

Numerical errors in computer simulations

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

Chemical Process Intensification group
Eindhoven University of Technology

Numerical Methods (6BER03), 2024-2025

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?

Enter:

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

(repeat until A30)

Example 2

Start Python

Investigate the result of `np.sin(1e40 * np.pi)`

Create a vector `v` containing the powers of 10, e.g., from 10^0 up to 10^{40} , and solve `np.sin(v * np.pi)`:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 v = np.logspace(0, 40, 41)
5 y = np.sin(v * np.pi)
6 plt.loglog(v, np.abs(y)) # Double log plot, values on y-axis must be positive.
7 plt.show()
```

Errors in computer simulations

In this lecture I will outline different numerical errors that can 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

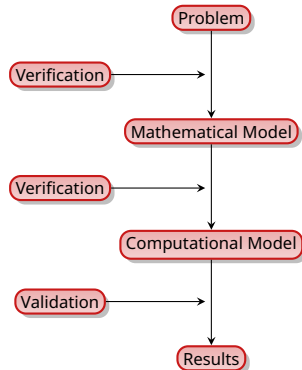
Verification and validation

Verification

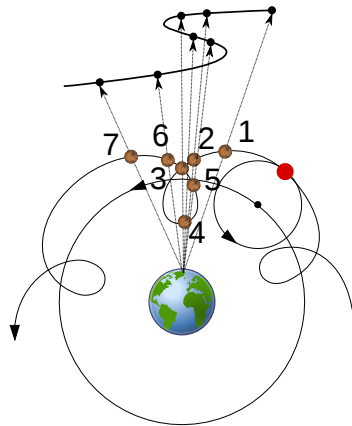
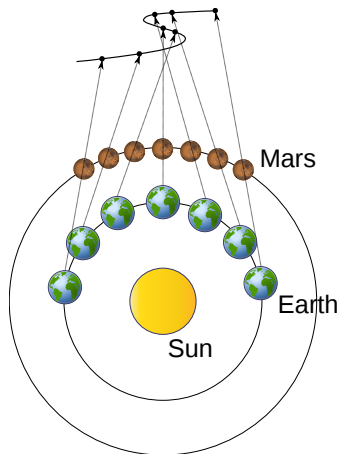
Verification is the process of mathematically and computationally assuring that the model computes the equations you intended to implement.

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 lecture I will outline different numerical errors that can 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

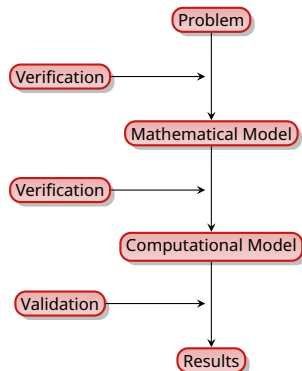
Verification and validation

Verification

Verification is the process of mathematically and computationally assuring that the model computes the equations you intended to implement.

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 lecture I will outline different numerical errors that can 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

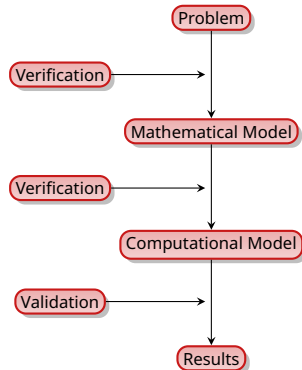
Verification and validation

Verification

Verification is the process of mathematically and computationally assuring that the model computes the equations you intended to implement.

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 lecture I will outline different numerical errors that can 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:

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

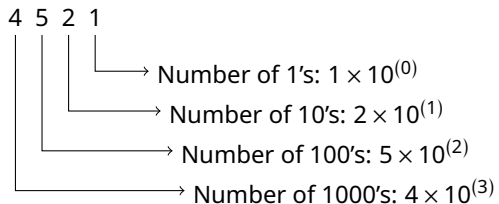
3 significant digits

Today's outline

- Introduction
- Roundoff and truncation errors
- Break errors
- Loss of digits
- (Un)stable methods
- Symbolic math
- Summary
- Introduction
- Matrix inversion
- Solving a linear system
- Towards larger systems
- 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

Exercise

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

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

- Excel:
 - Decimal: =DEC2BIN(214)
 - Octal: =DEC2OCT(214)
 - Hexadecimal: =DEC2HEX(214)
- Python:
 - Decimal: `bin(214)`
 - Octal: `oct(214)`
 - Hexadecimal: `hex(214)`

Arithmetic operations with binary numbers

Addition:

$0 + 0 = 0$
 $0 + 1 = 1$
 $1 + 0 = 1$
 $1 + 1 = 0$ (carry one)

$$\begin{array}{r}
 145 \\
 + \quad 23 \\
 \hline
 168
 \end{array}$$

$$\begin{array}{r}
 10010001 \\
 + 00010111 \\
