

NUMERICAL PROGRAMMING 1 CSE

Prof. Oliver Juge

Dr. Dawid Karrasch

Center for Mathematics – M3

TUM

„numerical programming“ :

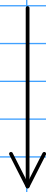
real world problem



mathematical problem



numerical algorithm



computer program

weather prediction, crash test, ...

equations, e.g. $f(x) = 0$

Newton's method

newton.m

1. Basics

1.1 Floating point arithmetic

$$\pi = 3.14159|2... \quad (\text{base } 10)$$

For each real number, we are only storing finitely many digits:

$$\frac{1}{10} = 0.1_{[10]} = 0.000|100|100..._{[2]}$$

$\uparrow \uparrow$ $\uparrow \uparrow \uparrow \uparrow$ \uparrow
 $1 \quad \frac{1}{10}$ $1 \quad \frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{8} \dots$ 2^{-n}

$$= 1.0 \cdot 10^{-1} = 1.100|100... \cdot 2^{-4}$$

floating point numbers

$$x = \underbrace{\pm}_{\text{sign}} \underbrace{d_1 d_2 d_3 \dots d_p}_{\text{mantissa}} \cdot \underbrace{b^e}_{\text{base}} \quad \text{exponent}$$

with digits $d_k \in \{0, 1, \dots, b-1\}$

normalization: „float“ the decimal point such that $d_1 \neq 0$ (and adjust the exponent appropriately)

Since we can only store a bounded range of integers, we need

$$e_{\min} \leq e \leq e_{\max}$$

For numbers with $e > e_{\max}$ we get an overflow error / luf
 $e < e_{\min}$ we get an underflow error / 0

IEEE 754 floating point numbers

double (precision): 64 bits = 8 bytes

s: 1	e: 11	f: 52
------	-------	-------

$$x = \begin{cases} (-1)^s 1.d_2 \dots d_p \cdot e^{e-1023} & \text{if } e \in 1:2046 \\ (-1)^s 0.d_2 \dots d_p \cdot e^{-1023} & \text{if } e = 0 \\ \text{luf} & \text{if } e = 2047 \text{ and } f = 0 \\ \text{NaN} \text{ „not a number“} & \text{if } e = 2047 \text{ and } f \neq 0 \end{cases}$$

Machine epsilon

is the smallest number ε such that (in floating point arithmetic)

$$1 + \varepsilon \neq 1$$

$$\begin{array}{rcl} 1.00 \dots 0 & \cdot 2^0 & = 1 \\ + 0.00 \dots 1 & \cdot 2^0 & = \varepsilon \\ \hline & & = 1.00 \dots 1 \cdot 2^0 \neq 1 \end{array}$$

double precision: $\varepsilon \approx 2 \cdot 10^{-16}$

Rounding

$a \in \mathbb{R}$, $a \neq 0$, then the floating point number \hat{a} which represents a satisfies

$$|a - \hat{a}| \leq \frac{\varepsilon}{2} |a|$$

\Leftrightarrow

$$\frac{|a - \hat{a}|}{|a|} \leq \frac{\varepsilon}{2}$$

relative error of \hat{a}

IEEE 754 arithmetic

satisfies the following :

for any arithmetic operation $\circ \in \{+, -, \cdot, /\}$ and any floating point numbers a, b :

$$\frac{|a \circ b - \widehat{a \circ b}|}{|a \circ b|} \leq \frac{\epsilon}{2}.$$

