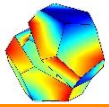


Theoretical and practical aspects of computer arithmetic

Siegfried M. Rump, Hamburg/Tokyo



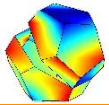
1/63



Back

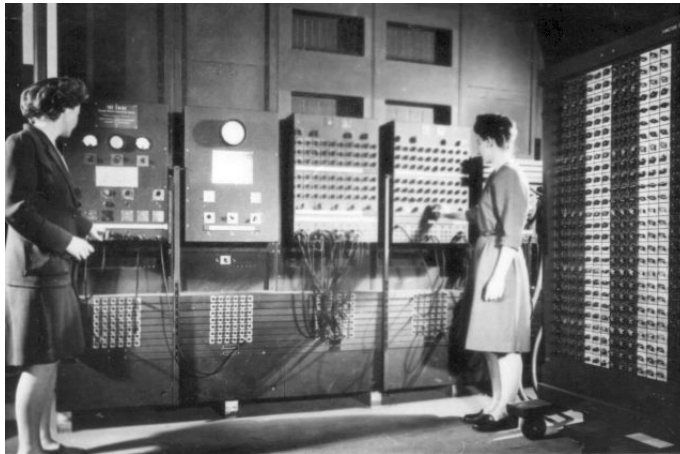
Close

The origin of floating-point



2/63

| Computer | base | arithmetic | method | Turing complete |
|-----------------|---------|----------------|------------|-----------------------|
| Zuse Z3 | binary | floating-point | relais | yes |
| Atanasoff-Berry | binary | fixed point | tubes | no [linsys $n < 30$] |
| Colossus | binary | fixed point | tubes | no [deciphering] |
| Mark I | decimal | fixed point | relais | yes |
| Eniac | decimal | fixed point | tubes | yes |
| Babbage | decimal | fixed point | mechanical | yes [not built] |



Back

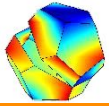
Close

The origin of error analysis I

Carl-Friedrich Gauß was fully aware of computational errors and developed a complete and rigorous error analysis



Based on his computations Ceres was rediscovered



3/63



Back

Close

The origin of error analysis II

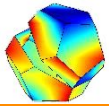
In their seminal paper

Numerical inverting of matrices of high order (1947)

John v. Neumann and Hermann Goldstine stated:

“Cholesky decomposition in 24-bit fixed point arithmetic may produce reliable results up to dimension $n \leq 9$.”

The analysis is correct but far too pessimistic

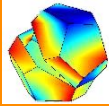


4/63



Back

Close



5/63

Limits of computer arithmetic

Let $A \subseteq \mathbb{R}$ with $|A| < \infty$.

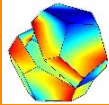
There is no isomorphism from \mathbb{R} to A .

There is no meaningful homomorphism respecting order relations.



Back

Close



Limits of computer arithmetic

Let $\mathbb{A} \subseteq \mathbb{R}$ with $|\mathbb{A}| < \infty$.

There is no isomorphism from \mathbb{R} to \mathbb{A} .

There is no meaningful homomorphism respecting order relations.

Under very general assumptions it can be shown that operations on \mathbb{A} cannot meet the law of associativity or distributivity.

That is due to the finiteness of \mathbb{A} .



Back

Close

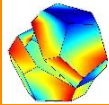
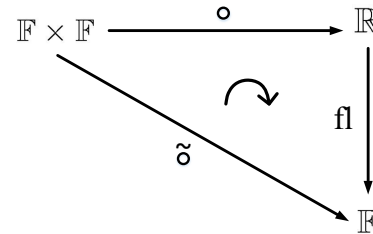
The IEEE 754 arithmetic standard 1984 - a closer look

$\pm 1.m_1m_2\dots m_k \cdot 2^e$ binary floating-point

\mathbb{F} set of floating-point numbers

Define a mapping (rounding) $\text{fl} : \mathbb{R} \rightarrow \mathbb{F}$

Operations $\tilde{\circ} : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$ are defined by
 $a \tilde{\circ} b := \text{fl}(a \circ b)$



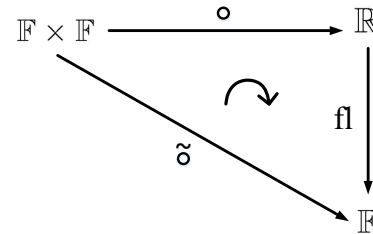
The IEEE 754 arithmetic standard 1984 - a closer look

$\pm 1.m_1m_2\dots m_k \cdot 2^e$ binary floating-point

\mathbb{F} set of floating-point numbers

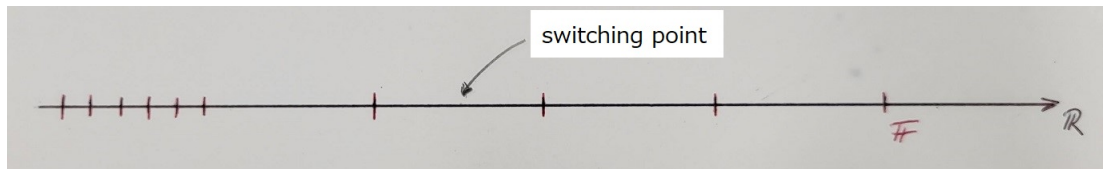
Define a mapping (rounding) $\text{fl} : \mathbb{R} \rightarrow \mathbb{F}$

Operations $\tilde{\circ} : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$ are defined by
 $a \tilde{\circ} b := \text{fl}(a \circ b)$

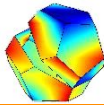


In rounding to nearest, the mapping fl_{\square} has minimal error:

$$x \in \mathbb{R} \Rightarrow |\text{fl}_{\square}(x) - x| = \min\{|f - x| : f \in \mathbb{F}\}$$



The results of arithmetic operations $\tilde{\circ}$ is best possible.



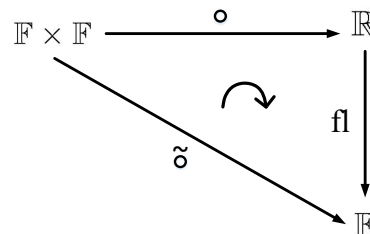
The IEEE 754 arithmetic standard 1984 - a closer look

$\pm 1.m_1m_2\dots m_k \cdot 2^e$ binary floating-point

\mathbb{F} set of floating-point numbers

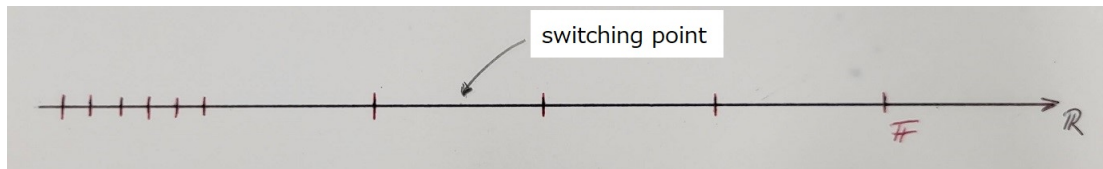
Define a mapping (rounding) $\text{fl} : \mathbb{R} \rightarrow \mathbb{F}$

Operations $\tilde{\circ} : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$ are defined by
 $a \tilde{\circ} b := \text{fl}(a \circ b)$



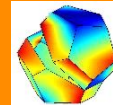
In rounding to nearest, the mapping fl_{\square} has minimal error:

$$x \in \mathbb{R} \Rightarrow |\text{fl}_{\square}(x) - x| = \min\{|f - x| : f \in \mathbb{F}\}$$



The results of arithmetic operations $\tilde{\circ}$ is best possible.

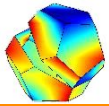
What means "best"?



The relative rounding error - switching points

First standard model $E_1(x) := \left| \frac{\text{fl}(x) - x}{x} \right|$ rel. err. w.r.t. x

Switching point: *arithmetic* mean of adjacent fl-pt numbers



7/63



Back

Close

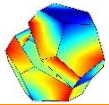
The relative rounding error - switching points

First standard model $E_1(x) := \left| \frac{\text{fl}(x) - x}{x} \right|$ rel. err. w.r.t. x

Switching point: *arithmetic* mean of adjacent fl-pt numbers

Second standard model $E_2(x) := \left| \frac{\text{fl}(x) - x}{\text{fl}(x)} \right|$ rel. err. w.r.t. $\text{fl}(x)$

Switching point: *harmonic* mean of adjacent fl-pt numbers



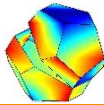
7/63



Back

Close

The relative rounding error - switching points



7/63

First standard model $E_1(x) := \left| \frac{\text{fl}(x) - x}{x} \right|$ rel. err. w.r.t. x

Switching point: *arithmetic* mean of adjacent fl-pt numbers

Second standard model $E_2(x) := \left| \frac{\text{fl}(x) - x}{\text{fl}(x)} \right|$ rel. err. w.r.t. $\text{fl}(x)$

Switching point: *harmonic* mean of adjacent fl-pt numbers

Minimize $\max\{E_1(x), E_2(x)\}$

Switching point: *geometric* mean of adjacent fl-pt numbers

S.M. Rump and M. Lange. On the Definition of Unit Roundoff.

BIT Numerical Mathematics, 56(1):309–317, 2015.

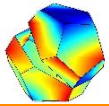


Back

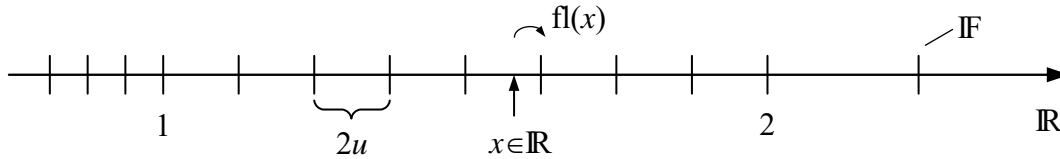
Close

The standard models for the relative rounding error

Rounding to nearest with relative rounding error unit \mathbf{u}



8/63



$$x \in [1, 2] : \quad |\text{fl}(x) - x| \leq \mathbf{u}$$

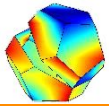


Back

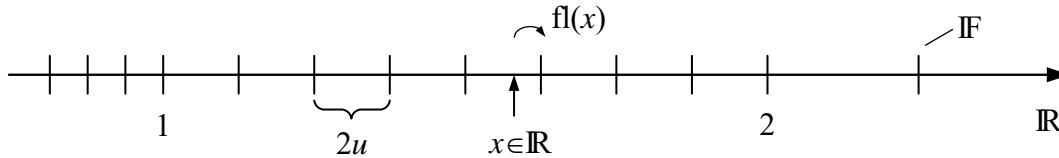
Close

The standard models for the relative rounding error

Rounding to nearest with relative rounding error unit \mathbf{u}



8/63



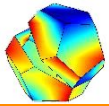
$$x \in [1, 2] : \quad |\text{fl}(x) - x| \leq \mathbf{u}$$

$$\text{relative rounding error} \quad E_2(x) := \left| \frac{\text{fl}(x) - x}{\text{fl}(x)} \right| \leq \frac{\mathbf{u}}{1} = \mathbf{u} \quad \text{w.r.t. } \text{fl}(x)$$

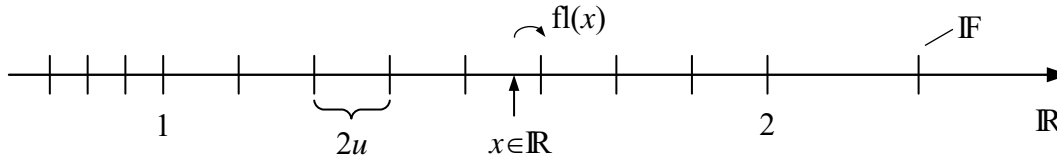
$$\Rightarrow \quad (1 + \varepsilon)\text{fl}(x) = x \quad |\varepsilon| \leq \mathbf{u}$$

