

The question is: What step size, or how many steps to use?

- ① Accuracy \Rightarrow need information on numerical error!
- ② Stability \Rightarrow need information on stability limits!



Reaction rate: $k = 1 \text{ s}^{-1}$



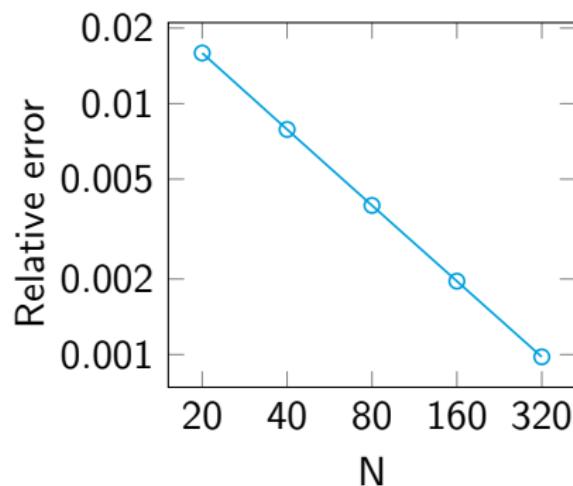
Reaction rate: $k = 50 \text{ s}^{-1}$

Accuracy

Comparison with analytical solution for $k = 1 \text{ s}^{-1}$:

$$c(t) = c_0 \exp(-kt) \Rightarrow \zeta = 1 - \exp(-kt) \Rightarrow \zeta_{\text{analytical}} = 0.864665$$

N	ζ	$\frac{\zeta_{\text{numerical}} - \zeta_{\text{analytical}}}{\zeta_{\text{analytical}}}$
20	0.878423	0.015912
40	0.871488	0.007891
80	0.868062	0.003929
160	0.866360	0.001961
320	0.865511	0.000979



Accuracy

For Euler's method: Error halves when the number of grid points is doubled, i.e. error is proportional to Δt : first order method.

Error estimate:

$$\left. \frac{dx}{dt} \right|_{t_0} = \frac{x(t_0 + \Delta t) - x(t_0)}{\Delta t} + \frac{1}{2} \left. \frac{d^2x}{dt^2} \right|_{t_0} (\Delta t) + \mathcal{O}(\Delta t)^2$$

$$\frac{x(t_0 + \Delta t) - x(t_0)}{\Delta t} = f(x_0, t_0) - \frac{1}{2} \left. \frac{d^2x}{dt^2} \right|_{t_0} (\Delta t) + \mathcal{O}(\Delta t)^2$$

Today's outline

① Introduction

② Euler's method

Forward Euler

③ Rates of convergence

④ Runge-Kutta methods

RK2 methods

RK4 method

⑤ Step size control

⑥ Solving ODEs in Matlab

Errors and convergence rate

L_2 norm (Euclidean norm)

$$\|\mathbf{v}\|_2 = \sqrt{v_1^2 + v_2^2 + \dots + v_n^2} = \sqrt{\sum_{i=1}^n v_i^2}$$

L_∞ norm (maximum norm)

$$\|\mathbf{v}\|_\infty = \max(|v_1|, \dots, |v_n|)$$

Absolute difference

$$\epsilon_{\text{abs}} = \left\| \mathbf{y}_{\text{numerical}} - \mathbf{y}_{\text{analytical}} \right\|_{2,\infty}$$

Relative difference

$$\epsilon_{\text{rel}} = \frac{\left\| \mathbf{y}_{\text{numerical}} - \mathbf{y}_{\text{analytical}} \right\|_{2,\infty}}{\left\| \mathbf{y}_{\text{analytical}} \right\|_{2,\infty}}$$

Errors and convergence rate

Convergence rate (or: order of convergence) r

$$\epsilon = \lim_{\Delta x \rightarrow 0} c(\Delta x)^r$$

- A first order method reduces the error by a factor 2 when increasing the number of steps by a factor 2
- A second order method reduces the error by a factor 4 when increasing the number of steps by a factor 2

Computing the rate of convergence

When the analytical solution is available, choose ① or ② for a particular number of grid points N :

- ① Compute the relative or absolute error vector $\bar{\epsilon}$. Take the norm to compute a single error value ϵ following:

- Based on L_1 -norm: $\epsilon = \frac{\|\bar{\epsilon}\|_1}{N}$
- Based on L_2 -norm: $\epsilon = \frac{\|\bar{\epsilon}\|_2}{\sqrt{N}}$
- Based on L_∞ -norm: $\epsilon = \|\bar{\epsilon}\|_\infty$

- ② Compute the relative or absolute error at a single indicative points (e.g. middle of domain, outlet).

Compare to calculations with different number of steps: $\epsilon_1 = c(\Delta x_1)^r$ and $\epsilon_2 = c(\Delta x_2)^r$ and solve for r :

$$\frac{\epsilon_2}{\epsilon_1} = \frac{c(\Delta x_2)^r}{c(\Delta x_1)^r} = \left(\frac{\Delta x_2}{\Delta x_1} \right)^r \Rightarrow \log \left(\frac{\epsilon_2}{\epsilon_1} \right) = \log \left(\frac{\Delta x_2}{\Delta x_1} \right)^r$$

$$\Rightarrow r = \frac{\log \left(\frac{\epsilon_2}{\epsilon_1} \right)}{\log \left(\frac{\Delta x_2}{\Delta x_1} \right)} = \frac{\log \left(\frac{\epsilon_2}{\epsilon_1} \right)}{\log \left(\frac{N_1}{N_2} \right)} \quad \text{in the limit of } \Delta x \rightarrow 0 \quad \text{or} \quad N \rightarrow \infty$$

Computing the rate of convergence

When the analytical solution is not available:

- ① Compute the solution with $N + 1$, N , $N - 1$ and $N - 2$ grid points
- ② Select a single indicative grid point (e.g. middle of domain, outlet) that lies at exactly the same position in each computation
- ③ Use the solution c at this grid point for various grid sizes to compute:

$$r = \frac{\log \frac{c_{N+1} - c_N}{c_N - c_{N-1}}}{\log \frac{c_N - c_{N-1}}{c_{N-1} - c_{N-2}}}$$

- ④ Alternative for simulations with $2N$, N and $\frac{N}{2}$ grid points:

$$r = \frac{\log \left| \frac{c_{2N} - c_N}{c_N - c_{\frac{N}{2}}} \right|}{\log \left| \frac{2N}{N} \right|}$$

Example: Euler's method — order of convergence

N	ζ	$\frac{\zeta_{\text{numerical}} - \zeta_{\text{analytical}}}{\zeta_{\text{analytical}}}$	$r = \frac{\log\left(\frac{\epsilon_j}{\epsilon_{j-1}}\right)}{\log\left(\frac{N_{j-1}}{N_j}\right)}$
20	0.878423	0.015912	—
40	0.871488	0.007891	1.011832
80	0.868062	0.003929	1.005969
160	0.866360	0.001961	1.002996
320	0.865511	0.000979	1.001500

⇒ Euler's method is a first order method (as we already knew from the truncation error analysis)

Wouldn't it be great to have a method that can give the answer using much less steps? ⇒ Higher order methods

Today's outline

① Introduction

② Euler's method

Forward Euler

③ Rates of convergence

④ Runge-Kutta methods

RK2 methods

RK4 method

⑤ Step size control

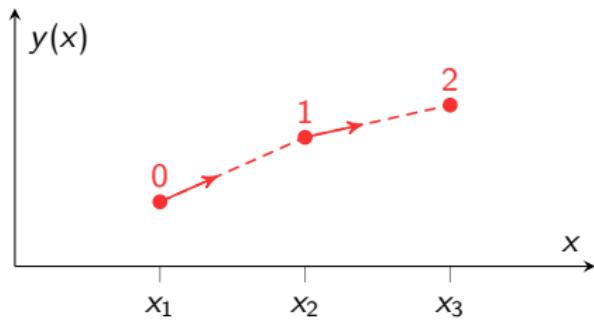
⑥ Solving ODEs in Matlab

Runge-Kutta methods

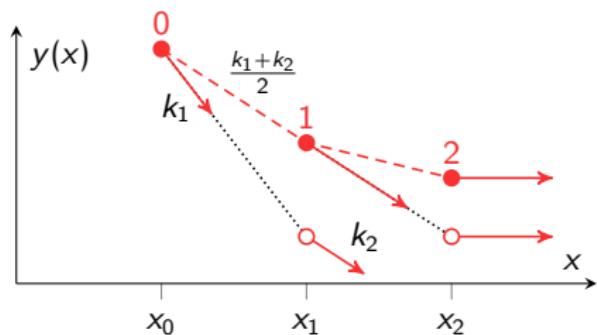
Propagate a solution by combining the information of several Euler-style steps (each involving one function evaluation) to match a Taylor series expansion up to some higher order.

Euler: $y_{i+1} = y_i + hf(x_i, y_i)$ with $h = \Delta x$, i.e.
slope = $k_1 = f(x_i, y_i)$.

Euler's method



RK2 method



Classical second order Runge-Kutta (RK2) method

This method is also called Heun's method, or improved Euler method:

- ① Approximate the slope at x_i : $k_1 = f(x_i, y_i)$
- ② Approximate the slope at x_{i+1} : $k_2 = f(x_{i+1}, y_{i+1})$ where we use Euler's method to approximate $y_{i+1} = y_i + hf(x_i, y_i) = y_i + hk_1$
- ③ Perform an Euler step with the average of the slopes:

$$y_{i+1} = y_i + h \frac{1}{2}(k_1 + k_2)$$

In pseudocode:

```
x = x0, y = y0
while x < xend do
    xi+1 = xi + h
    k1 = f(xi, yi)
    k2 = f(xi + h, yi + hk1)
    yi+1 = yi + h  $\frac{1}{2}$ (k1 + k2)
end while
```

Runge-Kutta methods — derivation

$$\frac{dy}{dx} = f(x, y(x))$$

Using Taylor series expansion: $y_{i+1} = y_i + h \left. \frac{dy}{dx} \right|_i + \frac{h^2}{2} \left. \frac{d^2y}{dx^2} \right|_i + \mathcal{O}(h^3)$

$$\left. \frac{dy}{dx} \right|_i = f(x_i, y_i) \equiv f_i$$

$$\left. \frac{d^2y}{dx^2} \right|_i = \left. \frac{d}{dx} f(x, y(x)) \right|_i = \left. \frac{\partial f}{\partial x} \right|_i + \left. \frac{\partial f}{\partial y} \right|_i \left. \frac{\partial y}{\partial x} \right|_i = \left. \frac{\partial f}{\partial x} \right|_i + \left. \frac{\partial f}{\partial y} \right|_i f_i \quad (\text{chain rule})$$

Substitution gives:

$$y_{i+1} = y_i + hf_i + \frac{h^2}{2} \left(\left. \frac{\partial f}{\partial x} \right|_i + \left. \frac{\partial f}{\partial y} \right|_i f_i \right) + \mathcal{O}(h^3)$$

$$y_{i+1} = y_i + \frac{h}{2} f_i + \frac{h}{2} \left(f_i + h \left. \frac{\partial f}{\partial x} \right|_i + hf_i \left. \frac{\partial f}{\partial y} \right|_i \right) + \mathcal{O}(h^3)$$

Runge-Kutta methods — derivation

Note multivariate Taylor expansion:

$$f(x_i + h, y_i + k) = f_i + h \frac{\partial f}{\partial x} \Big|_i + h \frac{\partial f}{\partial y} \Big|_i + \mathcal{O}(h^2)$$

$$\Rightarrow \frac{h}{2} \left(f_i + h \frac{\partial f}{\partial x} \Big|_i + hf_i \frac{\partial f}{\partial y} \Big|_i \right) = \frac{h}{2} f(x_i + h, y_i + hf_i) + \mathcal{O}(h^3)$$

Concluding:

$$y_{i+1} = y_i + \frac{h}{2} f_i + \frac{h}{2} f(x_i + h, y_i + hf_i) + \mathcal{O}(h^3)$$

Rewriting:

$$k_1 = f(x_i, y_i)$$

$$k_2 = f(x_i + h, y_i + hk_1)$$

$$\Rightarrow y_{i+1} = y_i + \frac{h}{2}(k_1 + k_2)$$

Runge-Kutta methods — derivation

Generalization: $y_{i+1} = y_i + h(b_1 k_1 + b_2 k_2) + \mathcal{O}(h^3)$

with $k_1 = f_i$, $k_2 = f(x_i + c_2 h, y_i + a_{2,1} h k_1)$

(Note that classical RK2: $b_1 = b_2 = \frac{1}{2}$ and $c_2 = a_{2,1} = 1.$)

Bivariate Taylor expansion:

$$f(x_i + c_2 h, y_i + a_{2,1} h k_1) = f_i + c_2 h \left. \frac{\partial f}{\partial x} \right|_i + a_{2,1} h k_1 \left. \frac{\partial f}{\partial y} \right|_i + \mathcal{O}(h^2)$$

$$\begin{aligned} y_{i+1} &= y_i + h(b_1 k_1 + b_2 k_2) + \mathcal{O}(h^3) \\ &= y_i + h [b_1 f_i + b_2 f(x_i + c_2 h, y_i + a_{2,1} h k_1)] + \mathcal{O}(h^3) \\ &= y_i + h \left[b_1 f_i + b_2 \left\{ f_i + c_2 h \left. \frac{\partial f}{\partial x} \right|_i + a_{2,1} h k_1 \left. \frac{\partial f}{\partial y} \right|_i + \mathcal{O}(h^2) \right\} \right] + \mathcal{O}(h^3) \\ &= y_i + h(b_1 + b_2)f_i + h^2 b_2 \left(c_2 \left. \frac{\partial f}{\partial x} \right|_i + a_{2,1} f_i \left. \frac{\partial f}{\partial y} \right|_i \right) + \mathcal{O}(h^3) \end{aligned}$$

Comparison with Taylor:

$$y_{i+1} = y_i + h f_i + \frac{h^2}{2} \left(\left. \frac{\partial f}{\partial x} \right|_i + \left. \frac{\partial f}{\partial y} \right|_i f_i \right) + \mathcal{O}(h^3)$$

Runge-Kutta methods — derivation

$$y_{i+1} = y_i + h(b_1 + b_2)f_i + h^2 b_2 \left(c_2 \frac{\partial f}{\partial x} \Big|_i + a_{2,1} f_i \frac{\partial f}{\partial y} \Big|_i \right) + \mathcal{O}(h^3)$$

$$y_{i+1} = y_i + hf_i + \frac{h^2}{2} \left(\frac{\partial f}{\partial x} \Big|_i + \frac{\partial f}{\partial y} \Big|_i f_i \right) + \mathcal{O}(h^3)$$

⇒ 3 eqns and 4 unknowns ⇒ multiple possibilities!