 \hline
 10110000
 \end{array}$$

Subtraction:

$0 - 0 = 0$
 $1 - 0 = 1$
 $1 - 1 = 0$
 $0 - 1 = 1$ (borrow one)

$$\begin{array}{r}
 145 \\
 - \quad 23 \\
 \hline
 122
 \end{array}$$

$$\begin{array}{r}
 10010001 \\
 - 00010111 \\
 \hline
 01111000
 \end{array}$$

- Multiplication and division are more expensive, and more elaborate

Exercise

Try the following commands in Python:

| Command | Result |
|---|--|
| <code>np.iinfo(np.int32).min</code> | <code>-2147483648</code> |
| <code>np.iinfo(np.int32).max</code> | <code>2147483647</code> |
| <code>i = np.int16(np.iinfo(np.int32).max)</code> | <code>i = 32767</code> |
| <code>type(i)</code> | <code><class 'numpy.int16'></code> |
| <code>i = i + 100</code> | <code>i = 32767</code> |
| <code>np.finfo(np.float64).max</code> | <code>1.7976931348623157e+308</code> |
| <code>f = 0.1</code> | |
| <code>type(f)</code> | <code><class 'float'></code> |
| <code>print(":.16f".format(f))</code> | <code>0.100000000000000000</code> |
| <code>print(":.20f".format(f))</code> | <code>0.100000000000000000555</code> |

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

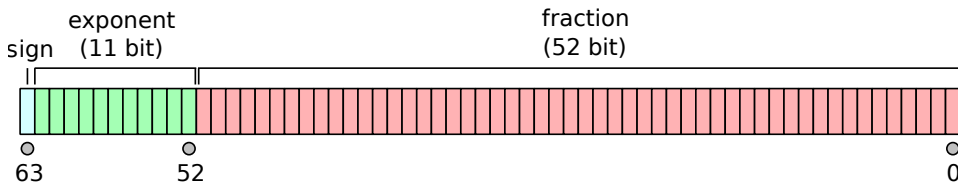


Image: [Wikimedia Commons CC by-SA](#)

Representation of real (floating point) numbers

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

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

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

Today's outline

- Introduction
- Roundoff and truncation errors
- **Break errors**
- Loss of digits
- (Un)stable methods
- Symbolic math
- Summary
- Introduction
- Matrix inversion
- Solving a linear system
- Towards larger systems
- 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

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

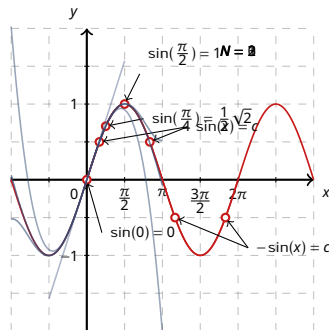
$$e^x = \sum_{n=0}^N \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots + \frac{x^N}{N!}$$

- 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

- Introduction
- Roundoff and truncation errors
- Break errors
- **Loss of digits**
- (Un)stable methods
- Symbolic math
- Summary
- **Introduction**
- Matrix inversion
- Solving a linear system
- Towards larger systems
- 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

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

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

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

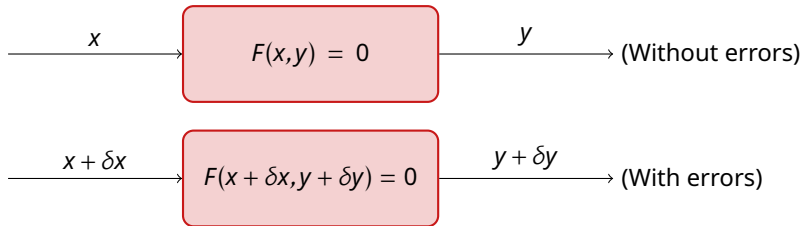
- The absolute error is small ($\approx 10^{-5}$), but the relative error is much bigger (0.028).
- Adding up the errors results in a loss of significant digits!

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: `round(term, 4)`
- 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:



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

Propagated error on the basis of Taylor expansion

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

Condition criterion, $C < 10$ error development small

Badly (ill) conditioned problems: Example

Solve the following linear system in Python 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