The standard models for the relative rounding error

Rounding to nearest with relative rounding error unit \mathbf{u}



8/63



$$x \in [1, 2] : \quad |\text{fl}(x) - x| \leq \mathbf{u}$$

$$\text{relative rounding error} \quad E_2(x) := \left| \frac{\text{fl}(x) - x}{\text{fl}(x)} \right| \leq \frac{\mathbf{u}}{1} = \mathbf{u} \quad \text{w.r.t. } \text{fl}(x)$$

$$\Rightarrow \quad (1 + \varepsilon)\text{fl}(x) = x \quad |\varepsilon| \leq \mathbf{u}$$

$$E_1(x) := \left| \frac{\text{fl}(x) - x}{x} \right| = \left| \frac{\varepsilon \text{fl}(x)}{(1 + \varepsilon)\text{fl}(x)} \right| = \left| \frac{\varepsilon}{1 + \varepsilon} \right| \leq \frac{\mathbf{u}}{1 + \mathbf{u}} \quad \text{w.r.t. } x$$

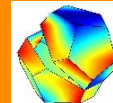
P.H. Sterbenz: Floating-Point Computations, Prentice-Hall, 1974



Back

Close

Optimal bounds of floating-point operations



9/63

| t | bound on $E_1(t)$ | bound on $E_2(t)$ |
|-------------|--|--|
| real number | $\frac{\mathbf{u}}{1+\mathbf{u}}$ | \mathbf{u} |
| $a \pm b$ | $\frac{\mathbf{u}}{1+\mathbf{u}}$ | \mathbf{u} |
| ab | $\frac{\mathbf{u}}{1+\mathbf{u}}$ | \mathbf{u} |
| a/b | $\begin{cases} \mathbf{u} - 2\mathbf{u}^2 & \text{if } \beta = 2, \\ \frac{\mathbf{u}}{1+\mathbf{u}} & \text{if } \beta > 2 \end{cases}$ | $\begin{cases} \frac{\mathbf{u}-2\mathbf{u}^2}{1+\mathbf{u}-2\mathbf{u}^2} & \text{if } \beta = 2, \\ \mathbf{u} & \text{if } \beta > 2 \end{cases}$ |
| \sqrt{a} | $1 - \frac{1}{\sqrt{1+2\mathbf{u}}}$ | $\sqrt{1 + 2\mathbf{u}} - 1$ |

The bounds are optimal for p -digit base- β IEEE-754 arithmetic under some mild conditions.

For example, multiplication in base $\beta = 2$ requires that

$2^p + 1$ is not a Fermat prime.

C.-P. Jeannerod and S.M. Rump. On relative errors of floating-point operations:

Optimal bounds and applications. *Mathematics of Computation*, 87:803–819, 2018.

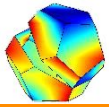


Composed operations: Classical Wilkinson-type error estimates

Summation $p_1 + p_2 + \dots + p_n$

recursive summation $\hat{s} := p_1$

$$\hat{s}_i := \hat{s}_{i-1} \tilde{+} p_i \quad \text{for } i \in \{2, \dots, n\}$$



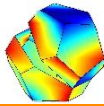
10/63



Back

Close

Composed operations: Classical Wilkinson-type error estimates



10/63

Summation $p_1 + p_2 + \dots + p_n$

recursive summation $\hat{s} := p_1$

$$\hat{s}_i := \hat{s}_{i-1} \tilde{+} p_i \quad \text{for } i \in \{2, \dots, n\}$$

... now “Epsilontik” starts

classical $\hat{s}_n = (\dots ((p_1 + p_2)(1 + \varepsilon_1) + p_3)(1 + \varepsilon_2) + \dots p_n)(1 + \varepsilon_{n-1})$

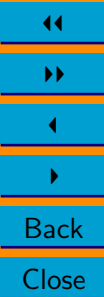
$$\Rightarrow \left| \hat{s}_n - \sum_{i=1}^n p_i \right| \leq ((1 + \mathbf{u})^{n-1} - 1) \sum_{i=1}^n |p_i| \leq \underbrace{\frac{(n-1)\mathbf{u}}{1 - (n-1)\mathbf{u}}}_{\gamma_{n-1}} \sum_{i=1}^n |p_i|$$

[provided that $(n-1)\mathbf{u} < 1$]

γ_{n-1}

Classical since the 1960's

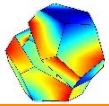
but not “nice”



Linearized bounds for composed operations !

$$[\text{R. 2012}] \quad \left| \hat{s} - \sum_{i=1}^n p_i \right| \leq (n-1)\mathbf{u} \sum_{i=1}^n |p_i|$$

no limit on n



11/63



Back

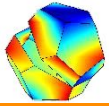
Close

Linearized bounds for composed operations !

$$[\text{R. 2012}] \quad \left| \hat{s} - \sum_{i=1}^n p_i \right| \leq (n-1)\mathbf{u} \sum_{i=1}^n |p_i|$$

no limit on n

... the race began



11/63



Back

Close

Linearized bounds for composed operations !

$$[\text{R. 2012}] \quad \left| \hat{s} - \sum_{i=1}^n p_i \right| \leq (n-1)\mathbf{u} \sum_{i=1}^n |p_i|$$

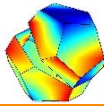
... the race began

no limit on n

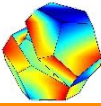
$$[\text{Jeannerod, R. 2013}] \quad \left| \hat{s} - \sum_{i=1}^n x_i \right| \leq n\mathbf{u} \sum_{i=1}^n |x_i|$$

- $x_i \in \mathbb{R}$
- summation of $\text{fl}(x_i)$ in floating-point
- any base $\beta \geq 2$
- any order of evaluation
- no limit on n

$$\text{Corollary} \quad |\hat{r} - a^T b| \leq n\mathbf{u} |a^T| |b| \quad \text{for } a, b \in \mathbb{F}^n$$



More linearized bounds for compound operations



12/63

[Graillat, Lefèvre, Muller 2015] power

$$|\hat{r} - a^{k+1}| \leq k\mathbf{u}|a^{k+1}| \quad \text{if } k \leq \sqrt{2^{1/3} - 1}\mathbf{u}^{-1/2} - 1$$

- base $\beta = 2$
- successive multiplication

[R., Bünger, Jeannerod 2015] products

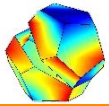
$$\left| \hat{r} - \prod_{i=0}^k x_i \right| \leq k\mathbf{u} \left| \prod_{i=0}^k x_i \right| \quad \text{for } x_i \in \mathbb{F}, \beta = 2, k < \mathbf{u}^{-1/2}$$

- any order of evaluation
- limit on k is mandatory
- $k < \mathbf{u}^{-1/2}$ cannot be replaced by $k < 12\mathbf{u}^{-1/2}$



Back

Close



[R., Bünger, Jeannerod 2015]

Horner's scheme

$$\left| \hat{r} - \sum_{i=0}^n a_i x^i \right| \leq 2n\mathbf{u} \sum_{i=0}^n |a_i x^i| \quad \text{if } n < \frac{1}{2} \left(\sqrt{\frac{\omega}{\beta}} \mathbf{u}^{-1/2} - 1 \right).$$

Classical

$$|\hat{r} - \|p\|_2| \leq ((1 + \mathbf{u})^{n/2+1} - 1) \|p\|_2 \quad \text{for } p \in \mathbb{F}^n$$

[Jeannerod, R. 2016]

$$|\hat{r} - \|p\|_2| \leq \left(\frac{n}{2} + 1\right) \mathbf{u} \|p\|_2$$

- any order of evaluation
- no restriction on n



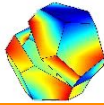
Back

Close

Linearized bounds for algorithms

Classical $\gamma_k := \frac{k\mathbf{u}}{1-k\mathbf{u}}$, $k\mathbf{u} < 1$

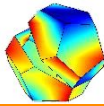
- $A \in \mathbb{F}^{m \times n}$, computed LU -factors \hat{L}, \hat{U} :
 $\hat{L}\hat{U} = A + \Delta A, \quad |\Delta A| \leq \gamma_n |\hat{L}| |\hat{U}|$
- $A \in \mathbb{F}^{n \times n}$, computed Cholesky factor \hat{R} :
 $\hat{R}^T \hat{R} = A + \Delta A, \quad |\Delta A| \leq \gamma_{n+1} |\hat{R}^T| |\hat{R}|$
- $T \in \mathbb{F}^{n \times n}$ triangular, $b \in \mathbb{F}^n$, $\hat{x} = T \backslash b$:
 $(T + \Delta T)\hat{x} = b, \quad |\Delta T| \leq \gamma_n |T|$



Linearized bounds for algorithms

Improved [R., Jeannerod (2015)]

no limit on n



15/63

- $A \in \mathbb{F}^{m \times n}$, computed LU -factors \hat{L}, \hat{U} :
 $\hat{L}\hat{U} = A + \Delta A, \quad |\Delta A| \leq n\mathbf{u}|\hat{L}||\hat{U}|$
- $A \in \mathbb{F}^{n \times n}$, computed Cholesky factor \hat{R} :
 $\hat{R}^T \hat{R} = A + \Delta A, \quad |\Delta A| \leq (n+1)\mathbf{u}|\hat{R}^T||\hat{R}|$
- $T \in \mathbb{F}^{n \times n}$ triangular, $b \in \mathbb{F}^n$, $\hat{x} = T \setminus b$:
 $(T + \Delta T)\hat{x} = b, \quad |\Delta T| \leq n\mathbf{u}|T|$



Back

Close

Towards a more general perspective

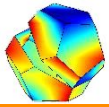
Up to now:

- We *actively* assumed base- β IEEE-754 conform arithmetic.
- Every result relied on that specific arithmetic.

Next:

- *Passively* identify sufficient assumptions to prove linearized bounds.

⇒ Understand “Machine numbers” \mathbb{M} as a subset of \mathbb{R}



16/63



Back

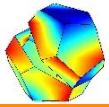
Close

An arithmetic on a general subset of \mathbb{R}

$\mathbb{M} \subseteq \mathbb{R}$, $\square : \mathbb{M} \times \mathbb{M} \rightarrow \mathbb{M}$ for $\circ \in \{+, -, \times, /\}$, also $\sqrt{\cdot}$.

$$x, y \in \mathbb{M}: \quad x \square y = (x \circ y)(1 + \delta) \quad |\delta| \leq \textit{eps}$$

for some constant *eps*. We do *not* assume a rounding function fl !



17/63



Back

Close

An arithmetic on a general subset of \mathbb{R}

$\mathbb{M} \subseteq \mathbb{R}$, $\square : \mathbb{M} \times \mathbb{M} \rightarrow \mathbb{M}$ for $\circ \in \{+, -, \times, /\}$, also $\sqrt{\cdot}$.

$$x, y \in \mathbb{M}: \quad x \square y = (x \circ y)(1 + \delta) \quad |\delta| \leq \textit{eps}$$

for some constant *eps*. We do *not* assume a rounding function fl !

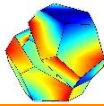
Much freedom:

- $x \circ y \in \mathbb{M} \quad \not\Rightarrow \quad x \square y = x \circ y$
- $a \circ b = c \circ d \quad \not\Rightarrow \quad a \square b = c \square d$

Example 3-digit decimal format, $p = 3$, $\textit{eps} = \frac{1}{2}\beta^{1-p} = 0.005$

$$x + y = 9.96$$

$$\Rightarrow x \square + y \in \{9.92, 9.93, 9.94, 9.95, 9.96, 9.97, 9.98, 9.99, 10.0\}$$



An arithmetic on a general subset of \mathbb{R}

$\mathbb{M} \subseteq \mathbb{R}$, $\square : \mathbb{M} \times \mathbb{M} \rightarrow \mathbb{M}$ for $\circ \in \{+, -, \times, /\}$, also $\sqrt{\cdot}$.

$$x, y \in \mathbb{M}: \quad x \square y = (x \circ y)(1 + \delta) \quad |\delta| \leq eps$$

for some constant *eps*. We do *not* assume a rounding function fl !