① Classical RK2:

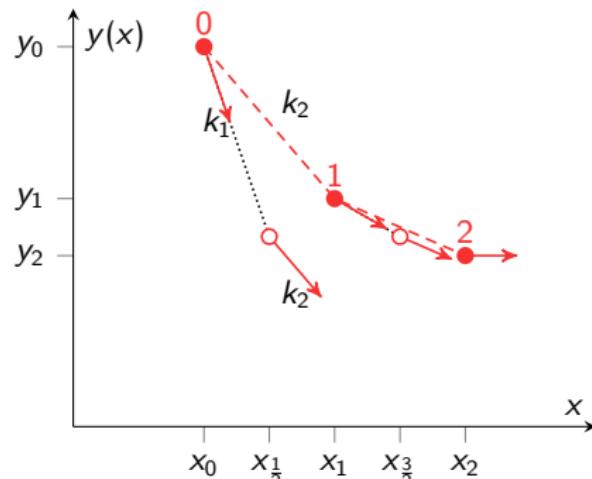
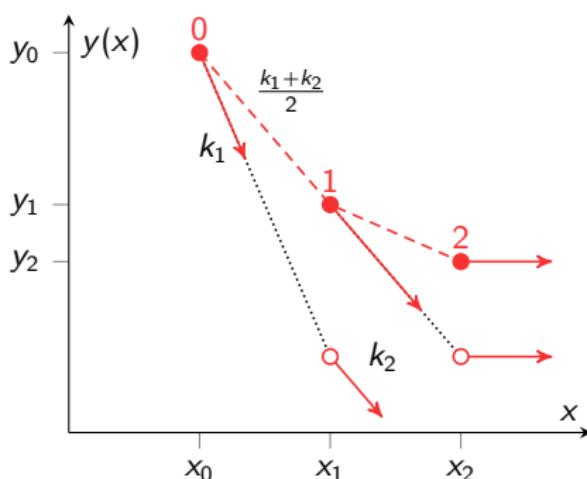
$$b_1 = b_2 = \frac{1}{2} \text{ and } c_2 = a_{2,1} = 1$$

② Midpoint rule (modified Euler):

$$b_1 = 0, b_2 = 1, c_2 = a_{2,1} = \frac{1}{2}$$

Second order Runge-Kutta methods

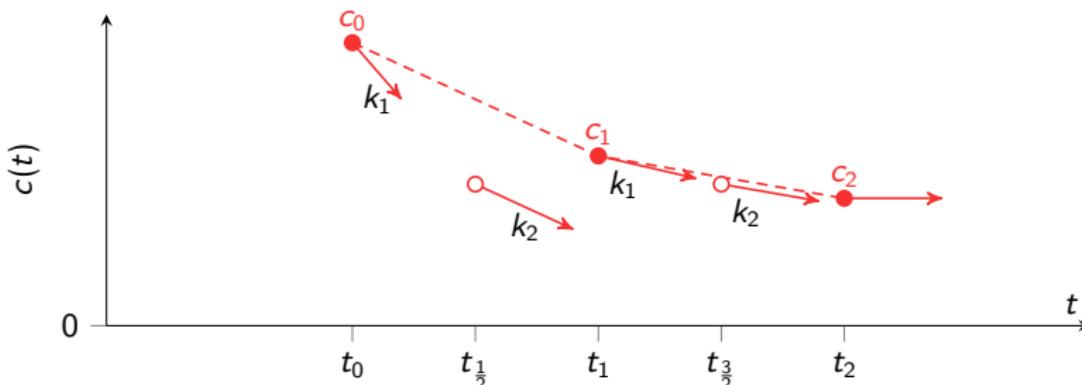
Classical RK2 method (= Heun's method, improved Euler method)	Explicit midpoint rule (modified Euler method)
$k_1 = f_i$	$k_1 = f_i$
$k_2 = f(x_i + h, y_i + hk_1)$	$k_2 = f(x_i + \frac{1}{2}h, y_i + \frac{1}{2}hk_1)$
$y_{i+1} = y_i + \frac{1}{2}h(k_1 + k_2)$	$y_{i+1} = y_i + hk_2$



Second order Runge-Kutta method — Example

First order reaction in a batch reactor: $\frac{dc}{dt} = -kc$ with
 $c(t = 0) = 1 \text{ mol m}^{-3}$, $k = 1 \text{ s}^{-1}$, $t_{\text{end}} = 2 \text{ s}$.

Time [s]	$C \text{ [mol m}^{-3}]$	$k_1 = hf(x_i, y_i)$	$k_2 = hf(x_i + \frac{1}{2}h, y_n + \frac{1}{2}k_1)$
0	1.00	$0.1 \cdot (-1 \cdot 1) = -0.1$	$0.1 \cdot (-1 \cdot (1 - 0.5 \cdot 0.1)) = -0.095$
0.1	$1 - 0.095 = 0.905$	$0.1 \cdot (-1 \cdot 0.905) = -0.0905$	$0.1 \cdot (-1 \cdot (0.905 - 0.5 \cdot 0.0905)) = -0.085975$
...
2	0.1358225	-0.0135822	-0.0129031



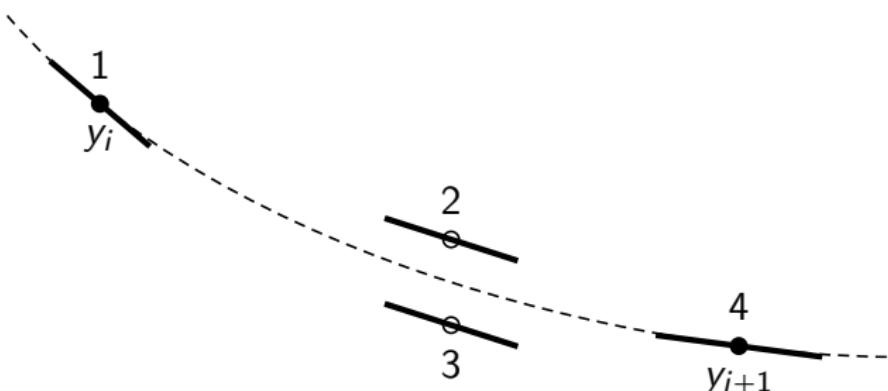
RK2 method — order of convergence

N	ζ	$\frac{\zeta_{\text{numerical}} - \zeta_{\text{analytical}}}{\zeta_{\text{analytical}}}$	$r = \frac{\log\left(\frac{\epsilon_i}{\epsilon_{i-1}}\right)}{\log\left(\frac{N_{i-1}}{N_i}\right)}$
20	0.864178	5.634×10^{-4}	—
40	0.864548	1.355×10^{-4}	2.056
80	0.864636	3.323×10^{-5}	2.028
160	0.864658	8.229×10^{-6}	2.014
320	0.864663	2.048×10^{-6}	2.007

⇒ RK2 is a second order method. Doubling the number of cells reduces the error by a factor 4!

Can we do even better?

RK4 method (classical fourth order Runge-Kutta method)



$$k_1 = f(x_i, y_i)$$

$$k_2 = f\left(x_i + \frac{1}{2}h, y_i + \frac{1}{2}hk_1\right)$$

$$k_3 = f\left(x_i + \frac{1}{2}h, y_i + \frac{1}{2}hk_2\right)$$

$$k_4 = f(x_i + h, y_i + hk_3)$$

$$y_{i+1} = y_i + h \left(\frac{1}{6}k_1 + \frac{1}{3}(k_2 + k_3) + \frac{1}{6}k_4 \right)$$

RK4 method — order of convergence

N	ζ	$\frac{\zeta_{\text{numerical}} - \zeta_{\text{analytical}}}{\zeta_{\text{analytical}}}$	$r = \frac{\log\left(\frac{\epsilon_i}{\epsilon_{i-1}}\right)}{\log\left(\frac{N_{i-1}}{N_i}\right)}$
20	0.864664472	2.836×10^{-7}	—
40	0.864664702	1.700×10^{-8}	4.060
80	0.864664716	1.040×10^{-9}	4.030
160	0.864664717	6.435×10^{-11}	4.015
320	0.864664717	4.001×10^{-12}	4.007

⇒ RK4 is a fourth order method: Doubling the number of cells reduces the error by a factor 16!

Can we do even better?

Today's outline

① Introduction

② Euler's method

Forward Euler

③ Rates of convergence

④ Runge-Kutta methods

RK2 methods

RK4 method

⑤ Step size control

⑥ Solving ODEs in Matlab

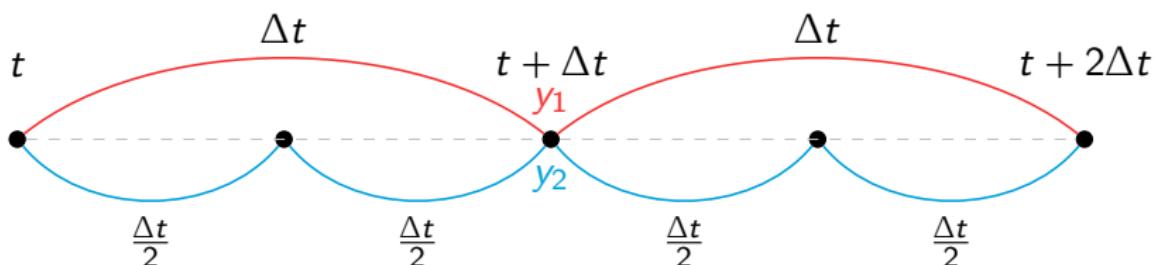
Adaptive step size control

The step size (be it either position, time or both (PDEs)) cannot be decreased indefinitely to favour a higher accuracy, since each additional grid point causes additional computation time. It may be wise to adapt the step size according to the computation requirements.

Globally two different approaches can be used:

- ① Step doubling: compare solutions when taking one full step or two consecutive halve steps
- ② Embedded methods: Compare solutions when using two approximations of different order

Adaptive step size control: step doubling



- RK4 with one large step of h : $y_{i+1} = y_1 + ch^5 + \mathcal{O}(h^6)$
- RK4 with two steps of $\frac{1}{2}h$: $y_{i+1} = y_2 + 2c(\frac{1}{2}h)^5 + \mathcal{O}(h^6)$

Adaptive step size control: step doubling

- Estimation of truncation error by comparing y_1 and y_2 :
$$\Delta = y_2 - y_1$$
- If Δ too large, reduce step size for accuracy
- If Δ too small, increase step size for efficiency.
- Ignoring higher order terms and solving for c :
$$\Delta = \frac{15}{16}ch^5 \Rightarrow ch^5 = \frac{16}{15}\Delta \Rightarrow y_{i+1} = y_2 + \frac{\Delta}{15} + \mathcal{O}(h^6)$$

(local Richardson extrapolation)

Note that when we specify a tolerance tol , we can estimate the maximum allowable step size as: $h_{\text{new}} = \alpha h_{\text{old}} \left| \frac{\text{tol}}{\Delta} \right|^{\frac{1}{5}}$ with α a safety factor (typically $\alpha = 0.9$).

Adaptive step size control: embedded methods

Use a special fourth and a fifth order Runge-Kutta method to approximate y_{i+1}

- The fourth order method is special because we want to use the same positions for the evaluation for computational efficiency.
- RK45 is the preferred method (minimum number of function evaluations) (this is built in Matlab as `ode45`).

Today's outline

① Introduction

② Euler's method

Forward Euler

③ Rates of convergence

④ Runge-Kutta methods

RK2 methods

RK4 method

⑤ Step size control

⑥ Solving ODEs in Matlab

Solving ODEs in Matlab

Matlab provides convenient procedures to solve (systems of) ODEs automatically.

The procedure is as follows:

- ① Create a function that specifies the ODE(s). Specifically, this function returns the $\frac{dy}{dx}$ value (vector).
- ② Initialise solver variables and settings (e.g. step size, initial conditions, tolerance), in a separate script
- ③ Call the ODE solver function, using a *function handle* to the ODE function described in point 1.
 - The ODE solver will return the vector for the independent variable, and a solution vector (matrix for systems of ODEs).

Solving ODEs in Matlab: example 1

We solve the system: $\frac{dx}{dt} = -k_1x + k_2$, $k_1 = 0.2$, $k_2 = 2.5$

- Create an anonymous function handle:

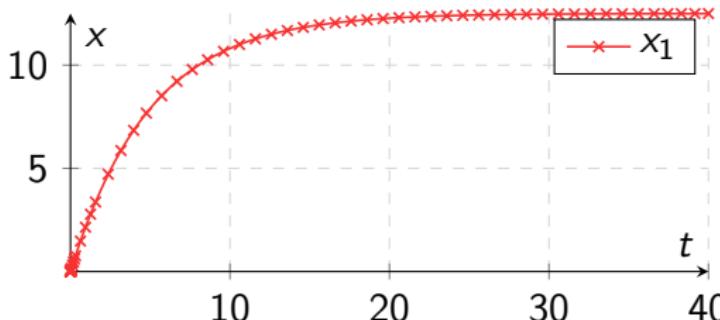
```
>> myEqn = @(t,x) (-0.2*x + 2.5)
```

- Solve with a call to

```
ode45(function_handle, timespan, initial_condition):
```

```
>> ode45(myEqn, [0 40], 0);
```

- By omitting the output of this function, the graph is automatically drawn.



Solving ODEs in Matlab: example 2

We solve the system:

$$\frac{dx}{dt} = \begin{cases} -\frac{k_1}{x^2} & t \leq 10 \\ \frac{k_2}{x} - \frac{k_1}{x^2} & t > 10 \end{cases} \quad \text{with } k_1 = 0.5, k_2 = 1, x(0) = 2$$

Create an ODE function

```
function [dxdt] = myEqnFunction(t,x)
k1 = 0.5;
k2 = 1;
dxdt = (t>10)*k2/x - k1/x^2;
```

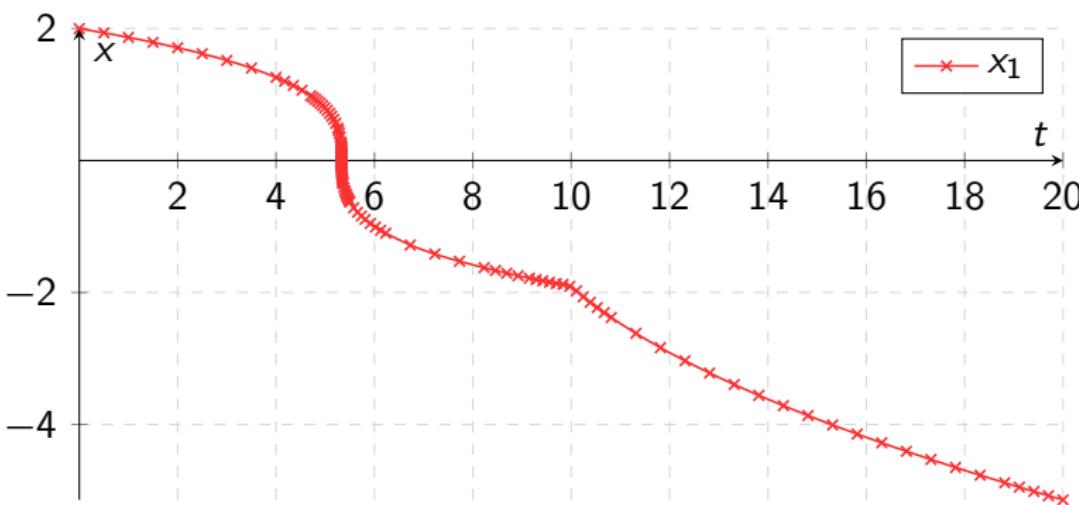
Create a solution script

```
x_init = 2; % Initial condition
tspan = [0 20]; % Time span
options = odeset('RelTol',1e-8,'AbsTol',1e-8);
[t,x] = ode45(@myEqnFunction,tspan,x_init,options);
```

Solving ODEs in Matlab: example 2

Plot the solution:

```
plot(t,x, 'r-x')
```



Note the refinement in regions where large changes occur.

Solving ODEs in Matlab: example

A few notes on working with `ode45` and other solvers. If we want to give additional arguments (e.g. `k1` and `k2`) to our ODE function, we can list them in the function line:

```
function [dxdt] = myEqn(t,x,k1,k2)
```

The additional arguments can now be set in the solver script by *adding them after the options*:

```
[t,x] = ode45(@myEqn,tspan,x_0,options,k1,k2);
```

- Of course, in the solver script, the variables do not have to be called `k1` and `k2`:

```
[t,x] = ode45(@myODE,tspan,x_0,options,q,u);
```

- These variables may be of any type (vectors, matrix, struct). Especially a struct is useful to carry many values in 1 variable.

Solving systems of ODEs in Matlab: example

You have noticed that the step size in t varied. This is because we have given just the begin and end times of our time span:

```
tspan = [0 10];
```

You can also solve at specific steps, by supplying all steps explicitly, e.g.:

```
tspan = linspace(0,10,101);
```

This example provides 101 explicit time steps between 0 and 10 seconds.

Note that the results are interpolated to these data points afterwards; you do not influence the efficiency and accuracy of the solver algorithm this way!

Ordinary differential equations

Martin van Sint Annaland, Ivo Roghair

m.v.sintannaland@tue.nl

Chemical Process Intensification,
Eindhoven University of Technology

Part II

Implicit methods, Systems of ODEs and Boundary Value Problems

Today's outline

⑦ Implicit methods

Backward Euler

Implicit midpoint method

⑧ Systems of ODEs

Solution methods for systems of ODEs

Solving systems of ODEs in Matlab

Stiff systems of ODEs

⑨ Boundary value problems

Shooting method

⑩ Conclusion

Problems with Euler's method: instability

Consider the ODE:

$$\frac{dy}{dx} = f(x, y(x)) \quad \text{with} \quad y(x=0) = y_0$$

First order approximation of derivative: $\frac{dy}{dx} = \frac{y_{i+1} - y_i}{\Delta x}$.

Where to evaluate the function f ?

- ① Evaluation at x_i : Explicit Euler method (forward Euler)
 - ② Evaluation at x_{i+1} : Implicit Euler method (backward Euler)

Problems with Euler's method: instability – forward Euler

Explicit Euler method (forward Euler):

- Use values at x_i :

$$\frac{y_{i+1} - y_i}{\Delta x} = f(x_i, y_i) \Rightarrow y_{i+1} = y_i + h f(x_i, y_i).$$
 - This is an explicit equation for y_{i+1} in terms of y_i .
 - It can give instabilities with large function values.

Consider the first order batch reactor:

$$\frac{dc}{dt} = -kc \Rightarrow c_{i+1} = c_i - k \textcolor{red}{c}_i \Delta t \Rightarrow \frac{c_{i+1}}{c_i} = 1 - k \Delta t$$

It follows that unphysical results are obtained for $k\Delta t \geq 1$!!