where $\widehat{a \circ b} = \widehat{a \circ b}$

1.2. Conditioning of problems

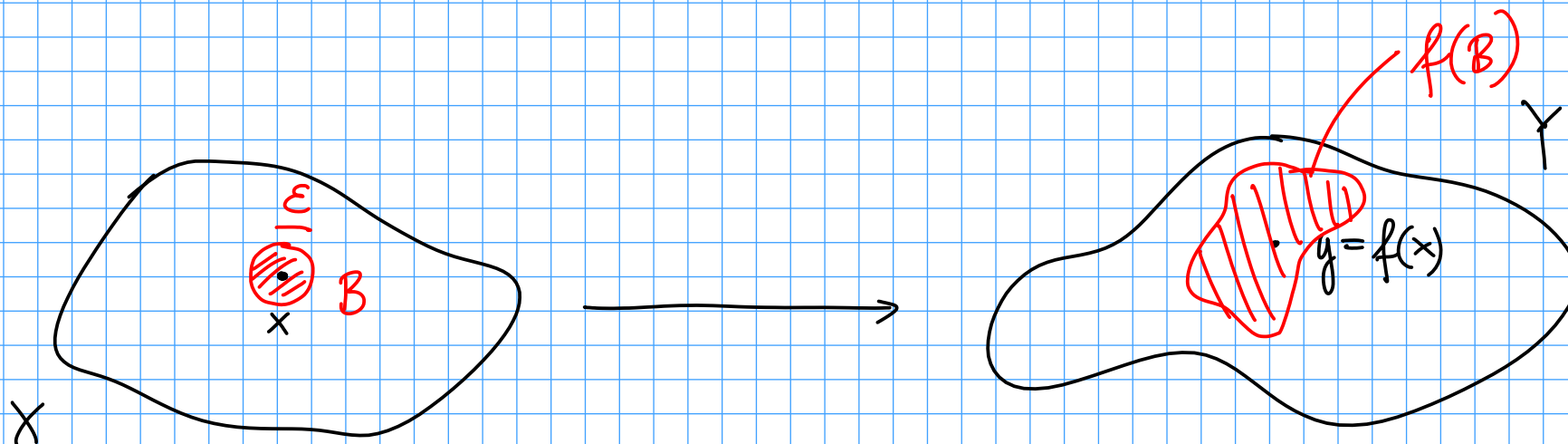
problem (in a mathematical sense): $f: X \rightarrow Y$
input output

example: find the zeros of some quadratic polynomial

$$x^2 - 2px - q$$

input: $(p, q) \in \mathbb{R}^2 = X$

output: $(x_1, x_2) \in \mathbb{R}^2 = Y$ s.t. $x_{1,2}^2 - 2px_{1,2} - q = 0$



condition number of f (at x) is the smallest number $K_f > 0$

st.

$$\|f(x) - f(\hat{x})\| \leq K_f \|x - \hat{x}\| \quad (*)$$

where $\|\cdot\|$ measures the error, typically this the relative error

$$\|x - \hat{x}\| = \frac{|x - \hat{x}|}{|x|}$$

example: $f: (x, y) \in \mathbb{R}^2 \mapsto x + y$

let's compute the condition of f :

$$\frac{|(x+y) - (\hat{x} + \hat{y})|}{|x+y|} = \frac{|(x-\hat{x}) + (y-\hat{y})|}{|x+y|}$$

$$= \frac{|\delta x + \delta y|}{|x+y|}$$

$$\leq \frac{|\delta x| + |\delta y|}{|x+y|}$$

$$\begin{aligned} x &= \hat{x} + \delta x \\ y &= \hat{y} + \delta y \end{aligned}$$

now use that $\frac{|x - \hat{x}|}{|x|} = \frac{|\delta x|}{|x|} \leq \varepsilon$, i.e. $|\delta x| \leq \varepsilon |x|$

likewise $|\delta y| \leq \varepsilon |y|$

ie we obtain $\frac{|\delta x| + |\delta y|}{|x+y|} \leq \frac{\varepsilon |x| + \varepsilon |y|}{|x+y|}$

$$= \underbrace{\frac{|x| + |y|}{|x+y|}}_{= K_f} \varepsilon = K_f \llbracket x - \hat{x} \rrbracket$$

Case 1: x and y have the same sign : $|x| + |y| = |x+y|$

$$\Rightarrow K_f = 1$$

Case 2: x and y have opposite sign : $|x+y| = ||x| - |y||$

So for $|x| \approx |y|$, $|x+y| \approx 0$ and K_f will be large !