Much freedom:

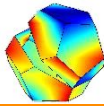
- $x \circ y \in \mathbb{M} \quad \not\Rightarrow \quad x \square y = x \circ y$
- $a \circ b = c \circ d \quad \not\Rightarrow \quad a \square b = c \square d$

Example 3-digit decimal format, $p = 3$, $eps = \frac{1}{2}\beta^{1-p} = 0.005$

$$x + y = 9.96$$

$$\Rightarrow x \square + y \in \{9.92, 9.93, 9.94, 9.95, 9.96, 9.97, 9.98, 9.99, 10.0\}$$

e.g. $9.90 \square + 0.06 = 10$ $9.91 \square + 0.05 = 9.92$



An arithmetic on a general subset of \mathbb{R}

$\mathbb{M} \subseteq \mathbb{R}$, $\square : \mathbb{M} \times \mathbb{M} \rightarrow \mathbb{M}$ for $\circ \in \{+, -, \times, /\}$, also $\sqrt{\cdot}$.

$$x, y \in \mathbb{M}: \quad x \square y = (x \circ y)(1 + \delta) \quad |\delta| \leq \textit{eps}$$

for some constant *eps*. We do *not* assume a rounding function fl !

Much freedom:

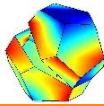
- $x \circ y \in \mathbb{M} \quad \not\Rightarrow \quad x \square y = x \circ y$ also $x \square y$ may change
- $a \circ b = c \circ d \quad \not\Rightarrow \quad a \square b = c \square d$

Example 3-digit decimal format, $p = 3$, $\mathbf{u} = \frac{1}{2}\beta^{1-p} = 0.005$

$$x + y = 9.96$$

$$\Rightarrow x \square + y \in \{9.92, 9.93, 9.94, 9.95, 9.96, 9.97, 9.98, 9.99, 10.0\}$$

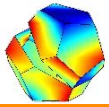
e.g. $9.90 \square + 0.06 = 10$ $9.91 \square + 0.05 = 9.92$ $9.91 \square + 0.05 = 9.96$



Linearized bounds: An even simplified exposition

$$\forall a, b \in \mathbb{M}: \quad |(a \boxplus b) - (a + b)| \leq \min(|a|, |b|)$$

Assumption A



19/63



Back

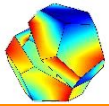
Close

Linearized bounds: An even simplified exposition

$$\forall a, b \in \mathbb{M}: \quad |(a \boxplus b) - (a + b)| \leq \min(|a|, |b|)$$

$$\text{Very weak: } |3 \boxplus 4 - (3 + 4)| \leq \min(3, 4) = 3$$

Assumption A



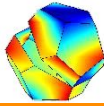
19/63



Back

Close

Linearized bounds: An even simplified exposition



19/63

$$\forall a, b \in \mathbb{M}: \quad |(a \boxplus b) - (a + b)| \leq \min(|a|, |b|) \quad \text{Assumption A}$$

$$\text{Very weak: } |3 \boxplus 4 - (3 + 4)| \leq \min(3, 4) = 3$$

$$\text{IEEE-754 } x \in \mathbb{R}: \quad |\text{fl}(x) - x| = \min\{|f - x| : f \in \mathbb{F}\} \quad \text{nearest}$$

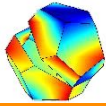
$$\begin{aligned} \Rightarrow \quad |a \boxplus b - (a + b)| &= |\text{fl}(a + b) - (a + b)| \\ &= \min(|f - (a + b)| : f \in \mathbb{F}) \\ &\leq \min(|a - (a + b)|, |b - (a + b)|) \\ &= \min(|a|, |b|) \end{aligned}$$



Back

Close

Linearized bounds: An even simplified exposition



19/63

$$\forall a, b \in \mathbb{M}: \quad |(a \boxplus b) - (a + b)| \leq \min(|a|, |b|) \quad \text{Assumption A}$$

$$\text{Very weak:} \quad |3 \boxplus 4 - (3 + 4)| \leq \min(3, 4) = 3$$

$$\text{IEEE-754} \quad x \in \mathbb{R}: \quad |\text{fl}(x) - x| = \min\{|f - x| : f \in \mathbb{F}\} \quad \text{nearest}$$

$$\begin{aligned} \Rightarrow \quad |a \boxplus b - (a + b)| &= |\text{fl}(a + b) - (a + b)| \\ &= \min(|f - (a + b)| : f \in \mathbb{F}) \\ &\leq \min(|a - (a + b)|, |b - (a + b)|) \\ &= \min(|a|, |b|) \end{aligned}$$

Not satisfied for rounding upwards:

$$1 \boxplus \mathbf{u}^2 = \text{succ}(1) = 1 + 2\mathbf{u} \quad \Rightarrow \quad 2\mathbf{u} - \mathbf{u}^2 \not\leq \min(1, \mathbf{u}^2) = \mathbf{u}^2$$



Back

Close

The linearized error estimate

Theorem. Let an arithmetic on \mathbb{M} with Assumption A be given.
For $p \in \mathbb{M}^n$ define

$$\hat{s}_1 := p_1; \quad \hat{s}_k = \hat{s}_{k-1} \boxplus p_k = (\hat{s}_{k-1} + p_k)(1 + \delta_k) \quad \text{for } 2 \leq k \leq n$$

with $|\delta_k| \leq eps$.

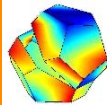
Then

$$\left| \hat{s}_n - \sum_{i=1}^n p_i \right| \leq \sum_{i=1}^n |\delta_i| \sum_{i=1}^n |p_i| \leq (n-1)eps \sum_{i=1}^n |p_i| \quad (*)$$

The result is true under much more general assumptions

E.g. (*) is true for directed rounding (not satisfying Assumption A)

M. Lange and S.M. Rump. Error estimates for the summation of real numbers with application to floating-point summation. BIT, 57:927–941, 2017.



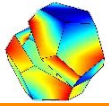
Optimal bounds for summation

Worst case $1 + \mathbf{u} + \mathbf{u} + \dots ?$

Mascarenhas 2016:

$$\beta = 2, \quad p \in \mathbb{F}^n, \quad n \leq \frac{1}{5} 2^{p-2} : \quad \left| \hat{s} - \sum_{i=1}^n p_i \right| \leq \frac{(n-1)\mathbf{u}}{1+(n-1)\mathbf{u}} \sum_{i=1}^n |p_i|$$

Proof uses some optimization and continuous mathematics



Optimal bounds for summation

Worst case $1 + \mathbf{u} + \mathbf{u} + \dots ?$

Mascarenhas 2016:

$$\beta = 2, \quad p \in \mathbb{F}^n, \quad n \leq \frac{1}{5} 2^{p-2} : \quad \left| \hat{s} - \sum_{i=1}^n p_i \right| \leq \frac{(n-1)\mathbf{u}}{1+(n-1)\mathbf{u}} \sum_{i=1}^n |p_i|$$

Proof uses some optimization and continuous mathematics

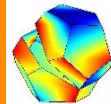
Theorem For an arithmetic on \mathbb{M} with Assumption A and $x \in \mathbb{M}^n$

$$\left| \hat{s} - \sum_{i=1}^n x_i \right| \leq \frac{\sum_{i=1}^{n-1} \xi_i}{1 + \sum_{i=1}^{n-1} \xi_i} \sum_{i=1}^n |x_i| \quad [\text{IEEE-754: } |\xi_i| \leq \mathbf{u}]$$

The estimate is sharp.

M. Lange and S.M. Rump. Sharp estimates for perturbation errors in summations.

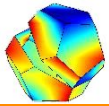
Math. of Comp., 88:349–368, 2019.



Error-free transformations

```
function [x,y] = TwoSum(a,b)
    x = a + b;
    z = x - a;
    y = ( a - (x-z) ) + (b-z);
```

Knuth 1969: $a, b \in \mathbb{F} \Rightarrow x + y = a + b$



22/63



Back

Close

Error-free transformations

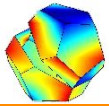
```
function [x,y] = TwoSum(a,b)
    x = a + b;
    z = x - a;
    y = ( a - (x-z) ) + (b-z);
```

Knuth 1969: $a, b \in \mathbb{F} \Rightarrow x + y = a + b$

```
function [x,y] = FastTwoSum(a,b)
    x = a + b;
    y = a - (x - b);
```

Dekker 1971: $a, b \in \mathbb{F}, |a| \geq |b| \Rightarrow x + y = a + b$

FastTwoSum with comparison often 2 times slower than **TwoSum**



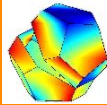
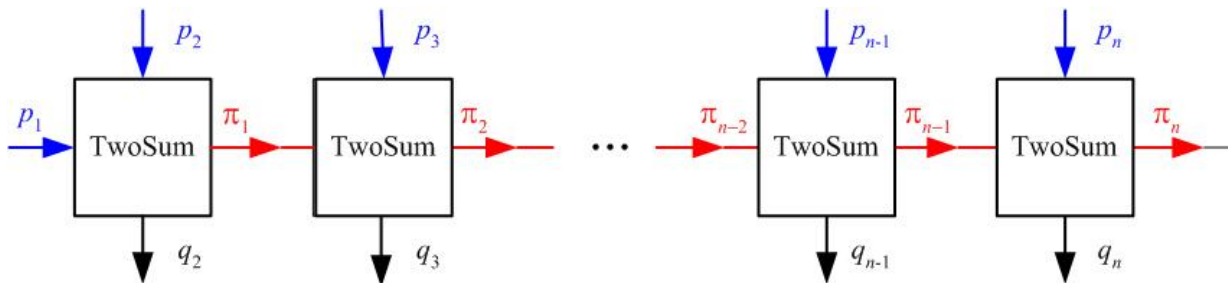
Error-free vector transformations

```
function p = VecSum(p)
```

```
  for i=2:n
```

```
    [p(i),p(i-1)] = TwoSum(p(i),p(i-1))
```

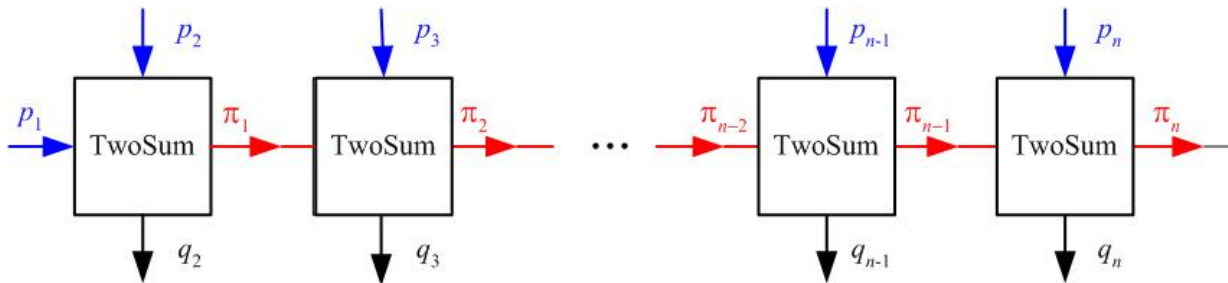
$$q = \text{VecSum}(p) \Rightarrow \sum q_i = \sum p_i, \quad q_n = \text{float}(\sum p_i)$$



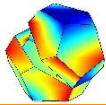
Error-free vector transformations

```
function p = VecSum(p)
    for i=2:n
        [p(i),p(i-1)] = TwoSum(p(i),p(i-1))
```

