# Numerical methods

Lectures:Doc. Mgr. Jozef Kristek, PhD.F1-207Exercises:Mgr. David GregorF1-204

### Contents

- 1. Introduction, syllabus, evaluation
- 2. Source and types of errors
- 3. Definition of errors
- 4. Rounding, propagation of errors
- 5. Representation of numbers
- 6. Conditionality of numerical problems and numerical stability of algorithms

#### real problem -> mathematical model

mathematical model - solved using computers

**Numerical mathematics** is a scientific discipline, which develop and analyze methods based on manipulations with numbers

#### **Numerical problem**

clear and unambiguous description of functional relation between the finite number of input and output data

#### Algorithm of the numerical problem

clear and unambiguous specification of finite sequence of operations. The operations uniquely assign *n*-tuple numbers of results to the *m*-tuple numbers from certain set of input data

Pre- a post-processing

# Contents

- 1. Introduction, syllabus, evaluation
- 2. Source and types of errors
- 3. Definition of errors
- 4. Rounding, propagation of errors
- 5. Representation of numbers
- 6. Conditionality of numerical problems and numerical stability of algorithms

#### Human errors

#### **Errors of mathematical model**

difference between solution of mathematical (idealized) problem and solution of real problem

Example: Estimation of the surface of the Earth using equation  $S = 4 \pi r^2$ 

**Errors of input data** due to inaccuracy of measurements

#### **Errors of numerical method**

comes from taking a numerical problem instead of mathematical problem

The estimation of this error is necessary part of solution of numerical problem

#### **Errors of numerical method**

comes from taking a numerical problem instead of mathematical problem

The estimation of this error is necessary part of solution of numerical problem

Example: Computation of the value of sin x for x=1 using the summation of the finite number of terms in Taylor expansion

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \frac{x^9}{9!} - \dots + (-1)^n \frac{x^{2n+1}}{(2n+1)!} + \dots$$

#### **Errors of numerical method**

comes from taking a numerical problem instead of mathematical problem

The estimation of this error is necessary part of solution of numerical problem

Example: Computation of the value of sin x for x=1 using the summation of the finite number of terms in Taylor expansion

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \frac{x^9}{9!} - \dots + (-1)^n \frac{x^{2n+1}}{(2n+1)!} + \dots$$

We know that summation of the first n terms of expansion gives the error in estimation at most 1/(2n+1)!

#### **Round-off errors**

Round-off errors can cumulate but also cancel

*Example: The number*  $\pi$ , result of 2/3

#### **Round-off errors**

Round-off errors can cumulate but also cancel

*Example: The number*  $\pi$ , result of 2/3

#### real problem – all kind of errors together



# Contents

- 1. Introduction, syllabus, evaluation
- 2. Source and types of errors

# 3. Definition of errors

- 4. Rounding, propagation of errors
- 5. Representation of numbers
- 6. Conditionality of numerical problems and numerical stability of algorithms

#### Definition of errors

Let x is the exact value of some number and  $\tilde{x}$  is its approximation

 $\Delta(x) = \tilde{x} - x$ 

we call absolute error of approximation

#### Definition of errors

Let x is the exact value of some number and  $\tilde{x}$  is its approximation

 $\Delta(x) = \tilde{x} - x$ 

#### we call absolute error of approximation

Relative error  $\frac{\Delta(x)}{x} = \frac{\tilde{x} - x}{x}$ 

#### **Estimation of errors**

Each non-negative number  $\mathcal{E}$  , for which  $\left|\Delta(x)\right| \leq \mathcal{E}$  i.e.  $\tilde{x} - \mathcal{E} \leq x \leq \tilde{x} + \mathcal{E}$ 

we call estimation of absolute error

#### **Estimation of errors**

Each non-negative number  $\mathcal{E}$  , for which  $\left|\Delta(x)\right| \leq \mathcal{E}$  i.e.  $\tilde{x} - \mathcal{E} \leq x \leq \tilde{x} + \mathcal{E}$ 

we call estimation of absolute error

Each non-negative number  $\delta$ , for which  $\left|\frac{\Delta(x)}{x}\right| \leq \delta$ 

we call estimation of relative error

#### **Estimation of errors**

Each non-negative number  $\mathcal{E}$  , for which  $\left|\Delta(x)\right| \leq \mathcal{E}$ i.e.  $\tilde{x} - \mathcal{E} \leq x \leq \tilde{x} + \mathcal{E}$ 