```
1 import numpy as np
2
3 # Double precision
4 A = np.array([[1.2969, 0.8648], [0.2161,
5     0.1441]], dtype=np.float64)
6 x = np.array([0.8642, 0.1440], dtype=np.float64)
7 y = np.linalg.solve(A, x)
8 print("y =")
9 print(y)
```

Single precision

```
1 import numpy as np
2
3 # Single precision
4 A = np.array([[1.2969, 0.8648], [0.2161,
5     0.1441]], dtype=np.float32)
6 x = np.array([0.8642, 0.1440], dtype=np.float32)
7 y = np.linalg.solve(A, x)
8 print("y =")
9 print(y)
```

Badly (ill) conditioned problems: Example

- Python does not automatically warn about bad condition number. You need to compute and check it manually using NumPy:

```
if np.linalg.cond(A) > threshold: raise("Ill conditioned error")
```
- The `cond` number is the condition number computed using NumPy.
- A small error in A (the matrix) results in a big error in the solution. This is called an ill-conditioned problem.

Today's outline

- Introduction
 - Roundoff and truncation errors
 - Break errors
 - Loss of digits
 - (Un)stable methods
 - Symbolic math
 - Summary
-
- Introduction
 - Matrix inversion
 - Solving a linear system
 - Towards larger systems
 - 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

Recurrent version

```

1 import numpy as np
2
3 def golden_mean_recurrent(Ntot):
4     # Initialize the series with the given
5     # initial conditions
6     y = np.zeros(Ntot)
7     y[0] = 1
8     y[1] = 2 / (1 + np.sqrt(5))
9
10    # Perform the recurrence to fill in the rest
11    # of the series
12    for n in range(2, Ntot):
13        y[n] = y[n-1] - y[n-2]
14    return y

```

Power law version

```
1 import numpy as np
2
3 def golden_mean_powerlaw(Ntot):
4     # Initialize the constant value
5     x = (1 + np.sqrt(5)) / 2
6
7     # Generate a range of values from 0 to Ntot
8     # and apply the power law
9     y = x ** -np.arange(0, Ntot + 1)
10
11     return y
```

- Compare the outcomes: `plot(goldenMeanPowerlaw(40)- goldenMeanRecurrent(40))`
- See what happens if you use single precision (uncomment the second line of both functions).

The Golden mean

| n | Recurrent | Power law |
|-----|------------------------|------------------------|
| 0 | 1.0000 | 1.0000 |
| 1 | 0.6180 | 0.6180 |
| 2 | 0.3820 | 0.3820 |
| 3 | 0.2361 | 0.2361 |
| ... | ... | ... |
| 37 | $1.714 \cdot 10^{-08}$ | $1.851 \cdot 10^{-08}$ |
| 38 | $1.366 \cdot 10^{-08}$ | $1.144 \cdot 10^{-08}$ |
| 39 | $3.485 \cdot 10^{-08}$ | $7.071 \cdot 10^{-09}$ |
| 40 | $1.017 \cdot 10^{-08}$ | $4.370 \cdot 10^{-09}$ |

- 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 (large sine series)

The `np.sin(1e40*np.pi)` result gives poor results because `1e40` has an error margin on the order of floating-point machine epsilon, which is roughly 1×10^{-16} in Python (double-precision floating-point format). In Python, as in many computing environments, the number of $2 \cdot \pi$ cycles is still much larger than $10^{40} \cdot 10^{-16}$. Also, π is not stored with an infinite number of digits, which further contributes to the imprecision.

Today's outline

- Introduction
- Roundoff and truncation errors
- Break errors
- Loss of digits
- (Un)stable methods
- **Symbolic math**
- Summary
- Introduction
- Matrix inversion
- **Solving a linear system**
- Towards larger systems
- Summary

Maple Well known, license available via TU/e

- h[Alpha] Web-based interface by Mathematica developer. Less powerful in

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

Matlab Symbolic math toolbox

Python SymPy

Symbolic math: simplify

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

```
1 from sympy import symbols, simplify
2
3 x = symbols('x') # Alt: from sympy.abc import x,y,z
4 f = (x - 1)*(x + 1)*(x**2 + 1) + 1
5 f_simplified = simplify(f)
6 print(f_simplified)
```

```
f_simplified = x**4
```


Symbolic math: integration and differentiation

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

```

1 from sympy import symbols, simplify, integrate, diff
2
3 x = symbols('x')
4 f = 1/(x**3+1)
5 my_f_int = integrate(f, x)
6 my_f_diff = diff(my_f_int, x)
7 my_f_diff_simplified = simplify(my_f_diff)
8 print(my_f_diff_simplified)

```

```
my_f_diff_simplified = 1/(x**3 + 1)
```

Symbolic math: exercises

Exercise 1

Simplify the following expression:

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

```
1 from sympy import symbols, trigsimp, tan
2
3 x = symbols('x')
4 expr = 2*tan(x)/(1 + tan(x)**2)
5 simplified_expr = trigsimp(expr)
```

Symbolic math: exercises

Exercise 2

Calculate the *value* of p :

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

```

1 from sympy import exp, sinh, integrate, symbols
2
3 x = symbols("x")
4 f = ((exp(x)- exp(-x))/sinh(x)).simplify()
5 p = integrate(f, (x, 0, 10))
6 print(p)

```

$p = 20$

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 Python.

Symbolic math function

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

```

1 from sympy import solve, symbols
2
3 x = symbols("x")
4 f = 3 / (x**2 + 3*x) - 2
5 solutions = solve(f, x)
6 print(solutions)
    
```

`solutions = [-3/2 + sqrt(15)/2, -sqrt(15)/2 - 3/2]`

Today's outline

- Introduction
- Roundoff and truncation errors
- Break errors
- Loss of digits
- (Un)stable methods
- Symbolic math
- **Summary**
- Introduction
- Matrix inversion
- Solving a linear system
- **Towards larger systems**

204/612

- 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.