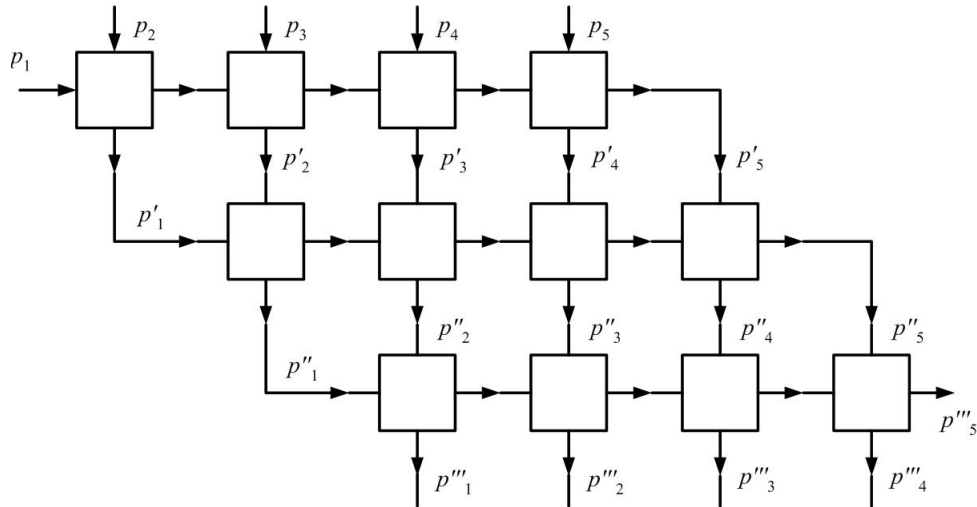
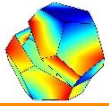
$$q = \text{VecSum}(p) \Rightarrow \sum q_i = \sum p_i, \quad q_n = \text{float}(\sum p_i)$$



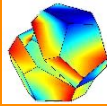
Error of $\text{sum}(p)$ of the order $[(n-1)\mathbf{u}]^2$



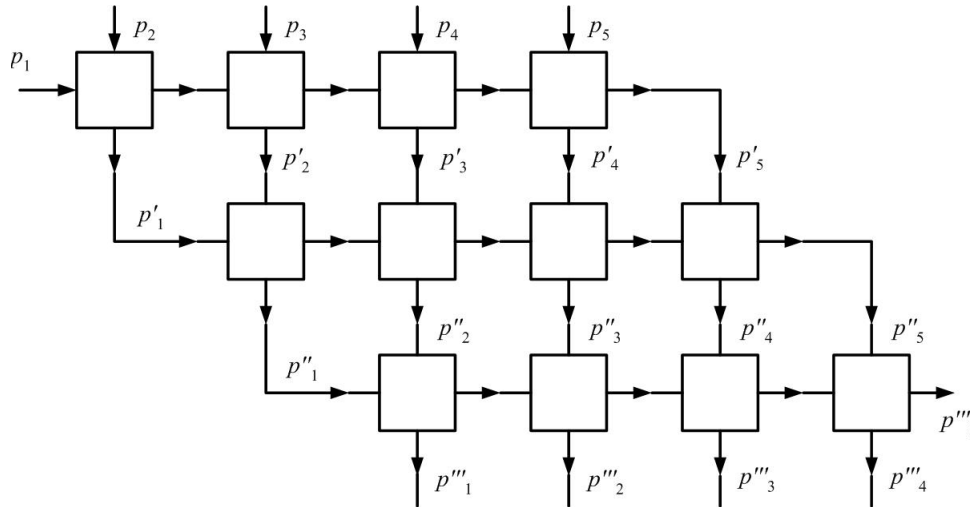
Iterated error-free vector transformations



Iterated error-free vector transformations



24/63



Error of $\mathbf{sum}(\mathbf{p})$ of the order $[(n - 1)\mathbf{u}]^{K+1}$ after K transformations

Similar routines for dot products, most important in numerical analysis

T. Ogita, S.M. Rump, and S. Oishi. Accurate sum and dot product. SIAM Journal on Scientific Computing (SISC), 26(6):1955–1988, 2005.



Back

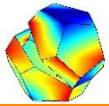
Close

The power of modern error analysis

John v. Neumann and Hermann Goldstine stated:

“Cholesky decomposition in 24-bit fixed point arithmetic may produce reliable results up to dimension $n \leq 9$.”

Theorem. Let $A \in \mathbb{F}^{n \times n}$ with $A^T = A$ be given, and let $B = A - D \in \mathbb{F}^{n \times n}$ for diagonal D with $D \geq 2\alpha I$ and $\alpha \geq \gamma_{n+1} \text{trace}(A) > 0$.



25/63



Back

Close

The power of modern error analysis

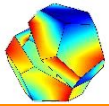
John v. Neumann and Hermann Goldstine stated:

“Cholesky decomposition in 24-bit fixed point arithmetic may produce reliable results up to dimension $n \leq 9$.”

Theorem. Let $A \in \mathbb{F}^{n \times n}$ with $A^T = A$ be given, and let $B = A - D \in \mathbb{F}^{n \times n}$ for diagonal D with $D \geq 2\alpha I$ and $\alpha \geq \gamma_{n+1} \text{trace}(A) > 0$.

If the *floating-point* Cholesky decomposition of B runs to completion, then A is symmetric positive definite, and for any $\tilde{x} \in \mathbb{R}^n$

$$\|A^{-1}b - \tilde{x}\|_2 \leq \alpha^{-1} \|A\tilde{x} - b\|_2.$$



25/63



Back

Close

The power of modern error analysis

John v. Neumann and Hermann Goldstine stated:

“Cholesky decomposition in 24-bit fixed point arithmetic may produce reliable results up to dimension $n \leq 9$.”

Theorem. Let $A \in \mathbb{F}^{n \times n}$ with $A^T = A$ be given, and let $B = A - D \in \mathbb{F}^{n \times n}$ for diagonal D with $D \geq 2\alpha I$ and $\alpha \geq \gamma_{n+1} \text{trace}(A) > 0$.

If the *floating-point* Cholesky decomposition of B runs to completion, then A is symmetric positive definite, and for any $\tilde{x} \in \mathbb{R}^n$

$$\|A^{-1}b - \tilde{x}\|_2 \leq \alpha^{-1} \|A\tilde{x} - b\|_2.$$

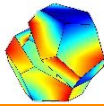
That approach works for dimensions n in the 10-thousands.

All operations are in ordinary floating-point arithmetic !

The analysis is based on properties of a symm. pos. def. matrix

S.M. Rump and T. Ogita. Super-fast validated solution of linear systems.

JCAM, 199(2):199–206, 2006.



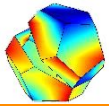
Towards solving general problems

What about general linear systems, nonlinear systems, global optimization, differential equations etc. ?

We may use interval arithmetic:

$$[a, b] \circ [c, d] := [\min x, \max x] \quad \text{for} \quad x \in \{a \circ c, a \circ d, b \circ c, b \circ d\}$$

On the computer we use directed roundings.



26/63



Back

Close

Towards solving general problems

What about general linear systems, nonlinear systems, global optimization, differential equations etc. ?

We may use interval arithmetic:

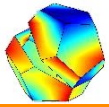
$$[a, b] \circ [c, d] := [\min x, \max x] \quad \text{for} \quad x \in \{a \circ c, a \circ d, b \circ c, b \circ d\}$$

On the computer we use directed roundings.

Fundamental inclusion property:

$$\forall a \in A, b \in B: \quad a \circ b \in A \circ B \quad \text{for interval quantities } A, B$$

Covers all elementary standard functions, erf, $\Gamma(x)$ etc. as well



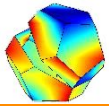
Towards solving general problems

Fundamental observation:

Replace in an algorithm all operations
by the corresponding interval operations.

If finished successfully, i.e., no division by a zero interval, then

- It is mathematically certain that the problem is solvable, and
- the computed results do contain the true solution.



27/63



Back

Close

Towards solving general problems

Fundamental observation:

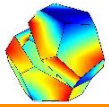
Replace in an algorithm all operations
by the corresponding interval operations.

If finished successfully, i.e., no division by a zero interval, then

- It is mathematically certain that the problem is solvable, and
- the computed results do contain the true solution.

This is called *naive* interval arithmetic

Why does interval arithmetic has a bad reputation?



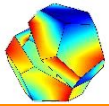
27/63



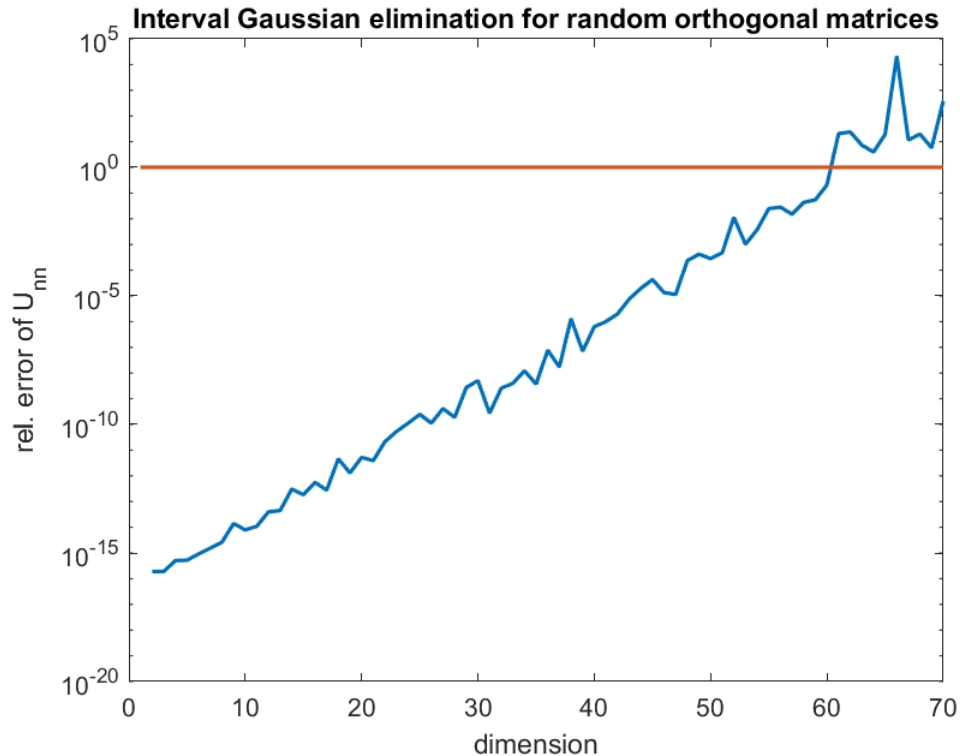
Back

Close

Naive interval arithmetic: Interval Gaussian elimination (IGA)



28/63



The matrices are perfectly well conditioned: $\text{cond}(A) = 1$



Back

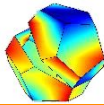
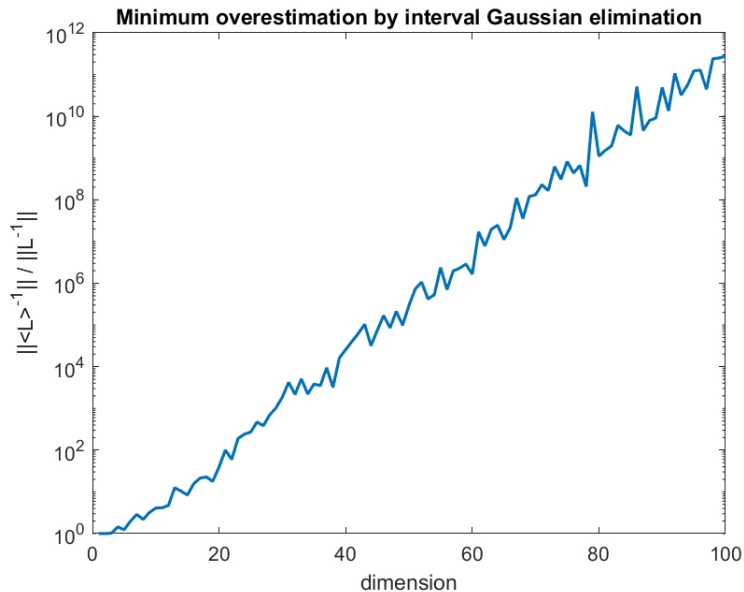
Close

Minimum overestimation for Interval Gaussian elimination (IGA)

Theorem [R., 2010] For $A \in \mathbb{R}^{n \times n}$ perform Gaussian elimination with total pivoting using *real* interval operations everywhere.