we call estimation of absolute error

Each non-negative number  $\delta$ , for which  $\left|\frac{\Delta(x)}{x}\right| \leq \delta$ 

we call estimation of relative error

Usually we write  

$$x = \tilde{x} \pm \varepsilon$$
  $\tilde{x} = x(1 \pm \delta)$ 

Now evaluate an error of value  $f(x_1, x_2, ..., x_n)$  of function f, if exact values  $x_i$ will be changed by approximate values  $\tilde{x}_i = x_i + \Delta x_i$ . Now evaluate an error of value  $f(x_1, x_2, ..., x_n)$  of function f, if exact values  $x_i$ will be changed by approximate values  $\tilde{x}_i = x_i + \Delta x_i$ .  $f(\tilde{x}) = f(\mathbf{x}) + \sum_{i=1}^n \Delta x_i \frac{\partial f(\mathbf{x})}{\partial x_i} + \frac{1}{2} \sum_{i, i=1}^n \Delta x_i \Delta x_j \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} + \cdots$ . Now evaluate an error of value  $f(x_1, x_2, ..., x_n)$  of function f, if exact values  $x_i$ will be changed by approximate values  $\tilde{x}_i = x_i + \Delta x_i$ .  $f(\tilde{x}) = f(\mathbf{x}) + \sum_{i=1}^n \Delta x_i \frac{\partial f(\mathbf{x})}{\partial x_i} + \frac{1}{2} \sum_{i,j=1}^n \Delta x_i \Delta x_j \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} + \cdots$ .

Assuming that  $\Delta x_i \Delta x_j$  are small

Now evaluate an error of value  $f(x_1, x_2, ..., x_n)$  of function f, if exact values  $x_i$ will be changed by approximate values  $\tilde{x}_i = x_i + \Delta x_i$ .

$$f(\tilde{x}) \approx f(\mathbf{x}) + \sum_{i=1}^{n} \Delta x_i \frac{\partial f(\mathbf{x})}{\partial x_i}$$

Assuming that  $\Delta x_i \Delta x_j$  are small

Now evaluate an error of value  $f(x_1, x_2, ..., x_n)$  of function f, if exact values  $x_i$ 

will be changed by approximate values  $\tilde{x}_i = x_i + \Delta x_i$ .

$$f(\tilde{x}) \approx f(\mathbf{x}) + \sum_{i=1}^{n} \Delta x_i \frac{\partial f(\mathbf{x})}{\partial x_i}$$

Assuming that  $\Delta x_i \Delta x_j$  are small, we can define absolute error as

$$\left|\Delta f(\mathbf{x})\right| \coloneqq \left|f(\tilde{\mathbf{x}}) - f(\mathbf{x})\right| \doteq \left|\sum_{i=1}^{n} \frac{\partial f(\mathbf{x})}{\partial x_{i}} \Delta x_{i}\right|$$

Now evaluate an error of value  $f(x_1, x_2, ..., x_n)$  of function f, if exact values  $x_i$ 

will be changed by approximate values  $\ \widetilde{x}_i = x_i + \Delta x_i$  .

$$f(\tilde{x}) \approx f(\mathbf{x}) + \sum_{i=1}^{n} \Delta x_i \frac{\partial f(\mathbf{x})}{\partial x_i}$$

Assuming that  $\Delta x_i \Delta x_j$  are small, we can define absolute error as

$$\left|\Delta f(\mathbf{x})\right| \coloneqq \left|f(\tilde{\mathbf{x}}) - f(\mathbf{x})\right| \doteq \left|\sum_{i=1}^{n} \frac{\partial f(\mathbf{x})}{\partial x_{i}} \Delta x_{i}\right| \le \sum_{i=1}^{n} \left|\frac{\partial f(\mathbf{x})}{\partial x_{i}}\right| \cdot \left|\Delta x_{i}\right| \quad (1)$$

#### Definition of errors

And for relative error

 $\left|\frac{\Delta f(\mathbf{x})}{f(\mathbf{x})}\right| \doteq \left|\sum_{i=1}^{n} \frac{x_{i}}{f(\mathbf{x})} \frac{\partial f(\mathbf{x})}{\partial x_{i}} \frac{\Delta x_{i}}{x_{i}}\right|$ 