Stability requirement

$$k\Delta t < 1$$

(but probably accuracy requirements are more stringent here!)

Problems with Euler's method: instability – backward Euler

Implicit Euler method (backward Euler):

- Use values at x_{i+1} :

$$\frac{y_{i+1} - y_i}{\Delta x} = f(x_{i+1}, y_{i+1}) \Rightarrow y_{i+1} = y_i + h f(x_{i+1}, y_{i+1}).$$
 - This is an implicit equation for y_{i+1} , because it also depends on terms of y_{i+1} .

Consider the first order batch reactor:

$$\frac{dc}{dt} = -kc \Rightarrow c_{i+1} = c_i - k \textcolor{red}{c_{i+1}} \Delta t \Rightarrow \frac{c_{i+1}}{c_i} = \frac{1}{1 + k \Delta t}$$

This equation does never give unphysical results!

The implicit Euler method is *unconditionally stable* (but maybe not very accurate or efficient).

Semi-implicit Euler method

Usually f is a non-linear function of y , so that linearization is required (recall Newton's method).

$$\frac{dy}{dx} = f(y) \Rightarrow y_{i+1} = y_i + hf(y_{i+1}) \quad \text{using} \quad f(y_{i+1}) = f(y_i) + \left. \frac{df}{dy} \right|_i (y_{i+1} - y_i) + \dots$$

$$\Rightarrow y_{i+1} = y_i + h \left[f(y_i) + \frac{df}{dy} \Big|_i (y_{i+1} - y_i) \right]$$

$$\Rightarrow \left(1 - h \frac{df}{dy} \Big|_i \right) y_{i+1} = \left(1 - h \frac{df}{dy} \Big|_i \right) y_i + hf(y_i)$$

$$\Rightarrow y_{i+1} = y_i + h \left(1 - h \frac{df}{dy} \Big|_i \right)^{-1} f(y_i)$$

For the case that $f(x, y(x))$ we could add the variable x as an additional variable $y_{n+1} = x$. Or add one fully implicit Euler step (which avoids the computation of $\frac{\partial f}{\partial x}$):

$$y_{i+1} = y_i + h f(x_{i+1}, y_{i+1}) \Rightarrow y_{i+1} = y_i + h \left(1 - h \frac{df}{dy} \Big|_i \right)^{-1} f(x_{i+1}, y_i)$$

Semi-implicit Euler method - example

Second order reaction in a batch reactor

$$\frac{dc}{dt} = -kc^2 \text{ with } c_0 = 1 \text{ mol m}^{-3}, k = 1 \text{ m}^3 \text{ mol}^{-1} \text{ s}^{-1}, t_{\text{end}} = 2 \text{ s}$$

Analytical solution: $c(t) = \frac{c_0}{1+k c_0 t}$

Define $f = -kc^2$, then $\frac{df}{dc} = -2kc \Rightarrow c_{i+1} = c_i - \frac{hkc_i^2}{1+2hkc_i}$.

N	ζ	$\frac{\zeta_{\text{numerical}} - \zeta_{\text{analytical}}}{\zeta_{\text{analytical}}}$	$r = \frac{\log\left(\frac{\epsilon_i}{\epsilon_{i-1}}\right)}{\log\left(\frac{N_{i-1}}{N_i}\right)}$
20	0.654066262	1.89×10^{-2}	—
40	0.660462687	9.31×10^{-3}	1.02220
80	0.663589561	4.62×10^{-3}	1.01162
160	0.665134433	2.30×10^{-3}	1.00594
320	0.665902142	1.15×10^{-3}	1.00300

Second order implicit method: Implicit midpoint method

Implicit midpoint rule (second order)	Explicit midpoint rule (modified Euler method)
$y_{i+1} = y_i + hf\left(x_i + \frac{1}{2}h, \frac{1}{2}(y_i + y_{i+1})\right)$	$y_{i+1} = y_i + hf(x_i + \frac{1}{2}h, y_i + \frac{1}{2}hk_1)$

in case $f(y)$ then:

$$f\left(\frac{1}{2}(y_i + y_{i+1})\right) = f_i + \frac{df}{dy}\Big|_i \left(\frac{1}{2}(y_i + y_{i+1}) - y_i\right) = f_i + \frac{1}{2} \frac{df}{dy}\Big|_i (y_{i+1} - y_i)$$

Implicit midpoint rule reduces to:

$$y_{i+1} = y_i + h f_i + \frac{h}{2} \left. \frac{df}{dy} \right|_i (y_{i+1} - y_i)$$

$$\Rightarrow \left(1 - \frac{h}{2} \left. \frac{df}{dy} \right|_i \right) y_{i+1} = \left(1 - \frac{h}{2} \left. \frac{df}{dy} \right|_i \right) y_i + h f_i$$

$$\Rightarrow y_{i+1} = y_i + h \left(1 - \frac{h}{2} \frac{df}{dy} \Big|_i \right)^{-1} f_i$$

Implicit midpoint method — example

Second order reaction in a batch reactor

$\frac{dc}{dt} = -kc^2$ with $c_0 = 1 \text{ mol m}^{-3}$, $k = 1 \text{ m}^3 \text{ mol}^{-1} \text{ s}^{-1}$, $t_{\text{end}} = 2 \text{ s}$
 (Analytical solution: $c(t) = \frac{c_0}{1 + k c_0 t}$).

Define $f = -kc^2$, then $\frac{df}{dc} = -2kc$.

Substitution:

$$c_{i+1} = c_i + h \left(1 - \frac{h}{2} \cdot (-2kc_i) \right)^{-1} \cdot (-kc_i^2)$$

$$= c_i - \frac{hkc_i^2}{1 + hkc_i} = \frac{c_i + hkc_i^2 - hkc_i^2}{1 + hkc_i} \Rightarrow c_{i+1} = \frac{c_i}{1 + hkc_i}$$

You will find that this method is exact for all step sizes h because of the quadratic source term!

Implicit midpoint method — example

Second order reaction in a batch reactor:

$$\frac{dc}{dt} = -kc^2 \text{ with } c_0 = 1 \text{ mol m}^{-3}, k = 1 \text{ m}^3 \text{ mol}^{-1} \text{ s}^{-1}, t_{\text{end}} = 2 \text{ s}$$

Analytical solution: $c(t) = \frac{c_0}{1+kc_0t}$

$$c_{i+1} = \frac{c_i}{1 + hkc_i}$$

N	ζ	$\frac{\zeta_{\text{numerical}} - \zeta_{\text{analytical}}}{\zeta_{\text{analytical}}}$	$r = \frac{\log\left(\frac{\epsilon_i}{\epsilon_{i-1}}\right)}{\log\left(\frac{N_{i-1}}{N_i}\right)}$
20	0.6666666667	1.665×10^{-16}	—
40	0.6666666667	0	—
80	0.6666666667	0	—
160	0.6666666667	0	—
320	0.6666666667	0	—

Implicit midpoint method — example

Third order reaction in a batch reactor: $\frac{dc}{dt} = -kc^3$

Analytical solution: $c(t) = \frac{c_0}{\sqrt{1+2kc_0^2 t}}$

$$c_{i+1} = c_i - \frac{h k c_i^3}{1 + \frac{3}{2} h k c_i^2}$$

N	ζ	$\frac{\zeta_{\text{numerical}} - \zeta_{\text{analytical}}}{\zeta_{\text{analytical}}}$	$r = \frac{\log\left(\frac{\epsilon_i}{\epsilon_{i-1}}\right)}{\log\left(\frac{N_{i-1}}{N_i}\right)}$
20	0.5526916174	1.71×10^{-4}	—
40	0.5527633731	4.17×10^{-5}	2.041
80	0.5527807304	1.03×10^{-5}	2.021
160	0.5527849965	2.55×10^{-6}	2.011
320	0.5527860538	6.34×10^{-7}	2.005

Today's outline

⑦ Implicit methods

Backward Euler

Implicit midpoint method

⑧ Systems of ODEs

Solution methods for systems of ODEs

Solving systems of ODEs in Matlab

Stiff systems of ODEs

⑨ Boundary value problems

Shooting method

⑩ Conclusion

Systems of ODEs

A system of ODEs is specified using vector notation:

$$\frac{dy}{dx} = \mathbf{f}(x, \mathbf{y}(x))$$

for

$$\frac{dy_1}{dx} = f_1(x, y_1(x), y_2(x)) \quad \text{or} \quad f_1(x, y_1, y_2)$$

$$\frac{dy_2}{dx} = f_2(x, y_1(x), y_2(x)) \quad \text{or} \quad f_2(x, y_1, y_2)$$

The solution techniques discussed before can also be used to solve systems of equations.

Systems of ODEs: Explicit methods

Forward Euler method

$$\mathbf{y}_{i+1} = \mathbf{y}_i + h\mathbf{f}(x_i, \mathbf{y}_i)$$

Improved Euler method (classical RK2)

$$y_{i+1} = y_i + \frac{h}{2}(k_1 + k_2) \quad \text{using} \quad \begin{aligned} k_1 &= f(x_i, y_i) \\ k_2 &= f(x_i + h, y_i + hk_1) \end{aligned}$$

Modified Euler method (midpoint rule)

$$y_{i+1} = y_i + h k_2 \quad \text{using} \quad \begin{aligned} k_1 &= f(x_i, y_i) \\ k_2 &= f\left(x_i + \frac{h}{2}, y_i + \frac{h}{2} k_1\right) \end{aligned}$$

Systems of ODEs: Explicit methods

Classical fourth order Runge-Kutta method (RK4)

$$\mathbf{y}_{i+1} = \mathbf{y}_i + h \left(\frac{\mathbf{k}_1}{6} + \frac{1}{3} (\mathbf{k}_2 + \mathbf{k}_3) + \frac{\mathbf{k}_4}{6} \right)$$

$$\mathbf{k}_1 = \mathbf{f}(x_i, \mathbf{y}_i)$$

$$\mathbf{k}_2 = \mathbf{f}\left(x_i + \frac{h}{2}, \mathbf{y}_i + \frac{h}{2}\mathbf{k}_1\right)$$

using

$$\mathbf{k}_3 = \mathbf{f}\left(x_i + \frac{h}{2}, \mathbf{y}_i + \frac{h}{2}\mathbf{k}_2\right)$$

$$\mathbf{k}_4 = \mathbf{f}(x_i + h, \mathbf{y}_i + h\mathbf{k}_3)$$

Solving systems of ODEs in Matlab

Solving systems of ODEs in Matlab is completely analogous to solving a single ODE:

- ① Create a function that specifies the ODEs. This function returns the $\frac{dy}{dx}$ vector.
- ② Initialise solver variables and settings (e.g. step size, initial conditions, tolerance), in a separate script. Initial conditions and tolerances should be given per-equation, i.e. as a vector.
- ③ Call the ODE solver function, using a *function handle* to the ODE function described in point 1.
 - The ODE solver will return the vector for the independent variable, and a solution matrix, with a column as the solution for each equation in the system.

Solving systems of ODEs in Matlab: example

We solve the system: $\frac{dx_1}{dt} = -x_1 - x_2, \quad \frac{dx_2}{dt} = x_1 - 2x_2$

Create an ODE function

```
function [dxdt] = myODEFunction(t,x)
dxdt(1) = -x(1) - x(2);
dxdt(2) = x(1) - 2*x(2);
dxdt=dxdt'; % Transpose to column vector
return
```

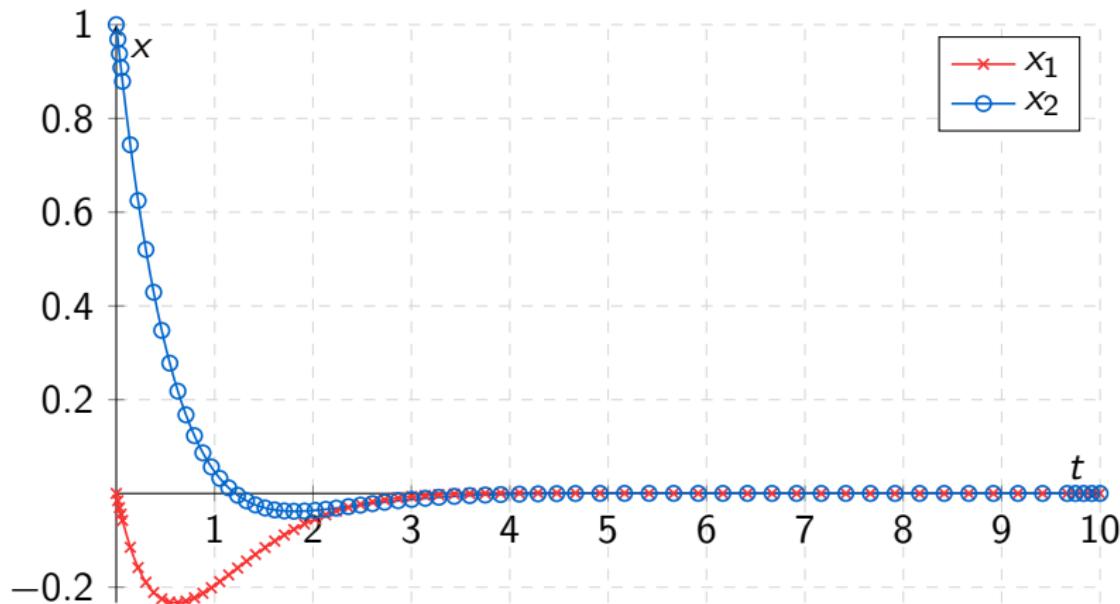
Create a solution script

```
x_init = [0 1]; % Initial conditions
tspan = [0 10]; % Time span
options = odeset('RelTol',1e-4,'AbsTol',[1e-4 1e-4]);
[t,x] = ode45(@myODEfunction,tspan,x_init,options);
```

Solving systems of ODEs in Matlab: example

Plot the solution:

```
plot(t,x(:,1),'r-x',t,x(:,2),'b-o')
```



Solving systems of ODEs in Matlab: repeated notes

A few notes on working with `ode45` and other solvers. If we want to give additional arguments (e.g. `a`, `b` and `c`) to our ODE function, we can list them in the function line:

```
function [dxdt] = myODE(t,x,a,b,c)
```

The additional arguments can now be set in the solver script by *adding them after the options*:

```
[t,x] = ode45(@myODE,tspan,x_0,options,a,b,c);
```

- Of course, in the solver script, the variables do not need to be called `a`, `b` and `c`:

```
[t,x] = ode45(@myODE,tspan,x_0,options,k1,phi,V);
```

- These variables may be of any type (vectors, matrix, struct). Especially a struct is useful to carry many values in 1 variable.

Systems of ODEs: Implicit methods

Backward Euler method

$$\mathbf{y}_{i+1} = \mathbf{y}_i + h \left(1 - h \frac{d\mathbf{f}}{d\mathbf{y}} \Big|_i \right)^{-1} \mathbf{f}(\mathbf{y}_i)$$

Implicit midpoint method

$$\mathbf{y}_{i+1} = \mathbf{y}_i + h \left(1 - \frac{h}{2} \frac{d\mathbf{f}}{d\mathbf{y}} \Big|_i \right)^{-1} \mathbf{f}(\mathbf{y}_i)$$

Stiff systems of ODEs

A system of ODEs can be stiff and require a different solution method. For example:

$$\frac{dc_1}{dt} = 998c_1 + 1998c_2 \quad \frac{dc_2}{dt} = -999c_1 - 1999c_2$$

with boundary conditions $c_1(t = 0) = 1$ and $c_2(t = 0) = 0$.

The analytical solution is:

$$c_1 = 2e^{-t} - e^{-1000t} \quad c_2 = -e^{-t} + e^{-1000t}$$

For the explicit method we require $\Delta t < 10^{-3}$ despite the fact that the term is completely negligible, but essential to keep stability.

The “disease” of stiff equations: we need to follow the solution on the shortest length scale to maintain stability of the integration, although accuracy requirements would allow a much larger time step.