If finished successfully, then elementwise

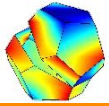
$$\text{rad}(U) \geq \text{upper triangle} (\langle L \rangle^{-1} \cdot \text{rad}(A))$$



The reason for the poor reputation of interval arithmetic

Historically, interval arithmetic was (at least) known to Gauss.

It was taught in German junior high schools from the mid 19th century.



30/63



Back

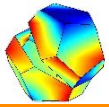
Close

The reason for the poor reputation of interval arithmetic

Historically, interval arithmetic was (at least) known to Gauss.

It was taught in German junior high schools from the mid 19th century.

It was *re-discovered* in the 1960's and *advocated as the holy grail* .



30/63



Back

Close

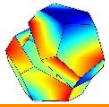
The reason for the poor reputation of interval arithmetic

Historically, interval arithmetic was (at least) known to Gauss.

It was taught in German junior high schools from the mid 19th century.

It was *re-discovered* in the 1960's and *advocated as the holy grail* .

The problem is not the tool [interval arithmetic],
but the way it was used :



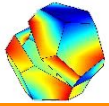
30/63



Back

Close

Tools — should be used appropriately I



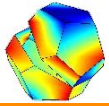
31/63



Back

Close

Tools — should be used appropriately II



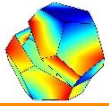
32/63



Back

Close

Is interval arithmetic of any use?



33/63

The (unique) advantage of interval arithmetic is to compute bounds for the range of a function over some domain.

The bounds may overestimate the true range, but they are always mathematically true.

A Matlab example ...



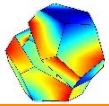
Back

Close

How to fight overestimation of interval arithmetic

A verification method should:

- use floating-point arithmetic wherever possible
- try to avoid the dependency problem
- try to scale intervals by a small number



34/63



Back

Close

How to fight overestimation of interval arithmetic

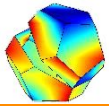
A verification method should:

- use floating-point arithmetic wherever possible
- try to avoid the dependency problem
- try to scale intervals by a small number

THEOREM. Let $A, R \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$. If for given $X \in \mathbb{IR}^n$

$$Rb + (I - RA)X \subseteq \text{int}(X)$$

then A is nonsingular and $A^{-1}b \in X$.



How to fight overestimation of interval arithmetic

A verification method should:

- use floating-point arithmetic wherever possible
- try to avoid the dependency problem
- try to scale intervals by a small number

THEOREM. Let $A, R \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$. If for given $X \in \mathbb{IR}^n$

$$Rb + (I - RA)X \subseteq \text{int}(X)$$

then A is nonsingular and $A^{-1}b \in X$.

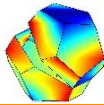
PROOF. Define $f(x) := Rb + (I - RA)x$. Then

$$\forall x \in X : f(x) \in X \quad \Rightarrow \quad \exists \hat{x} \in X : f(\hat{x}) = \hat{x} = Rb + \hat{x} - RA\hat{x}$$

by Brouwer's fixed point Theorem.

Inclusion in $\text{int}(X)$ implies R, A to be non-singular.

S.M. Rump. Kleine Fehlerschranken bei Matrixproblemen. PhD thesis, Univ. Karlsruhe, 1980.



How to fight overestimation of interval arithmetic

A verification method should:

- use floating-point arithmetic wherever possible
- try to avoid the dependency problem
- try to scale intervals by a small number

THEOREM. Let $A, R \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$. If for given $X \in \mathbb{IR}^n$

$$Rb + (I - RA)X \subseteq \text{int}(X) \quad \text{do NOT use } X + R(b - AX)$$

then A is nonsingular and $A^{-1}b \in X$.

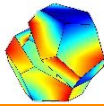
PROOF. Define $f(x) := Rb + (I - RA)x$. Then

$$\forall x \in X : f(x) \in X \quad \Rightarrow \quad \exists \hat{x} \in X : f(\hat{x}) = \hat{x} = Rb + \hat{x} - RA\hat{x}$$

by Brouwer's fixed point Theorem.

Inclusion in $\text{int}(X)$ implies R, A to be non-singular.

S.M. Rump. Kleine Fehlerschranken bei Matrixproblemen. PhD thesis, Univ. Karlsruhe, 1980.



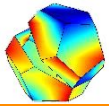
A verification method for systems of nonlinear equations

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $f \in \mathcal{C}^1$, $R \in \mathbb{R}^{n \times n}$, $\tilde{x} \in \mathbb{R}^n$, $X \in \mathbb{IR}^n$. If $\tilde{x} \in X$ and

$$(*) \quad -Rf(\tilde{x}) + (I - RJ_f(X))X \subseteq \text{int}(X),$$

then there exists a unique root \hat{x} of $f(x) = 0$ in $\tilde{x} + X$.

Verify (*) using interval arithmetic and algorithmic differentiation.



36/63



Back

Close

A verification method for systems of nonlinear equations

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $f \in C^1$, $R \in \mathbb{R}^{n \times n}$, $\tilde{x} \in \mathbb{R}^n$, $X \in \mathbb{IR}^n$. If $\tilde{x} \in X$ and

$$(*) \quad -Rf(\tilde{x}) + (I - RJ_f(X))X \subseteq \text{int}(X),$$

then there exists a unique root \hat{x} of $f(x) = 0$ in $\tilde{x} + X$.

Verify (*) using interval arithmetic and algorithmic differentiation.

Rationale, i.e., why is it working well:

$f(\tilde{x}) \approx 0$, $R \approx \frac{\partial f}{\partial x}(\tilde{x})^{-1}$ ensured by good fl-pt approximations

The error w.r.t. to the approximate solution \tilde{x} is included

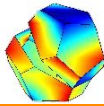
The product $(I - RJ_f(X))X$ is small in magnitude.

There is a dichotomy:

Either mathematically rigorous inclusion of the solution

or no result (error message)

S.M. Rump. Solving Algebraic Problems with High Accuracy. Habilitation, Acad. Press 1983.



36/63



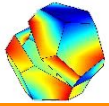
Back

Close

Global optimization in n dimensions

Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, minimize $f(x)$ over a box,
possibly subject to constraints

The main problem: To discard sub-boxes.



37/63



Back

Close

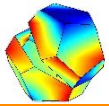
Global optimization in n dimensions

Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, minimize $f(x)$ over a box,
possibly subject to constraints

The main problem: To discard sub-boxes.

This is basically outside the scope of (purely) numerical algorithms.

Even if Lipschitz constants are known, rounding errors may have
disastrous effects.



37/63



Back

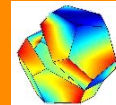
Close

Global minimization — Exclusion regions I

(1) Necessarily $\frac{\partial f}{\partial x}(\hat{x}) = 0 \rightarrow$

If $0 \notin \left[\frac{\partial f}{\partial x}(Y) \right]_i$ for some $1 \leq i \leq n$
and $Y \subseteq \text{int}(X)$

then Y can be discarded.

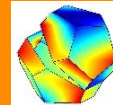


38/63



Back

Close



(1) Necessarily $\frac{\partial f}{\partial x}(\hat{x}) = 0 \rightarrow$

If $0 \notin \left[\frac{\partial f}{\partial x}(Y) \right]_i$ for some $1 \leq i \leq n$
and $Y \subseteq \text{int}(X)$

then Y can be discarded.

(2) Dimension reduction

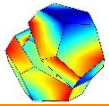
If $0 \notin \left[\frac{\partial f}{\partial x}(Y) \right]_i$ but $Y_i \cap \partial X_i \neq \emptyset$

then Y_i can be replaced by corresponding ∂X_i .

Exclusion regions II — The expansion principle (Jansson)

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f \in C^1$ be given.

For a given box X , our verification methods can prove that there is exactly one stationary point of f in X .



39/63



Back

Close

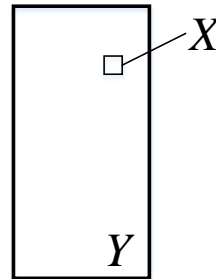
Exclusion regions II — The expansion principle (Jansson)

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f \in \mathcal{C}^1$ be given.

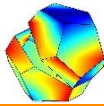
For a given box X , our verification methods can prove that there is exactly one stationary point of f in X .

Intentionally widen X into $Y \supseteq X$ and suppose that Y as well contains exactly one stationary point.

Then f has no minimum in $Y \setminus X$



C. Jansson. On Self-Validating Methods for Optimization Problems. In J. Herzberger (ed.) Topics in Validated Computations - Studies in Computational Mathematics 5, 381–438, North-Holland, Amsterdam, 1994.



39/63



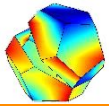
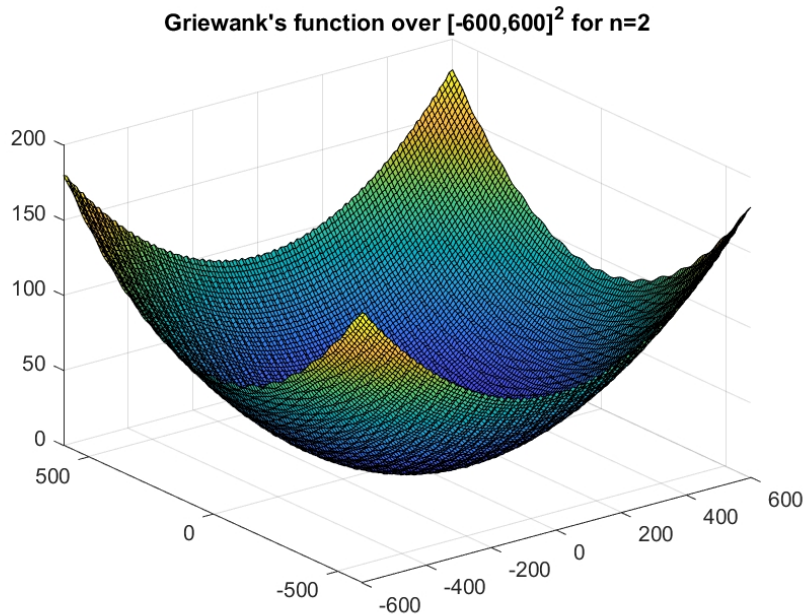
Back

Close

A famous test function in the global optimization community

Minimize Griewank's function $G : \mathbb{R}^n \rightarrow \mathbb{R}$ on $X = [-600, 600]^n$