#### Definition of errors

And for relative error

$$\frac{\left|\Delta f\left(\mathbf{x}\right)\right|}{f\left(\mathbf{x}\right)} \doteq \left|\sum_{i=1}^{n} \frac{x_{i}}{f\left(\mathbf{x}\right)} \frac{\partial f\left(\mathbf{x}\right)}{\partial x_{i}} \frac{\Delta x_{i}}{x_{i}}\right| \leq \sum_{i=1}^{n} \left|\frac{x_{i}}{f\left(\mathbf{x}\right)} \frac{\partial f\left(\mathbf{x}\right)}{\partial x_{i}}\right| \cdot \left|\frac{\Delta x_{i}}{x_{i}}\right|. \quad (2)$$

In practice the function values and values of derivatives are evaluated at  $\widetilde{x}.$ 

#### Error of basic arithmetic operations

Let  $f(x, y) = x \pm y$ .

Let 
$$f(x, y) = x \pm y$$
.

# Using eqs. (1) and (2) we obtain absolute and relative error of **addition** and **subtraction**

$$\Delta(x \pm y) \le \Delta x + \Delta y \qquad \qquad \left| \frac{\Delta(x \pm y)}{x \pm y} \right| \le \left| \frac{x}{x \pm y} \right| \left| \frac{\Delta x}{x} \right| + \left| \frac{y}{x \pm y} \right| \left| \frac{\Delta y}{y} \right|$$

Let 
$$f(x, y) = x \pm y$$
.

# Using eqs. (1) and (2) we obtain absolute and relative error of **addition** and **subtraction**

$$\Delta(x \pm y) \le \Delta x + \Delta y \qquad \qquad \left| \frac{\Delta(x \pm y)}{x \pm y} \right| \le \left| \frac{x}{x \pm y} \right| \left| \frac{\Delta x}{x} \right| + \left| \frac{y}{x \pm y} \right| \left| \frac{\Delta y}{y} \right|$$

Relative error of addition or subtraction could be significantly larger then relative errors of each operand in case when  $|x \pm y|$  is significantly smaller than |x| or |y|.

Let 
$$f(x, y) = xy$$
.

Then absolute and relative error of **multiplication** 

$$\Delta(xy) \le |y| \Delta x + |x| \Delta y \qquad \qquad \left| \frac{\Delta(xy)}{xy} \right| \le \frac{\Delta x}{|x|} + \frac{\Delta y}{|y|}$$

Let 
$$f(x, y) = x / y$$
.

#### Then absolute and relative error of **division**

$$\Delta\left(\frac{x}{y}\right) \le \left|\frac{1}{y}\right| \Delta x + \left|\frac{x}{y^2}\right| \Delta y$$

$$\left|\frac{\Delta(x / y)}{x / y}\right| \leq \frac{\Delta x}{|x|} + \frac{\Delta y}{|y|}$$

# Contents

- 1. Introduction, syllabus, evaluation
- 2. Source and types of errors
- 3. Definition of errors
- 4. Rounding, propagation of errors
- 5. Representation of numbers
- 6. Conditionality of numerical problems and numerical stability of algorithms

Let  $\tilde{x}$  is approximation of x written in decimal representation

$$\tilde{x} = \pm \left[ d_1 \cdot 10^e + d_2 \cdot 10^{e-1} + \dots + d_k \cdot 10^{e+1-k} + \dots \right], \quad d_1 \neq 0.$$

Let  $\tilde{x}$  is approximation of x written in decimal representation  $\tilde{x} = \pm \left[ d_1 \cdot 10^e + d_2 \cdot 10^{e-1} + \dots + d_k \cdot 10^{e+1-k} + \dots \right], \quad d_1 \neq 0.$ 

We say that *k-th* decimal digit  $d_k$  is significant if

$$\left|x - \tilde{x}\right| \le 0, 5 \cdot 10^{e+1-k} \tag{3}$$

i.e. if  $\tilde{x}$  differs from x

at most of 5 units of order of subsequent digit.

Let  $\tilde{x}$  is approximation of x written in decimal representation  $\tilde{x} = \pm \left[ d_1 \cdot 10^e + d_2 \cdot 10^{e-1} + \dots + d_k \cdot 10^{e+1-k} + \dots \right], \quad d_1 \neq 0.$ 

We say that *k***-th** decimal digit  $d_k$  is significant if

$$\left|x - \tilde{x}\right| \le 0, 5 \cdot 10^{e+1-k} \tag{3}$$

i.e. if  $\tilde{x}$  differs from x

at most of 5 units of order of subsequent digit. If inequality (3) holds for  $k \le p$ , but not for k = p + 1, we say, that  $\tilde{x}$  has **p** significant digits and is correctly rounded value of the number xto the p significant digits.