Demonstration with example

Forward Euler (explicit)

$$\frac{c_{1,i+1} - c_{1,i}}{dt} = 998c_{1,i} + 1998c_{2,i}$$

$$\frac{c_{2,i+1} - c_{2,i}}{dt} = -999c_{1,i} - 1999c_{2,i}$$

$$\Rightarrow \begin{aligned} c_{1,i+1} &= (1 + 998\Delta t) c_{1,i} + 1998\Delta t c_{2,i} \\ c_{2,i+1} &= -999\Delta t c_{1,i} + (1 - 1999\Delta t) c_{2,i} \end{aligned}$$

Demonstration with example

Backward Euler (implicit)

$$\frac{dc_{1,i+1} - c_{1,i+1}}{dt} = 998c_{1,i+1} + 1998c_{2,i+1}$$

$$\frac{dc_{2,i+1} - c_{2,i+1}}{dt} = -999c_{1,i+1} - 1999c_{2,i+1}$$

$$\Rightarrow \begin{aligned} (1 - 998\Delta t) c_{1,i+1} - 1998\Delta t c_{2,i} &= c_{1,i} \\ 999\Delta t c_{1,i+1} + (1 + 999\Delta t) c_{2,i+1} &= c_{2,i} \end{aligned}$$

$$A\mathbf{c}_{i+1} = \mathbf{c}_i \text{ with } A = \begin{pmatrix} 1 - 998\Delta t & -1998\Delta t \\ 999\Delta t & 1 + 999\Delta t \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} c_{1,i} \\ c_{2,i} \end{pmatrix}$$

Demonstration with example

Backward Euler (implicit) $A\mathbf{c}_{i+1} = \mathbf{c}_i$ with

$$A = \begin{pmatrix} 1 - 998\Delta t & -1998\Delta t \\ 999\Delta t & 1 + 1999\Delta t \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} c_{1,i} \\ c_{2,i} \end{pmatrix}$$

Cramers rule:

$$c_{1,i+1} = \frac{\begin{vmatrix} c_{1,i} & -1998\Delta t \\ c_{2,i} & 1 + 1999\Delta t \end{vmatrix}}{\det |A|} = \frac{(1+1999\Delta t)c_{1,i} + 1998\Delta t c_{2,i}}{(1-998\Delta t)(1+1999\Delta t) + 1998 \cdot 999 \Delta t^2}$$

$$c_{2,i+1} = \frac{\begin{vmatrix} 1 - 998\Delta t & c_{1,i} \\ 999\Delta t & c_{2,i} \end{vmatrix}}{\det |A|} = \frac{-999\Delta t c_{1,i} + (1-998\Delta t) c_{2,i}}{(1-998\Delta t)(1+1999\Delta t) + 1998 \cdot 999 \Delta t^2}$$

Forward Euler: $\Delta t \leq 0.001$ for stability

Backward Euler: always stable, even for $\Delta t > 100$ (but then not very accurate!)

Demonstration with example

Cure for stiff problems: use implicit methods! To find out whether your system is stiff: check whether one of the eigenvalues have an imaginary part

Implicit methods in Matlab

Matlab offers a stabilized solver, `ode15s`, for stiff problems.

$$\frac{dc_1}{dt} = 998c_1 + 1998c_2 \quad \frac{dc_2}{dt} = -999c_1 - 1999c_2, \quad c_1(0) = 1, \quad c_2(0) = 0$$

- Create the ode function

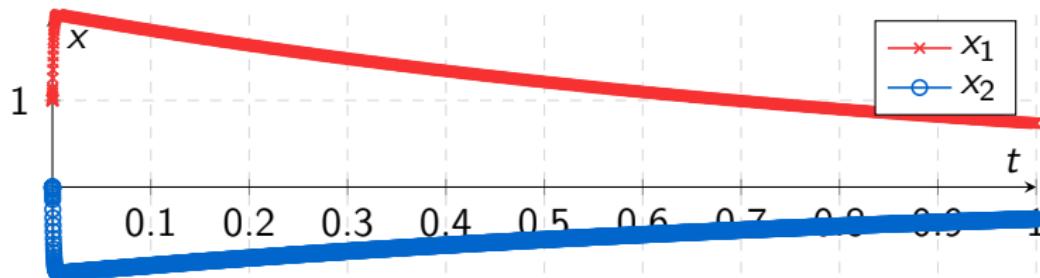
```
function [dcdt] = stiff_ode(t,c)
dcdt = zeros(2,1); % Pre-allocation
dcdt(1) = 998 * c(1) + 1998*c(2);
dcdt(2) = -999 * c(1) - 1999*c(2);
return
```

- Compare the resolution of the solutions

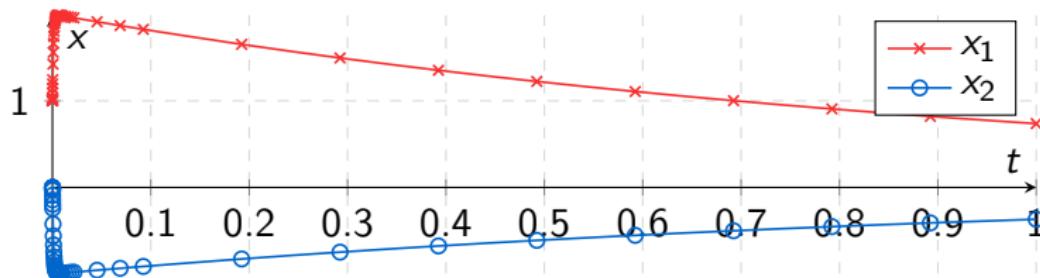
```
subplot(2,1,1);
ode45(@stiff_ode, [0 1], [1 0]);
subplot(2,1,2);
ode15s(@stiff_ode, [0 1], [1 0]);
```

Implicit methods in Matlab

ode45



ode15s



The explicit solver requires 1245 data points (default settings), the implicit solver requires just 48!

Today's outline

⑦ Implicit methods

Backward Euler

Implicit midpoint method

⑧ Systems of ODEs

Solution methods for systems of ODEs

Solving systems of ODEs in Matlab

Stiff systems of ODEs

⑨ Boundary value problems

Shooting method

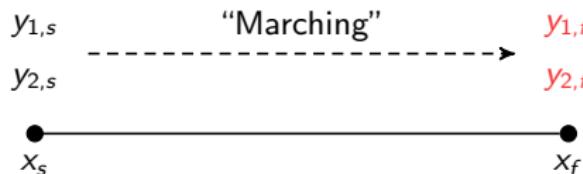
⑩ Conclusion

Importance of boundary conditions

The nature of boundary conditions determines the appropriate numerical method. Classification into 2 main categories:

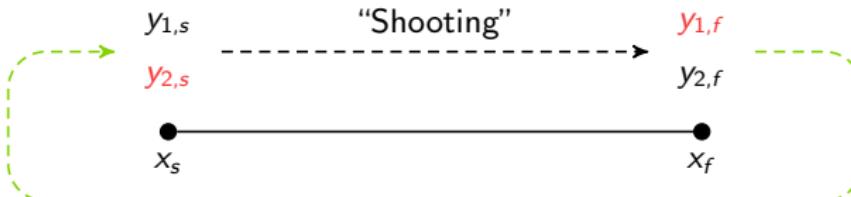
- *Initial value problems (IVP)*

We know the values of all y_i at some starting position x_s , and it is desired to find the values of y_i at some final point x_f .



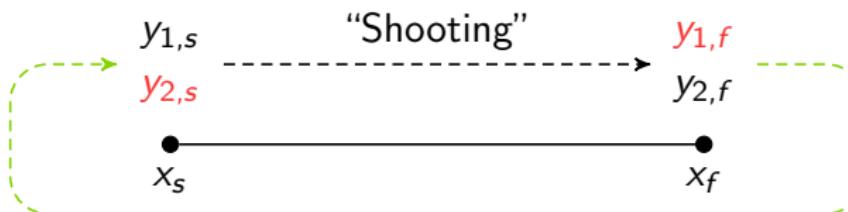
- *Boundary value problems (BVP)*

Boundary conditions are specified at more than one x . Typically, some of the BC are specified at x_s and the remainder at x_f .



Shooting method

How to solve a BVP using the shooting method:



- Define the system of ODEs
- Provide an initial guess for the unknown boundary condition
- Solve the system and compare the resulting boundary condition to the expected value
- Adjust the guessed boundary value, and solve again. Repeat until convergence.
 - Of course, you can subtract the expected value from the computed value at the boundary, and use a non-linear root finding method

BVP: example in Excel

Consider a chemical reaction in a liquid film layer of thickness δ :

$$\mathcal{D} \frac{d^2 c}{dx^2} = k_R c \text{ with } \begin{aligned} c(x=0) &= C_{A,i,L} = 1 && \text{(interface concentration)} \\ c(x=\delta) &= 0 && \text{(bulk concentration)} \end{aligned}$$

Question: compute the concentration profile in the film layer.

Step 1: Define the system of ODEs

This second-order ODE can be rewritten as a system of first-order ODEs, if we define the flux q as:

$$q = -\mathcal{D} \frac{dc}{dx}$$

Now, we find:

$$\frac{dc}{dx} = -\frac{1}{\mathcal{D}} q$$

$$\frac{dq}{dx} = -k_R c$$

BVP: example in Excel

Solving the two first-order ODEs in Excel. First, the cells with constants:

	A	B	C
1	CAiL	1	mol/m3
2	D	1e-8	m2/s
3	kR	10	1/s
4	delta	1e-4	m
5	N	100	
6	dx	=B4/B5	

$$\frac{dc}{dx} = -\frac{1}{D}q$$

$$\frac{dq}{dx} = -k_R c$$

Now, we program the forward Euler (explicit) schemes for c and q below:

	A	B	C
10	x	c	q
11	0	=B1	10
12	=A11+\$B\$6	=B11+\$B\$6*(-1/\$B\$2*C11)	=C11+\$B\$6*(-\$B\$3*B11)
13	=A12+\$B\$6	=B12+\$B\$6*(-1/\$B\$2*C12)	=C12+\$B\$6*(-\$B\$3*B12)
...
111	=A110+\$B\$6	=B110+\$B\$6*(-1/\$B\$2*C110)	=C110+\$B\$6*(-\$B\$3*B110)

BVP: example in Excel

- We now have profiles for c and q as a function of position x .
- The concentration $c(x = \delta)$ depends (eventually) on the boundary condition at the interface $q(x = 0)$
- We can use the solver to change $q(x = 0)$ such that the concentration at the bulk meets our requirement:
 $c(x = \delta) = 0$

BVP: example in Matlab

We first program the system of ODEs in a separate function:

$$\frac{dc}{dx} = -\frac{1}{D}q$$
$$\frac{dq}{dx} = -k_R c$$

```
function [dxdt] = BVPODE(t,x,ps)
dxdt(1)=-1/ps.D*x(2);
dxdt(2)=-ps.kR*x(1);
dxdt=dxdt';
return
```

Note that we pass a variable (type: struct) that contains required parameters: ps.

BVP: example in Matlab

The ODE function is solved via `ode45`, after setting a number of initial and boundary conditions:

```
function f = RunBVP(bcq,ps)
[x,cq] = ode45(@BVPODE,[0 ps.delta],[1 bcq],[],ps);
f = cq(end,1) - 0;
plotyy(x,cq(:,1),x,cq(:,2));
return;
```

Note the following:

- We use the interval $0 \leq x \leq \delta$
- Boundary conditions are given as: $c(x = 0) = 1$ and $q(x = 0) = bcq$, which is given as an argument to the function (i.e. changeable from ‘outside’!)
- The function returns `f`, the difference between the computed and desired concentration at $x = \delta$.

BVP: example in Matlab

Finally, we should solve the system so that we obtain the right boundary condition $q = \text{bcq}$ such that $c(x = \delta) = 0$. We can use the built-in function `fzero` to do this

```
% Parameter definition
ps.D=1e-8;
ps.kR=10;
ps.delta=1e-4;

% Solve for flux boundary condition (initial guess: 0)
opt = optimset('Display','iter');
flux = fzero(@RunBVP,0,opt,ps);
```

BVP example: analytical solution

Compare with the analytical solution:

$$q = k_L E_A C_{A,i,L} \quad \text{with}$$

$$E_A = \frac{Ha}{\tanh Ha} \quad (\text{Enhancement factor})$$

$$Ha = \frac{\sqrt{k_R D}}{k_L} \quad (\text{Hatta number})$$

$$k_L = \frac{D}{\delta} \quad (\text{mass transfer coefficient})$$

Today's outline

⑦ Implicit methods

Backward Euler

Implicit midpoint method

⑧ Systems of ODEs

Solution methods for systems of ODEs

Solving systems of ODEs in Matlab

Stiff systems of ODEs

⑨ Boundary value problems

Shooting method

⑩ Conclusion

Other methods

Other explicit methods:

- Burlisch-Stoer method (Richardson extrapolation + modified midpoint method)

Other implicit methods:

- Rosenbrock methods (higher order implicit Runge-Kutta methods)
- Predictor-corrector methods

Summary

- Several solution methods and their derivation were discussed:
 - Explicit solution methods: Euler, Improved Euler, Midpoint method, RK45
 - Implicit methods: Implicit Euler and Implicit midpoint method
 - A few examples of their spreadsheet implementation were shown
- We have paid attention to accuracy and instability, rate of convergence and step size
- Systems of ODEs can be solved by the same algorithms. Stiff problems should be treated with care.
- An example of solving ODEs with Matlab was demonstrated.

Partial differential equations

Prof.dr.ir. Martin van Sint Annaland, Dr.ir. Ivo Roghair

m.v.sintannaland@tue.nl

Chemical Process Intensification,
Eindhoven University of Technology

Today's outline

① Introduction

② Instationary diffusion equation

Discretization

Solving the diffusion equation

Non-linear source terms

③ Convection

Discretization

Central difference scheme

Upwind scheme

④ Conclusions

Other methods

Summary

Overview

Main question

How to solve parabolic PDEs like:

$$\frac{\partial c}{\partial t} = \mathcal{D} \frac{\partial^2 c}{\partial x^2} - u \frac{\partial c}{\partial x} + R$$

$$t = 0; 0 \leq x \leq \ell \Rightarrow c = c_0$$

with $t > 0; x = 0 \Rightarrow -\mathcal{D} \frac{\partial c}{\partial x} + uc = u_{\text{in}} c_{\text{in}}$

$$t > 0; x = \ell \Rightarrow \frac{\partial c}{\partial x} = 0$$

accurately and efficiently?

What is a PDE?

Partial differential equation

An equation containing a function and their derivatives to multiple independent variables.

Order of PDE

The highest derivative appearing in the PDE

General second order ODE:

$$A \frac{\partial^2 f}{\partial x^2} + B \frac{\partial^2 f}{\partial x \partial y} + C \frac{\partial^2 f}{\partial y^2} + D \frac{\partial f}{\partial x} + E \frac{\partial f}{\partial y} + Ff = G$$

- Linear equation: Coefficients A, B, \dots, G do not depend on x and y .
- Non-linear equation: Coefficients A, B, \dots, G are a function of x and y .

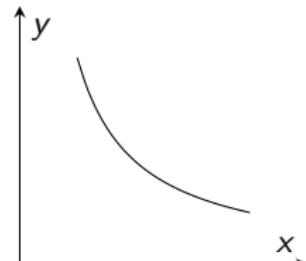
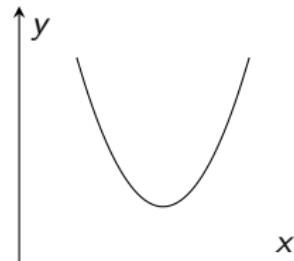
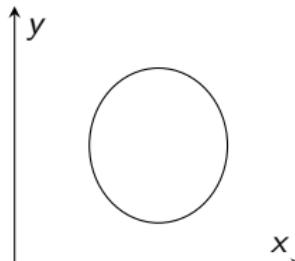
Classification of PDE's

$$A \frac{\partial^2 f}{\partial x^2} + B \frac{\partial^2 f}{\partial x \partial y} + C \frac{\partial^2 f}{\partial y^2} + D \frac{\partial f}{\partial x} + E \frac{\partial f}{\partial y} + Ff = G$$

The discriminant Δ of a quadratic polynomial is computed as (note: only the higher order coefficients are important):

$$\Delta = B^2 - 4AC$$

- $\Delta < 0 \Rightarrow$ Elliptic equation
(e.g. Laplace equation for stationary diffusion in 2D)
- $\Delta = 0 \Rightarrow$ Parabolic equation
(e.g. instationary heat penetration in 1D)
- $\Delta > 0 \Rightarrow$ Hyperbolic equation
(e.g. wave equation)

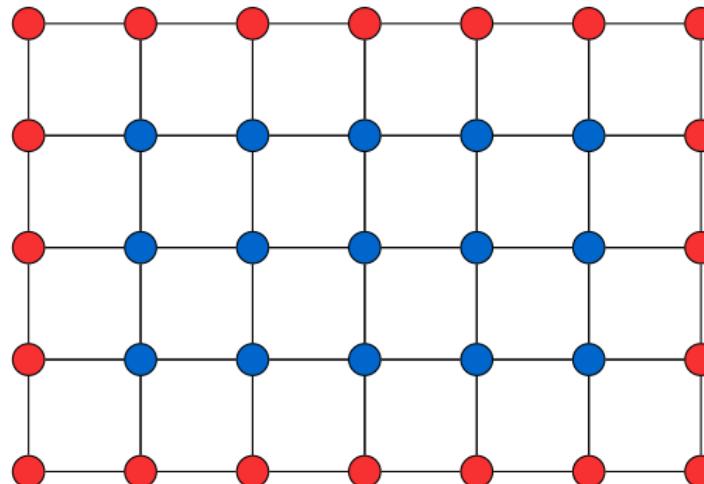


Importance of classification

Different PDE types require different solution techniques because of the difference in range of influence:

- *Characteristics*
Curves in xy -domain along with signal propagation takes place
- *Domain of dependence of point P*
points in xy -domain which influence the value of f in point P
- *Range of influence of point P*
points in xy -domain which are influenced by the value of f in point P

Example elliptic PDE (boundary value problems: BVP)



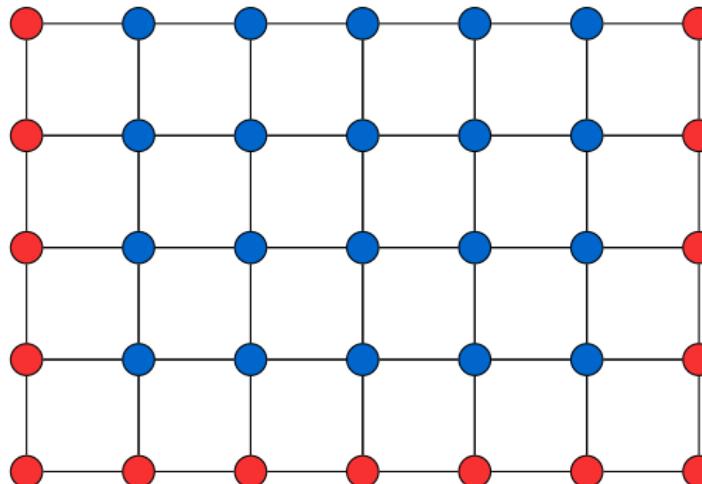
- Grid point at which dependent variable has to be computed
- Grid point at which boundary condition is specified

Typical example: Poisson equation

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} = f(x, y)$$

Efficiency (memory requirements, CPU time) of the numerical method is of crucial importance.