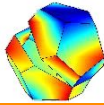
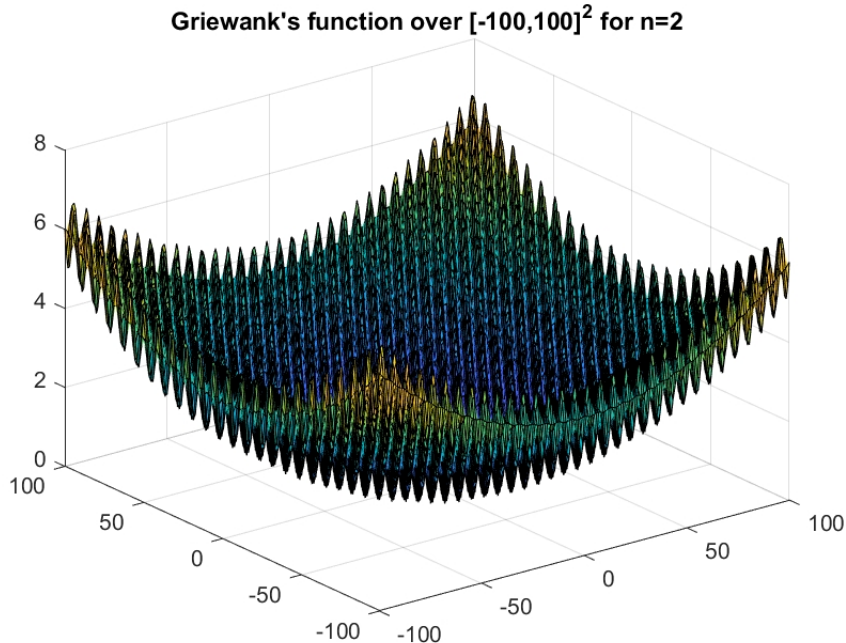
$$G(x) = 1 + \frac{1}{4000} \sum_{i=1}^n x_i^2 - \prod_{i=1}^n \cos\left(\frac{x_i}{\sqrt{i}}\right)$$



A famous test function

Minimize Griewank's function $G : \mathbb{R}^n \rightarrow \mathbb{R}$ on $X = [-600, 600]^n$

$$G(x) = 1 + \frac{1}{4000} \sum_{i=1}^n x_i^2 - \prod_{i=1}^n \cos\left(\frac{x_i}{\sqrt{i}}\right)$$



A famous test function

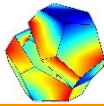
Minimize Griewank's function $G : \mathbb{R}^n \rightarrow \mathbb{R}$ on $X = [-600, 600]^n$

$$G(x) = 1 + \frac{1}{4000} \sum_{i=1}^n x_i^2 - \prod_{i=1}^n \cos\left(\frac{x_i}{\sqrt{i}}\right)$$

Timing [sec]

| n | $\#\nabla G(x) = 0$ | Montanher's | Csendes' | INTLAB |
|-----|---------------------|-------------|----------|--------|
| | | intsolver | GOP | |
| 5 | $\sim 10^{13}$ | 307*) | 229 | 0.6 |
| 10 | $\sim 10^{25}$ | | | 1.7 |
| 20 | $\sim 10^{51}$ | | | 5.2 |
| 30 | $\sim 10^{77}$ | | | 10.5 |
| 40 | $\sim 10^{103}$ | | | 17.9 |
| 50 | $\sim 10^{129}$ | | | 28.1 |

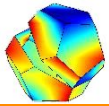
*) verification failed



INTLAB - the Matlab/Octave toolbox for Reliable Computing

- developing since 1998, >2000 routines, >70kLOC pure Matlab
- rigorous input and output
- Real and complex interval arithmetic and standard functions
- affine and Taylor arithmetic
- dense and sparse linear systems
- systems of nonlinear equations
- global optimization
- algorithmic differentiation, gradients, Hessians, Taylor series, slopes
- finding all roots of a nonlinear system
- Galois field toolbox
- etc.

<https://www.tuhh.de/ti3/rump/intlab/>



43/63



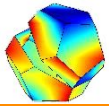
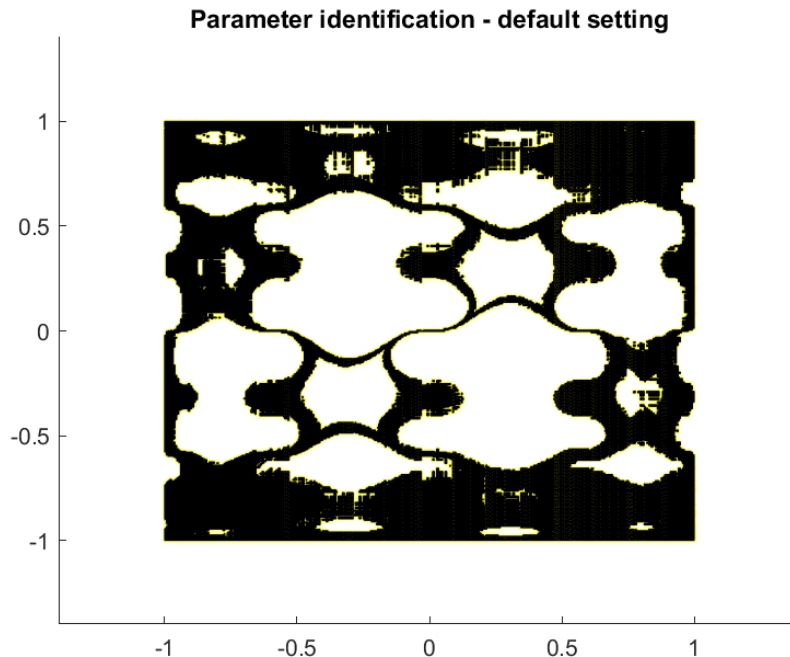
Back

Close

Parameter identification - ordinary interval arithmetic

$$f = -(5y - 20y^2 + 16y^5)^6 + (-(5x - 20x^3 + 16x^5)^3 + 5y^2 - 20y^3 + 16y^5)^2$$

```
X = infsup(-1,1)*ones(2,1);  
verifynlssparam(f,0,X)  
verifynlssparamset('Display','~');
```



44/63



Back

Close

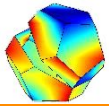
Interval arithmetic is no panacea — The wrapping effect

Ordinary interval arithmetic (executable INTLAB code):

```
A = infsup(1,3); intDiff = A-A
```

```
intval intDiff =
```

```
[ -2.0000, 2.0000]
```



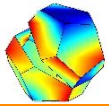
45/63



Back

Close

Interval arithmetic is no panacea — The wrapping effect



45/63

Ordinary interval arithmetic (executable INTLAB code):

```
A = infsup(1,3); intDiff = A-A
```

```
intval intDiff =  
[ -2.0000, 2.0000]
```

Affine arithmetic (executable INTLAB code):

```
B = affari(infsup(1,3)); affDiff = B-B
```

```
affari affDiff =  
[ 0.0000, 0.0000]
```



Back

Close

Affine arithmetic - References

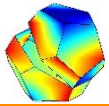
Ref.: Andrack, Comba, Stolfi 1994

Figueiredo/Stolfi, monograph 1997

Kashiwagi, monograph 2005

Stolfi, reference implementation 2007

R./Kashiwagi, Improvements of affine arithmetic, IEICE, 2015



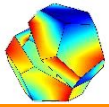
46/63



Back

Close

Affine arithmetic - References



46/63

Ref.: Andrack, Comba, Stolfi 1994

Figueiredo/Stolfi, monograph 1997

Kashiwagi, monograph 2005

Stolfi, reference implementation 2007

R./Kashiwagi, Improvements of affine arithmetic, IEICE, 2015

Representation of affine quantities:

$$C := \langle c; \gamma \rangle = \left\{ c + \sum_{i=1}^k \gamma_i \varepsilon_i : \varepsilon \in \mathcal{E}^k \right\} \quad \text{with } \mathcal{E} := [-1, 1]$$

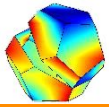
All ε_i vary independently in \mathcal{E} .



Back

Close

Affine arithmetic - References



46/63

Ref.: Andrack, Comba, Stolfi 1994

Figueiredo/Stolfi, monograph 1997

Kashiwagi, monograph 2005

Stolfi, reference implementation 2007

R./Kashiwagi, Improvements of affine arithmetic, IEICE, 2015

Representation of affine quantities:

$$C := \langle c; \gamma \rangle = \left\{ c + \sum_{i=1}^k \gamma_i \varepsilon_i : \varepsilon \in \mathcal{E}^k \right\} \quad \text{with } \mathcal{E} := [-1, 1]$$

All ε_i vary independently in \mathcal{E} .

For example,

`A = affari(infsup(1,3)); B = affari(infsup(-2,4));`

implies $A := \langle 2; 1 \rangle$ and $B := \langle 1; 0, 3 \rangle$



Back

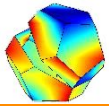
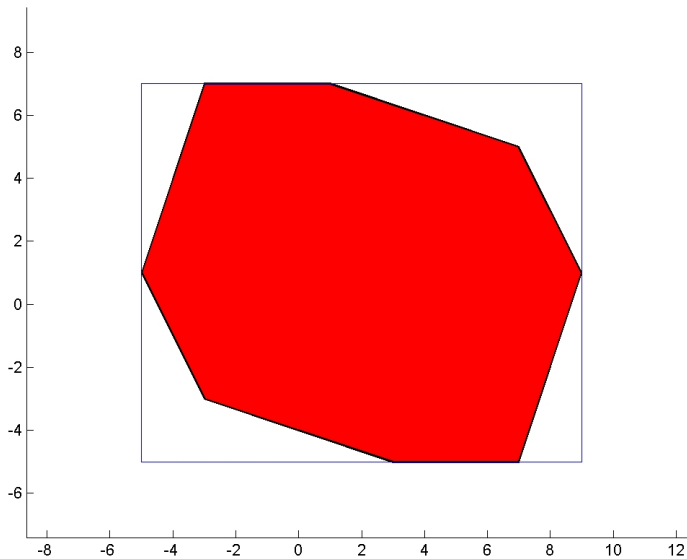
Close

Affine arithmetic

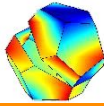
Example: $C := \langle 2; 1, -2, 3, -1 \rangle$

$D := \langle 1; 3, 0, -1, 2 \rangle$

$$C \times D = \left\{ \begin{pmatrix} 2 + \varepsilon_1 - 2\varepsilon_2 + 3\varepsilon_3 - \varepsilon_4 \\ 1 + 3\varepsilon_1 - \varepsilon_3 + 2\varepsilon_4 \end{pmatrix} : \varepsilon_i \in [-1, 1] \right\}$$



Affine arithmetic improvements



48/63

Given $C := \langle c; \gamma \rangle \stackrel{\wedge}{=} c + \sum \gamma_i \varepsilon_i \Rightarrow f(C)?$

Def. f is represented by $[[p, q, \Delta]]$ on $[a, b]$

$$\text{s.t. } \forall x \in \mathbf{X}: |px + q - f(x)| \leq \Delta$$

\Rightarrow Determine p, q, Δ for given f and given $[a, b]$

\Rightarrow Use $x \in [a, b] \Leftrightarrow x = c + \sum \gamma_i \varepsilon_i$ for $|\varepsilon_i| \leq 1$

$$\Rightarrow |pc + q + \sum p\gamma_i \varepsilon_i - f(x)| \leq \Delta$$

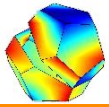
i.e. $\langle pc + q; p\gamma, \Delta \rangle$ represents $f(C)$ on $[a, b]$



Back

Close

Functions in the affine toolbox



49/63

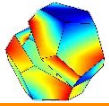
sqrt, sqr,
exp, log, log2, log10, power,
sin, cos, tan, cot, sec, csc,
asin, acos, atan, acot, asec, acsc,
sinh, cosh, tanh, coth,
asinh, acosh, atanh, acoth,
erf, erfc.