#### We say that *k***-th decimal place is significant** if

$$\left|x - \tilde{x}\right| \le 0, 5 \cdot 10^{-k} \tag{4}$$

i.e. if  $\tilde{x}$  differs from x

at most of 5 units of order of subsequent decimal place.

#### We say that *k*-th decimal place is significant if

$$\left|x - \tilde{x}\right| \le 0, 5 \cdot 10^{-k} \tag{4}$$

i.e. if  $\tilde{x}$  differs from x

at most of 5 units of order of subsequent decimal place. If inequality (4) hold for  $k \le p$  but not for k = p+1, we say that  $\tilde{x}$  has **p** significant decimal places.

#### Several examples

x	$\widetilde{x}$	<pre># of significant</pre>	<pre># of significant decimal places</pre>
374	380		
-27.6473	-27.598		
100.002	99.9973		
99.9973	100.002		
-0.003728	-0.0041		
1.841*10 <sup>-6</sup>	2.5*10 <sup>-6</sup>		
X	$\widetilde{x}$	<pre># of significant</pre>	<pre># of significant decimal places</pre>
------------	----------------------	-----------------------------	--
374	380	1	-
-27.6473	-27.598		
100.002	99.9973		
99.9973	100.002		
-0.003728	-0.0041		
1.841*10-6	2.5*10 <sup>-6</sup>		

x	$\widetilde{x}$	# of significant digits	<pre># of significant decimal places</pre>
374	380	1	-
-27.6473	-27.598	3	1
100.002	99.9973		
99.9973	100.002		
-0.003728	-0.0041		
1.841*10 <sup>-6</sup>	2.5*10 <sup>-6</sup>		

x	$\widetilde{x}$	# of significant digits	<pre># of significant decimal places</pre>
374	380	1	-
-27.6473	-27.598	3	1
100.002	99.9973	4	2
99.9973	100.002		
-0.003728	-0.0041		
1.841*10 <sup>-6</sup>	2.5*10 <sup>-6</sup>		

x	$\widetilde{x}$	# of significant digits	<pre># of significant decimal places</pre>
374	380	1	-
-27.6473	-27.598	3	1
100.002	99.9973	4	2
99.9973	100.002	5	2
-0.003728	-0.0041		
1.841*10-6	2.5*10 <sup>-6</sup>		

x	$\widetilde{x}$	<pre># of significant</pre>	<pre># of significant decimal places</pre>
374	380	1	-
-27.6473	-27.598	3	1
100.002	99.9973	4	2
99.9973	100.002	5	2
-0.003728	-0.0041	1	3
1.841*10-6	2.5*10 <sup>-6</sup>		

X	$\widetilde{x}$	<pre># of significant</pre>	<pre># of significant decimal places</pre>
374	380	1	-
-27.6473	-27.598	3	1
100.002	99.9973	4	2
99.9973	100.002	5	2
-0.003728	-0.0041	1	3
1.841*10-6	2.5*10-6	0	5

Example:

 $x = 4.998949 \cdot 10^{1}, \quad \tilde{x} = 4.999 \cdot 10^{1}, \quad |\Delta x| = 5.10 \cdot 10^{-4}, \quad \left|\frac{\Delta x}{x}\right| \doteq 1.020 \cdot 10^{-5},$  $y = 5.001848 \cdot 10^{1}, \quad \tilde{y} = 5.002 \cdot 10^{1}, \quad |\Delta y| = 1.52 \cdot 10^{-4}, \quad \left|\frac{\Delta y}{y}\right| \doteq 3.039 \cdot 10^{-5}$ 

Example:

 $x = 4.998949 \cdot 10^{1}, \quad \tilde{x} = 4.999 \cdot 10^{1}, \quad |\Delta x| = 5.10 \cdot 10^{-4}, \quad \left|\frac{\Delta x}{x}\right| \doteq 1.020 \cdot 10^{-5},$  $y = 5.001848 \cdot 10^{1}, \quad \tilde{y} = 5.002 \cdot 10^{1}, \quad |\Delta y| = 1.52 \cdot 10^{-4}, \quad \left|\frac{\Delta y}{y}\right| \doteq 3.039 \cdot 10^{-5}$ 

then for subtractions z = y - x,  $\tilde{z} = \tilde{y} - \tilde{x}$  we get

$$z = 2.899 \cdot 10^{-2}, \quad \tilde{z} = 3 \cdot 10^{-2}, \quad |\Delta z| = 1.01 \cdot 10^{-3}, \quad \left|\frac{\Delta z}{z}\right| \doteq 3.484 \cdot 10^{-2}$$