Example parabolic PDE (initial value problem: IVP)



- Grid point at which dependent variable has to be computed
- Grid point at which boundary condition is specified

Typical example: Poisson equation

$$\frac{\partial c}{\partial t} + u \frac{\partial c}{\partial x} = \mathcal{D} \frac{\partial^2 c}{\partial x^2} + R$$

Stability (in numerical sense) of the numerical method is of crucial importance.

Boundary conditions

- Dirichlet or fixed condition: prescribed value of f at boundary

$$f = f_0 \quad f_0 \text{ is a known function}$$

- Neumann condition: prescribed value of derivative of f at boundary

$$\frac{\partial f}{\partial n} = q \quad q \text{ is a known function}$$

- Mixed or Robin condition: relation between f and $\frac{\partial f}{\partial n}$ at boundary

$$a \frac{\partial f}{\partial n} + bf = c \quad a, b \text{ and } c \text{ are known functions}$$

Numerical solution method

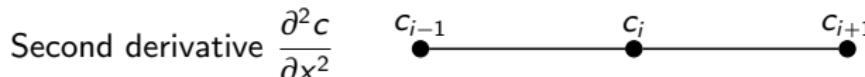
Finite differences (method of lines, MOL):

- ① Discretize spatial domain in discrete grid points
- ② Find suitable approximation for the spatial derivatives
- ③ Substitute approximations in PDE, which gives a system of ODE's, one for every grid points
- ④ Advance in time with a suitable ODE solver

Alternative methods: collocation, Galerkin or Finite Element methods

Instationary diffusion equation (Fick's second law)

$$\frac{\partial c}{\partial t} = \mathcal{D} \frac{\partial^2 c}{\partial x^2}, \quad \text{with} \quad \begin{aligned} t = 0; 0 \leq x \leq \ell &\Rightarrow c = c_0 \\ t > 0; x = 0 &\Rightarrow c = c_L \\ t > 0; x = \ell &\Rightarrow c = c_R \end{aligned}$$



$$c_{i+1} = c_i + \left. \frac{\partial c}{\partial x} \right|_i \Delta x + \frac{1}{2} \left. \frac{\partial^2 c}{\partial x^2} \right|_i \Delta x^2 + \frac{1}{6} \left. \frac{\partial^3 c}{\partial x^3} \right|_i \Delta x^3 + \dots$$

$$c_{i-1} = c_i - \left. \frac{\partial c}{\partial x} \right|_i \Delta x + \frac{1}{2} \left. \frac{\partial^2 c}{\partial x^2} \right|_i \Delta x^2 - \frac{1}{6} \left. \frac{\partial^3 c}{\partial x^3} \right|_i \Delta x^3 + \dots$$

$$c_{i+1} + c_{i-1} = 2c_i + \left. \frac{\partial^2 c}{\partial x^2} \right|_i \Delta x^2 + \mathcal{O}(\Delta x^4)$$

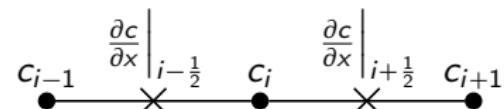
$$\Rightarrow \left. \frac{\partial^2 c}{\partial x^2} \right|_i = \frac{c_{i+1} - 2c_i + c_{i-1}}{\Delta x^2} + \mathcal{O}(\Delta x^2)$$

Due to symmetric discretization: second order (central discretization).

Instationary diffusion equation (Fick's second law)

An alternative discretization:

$$\frac{\partial^2 c}{\partial x^2} \Big|_i = \frac{\frac{\partial c}{\partial x} \Big|_{i+\frac{1}{2}} - \frac{\partial c}{\partial x} \Big|_{i-\frac{1}{2}}}{\Delta x} + \mathcal{O}(\Delta x^2)$$



$$= \frac{\frac{c_{i+1} - c_i}{\Delta x} - \frac{c_i - c_{i-1}}{\Delta x}}{\Delta x} = \frac{c_{i+1} - 2c_i + c_{i-1}}{\Delta x^2}$$

This is convenient for the derivation of $\frac{\partial}{\partial x} \left(\mathcal{D} \frac{\partial c}{\partial x} \right)$:

$$\begin{aligned} \frac{\partial}{\partial x} \left(\mathcal{D} \frac{\partial c}{\partial x} \right) &= \frac{\mathcal{D}_{i+\frac{1}{2}} \frac{\partial c}{\partial x} \Big|_{i+\frac{1}{2}} - \mathcal{D}_{i-\frac{1}{2}} \frac{\partial c}{\partial x} \Big|_{i-\frac{1}{2}}}{\Delta x} = \frac{\mathcal{D}_{i+\frac{1}{2}} \frac{c_{i+1} - c_i}{\Delta x} - \mathcal{D}_{i-\frac{1}{2}} \frac{c_i - c_{i-1}}{\Delta x}}{\Delta x} \\ &= \frac{\mathcal{D}_{i+\frac{1}{2}} c_{i+1} - (\mathcal{D}_{i+\frac{1}{2}} + \mathcal{D}_{i-\frac{1}{2}}) c_i + \mathcal{D}_{i-\frac{1}{2}} c_{i-1}}{(\Delta x)^2} \end{aligned}$$

Instationary diffusion equation (Fick's second law)

$$\frac{\partial^2 f}{\partial x^2} \quad i-1 \quad i-\frac{1}{2} \quad i \quad i+\frac{1}{2} \quad i+1$$

$$f_{i+\frac{1}{2}} = f_i + \frac{1}{2} \Delta x \left. \frac{\partial f}{\partial x} \right|_i \Delta x + \frac{1}{2} \left(\frac{1}{2} \Delta x \right)^2 \left. \frac{\partial^2 f}{\partial x^2} \right|_i + \mathcal{O}(\Delta x^3)$$

$$f_{i-\frac{1}{2}} = f_i - \frac{1}{2} \Delta x \left. \frac{\partial f}{\partial x} \right|_i \Delta x + \frac{1}{2} \left(\frac{1}{2} \Delta x \right)^2 \left. \frac{\partial^2 f}{\partial x^2} \right|_i + \mathcal{O}(\Delta x^3)$$

$$f_{i+\frac{1}{2}} - f_{i-\frac{1}{2}} = \Delta x \frac{\partial f}{\partial x} + \mathcal{O}(\Delta x^3)$$

$$\Rightarrow \left. \frac{\partial f}{\partial x} \right|_i = \frac{f_{i+\frac{1}{2}} - f_{i-\frac{1}{2}}}{\Delta x} + \mathcal{O}(\Delta x^2)$$

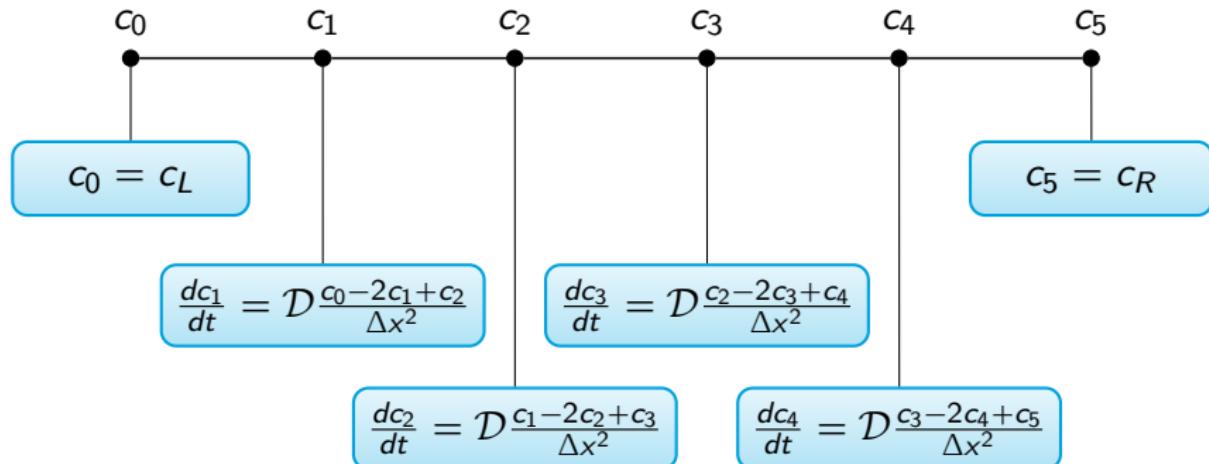
Symmetric discretization yields second order!

Instationary diffusion equation: spatial discretization

Substitution of spatial derivatives yields:

$$\frac{dc_i}{dt} = \mathcal{D} \frac{c_{i-1} - 2c_i + c_{i+1}}{\Delta x^2} \quad \text{for } i = 0, \dots, N$$

For example, using 6 (ridiculously low number!) grid points:



Instationary diffusion equation: boundary conditions

Two options:

- ① Keep boundary conditions as additional equations:

$$c_0 = c_L, \frac{dc_1}{dt} = \mathcal{D} \frac{c_0 - 2c_1 + c_2}{\Delta x^2}, \frac{dc_2}{dt} = \mathcal{D} \frac{c_1 - 2c_2 + c_3}{\Delta x^2},$$
$$\frac{dc_3}{dt} = \mathcal{D} \frac{c_2 - 2c_3 + c_4}{\Delta x^2}, \frac{dc_4}{dt} = \mathcal{D} \frac{c_3 - 2c_4 + c_5}{\Delta x^2}, c_5 = c_R$$

- ② Substitute boundary conditions to reduce number of equations:

$$\frac{dc_1}{dt} = \mathcal{D} \frac{c_L - 2c_1 + c_2}{\Delta x^2}, \frac{dc_2}{dt} = \mathcal{D} \frac{c_1 - 2c_2 + c_3}{\Delta x^2},$$
$$\frac{dc_3}{dt} = \mathcal{D} \frac{c_2 - 2c_3 + c_4}{\Delta x^2}, \frac{dc_4}{dt} = \mathcal{D} \frac{c_3 - 2c_4 + c_R}{\Delta x^2}$$

Instationary diffusion equation: temporal discretization

$$\frac{dc_i}{dt} = \mathcal{D} \frac{c_{i-1} - 2c_i + c_{i+1}}{\Delta x^2}$$

Time discretization: forward Euler (explicit)

$$\frac{c_i^{n+1} - c_i^n}{\Delta t} = \mathcal{D} \frac{c_{i-1}^n - 2c_i^n + c_{i+1}^n}{\Delta x^2}$$

$$\Rightarrow c_i^{n+1} = Foc_{i-1}^n + (1 - 2Fo)c_i^n + Foc_{i+1}^n \quad \text{with } Fo = \frac{\mathcal{D}\Delta t}{\Delta x^2}$$

Straightforward updating (explicit equation), simple to implement in a program but stability constraint $Fo = \frac{\mathcal{D}\Delta t}{\Delta x^2} < \frac{1}{2}$!

Small $\Delta x \Rightarrow$ small $\Delta t \Rightarrow$ patience required ☺

Instationary diffusion equation: temporal discretization

$$\frac{dc_i}{dt} = \mathcal{D} \frac{c_{i-1} - 2c_i + c_{i+1}}{\Delta x^2}$$

Time discretization: backward Euler (implicit)

$$\frac{c_i^{n+1} - c_i^n}{\Delta t} = \mathcal{D} \frac{c_{i-1}^{n+1} - 2c_i^{n+1} + c_{i+1}^{n+1}}{\Delta x^2}$$

$$\Rightarrow -Foc_{i-1}^{n+1} + (1+2Fo)c_i^{n+1} - Foc_{i+1}^{n+1} = c_i^n \quad \text{with } Fo = \frac{\mathcal{D}\Delta t}{\Delta x^2}$$

Requires the solution of a system of linear equations, but no stability constraints!

Note: extension to higher order schemes (with time step adaptation) straightforward.
Often second or third order optimal, because for each Euler-like step in the additional order an often large system needs to be solved (not treated in this course).

Solving the instationary diffusion equation: example

Solve the diffusion problem using explicit discretization:

$$\frac{\partial c_i}{\partial t} = \mathcal{D} \frac{\partial^2 c}{\partial x^2} \quad \text{with} \quad \begin{aligned} 0 &\leq x \leq \delta, \delta = 5 \cdot 10^{-3} \text{ m} \\ \delta/\Delta x &= 100 \text{ grid cells} \\ \mathcal{D} &= 1 \cdot 10^{-8} \text{ m}^2 \text{s}^{-1} \\ t_{\text{end}} &= 5000 \text{ s} \\ c_L &= 1 \text{ mol m}^{-3} \quad c_R = 0 \text{ mol m}^{-3} \end{aligned}$$

$$c_i^{n+1} = Foc_{i-1}^n + (1 - 2Fo)c_i^n + Foc_{i+1}^n \quad \text{with } Fo = \frac{\mathcal{D}\Delta t}{\Delta x^2}$$

- ① Initialise variables
- ② Compute time step so that $Fo \leq \frac{1}{2} \Rightarrow \Delta t = 0.125\text{s}!$
- ③ Compute 40000 time steps times 100 grid nodes!
- ④ Store solution

Solving the instationary diffusion equation: example

Initialise the variables and matrices:

```
Nx = 100; % Nc grid points
Nt = 40000; % Nt time steps
D = 1e-8; % m/s
c_L = 1.0; c_R = 0; % mol/m3
t_end = 5000.0; % s
x_end = 5e-3; % m

% Time step and grid size
dt = t_end/Nt;
dx = x_end/Nx;

% Fourier number
Fo=D*dt/dx/dx

% Initial matrices for solutions (Nx times Nt)
c = zeros(Nt+1,Nx+1); % All concentrations are zero
c(:,1) = c_L; % Concentration at left side
c(:,Nx+1) = c_R; % Concentration at right side

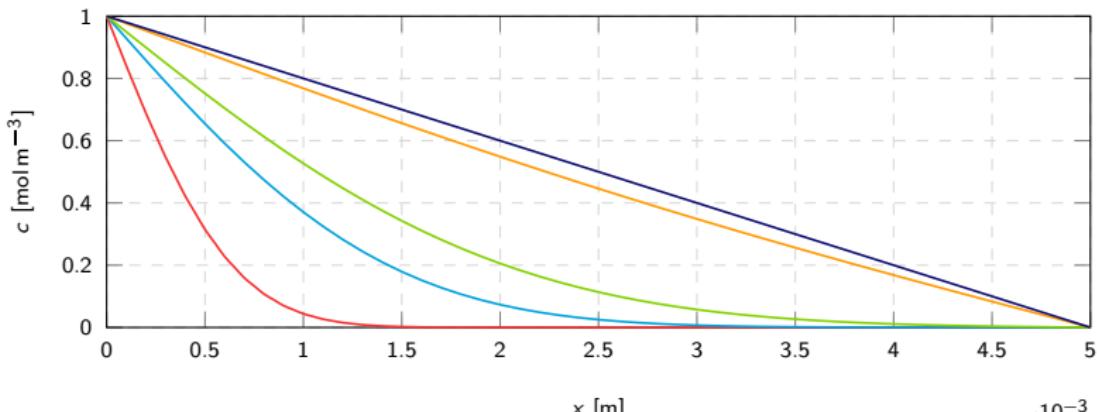
% Grid node and time step positions
x = linspace(0,x_end,Nx+1);
```

Solving the instationary diffusion equation: example

Compute the solution (nested time-and-grid loop):

```
for n = 1:Nt % time loop
    for i = 2:Nx % Nested loop for grid nodes
        c(n+1,i) = Fo*c(n,i-1) + (1-2*Fo)*c(n,i) + ...
                    Fo*c(n,i+1);
    end
end
```

Plotting the solution at $t = \{12.5, 62.5, 125, 625, 5000\}$ s.