Back

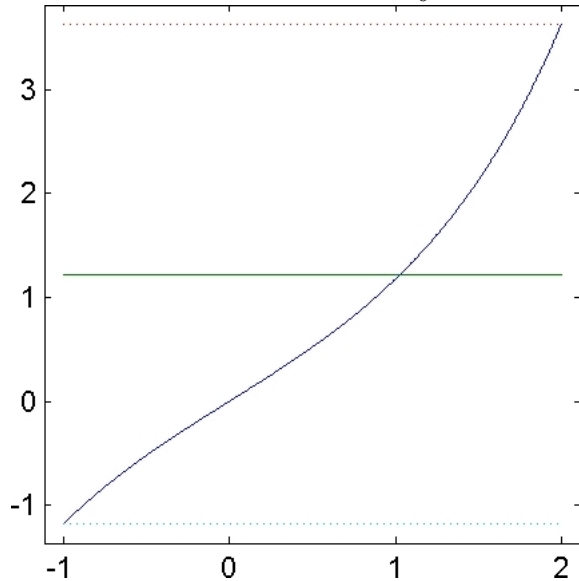
Close

Chebyshev representation

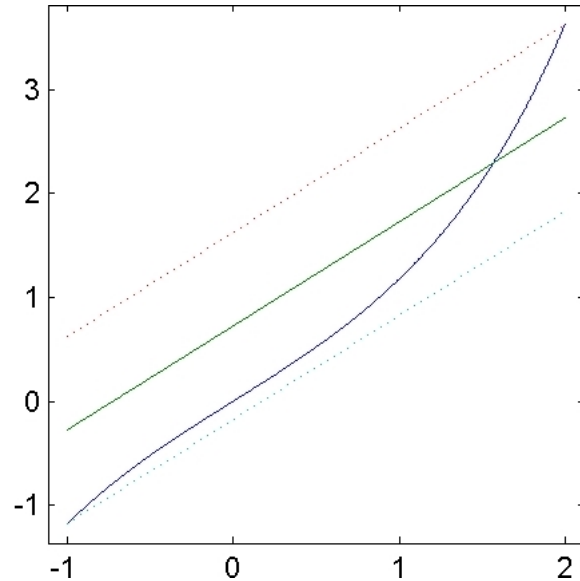


50/63

traditionally



new



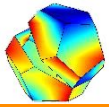
Min-Range representation of $\sinh(x)$ on $[-1, 2]$



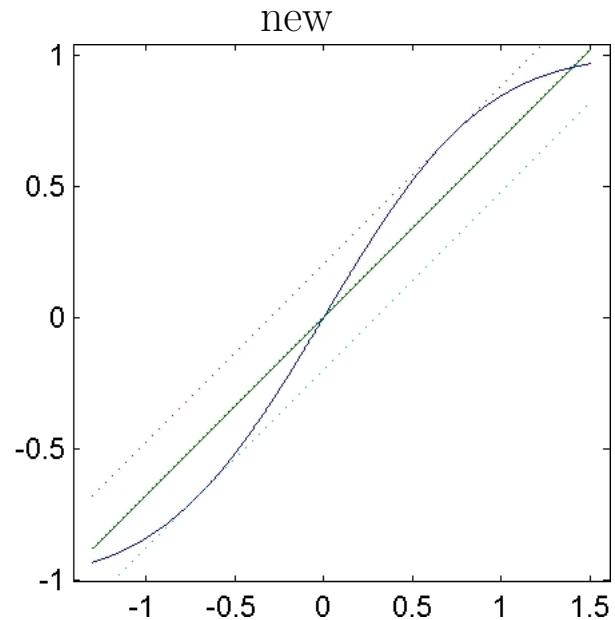
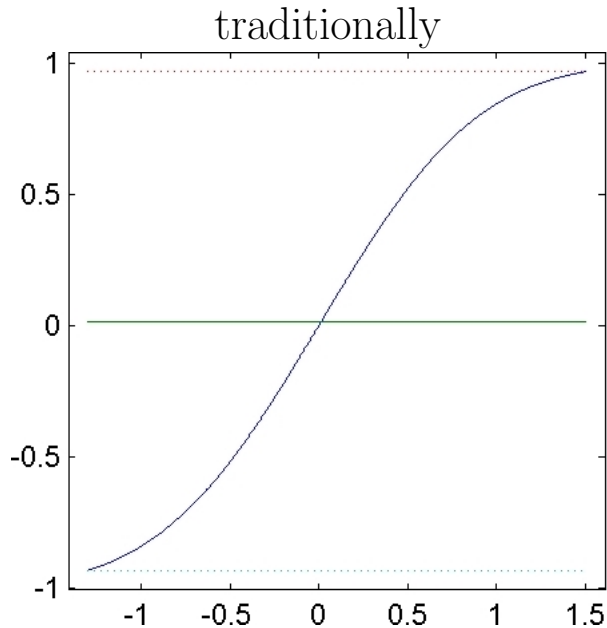
Back

Close

Chebyshev representation



51/63



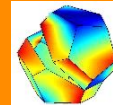
Chebyshev representation of $\text{erf}(x)$ on $[-1.3, 1.5]$



Back

Close

Parameter identification - ordinary interval arithmetic

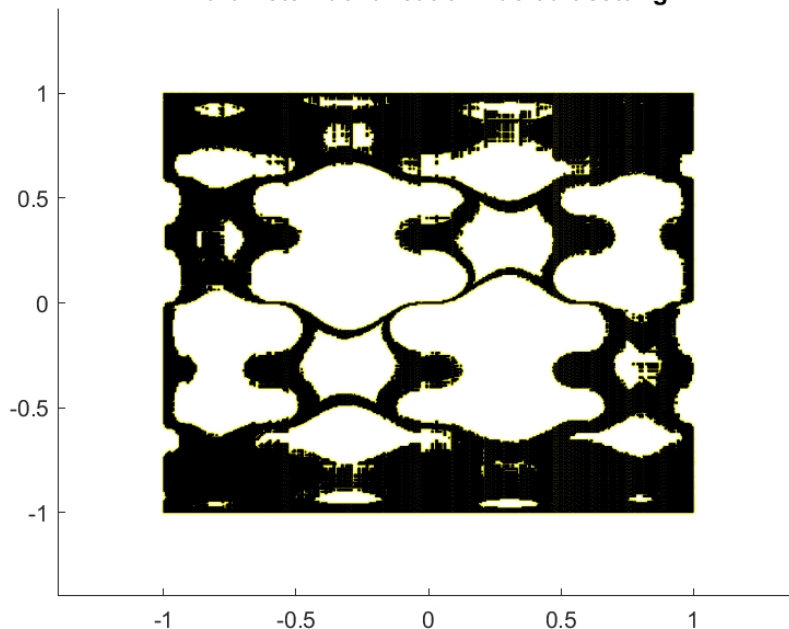


52/63

$$f = -(5y - 20y^2 + 16y^5)^6 + (-(5x - 20x^3 + 16x^5)^3 + 5y^2 - 20y^3 + 16y^5)^2$$

```
X = infsup(-1,1)*ones(2,1);  
verifynlssparam(f,0,X)  
verifynlssparamset('Display','~');
```

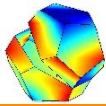
Parameter identification - default setting



Back

Close

Parameter identification - using affine arithmetic



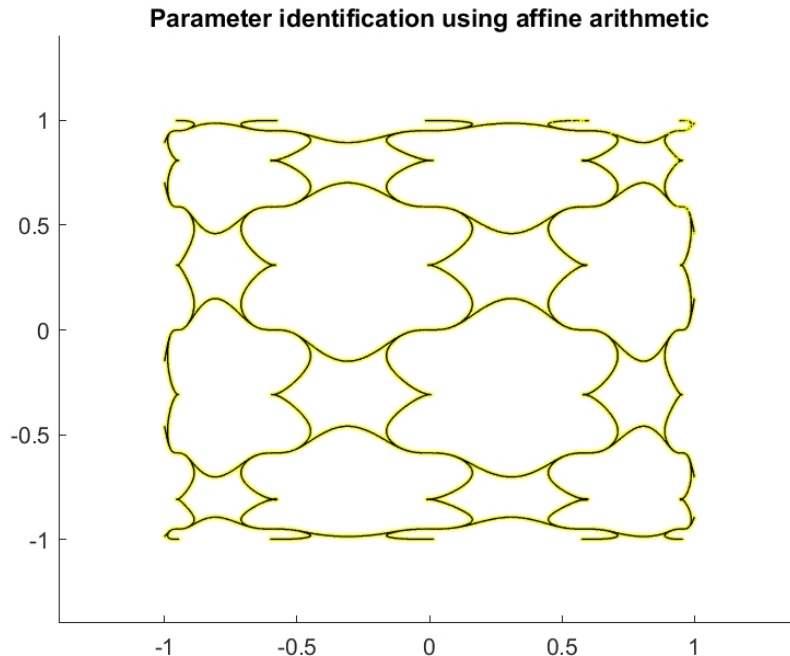
53/63

$$f = -(5y - 20y^2 + 16y^5)^6 + (-(5x - 20x^3 + 16x^5)^3 + 5y^2 - 20y^3 + 16y^5)^2$$

```
X = infsup(-1,1)*ones(2,1);
```

```
verifynlssparam(f,0,X)
```

```
verifynlssparamset('Display','~','Method','affari'));
```



Back

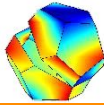
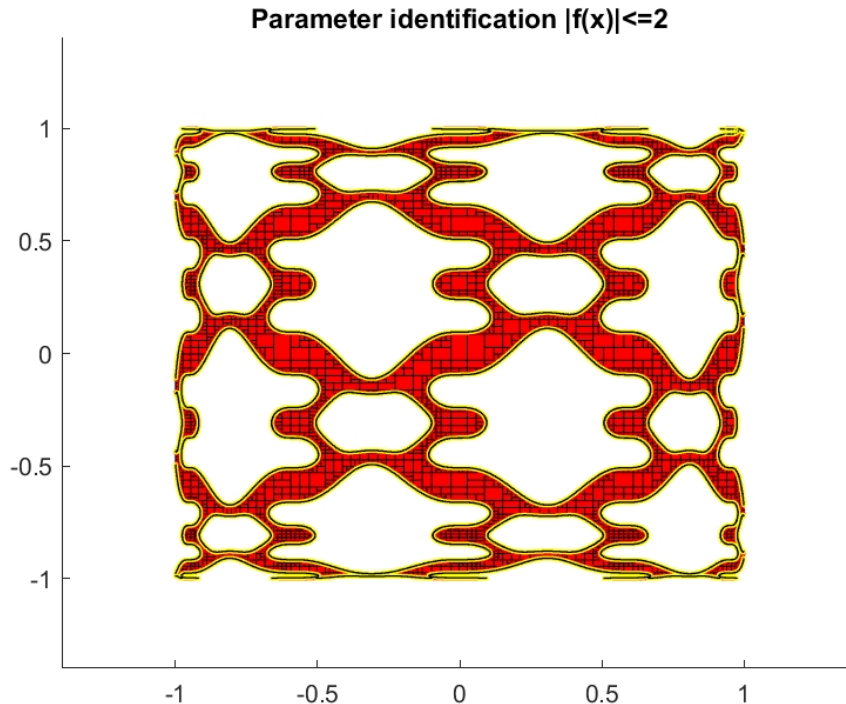
Close

Parameter identification - affine arithmetic, nontrivial interior

$$f = -(5y - 20y^2 + 16y^5)^6 + (-(5x - 20x^3 + 16x^5)^3 + 5y^2 - 20y^3 + 16y^5)^2$$

```
verifynlssparam(f, infsup(-0.2, 0.2), X, ...
```

```
verifynlssparamset('Display', '~', 'Method', 'affari'));
```



54/63



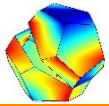
Back

Close

Limits of verification methods and of arithmetic

Verification methods can only solve well-posed problems

- No real interval can be verified to contain a root of $x^2 + 2x + 1$
- A complex interval can be verified to contain 2 roots



55/63



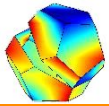
Back

Close

Limits of verification methods and of arithmetic

Verification methods can only solve well-posed problems

- No real interval can be verified to contain a root of $x^2 + 2x + 1$
- A complex interval can be verified to contain 2 roots
- We cannot verify that a matrix is singular
- We can verify that a matrix is nonsingular [even for $\text{cond}(A) > 10^{100}$]