Example:

 $x = 4.998949 \cdot 10^{1}, \quad \tilde{x} = 4.999 \cdot 10^{1}, \quad |\Delta x| = 5.10 \cdot 10^{-4}, \quad \left|\frac{\Delta x}{x}\right| \doteq 1.020 \cdot 10^{-5},$  $y = 5.001848 \cdot 10^{1}, \quad \tilde{y} = 5.002 \cdot 10^{1}, \quad |\Delta y| = 1.52 \cdot 10^{-4}, \quad \left|\frac{\Delta y}{y}\right| \doteq 3.039 \cdot 10^{-5}$ 

then for subtractions z = y - x,  $\tilde{z} = \tilde{y} - \tilde{x}$  we get

$$z = 2.899 \cdot 10^{-2}, \quad \tilde{z} = 3 \cdot 10^{-2}, \quad |\Delta z| = 1.01 \cdot 10^{-3}, \quad \left|\frac{\Delta z}{z}\right| \doteq 3.484 \cdot 10^{-2}$$

therefore  $\tilde{z}$  has one significant digit while  $\tilde{x}$  and  $\tilde{y}$  have four significant digits.

#### Example:

 $x = 1.3262 \pm 5 \cdot 10^{-5}$ ,  $y = -6.5347 \pm 5 \cdot 10^{-5}$ ,  $z = 13.235 \pm 5 \cdot 10^{-4}$ 

Find the approximation of

f = xy / z,

absolute and relative errors and the number of significant digits of result.

# Contents

- 1. Introduction, syllabus, evaluation
- 2. Source and types of errors
- 3. Definition of errors
- 4. Rounding, propagation of errors
- 5. Representation of numbers
- 6. Conditionality of numerical problems and numerical stability of algorithms

Real numbers in computers are represented in the **floating point format.** 

Basic idea is similar to the **semilogarithmic notation** (i.e. 2.457\*10<sup>5</sup>)

System of normalized floating point numbers  ${\mathcal F}$ 

is characterized by 4 integer numbers:

$$\begin{array}{ccc} \beta & \text{base} & \left(\beta \geq 2\right) \\ p & \text{precision} \left(p \geq 1\right) \\ \left[e_{\min}, e_{\max}\right] \text{ exponent range} \left(e_{\min} < 0 < e_{\max}\right) \end{array}$$

$$\begin{array}{ccc} \beta & \text{base} & \left(\beta \geq 2\right) \\ p & \text{precision} \left(p \geq 1\right) \\ \left[e_{\min}, e_{\max}\right] \text{ exponent range} \left(e_{\min} < 0 < e_{\max}\right) \end{array}$$

Each number  $x \in \mathcal{F}$  has form of  $x = \pm m \cdot \beta^e$ , where  $m = d_1 + \frac{d_2}{\beta} + \frac{d_3}{\beta^2} + \dots + \frac{d_p}{\beta^{p-1}}$ 

$$\begin{array}{ccc} \beta & \text{base} & \left(\beta \geq 2\right) \\ p & \text{precision} \left(p \geq 1\right) \\ \left[e_{\min}, e_{\max}\right] \text{ exponent range} \left(e_{\min} < 0 < e_{\max}\right) \end{array}$$

Each number  $x \in \mathcal{F}$  has form of  $x = \pm m \cdot \beta^e$ , where  $m = d_1 + \frac{d_2}{\beta} + \frac{d_3}{\beta^2} + \dots + \frac{d_p}{\beta^{p-1}}$ 

*m* is **normalized mantissa** (or **significand**),  $d_i \in \{0, 1, ..., \beta - 1\}, i = 1, 2, ..., p$  are digits of mantissa, *p* is the number of digits of mantissa and  $e \in \langle e_{\min}, e_{\max} \rangle$  is integer **exponent.** 

$$\begin{array}{ccc} \beta & \text{base} & \left(\beta \geq 2\right) \\ p & \text{precision} \left(p \geq 1\right) \\ \left[e_{\min}, e_{\max}\right] \text{ exponent range} \left(e_{\min} < 0 < e_{\max}\right) \end{array}$$