Solving the diffusion equation implicitly

Linear system $Ax = b$ from $-Foc_{i-1}^{n+1} + (1 + 2Fo)c_i^{n+1} - Foc_{i+1}^{n+1} = c_i^n$

$$\begin{pmatrix} 1 & 0 & 0 & 0 & \cdots & 0 \\ -Fo & (1 + 2Fo) & -Fo & 0 & \cdots & 0 \\ 0 & -Fo & (1 + 2Fo) & -Fo & \cdots & 0 \\ 0 & 0 & -Fo & (1 + 2Fo) & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 \end{pmatrix} \cdot \begin{pmatrix} c_0^{n+1} \\ c_1^{n+1} \\ c_2^{n+1} \\ c_3^{n+1} \\ \vdots \\ c_m^{n+1} \end{pmatrix} = \begin{pmatrix} c_0^n \\ c_1^n \\ c_2^n \\ c_3^n \\ \vdots \\ c_m^n \end{pmatrix}$$

$$1 \times c_0^{n+1} = c_0^n \text{ (boundary condition)}$$

$$-Foc_0^{n+1} + (1 + 2Fo)c_1^{n+1} - Foc_2^{n+1} = c_1^n$$

$$-Foc_1^{n+1} + (1 + 2Fo)c_2^{n+1} - Foc_3^{n+1} = c_2^n$$

$$-Foc_2^{n+1} + (1 + 2Fo)c_3^{n+1} - Foc_4^{n+1} = c_3^n$$

$$1 \times c_m^{n+1} = c_m^n \text{ (boundary condition)}$$

Solving the diffusion equation implicitly in Matlab

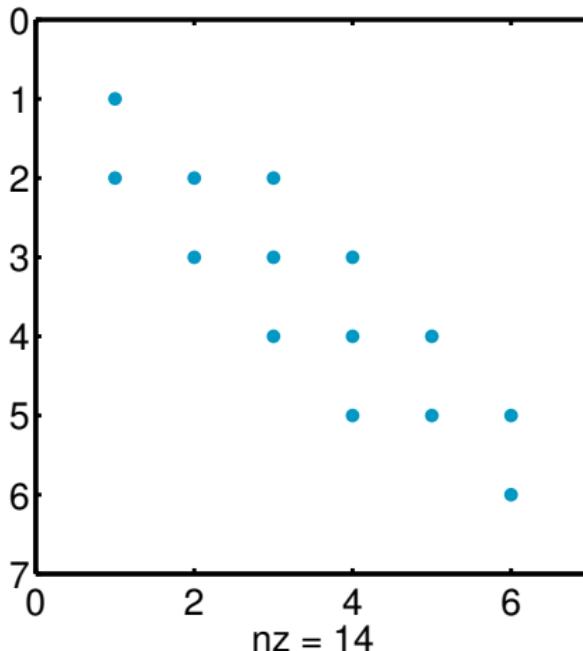
To solve the linear system, we need to define matrix A . It is clear that storing many zeros is not efficient in terms of memory. We use a *sparse matrix* format:

```
% Bands in matrix (internal cells)
A = sparse(Nx+1,Nx+1);
for i=2:Nx
    A(i,i-1) = -Fo;
    A(i,i) = (1+2*Fo);
    A(i,i+1) = -Fo;
end

% Set boundary cells, independent on neighbors:
A(1,1) = 1; % Left
A(Nx+1,Nx+1) = 1; % Right
```

Solving the diffusion equation implicitly in Matlab

The command `spy(A)` shows a figure with the non-zero positions.



Solving the diffusion equation implicitly in Matlab

The concentration matrix is initialised and the boundary conditions are set as follows:

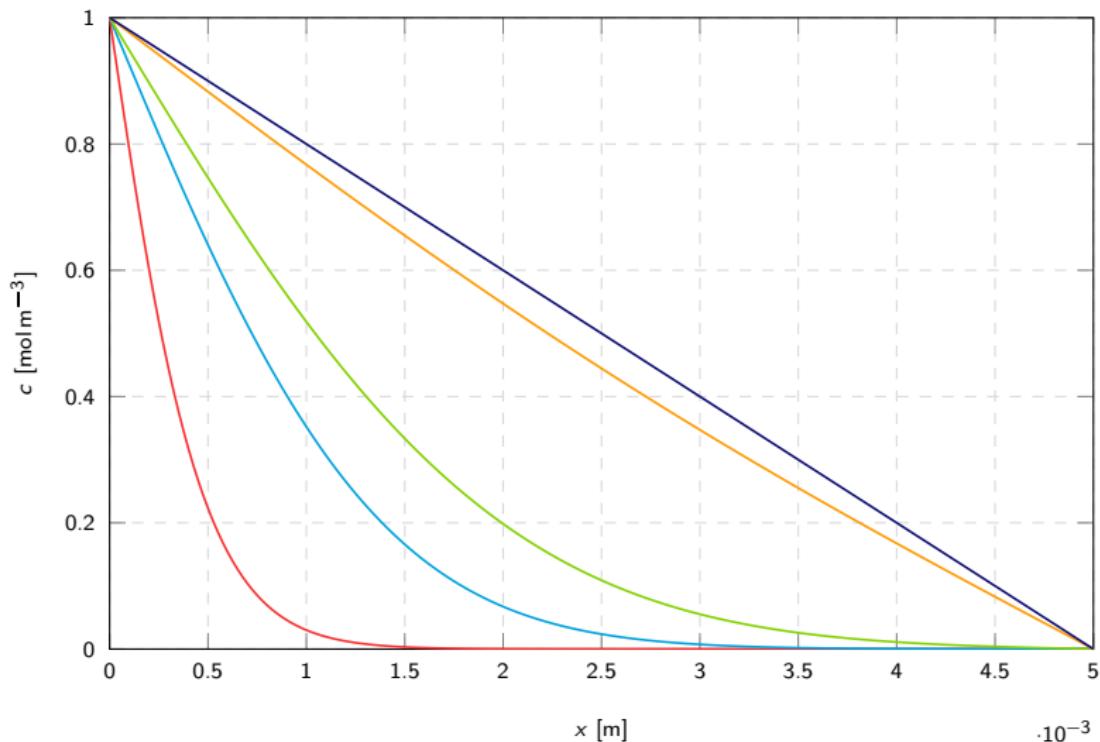
```
% Initial matrices for solutions (Nx times Nt)
c = zeros(Nt+1,Nx+1); % All concentrations are zero
c(:,1) = c_L;          % Concentration at left side
c(:,Nx+1) = c_R;       % Concentration at right side
```

The right hand side vector (b) can now be set during the time-loop:

```
for n = 1:Nt-1           % time loop
    b = c(n,:)';
    solX = A\b;
    c(n+1,:) = solX;      % Store solution each time step
end
```

Solving the diffusion equation implicitly in Matlab

Plotting the solution at $t = \{12.5, 62.5, 125, 625, 5000\}$ s.



About explicit vs. implicit solutions

- Explicit solution:
 - Easy to implement
 - Very small time steps required.
 - This problem took about 0.5 s.
- Implicit solution:
 - Harder to implement, needs sparse matrix solver
 - No stability constraint
 - This problem took about 0.05 s
- The difference will become much larger for systems with e.g. more grid nodes!

Extension with non-linear source terms

$$\frac{\partial c}{\partial t} = \mathcal{D} \frac{\partial^2 c}{\partial x^2} + R(c) \quad \text{with} \quad \begin{aligned} t = 0; 0 \leq x \leq \ell &\Rightarrow c = c_0 \\ t > 0; x = 0 &\Rightarrow c = c_L \\ t > 0; x = \ell &\Rightarrow c = c_R \end{aligned}$$

- Forward Euler (explicit): simply add to right-hand side

$$\begin{aligned} \frac{c_i^{n+1} - c_i^n}{\Delta t} &= \mathcal{D} \frac{c_{i-1}^n - 2c_i^n + c_{i+1}^n}{\Delta x^2} + R(c_i^n) \\ \Rightarrow c_i^{n+1} &= Foc_{i-1}^n + (1 - 2Fo)c_i^n + Foc_{i+1}^n + R_i^n \Delta t \end{aligned}$$

- Backward Euler (implicit): linearization required

$$\begin{aligned} R(c_i^{n+1}) &= R(c_i^n) + \left. \frac{dR}{dc} \right|_i^n (c_i^{n+1} - c_i^n) \\ \frac{c_i^{n+1} - c_i^n}{\Delta t} &= \mathcal{D} \frac{c_{i-1}^{n+1} - 2c_i^{n+1} + c_{i+1}^{n+1}}{\Delta x^2} + R(c_i^{n+1}) \end{aligned}$$

$$\Rightarrow -Foc_{i-1}^{n+1} + (1 + 2Fo - \left. \frac{dR}{dc} \right|_i^n \Delta t) c_i^{n+1} - Foc_{i+1}^{n+1} = c_i^n + \left(R_i^n - \left. \frac{dR}{dc} \right|_i^n c_i^n \right) \Delta t$$

Extension with convection terms

$$\frac{\partial c}{\partial t} = \mathcal{D} \frac{\partial^2 c}{\partial x^2} - u \frac{\partial c}{\partial x} + R$$

Discretization of first derivative $\frac{dc}{dx}$,
looks simple but is numerical headache!

Central discretization:

$$\frac{dc}{dx} = \frac{c_{i+1} - c_{i-1}}{2\Delta x}$$

⇒ simple and easy, too bad it doesn't work: yields unstable solutions if convection dominated.

Central difference scheme of 1st derivative

Unsteady convection:

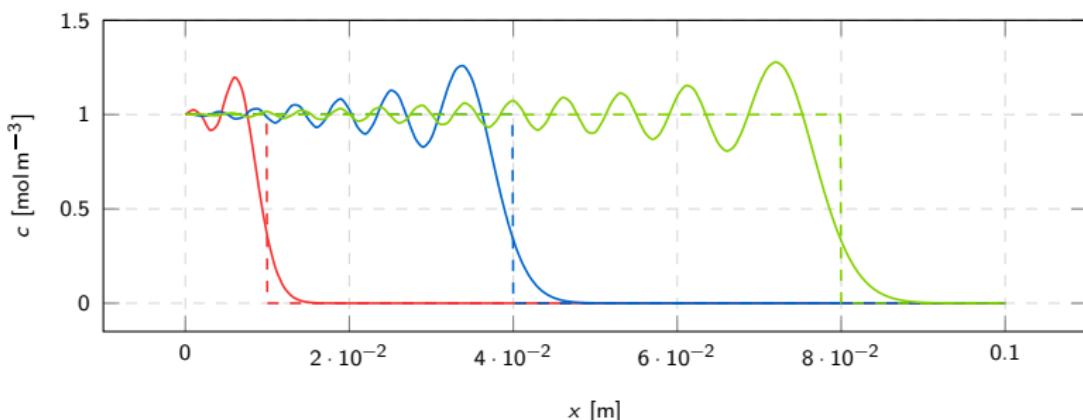
$$\frac{\partial c}{\partial t} = -u \frac{\partial c}{\partial x}$$

Central difference for first derivative:

$$\frac{dc}{dx} = \frac{c_{i+1} - c_{i-1}}{2\Delta x}$$

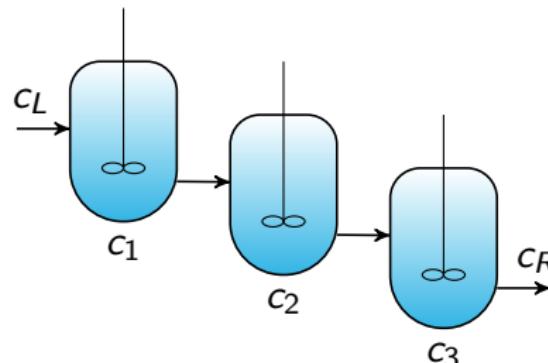
Forward Euler discretization of temporal and spatial domain:

$$\frac{c_i^{n+1} - c_i^n}{\Delta t} = -u \frac{c_{i+1} - c_{i-1}}{2\Delta x} \Rightarrow c_i^{n+1} = c_i^n - u \frac{c_{i+1}^n - c_{i-1}^n}{2\Delta x} \Delta t$$



Extension with convection terms

Solution: upwind discretization, like CSTR's in series:



First order upwind: $-u \frac{dc}{dx} \Big|_i = \begin{cases} -u \frac{c_i - c_{i-1}}{\Delta x} & \text{if } u \geq 0 \\ -u \frac{c_{i+1} - c_i}{\Delta x} & \text{if } u < 0 \end{cases}$

Stable if $Co = \frac{u\Delta t}{\Delta x} < 1$ (with Co the Courant number). However, only 1st order accurate (large smearing of concentration fronts). Higher order upwind requires TVD schemes (trick of the trade)...

First order upwind scheme of 1st derivative

Unsteady convection:

$$\frac{\partial c}{\partial t} = -u \frac{\partial c}{\partial x}$$

Upwind scheme for first derivative:

$$-u \frac{dc}{dx} \Big|_i = \begin{cases} -u \frac{c_i - c_{i-1}}{\Delta x} & \text{if } u \geq 0 \\ -u \frac{c_{i+1} - c_i}{\Delta x} & \text{if } u < 0 \end{cases}$$

Forward Euler discretization of temporal and spatial domain:

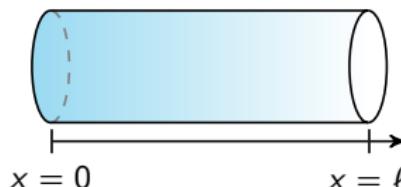
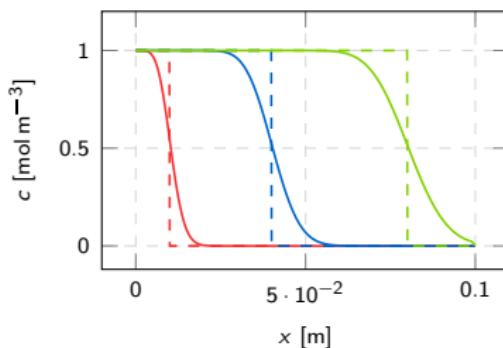
$$\frac{c_i^{n+1} - c_i^n}{\Delta t} = -u \frac{c_{i+1} - c_{i-1}}{2\Delta x}$$
$$\Rightarrow c_i^{n+1} = \begin{cases} c_i^n - u \frac{c_i - c_{i-1}}{\Delta x} & \text{if } u \geq 0 \\ c_i^n - u \frac{c_{i+1} - c_i}{\Delta x} & \text{if } u < 0 \end{cases}$$

Upwind scheme: example

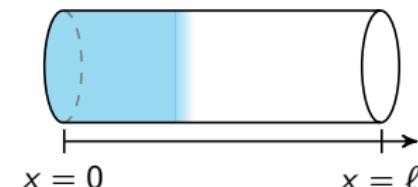
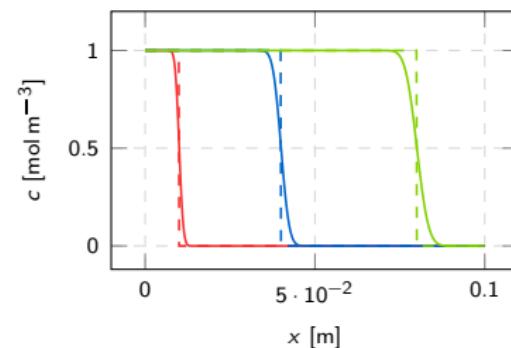
Unsteady convection through a pipe:

$$\frac{\partial c}{\partial t} = -u \frac{\partial c}{\partial x} \quad \text{with} \quad u = 0.1 \text{ m s}^{-1} \Rightarrow c_i^{n+1} = c_i^n - u \frac{c_i^n - c_{i-1}^n}{\Delta x} \Delta t$$

Using 100 grid cells



Using 1000 grid cells



Central difference and upwind in Matlab

The results from the previous slides were computed using this script:

```
Nx = 1000; % Nc grid points
Nt = 10000; % Nt time steps
u = 0.001; % m/s
c_in = 1.0; % mol/m3
t_end = 100.0; % s
x_end = 0.1; % m

% Time step and grid size
dt = t_end/Nt; dx = x_end/Nx;

% Courant number
Co=u*dt/dx

% Initial matrices for solutions (Nx times Nt)
c1 = zeros(Nt+1,Nx+1); % All concentrations are zero
c1(:,1) = c_in; % Concentration at inlet (all time steps)
an = c1; c2 = c1; % Analytical and upwind solution

% Grid node and time step positions
x = linspace(0,x_end,Nx+1);
t = linspace(0,t_end,Nt+1);
```

Central difference and upwind in Matlab

(continued)

```
for n = 1:Nt % time loop
    for i = 2:Nx % Nested loop for grid nodes
        % Central difference
        c1(n+1,i) = c1(n,i) - u*((c1(n,i+1) - ...
            c1(n,i-1))/(2*dx))*dt;
        % Upwind
        c2(n+1,i) = c2(n,i) - u*((c2(n,i) - ...
            c2(n,i-1))/(dx))*dt;
        % Analytical
        an(n+1,i) = (x(i) < u*t(n+1))*c_in;
    end
end
```

Extension to systems of PDE's

- Explicit methods: straightforward extension
- Implicit methods: yields block-tridiagonal matrix (note ordering of equations: all variables per grid cell)

Extension to 2D or 3D systems

Spatial discretization in 2 directions — different methods available:

- Explicit
- Fully implicit
 - 1D gives tri-diagonal matrix
 - 2D gives penta-diagonal matrix
 - 3D gives hepta-diagonal matrix

Use of dedicated matrix solvers (e.g. ICCG, multigrid, ...)