55/63



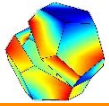
Back

Close

Limits of verification methods and of arithmetic

Verification methods can only solve well-posed problems

- No real interval can be verified to contain a root of $x^2 + 2x + 1$
- A complex interval can be verified to contain 2 roots
- We cannot verify that a matrix is singular
- We can verify that a matrix is nonsingular [even for $\text{cond}(A) > 10^{100}$]
- An eigenvector to a double eigenvalue cannot be included



55/63



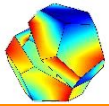
Back

Close

Limits of verification methods and of arithmetic

Verification methods can only solve well-posed problems

- No real interval can be verified to contain a root of $x^2 + 2x + 1$
- A complex interval can be verified to contain 2 roots
- We cannot verify that a matrix is singular
- We can verify that a matrix is nonsingular [even for $\text{cond}(A) > 10^{100}$]
- An eigenvector to a double eigenvalue cannot be included
- $\sin(4 \operatorname{atan}(1)) \geq 0$ cannot be decided in any precision



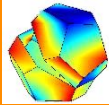
55/63



Back

Close

Fighting overestimation for verified ODE-solvers



56/63

The van der Pol equation

$$y'' - c(1 - y^2)y' + y = 0 \quad \text{for some scalar } c > 0$$

rewritten into a system of first order ODEs

$$\begin{aligned} y_1' &= y_2, \\ y_2' &= c(1 - y_1^2)y_2 - y_1 \end{aligned}$$

$$\text{Initial conditions } y_0 \in \begin{pmatrix} 3 \\ -3 \end{pmatrix} \pm 0.001$$



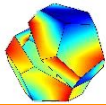
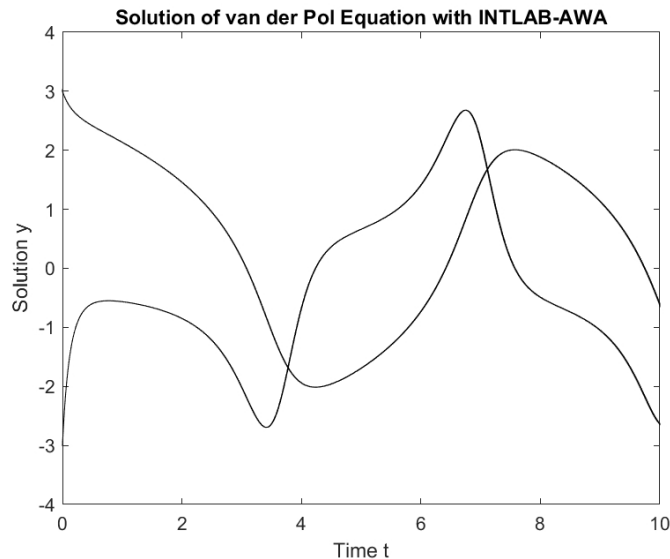
Back

Close

Solution of the van der Pol equation for $t \in [0, 10]$

```
function J = vdp_jac(t,y)
J = typeadjust([0,1;0,0],y);
J(2,1) = -2.*y(1).*y(2) - 1;
J(2,2) = 1 - sqr(y(1));
```

```
[T,Y] = awa(@vdp_fun,@vdp_jac,[0,10],midrad([3;-3],1e-3));
plot(T,Y)
```



57/63



Back

Close

Taylor models, implemented by Florian Bünger

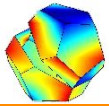
K. Makino, Rigorous analysis of nonlinear motion in particle accelerators, Dissertation at Michigan State University, 1998

K. Makino and M. Berz, Suppression of the wrapping effect by Taylor model - based validated integrators, MSU HEP Report 40910, 2003

A. Neumaier, Taylor Forms – Use and Limits, Reliable Computing 9, pp. 43-79, 2003

F. Bünger, Shrink wrapping for Taylor models revisited, Numerical Algorithms 78(4), pp. 1001-1017, 2018

F. Bünger, Reducing the truncation error in Taylor model multiplication, accepted for publication, 2023



Definition of Taylor models

$$p(x) = \sum_{a, |a| \leq d} p_a x^a, \quad |a| := a_1 + \dots + a_n, \quad x^a := x_1^{a_1} \dots x_n^{a_n},$$

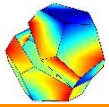
on $D = [u_1, v_1] \times \dots \times [u_n, v_n]$

An inclusion of $p(D - c) + E = \{p(x - c) + e \mid x \in D, e \in E\}$ is computed

Mostly the standard domain $D_s := [-1, 1]^n$, $cs := (0, \dots, 0)$ is used

In contrast to affine arithmetic, Taylor models need not to be convex

Options QR preconditioning, shrink wrapping,



59/63

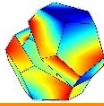


Back

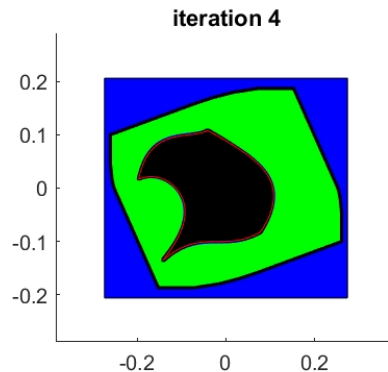
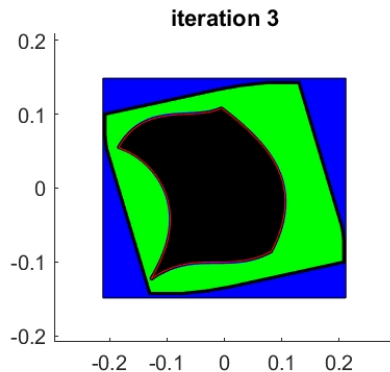
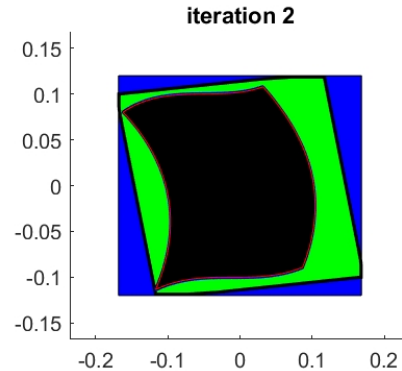
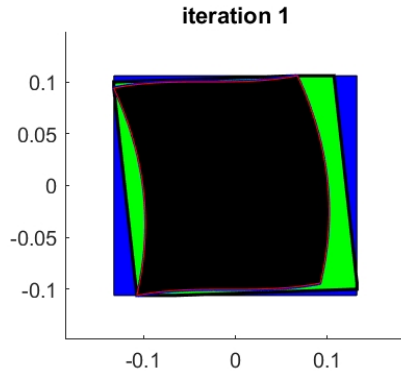
Close

Taylor models - An example

For $f(x, y) = (x - y(0.125 + 2y), y + 6x^3)$ compute the iterated image of $B := [-0.1, 0.1]^2$, i.e. $f(f(\dots f(B)\dots))$



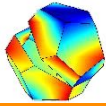
60/63



Back

Close

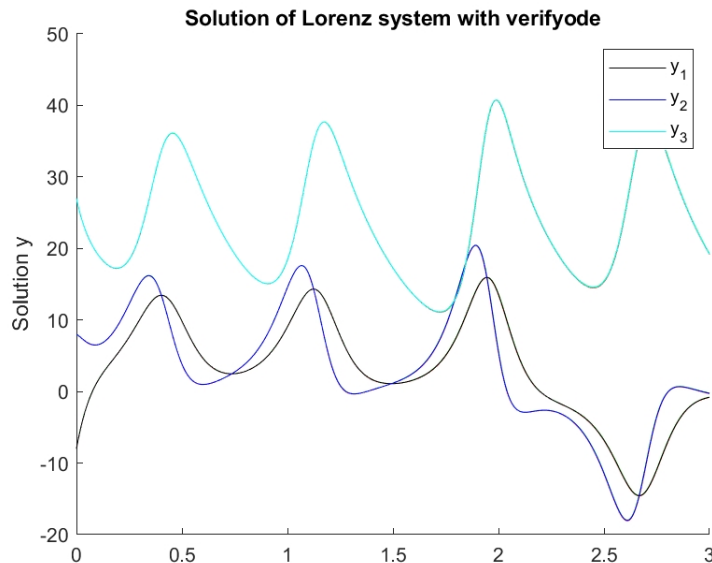
Taylor models - The Lorenz system



61/63

$$\begin{aligned}y_1' &= \sigma(y_2 - y_1) \\y_2' &= (\rho - y_3)y_1 - y_2 \quad \text{for } \sigma = 10, \rho = 28, \beta = 8/3 \\y_3' &= y_1y_2 - \beta y_3\end{aligned}$$

$$y_0 \in [-8.001, -7.998] \times [7.998, 8.001] \times [26.998, 27.001]$$



Back

Close

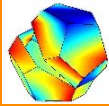
Summary

- Error standard models in numerical analysis
- Optimal bounds for IEEE-754 $+$, $-$, \times , $/$, $\sqrt{\cdot}$
- Weak sufficient assumptions for linearized bounds
- Optimal bounds for the error of summation
- error-free transformations
- Provably mathematical correct results
- Global optimization
- Affine arithmetic
- Taylor models

On verification methods:

S.M. Rump. Verification methods: Rigorous results using floating-point arithmetic.

Acta Numerica, 19:287–449, 2010.



62/63



Back

Close

An application of affine arithmetic — Julia sets

$$z_0, c \in \mathbb{C} : \quad z_{k+1} := z_k^2 + c \quad \text{for } k \geq 1$$

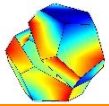
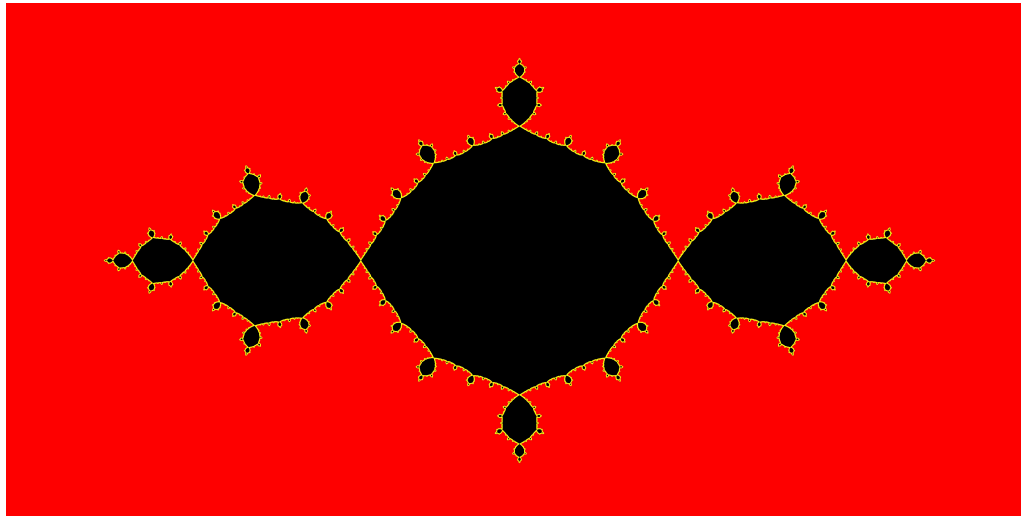
Given $c \in \mathbb{C}$, for which $z_0 \in \mathbb{C}$ is ∞ point of attraction?

Divergent for $\min\{|\operatorname{Re} c|, |\operatorname{Im} c|\} \geq 2$.

Color code red ∞ point of attraction

black iteration bounded for all $k \geq 1$

yellow don't know



63/63



Back

Close