Each number  $x \in \mathcal{F}$  has form of  $x = \pm m \cdot \beta^e$ , where  $m = d_1 + \frac{d_2}{\beta} + \frac{d_3}{\beta^2} + \dots + \frac{d_p}{\beta^{p-1}}$ 

*m* is **normalized mantissa** (or **significand**),  $d_i \in \{0, 1, ..., \beta - 1\}, i = 1, 2, ..., p$  are digits of mantissa, *p* is the number of digits of mantissa and  $e \in \langle e_{\min}, e_{\max} \rangle$  is integer **exponent.** 

Normalized mantissa means that for  $x \neq 0$  is  $d_1 \geq 1$ .

 $\begin{array}{ll} \beta = 2 & \mbox{binary numeral system} \\ \beta = 16 & \mbox{hexadecimal system} \\ \beta = 8 & \mbox{octal system} \\ \beta = 10 & \mbox{decimal system} \end{array}$ 

 $\begin{array}{ll} \beta = 2 & \mbox{binary numeral system} \\ \beta = 16 & \mbox{hexadecimal system} \\ \beta = 8 & \mbox{octal system} \\ \beta = 10 & \mbox{decimal system} \end{array}$ 

Set  ${\mathcal F}$  of floating point numbers is a finite set, count of number is

 $\begin{array}{ll} \beta = 2 & \mbox{binary numeral system} \\ \beta = 16 & \mbox{hexadecimal system} \\ \beta = 8 & \mbox{octal system} \\ \beta = 10 & \mbox{decimal system} \end{array}$ 

Set  ${\mathcal F}$  of floating point numbers is a finite set, count of number is

2

two signs,

$$eta=2$$
 binary numeral system  
 $eta=16$  hexadecimal system  
 $eta=8$  octal system  
 $eta=10$  decimal system

Set  ${\mathcal F}$  of floating point numbers is a finite set, count of number is

$$2(\beta-1)$$

two signs,  $\beta\!-\!1$  options for the first digit of mantissa,

$$eta=2$$
 binary numeral system  
 $eta=16$  hexadecimal system  
 $eta=8$  octal system  
 $eta=10$  decimal system

Set  ${\mathcal F}$  of floating point numbers is a finite set, count of number is

$$2(\beta-1)\beta^{p-1}$$

two signs,  $\beta - 1$  options for the first digit of mantissa,  $\beta$  options for other p - 1 digits of mantissa,

$$eta = 2$$
 binary numeral system  
 $eta = 16$  hexadecimal system  
 $eta = 8$  octal system  
 $eta = 10$  decimal system

Set  ${\mathcal F}$  of floating point numbers is a finite set, count of number is

$$2(\beta-1)\beta^{p-1}(e_{\max}-e_{\min}+1)$$

two signs,

 $\beta-1$  options for the first digit of mantissa,  $\beta$  options for other p-1 digits of mantissa,  $e_{\rm max}-e_{\rm min}+1$  possible values of exponent

$$eta=2$$
 binary numeral system  
 $eta=16$  hexadecimal system  
 $eta=8$  octal system  
 $eta=10$  decimal system

Set  ${\mathcal F}$  of floating point numbers is a finite set, count of number is

$$2(\beta - 1)\beta^{p-1}(e_{\max} - e_{\min} + 1) + 1$$

two signs,

 $\beta-1$  options for the first digit of mantissa,  $\beta\,$  options for other p-1 digits of mantissa,  $e_{\rm max}-e_{\rm min}+1$  possible values of exponent and one zero

The smallest positive number in  ${\mathscr F}$  is *UFL* (UnderFlow Level)

$$UFL = \beta^{e_{\min}},$$

which has the first digit of mantissa equal to one, others zero and exponent the smallest one.

The largest positive number in  $\mathcal{F}$  is *OFL* (OverFlow Level)  $OFL = \left(\beta - \beta^{1-p}\right)\beta^U$ ,

which has all digits of mantissa equal to  $~\beta\!-\!1~$  and exponent is the largest one.