- Alternating direction implicit (ADI)
 - Per direction implicit, but still overall unconditionally stable

Further extensions for parabolic PDEs

- Higher order temporal discretization (multi-step) with time step adaptation
- Non-uniform grids with automatic grid adaptation
- Higher-order discretization methods, especially higher order TVD (flux delimited) schemes for convective fluxes (e.g. WENO schemes)
- Higher-order finite volume schemes (Riemann solvers)

Summary

- Several classes of PDEs were introduced
 - Elliptic, Parabolic, Hyperbolic PDEs
- Diffusion equation: discretization of temporal and spatial domain was discussed
 - Solutions of the diffusion equation using explicit and implicit methods
 - How to add non-linear source terms
- Convection: upwind vs. central difference schemes

Curve fitting, data regression and optimization

Ivo Roghair, Martin van Sint Annaland

Chemical Process Intensification,
Eindhoven University of Technology

Today's outline

- ① Introduction
 - ② Curve fitting
 - ③ Regression
 - ④ Fitting numerical models
 - ⑤ Optimization
 - ⑥ Linear programming
 - ⑦ Summary

Today's outline

① Introduction

② Curve fitting

③ Regression

④ Fitting numerical models

⑤ Optimization

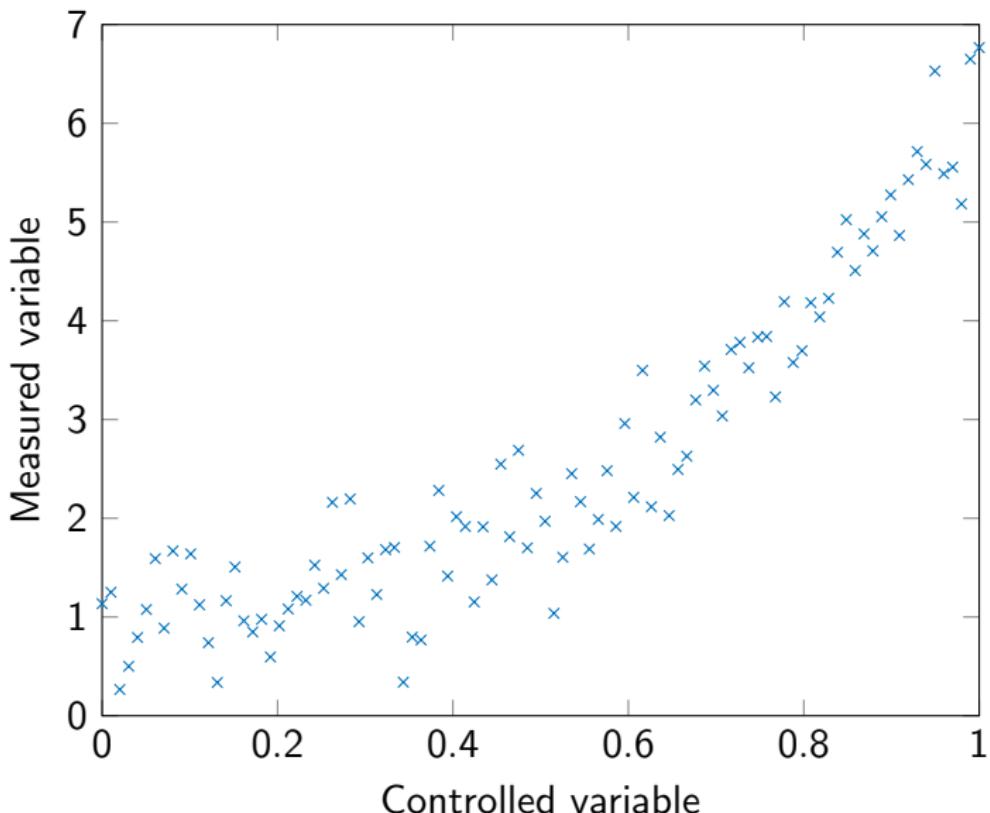
⑥ Linear programming

⑦ Summary

Overview

- We are going to fit measurements to models today
- You will also learn what R^2 actually means
- We get introduced to constrained and unconstrained optimization.
- We will use the simplex method to solve linear programming problems (LP)

Fitting models to data



How to fit a model to the data?

We would like to fit the following model to the data:

$$\hat{y} = a_1 + a_2x + a_3x^2 + a_4x^3$$

How to fit a model to the data?

We would like to fit the following model to the data:

$$\hat{y} = a_1 + a_2x + a_3x^2 + a_4x^3$$

First step: If we have N data points, we could write the model as the product of a matrix and a vector:

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ 1 & x_3 & x_3^2 & x_3^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}$$

$$\hat{y} = Xa$$

X is called the design matrix and a are the fit parameters.

Residuals

Second step: work out the residuals for each data point:

$$d_i = (y_i - \hat{y}_i)$$

Residuals

Second step: work out the residuals for each data point:

$$d_i = (y_i - \hat{y}_i)$$

Third step: work out the sum of squares of the residuals:

$$\text{SSE} = \sum_i d_i^2 = \sum_i (y_i - \hat{y}_i)$$

This can be written using the dot-product operation:

$$\text{SSE} = \sum_i d_i^2 = d \cdot d = d^T \cdot d = (y_i - \hat{y}_i)^T \cdot (y_i - \hat{y}_i)$$

Minimizing the sum of squares

Choose the parameter vector such that the sum of squares of the residuals is minimized; the partial derivative with respect to each parameter should be set to zero:

$$\frac{\partial}{\partial a_i} \left[\left(y - (Xa)^T \right) (y - Xa) \right]$$

Minimizing the sum of squares

Choose the parameter vector such that the sum of squares of the residuals is minimized; the partial derivative with respect to each parameter should be set to zero:

$$\frac{\partial}{\partial a_i} \left[(y - (Xa)^T) (y - Xa) \right]$$

In Matlab, we can solve our linear system $\hat{Y} = Xa$ simply by running `a = X\y`.

Minimizing the sum of squares

Choose the parameter vector such that the sum of squares of the residuals is minimized; the partial derivative with respect to each parameter should be set to zero:

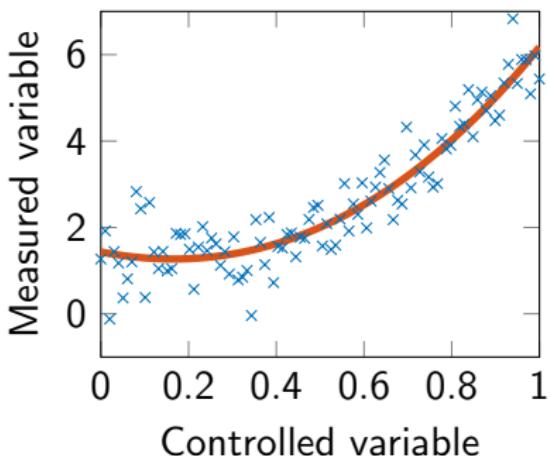
$$\frac{\partial}{\partial a_i} \left[(y - (Xa)^T) (y - Xa) \right]$$

In Matlab, we can solve our linear system $\hat{Y} = Xa$ simply by running $a = X\backslash y$.

- If there are more data points ($N > 4$), we can write an analogue, but maybe a consistent solution does not exist (the system is over specified).
- However, matlab will find values for the vector a which minimize $\|y - aX\|^2$, so i.e. a least squares fit.

Fitting our problem: Matlab solver

```
N=length(x);  
X(:,1) = ones(N,1);  
X(:,2) = x;  
X(:,3) = x.^2;  
X(:,4) = x.^3;  
A = X\y;
```



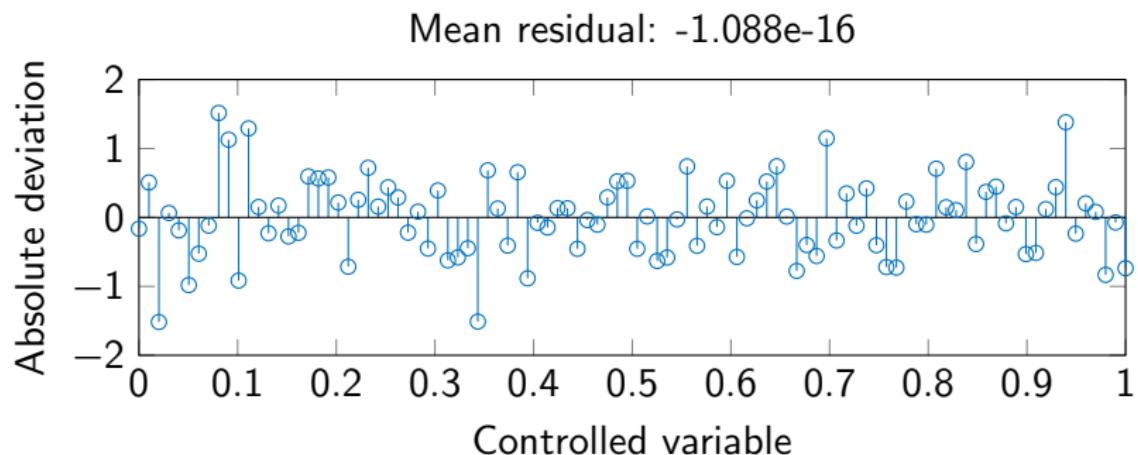
Fitting our problem: Matlab curve-fitting toolbox

- Start the toolbox: `cftool`
- Choose the dataset (x and y data)
- Choose the interpolant type (polynomial, exponential, ..., custom)
- Get the coefficients (save to workspace or write them down)

Fitting our problem: Excel

- Create a column with the independent and dependent data series (x and y)
- Create a column that computes \hat{y} , keeping the coefficients as separate cells
- Compute the sum of squares of the residuals (another column, sum the results)
- Use the solver to minimize this sum, modifying the coefficient cells
- Note: regression in Excel + display equation is dangerous if you choose the 'line' plot (use scatter if you can)

How good is the model?



- For a model to make sense the data points should be scattered randomly around the model predictions, the mean of the residuals d should be zero: $d_i = (y_i - \hat{y}_i)$
- It's always good to check if the residuals are not correlated with the measured values, if that is the case, it can indicate that your model is wrong.

Regression coefficients

- Variance measured in the data (y) is:

$$\sigma_y^2 = \frac{1}{N} \sum_i (y_i - \bar{y})^2$$

- Variance of the residuals is:

$$\sigma_{\text{error}}^2 = \frac{1}{N} \sum_i (d_i)^2$$

- Variance in the model is:

$$\sigma_{\text{model}}^2 = \frac{1}{N} \sum_i (\hat{y}_i - \bar{\hat{y}})^2$$

Regression coefficients

Given that the error is uncorrelated we can state that:

$$\sigma_y^2 = \sigma_{\text{error}}^2 + \sigma_{\text{model}}^2$$

$$R^2 = \frac{\sigma_{\text{model}}^2}{\sigma_y^2} = 1 - \frac{\sigma_{\text{error}}^2}{\sigma_y^2}$$

$$R^2 = 1 - \frac{\text{SSE}}{\text{SST}}$$

- SSE: Sum of errors squared
- SST: Total sum of squares (model)
- SSR: Sum of squares (data)

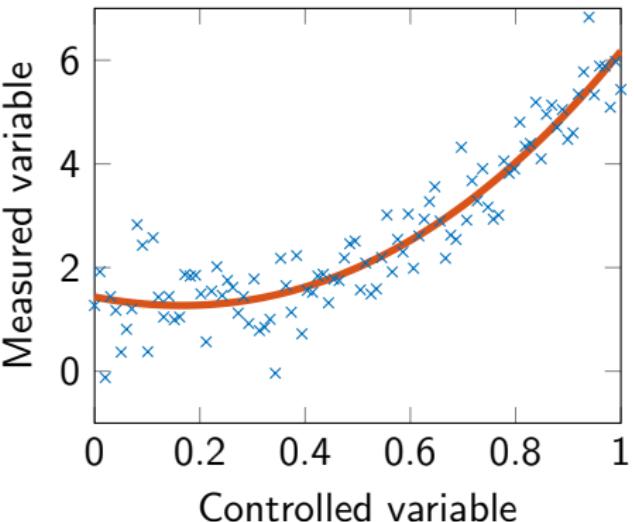
Regression coefficients

- An uncorrelated error ($\bar{d} \rightarrow 0$) means that SSE, SST and SSR will have χ^2 -distributions and the ratios will have an F-distribution. If SSR/SSE is large, the model is good!
- There is a chance that the model is rubbish, but that SSR/SSE will yield a good value, Analysis of Variance (ANOVA) will be a good tool to calculate the probability of such a thing happening!

Back to the example

The statistics:

	Value
N	100
SSE	32.042
SST	896.907
SSR	928.950
R^2	0.964



Dynamic fitting of non-linear equations: lsqnonlin

You may encounter situations where the model data is slightly more complicated to obtain (e.g. a numerical model based on ODEs where coefficients are unknown), or you want to perform fitting of multiple functions/coefficients, or just want to automate things via scripts. Matlab's Optimization toolbox gives access to a powerful function `lsqnonlin`, least-squares non-linear optimization.

General use of lsqnonlin

```
k = lsqnonlin(fun,k0,lb,ub,options)
```

- `fun` is a function handle to the fit criterium (e.g. `@myFitCrit`). The fit criterium function `myFitCrit` should return the residuals vector, e.g. $d_i = (y_i - \hat{y}_i)$. Here, y_i would again be the measurement data and \hat{y} the solution computed by a model.
- `k0` is the initial guess for the fitting coefficient (or: array of initial guesses when fitting multiple coefficients)
- `lb` and `ub` are the lower and upper boundaries for `k0`. These should both be the size of the `k0`-array.
- `options` are some fitting options, for more fine-grained control on the fit procedure. Use e.g.
`options = optimset(TolX ,1.0E-6, MaxFunEvals ,1000);` to create an options object, or leave it empty (`options = []`).