Remark: In the output, one typically loses $\log_{10} K_f$ decimal places of accuracy, i.e. $K_f = 10^{10} \Rightarrow$ we lose 10 digit of accuracy

example: loss of accuracy in addition:

$$\begin{aligned} \rightarrow x &= 1.23467 \cdot 10^0 \\ y &= -1.23456 \cdot 10^0 \\ x+y &= 0.00011 \cdot 10^0 \\ &= 1.1 \text{???} \cdot 10^{-4} \end{aligned}$$

cancellation of leading digits

$$\begin{aligned} K_f &= \frac{|x|+|y|}{|x+y|} \\ &\approx \frac{2}{10^{-4}} \approx 10^4 \end{aligned}$$

other arithmetic operation:

\otimes	\cdot	$/$	$\sqrt{\quad}$
K	2	2	$\frac{1}{2}$

More generally: $f: \mathbb{R}^n \rightarrow \mathbb{R}$ differentiable, $x = (x_1, \dots, x_n)^T \mapsto y$

$$K_f = \frac{\langle |\nabla f(x)|, |x| \rangle}{|f(x)|} = \frac{\sum_{i=1}^n |\partial_i f(x)| |x_i|}{|f(x)|}$$

(derive by linearization: $f(x+\delta x) - f(x) = \langle \nabla f(x), \delta x \rangle$)

1.3 Stability of algorithms

problem: $x \mapsto y = f(x)$

expected error in y : $\|y - \hat{y}\| \leq K_f \|x - \hat{x}\|$

algorithm: $f_n \circ \dots \circ f_2 \circ f_1$, $f_1, f_2, \dots \in \{+, -, \cdot, /, \sqrt{}\}$

goal: $= f$

this sequence of elementary arithmetic operations for evaluating f is not unique!

implementation: $\hat{f} = \hat{f}_n \circ \dots \circ \hat{f}_2 \circ \hat{f}_1 = fl \circ f_n \circ \dots \circ fl \circ f_2 \circ fl \circ f_1$
 \uparrow rounding from x to \hat{x}

simple case: $f = f_2 \circ f_1$
 $\hat{f} = f_2 \circ fl \circ f_1$

if f_2 has a large condition number, then the error in the output will be large

stable algorithm: $\|f(x) - \hat{f}(x)\| = O(K_f \cdot \epsilon)$
 $= C K_f \epsilon$
C constant "not too large"

unstable algorithm: $\|f(x) - \hat{f}(x)\| \gg K_f \epsilon$

example: problem: compute the zeros of a quadratic polynomial
 $x^2 - 2px - q$

input: $(p, q) \in \mathbb{R}^2$

output: $(x_1, x_2) \in \mathbb{R}^2$

problem is well conditioned

school algorithm: $x_1 = p - \sqrt{p^2 + q}$
 $x_2 = p + \sqrt{p^2 + q}$

observation:

$$(p, q) \xrightarrow{f_1} \begin{pmatrix} r = \sqrt{p^2 + q} \\ p \end{pmatrix} \xrightarrow{f_2} p - r = x_1 \quad \times$$

$$\xrightarrow{f_2} p + r = x_2 \quad \checkmark$$

K_{f_2} is large (in our case), i.e. this algorithm is *unstable*

better algorithm:

$$x^2 - 2px - q = (x - x_1)(x - x_2) = x^2 - (x_1 + x_2)x + x_1x_2$$

i.e. $-q = x_1x_2 \Rightarrow x_1 = -\frac{q}{x_2}$

$$(p, q) \xrightarrow{f_1} \begin{pmatrix} r = \sqrt{p^2 + q} \\ q \\ p \end{pmatrix} \xrightarrow{f_2} \begin{pmatrix} x_2 = p + r \\ q \end{pmatrix} \xrightarrow{f_3} \begin{pmatrix} x_2 \\ x_1 = -\frac{q}{x_2} \end{pmatrix}$$