The gaps between adjacent numbers **scale** with the size of the numbers

The gaps between adjacent numbers scale with the size of the numbers

Relative resolution is given by machine epsilon  $\varepsilon_{\rm machine} = 0.5 \, \beta^{1-p}$ 

For all *x*, there exists a floating point x'such that  $|x'-x| \le \varepsilon_{\text{machine}} |x|$  The gaps between adjacent numbers scale with the size of the numbers

Relative resolution is given by machine epsilon  $\varepsilon_{\rm machine} = 0.5 \, \beta^{1-p}$ 

For all x, there exists a floating point x'such that  $|x'-x| \le \varepsilon_{\text{machine}} |x|$ 

Example:  $\beta = 2$ , p = 3,  $e_{\min} = -1$ ,  $e_{\max} = 2$ 

The gaps between adjacent numbers scale with the size of the numbers

Relative resolution is given by machine epsilon  $\varepsilon_{\rm machine} = 0.5 \ \beta^{1\text{-}p}$ 

For all x, there exists a floating point x'such that  $|x'-x| \le \varepsilon_{\text{machine}} |x|$ 

- $\pm \infty$  is returned when an operation overflows
- $x/\pm \infty = 0$  for any number x,
- $x/0 = \pm \infty$  for any nonzero number x
- Operations with infinity are defined as limits, e.g.

$$4 - \infty = \lim_{x \to \infty} 4 - x = -\infty$$

 NaN (Not a Number) is returned when the an operation has no welldefined finite or infinite result

Examples:  $\infty - \infty$ ,  $\infty / \infty$ , 0 / 0, NaN  $\odot x$ 

#### **Denormalized Numbers**

• With normalized significand there is a "gap" between 0 and  $\beta^{emin}$ 



#### **Denormalized Numbers**

- With normalized significand there is a "gap" between 0 and  $\beta^{emin}$
- Solution: Allow non-normalized significand when the exponent is  $e_{\min}$
- This gradual underflow garantees that

$$x = y \Leftrightarrow x - y = \mathbf{0}$$



# **IEEE Single Precision**

• 1 sign bit, 8 exponent bits, 23 significand bits:

0	00000000	000000000000000000000000000000000000000
S	Е	М

• Represented number:

$$(-1)^S \times 1.M \times 2^{E-127}$$

• Special cases:

	E = 0	$0 \ < \ E \ < \ 255$	E = 255
M = 0	$\pm 0$	Powers of 2	$\pm\infty$
$M \neq 0$	Denormalized	Ordinary numbers	NaN

# **IEEE Single Precision, Examples**

S	E	Μ	Quantity
0	11111111	00000100000000000000000	NaN
1	11111111	00100010001001010101010	NaN
0	11111111	000000000000000000000000000000000000000	$\infty$
0	10000001	101000000000000000000000000000000000000	$+1 \cdot 2^{129 - 127} \cdot 1.101 = 6.5$
0	10000000	000000000000000000000000000000000000000	$+1 \cdot 2^{128 - 127} \cdot 1.0 = 2$
0	00000001	000000000000000000000000000000000000000	$+1 \cdot 2^{1-127} \cdot 1.0 = 2^{-126}$
0	00000000	100000000000000000000000000000000000000	$+1 \cdot 2^{-126} \cdot 0.1 = 2^{-127}$
0	00000000	000000000000000000000000000000000000000	$+1 \cdot 2^{-126} \cdot 2^{-23} = 2^{-149}$
0	00000000	000000000000000000000000000000000000000	0
1	00000000	000000000000000000000000000000000000000	-0
1	10000001	101000000000000000000000000000000000000	$-1 \cdot 2^{129 - 127} \cdot 1.101 = -6.5$
1	11111111	000000000000000000000000000000000000000	$-\infty$

# **IEEE Floating Point Data Types**

	Single precision	Double precision
Significand size ( $p$ )	24 bits	53 bits
Exponent size	8 bits	11
Total size	32 bits	64 bits
$e_{\max}$	+127	+1023
$e_{\min}$	-126	-1022
Smallest normalized	$2^{-126} \approx 10^{-38}$	$2^{-1022} \approx 10^{-308}$
Largest normalized	$2^{127} \approx 10^{38}$	$2^{1023} \approx 10^{308}$
$\epsilon_{\mathrm{machine}}$	$2^{-24} \approx 6 \cdots 10^{-8}$	$2^{-53} \approx 10^{-16}$

# **Floating Point Arithmetic**

 ${\buildrel \ }$  Define  ${\rm fl}(x)$  as the closest floating point approximation to x