Example use of lsqnonlin

We have experimental data stored in the file `tudataset1.mat`, containing T and U data. We want to fit a model with coefficients k_1 and k_2 with the following structure:

$$\frac{du}{dt} = -k_1 u + k_2$$

Example use of lsqnonlin

We have experimental data stored in the file `tudataset1.mat`, containing T and U data. We want to fit a model with coefficients k_1 and k_2 with the following structure:

$$\frac{du}{dt} = -k_1 u + k_2$$

- First, we need to create a function that contains the ODE:

```
function dudt = simpleode(t,u,k);
dudt = -k(1)*u + k(2);
```

Note that we supply a vector `k`, containing both coefficients for fitting

Example use of lsqnonlin

We have experimental data stored in the file `tudataset1.mat`, containing T and U data. We want to fit a model with coefficients k_1 and k_2 with the following structure:

$$\frac{du}{dt} = -k_1 u + k_2$$

- First, we need to create a function that contains the ODE:

```
function dudt = simpleode(t,u,k);
dudt = -k(1)*u + k(2);
```

Note that we supply a vector `k`, containing both coefficients for fitting

- We create a fit criterium function:

```
function err = fitcrit(ke,T,U,U0)
[t,ue] = ode45(@simpleode,T,U0,[],ke);
err = (ue-U);
```

Example use of lsqnonlin

Now let's make a script that uses `lsqnonlin` to yield k-values fitted to our dataset:

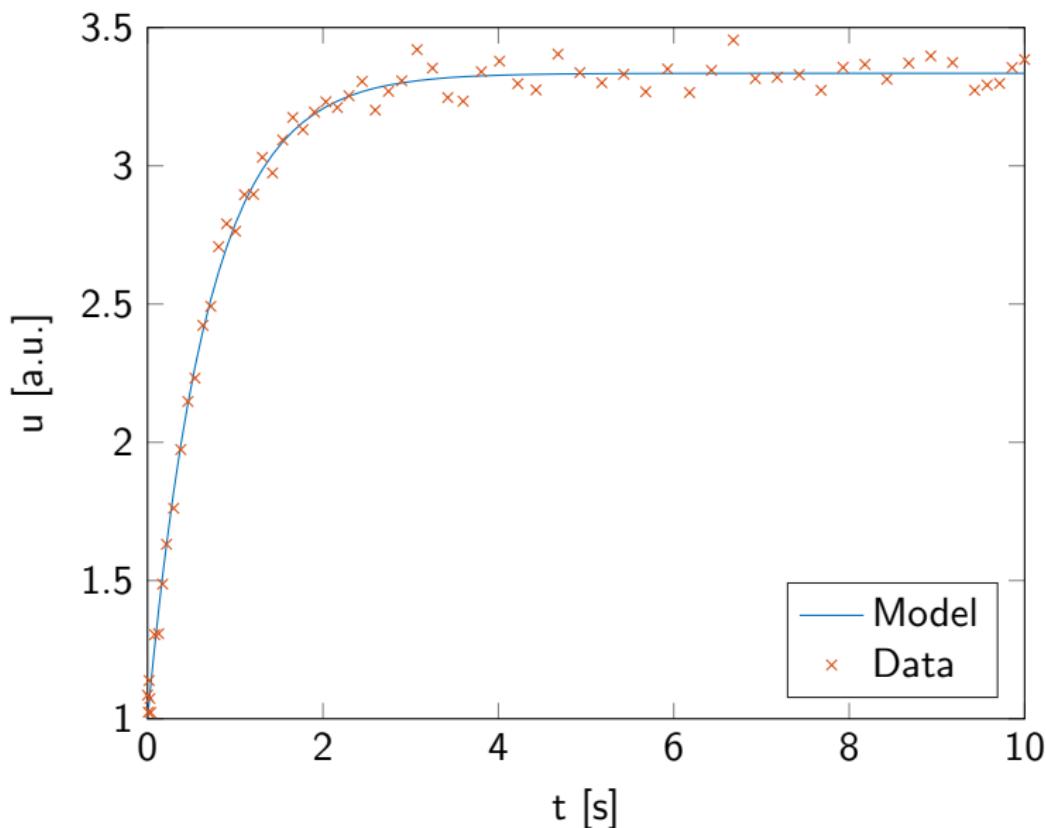
Example use of lsqnonlin

Now let's make a script that uses `lsqnonlin` to yield k-values fitted to our dataset:

```
% initial value
U0 = 1.00;
% initial guesses for model parameters
k0 = [1.00 1.00];
% lower and upper bounds for model parameters
LB = [0.00 0.00];
UB = [Inf Inf];
% Perform nonlinear least squares fit
options = optimset('TolX',1.0E-6,'MaxFunEvals',1000);
[ke,RESNORM,RESIDUAL,EXITFLAG,OUTPUT,LAMBDA,JACOBIAN]
    = lsqnonlin(@fitcrit,k0,LB,UB,options,T,U,U0);
```

Our fitted coefficients are stored in `ke`. Note that we get a lot more data back that allows to check the fitting results in more detail.

Example use of lsqnonlin



Today's outline

- ① Introduction
- ② Curve fitting
- ③ Regression
- ④ Fitting numerical models
- ⑤ Optimization
- ⑥ Linear programming
- ⑦ Summary

What is optimization?

Optimization is minimization or maximization of an objective function (also called a performance index or goal function) that may be subject to certain constraints.

- $\min f(x)$: Goal function
- $g(x) = 0$: Equality constraints
- $h(x) \geq 0$: Inequality constraints

Optimization Spectrum

Problem	Method	Solvers
LP	Simplex method	Linprog (Matlab)
	Barrier methods	CPLEX (GAMS, AIMMS, AMPL, OPB)
NLP	Lagrange multiplier method	Fminsearch/fmincon (Matlab)
	Successive linear programming	MINOS (GAMS, AMPL)
	Quadratic programming	CONOPT (GAMS)
MIP	Branch and bound	
MILP	Dynamic programming	Bintprog (Matlab)
MINLP	Generalized Benders decomposition	DICOPT (GAMS)
	Outer approximation method	BARON (GAMS)
MIQP	Disjunctive programming	

Factors of concern

- Continuity of the functions
- Convexity of the functions
- Global versus local optima
- Constrained versus unconstrained optima

Linear programming

In linear programming the objective function and the constraints are linear functions.

Linear programming

In linear programming the objective function and the constraints are linear functions.

For example:

$$\max z = f(x_1, x_2) = 40x_1 + 88x_2$$

s.t. (subject to)

$$2x_1 + 8x_2 \leq 60$$

$$5x_1 + 2x_2 \leq 60$$

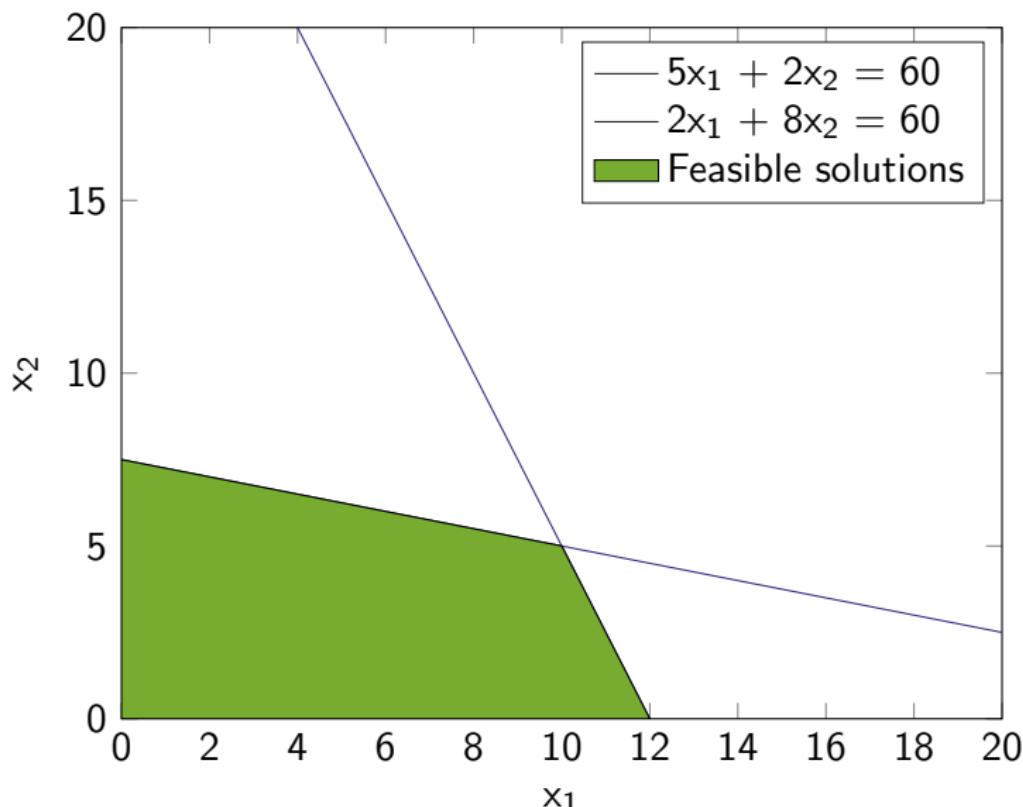
$$x_1 \geq 0$$

$$x_2 \geq 0$$

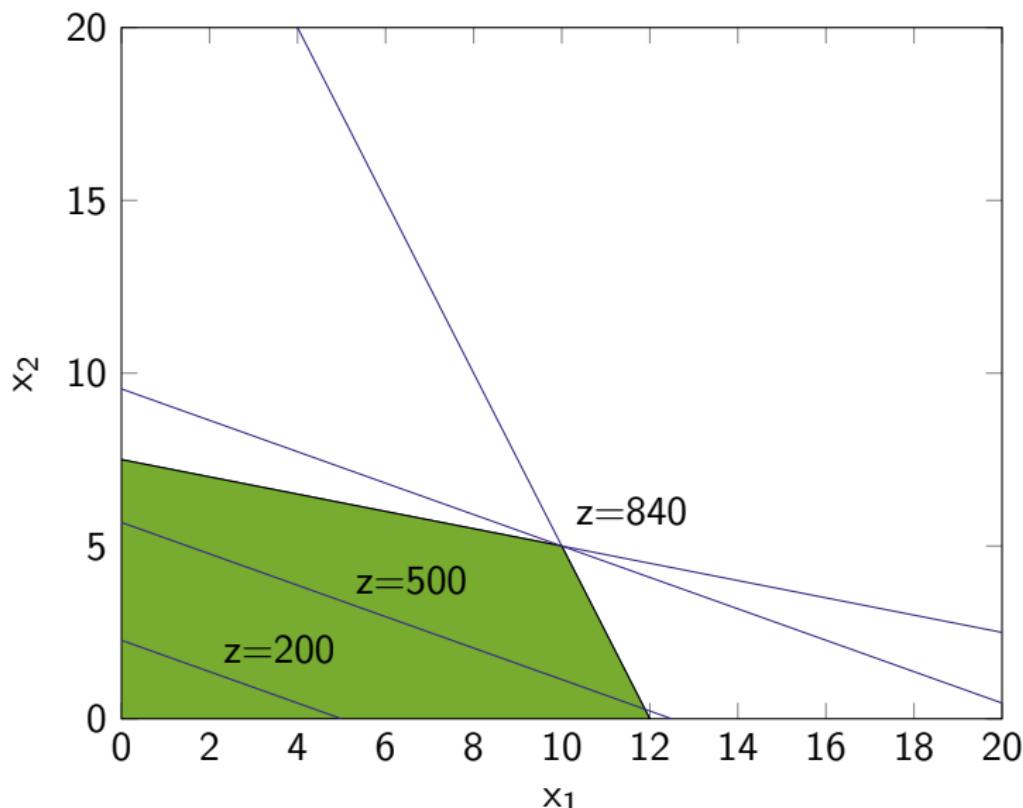
If the constraints are satisfied, but the objective function is not maximized/minimized we speak of a feasible solution.

If also the objective function is maximized/minimized, we speak of an optimal solution!

Plotting the constraints



Plotting the constraints



Normal form of an LP problem

$$\max z = f(x_1, x_2) = 40x_1 + 88x_2$$

s.t.

$$2x_1 + 8x_2 \leq 60$$

$$5x_1 + 2x_2 \leq 60$$

$$x_1 \geq 0$$

$$x_2 \geq 0$$

Normal form of an LP problem

$$\max z = f(x_1, x_2) = 40x_1 + 88x_2$$

s.t.

$$2x_1 + 8x_2 \leq 60$$

$$5x_1 + 2x_2 \leq 60$$

$$x_1 \geq 0$$

$$x_2 \geq 0$$

$$\max z = f(x) = 40x_1 + 88x_2$$

s.t.

$$2x_1 + 8x_2 + x_3 = 60$$

$$5x_1 + 2x_2 + x_4 = 60$$

$$x_i \geq 0 \quad i \in \{1, 2, 3, 4\}$$

x_3 and x_4 are called slack variables, they are non auxiliary variables introduced for the purpose of converting inequalities into equalities

The simplex method

We can formulate our earlier example to the normal form and consider it as the following augmented matrix with

$T_0 = [z \quad x_1 \quad x_2 \quad x_3 \quad x_4 \quad b]:$

$$T_0 = \begin{bmatrix} 1 & -40 & -88 & 0 & 0 & 0 \\ 0 & 2 & 8 & 1 & 0 & 60 \\ 0 & 5 & 2 & 0 & 1 & 60 \end{bmatrix}$$

This matrix is called the (initial) simplex table. Each simplex table has two kinds of variables, the basic variables (columns having only one nonzero entry) and the nonbasic variables

The simplex method

$$T_0 = \begin{bmatrix} 1 & -40 & -88 & 0 & 0 & 0 \\ 0 & 2 & 8 & 1 & 0 & 60 \\ 0 & 5 & 2 & 0 & 1 & 60 \end{bmatrix}$$

Every simplex table has a feasible solution. It can be obtained by setting the nonbasic variables to zero: $x_1 = 0$, $x_2 = 0$, $x_3 = 60/1$, $x_4 = 60/1$, $z = 0$.

The optimal solution?

- The optimal solution is now obtained stepwise by pivoting in such way that z reaches a maximum.
- The big question is, how to choose your pivot equation ...

Step 1: Selection of the pivot column

Select as the column of the pivot, the first column with a negative entry in Row 1. In our example, that's column 2 (-40)

$$T_0 = \begin{bmatrix} 1 & -40 & -88 & 0 & 0 & 0 \\ 0 & 2 & 8 & 1 & 0 & 60 \\ 0 & 5 & 2 & 0 & 1 & 60 \end{bmatrix}$$

Step 2: Selection of the pivot row

Divide the right sides by the corresponding column entries of the selected pivot column. In our example that is $60/2 = 30$ and $60/5 = 12$.

$$T_0 = \begin{bmatrix} 1 & -40 & -88 & 0 & 0 & 0 \\ 0 & 2 & 8 & 1 & 0 & 60 \\ 0 & 5 & 2 & 0 & 1 & 60 \end{bmatrix}$$

Take as the pivot equation the equation that gives the smallest quotient, so $60/5$.

Step 3: Elimination by row operations

- Row 1 = Row 1 + 8 * Row 3
- Row 2 = Row 2 - 0.4 * Row 3

$$T_1 = \begin{bmatrix} 1 & 0 & -72 & 0 & 8 & 480 \\ 0 & 0 & 7.2 & 1 & -0.4 & 36 \\ 0 & 5 & 2 & 0 & 1 & 60 \end{bmatrix}$$

The basic variables are now x_1 , x_3 and the nonbasic variables are x_2 , x_4 . Setting the nonbasic variables to zero will give a new feasible solution: $x_1 = 60/5$, $x_2 = 0$, $x_3 = 36/1$, $x_4 = 0$, $z = 480$.

The simplex method

- We moved from $z = 0$ to $z = 480$. The reason for the increase is because we eliminated a negative term from the equation, so: elimination should only be applied to negative entries in Row 1, but no others.
- Although we found a feasible solution, we did not find the optimal solution yet (the entry of -72 in our simplex table)
→ repeat step 1 to 3.

The simplex method

Another iteration is required:

- Step 1: Select column 3
- Step 2: $36/7.2 = 5$ and $60/2 = 30 \longrightarrow$ select 7.2 as the pivot
- Elimination by row operations:
 - Row 1 = Row 1 + 10*Row 2
 - Row 3 = Row 3 - (2/7.2)*Row 2

$$T_2 = \begin{bmatrix} 1 & 0 & 0 & 10 & 4 & 840 \\ 0 & 0 & 7.2 & 1 & -0.4 & 36 \\ 0 & 5 & 0 & -1/36 & 1/0.9 & 50 \end{bmatrix}$$

- The basic feasible solution: $x_1 = 50/5$, $x_2 = 36/7.2$, $x_3 = 0$, $x_4 = 0$, $z = 840$ (no more negative entries: so this solution is also the optimal solution)

Using Matlab for LP problems

We are going to solve the following LP problem:

$$\min f(x) = -5x_1 - 4x_2 - 6x_3$$

s.t.

$$x_1 - x_2 + x_3 \leq 20$$

$$3x_1 + 2x_2 + 4x_3 \leq 42$$

$$3x_1 + 2x_2 \leq 30$$

$$x_1 \geq 0$$

$$x_2 \geq 0$$

$$x_3 \geq 0$$

Using the function `linprog`:

```
f = [-5; -4; -6]
A = [1 -1 1; 3 2 4; 3 2 0];
b = [20; 42; 30];
lb = zeros(3,1);
[x,fval,exitflag,output,
lambda]
= linprog(f,A,b,[],[],lb);
```

Gives:

```
x = 0.00 15.00 3.00
lambda.ineqlin = 0 1.50
                0.50
lambda.lower = 1.00 0 0
```

Summary

- Curve fitting: Manual procedures for polynomial fitting in Matlab and Excel
- Curve fitting: Matlab's curve-fitting toolbox
- Curve fitting: Matlab's non-linear least-squares solver `lsqnonlin`
- Optimization: An introduction to the Simplex method
- Optimization: Use of the `linprog` solver