 $_{\bullet}\,$  By the definition of  $\epsilon_{machine},$  we have for the relative error:

For all  $x \in \mathbb{R}$ , there exists  $\epsilon$  with  $|\epsilon| \le \epsilon_{\text{machine}}$ such that  $fl(x) = x(1 + \epsilon)$ 

- The result of an operation  $\circledast$  using floating point numbers is  $\mathrm{fl}(a \circledast b)$
- If  $fl(a \circledast b)$  is the nearest floating point number to  $a \circledast b$ , the arithmetic *rounds correctly* (IEEE does), which leads to the following property:

For all floating point x, y, there exists  $\epsilon$  with  $|\epsilon| \le \epsilon_{\text{machine}}$  such that x $\circledast y = (x * y)(1 + \epsilon)$
## Contents

- 1. Introduction, syllabus, evaluation
- 2. Source and types of errors
- 3. Definition of errors
- 4. Rounding, propagation of errors
- 5. Representation of numbers
- 6. Conditionality of numerical problems and numerical stability of algorithms

We have to investigate the effects of small changes in input data and rounding on the result.

We have to investigate the effects of small changes in input data and rounding on the result.

Mathematical problem can be treated as mapping y = f(x), which each data from x the set D of input data assigns the result y from the set R of output data.

We have to investigate the effects of small changes in input data and rounding on the result.

Mathematical problem can be treated as mapping y = f(x), which each data from x the set D of input data assigns the result y from the set R of output data.

We say that the mathematical problem  $y = f(x), x \in D, y \in R,$  is **correct**, when

1. For each input  $x \in D$  exists only one result  $y \in R$ , 2. this result continuously depends on input data, i.e. if  $x \to a$ , then  $f(x) \to f(a)$ .

We say that the correct problem is **well-conditioned**, if small change in input data will cause small change of result.

We say that the correct problem is **well-conditioned**, if small change in input data will cause small change of result.

Condition number is defined as

 $C_p = \frac{|\text{relative error of input}|}{|\text{relative error of output}|}$ 

We say that the correct problem is **well-conditioned**, if small change in input data will cause small change of result.

Condition number is defined as

 $C_p = \frac{|\text{relative error of input}|}{|\text{relative error of output}|}$ 

If  $C_p \approx 1$ , the problem is well-conditioned.

For large  $C_p$  (>100) the problem is ill-conditioned.

## We say that the **algorithm is well-conditioned**,

if it is less sensitive to errors in input data.

If the effect of round-off errors on result is small then we say about **numerically stable algorithm.** 

Well-conditioned and numerically stable algorithm is called **stable**.

## Example:

## Estimate the condition number of the problem:

evaluate the functional value of (differentiable) function

y = f(x)

show on example function  $f(x) = \tan x$ 

Examples:

1. roots of quadratic equation  $x^2 - 2bx + c = 0$ 

2. Evaluation of integral 
$$E_n = \int_0^1 x^n e^{x-1} dx$$
  $n = 1, 2, ...$ 

$$E_n = \int_0^1 x^n e^{x-1} \, \mathrm{d}x \qquad n = 1, 2, \dots$$

$$\int_0^1 x^n e^{x-1} \, \mathrm{d}x = [x^n e^{x-1}]_0^1 - \int_0^1 n x^{n-1} e^{x-1} \, \mathrm{d}x \,,$$

$$E_n=1-nE_{n-1}.$$

$$E_1 = 1/e$$
,  $E_n = 1 - nE_{n-1}$ ,  $n = 2, 3, ...$ 

$$E_{n-1}=\frac{1-E_n}{n}\,,$$

$$E_{n-1}=\frac{1-E_n}{n}\,,$$

$$E_n = \int_0^1 x^n e^{x-1} \, \mathrm{d}x \le \int_0^1 x^n \, \mathrm{d}x = \frac{1}{n+1}$$

$$E_n 
ightarrow 0$$
 for  $n 
ightarrow \infty$ 

$$E_{n-1}=\frac{1-E_n}{n}\,,$$

$$E_n = \int_0^1 x^n e^{x-1} \, \mathrm{d}x \le \int_0^1 x^n \, \mathrm{d}x = \frac{1}{n+1}$$

$$E_n 
ightarrow 0$$
 for  $n 
ightarrow \infty$ 

$$E_N = 0$$
,  $E_{n-1} = \frac{1-E_n}{n}$ ,  $n = N, N - 1, ...,$