

Analisi Numerica
Error analysis

Università di **Perugia**, **Italia**

Approximating real numbers

Approximating real numbers

Real numbers: limits of sequences in $\mathbb{Q} \rightarrow$ can be approximated by rational numbers.

Approximating real numbers

Real numbers: limits of sequences in $\mathbb{Q} \rightarrow$ can be approximated by rational numbers.

Trivia: diophantine numbers are $x \in \mathbb{R}$ for which there exists $c > 0$ such that for any rational number p/q ,

$$\left| x - \frac{p}{q} \right| \geq \frac{c}{q^2}.$$

Approximating real numbers

Real numbers: limits of sequences in $\mathbb{Q} \rightarrow$ can be approximated by rational numbers.

Trivia: diophantine numbers are $x \in \mathbb{R}$ for which there exists $c > 0$ such that for any rational number p/q ,

$$\left| x - \frac{p}{q} \right| \geq \frac{c}{q^2}.$$

Rational numbers are infinitely many \Rightarrow we can only deal with finite sets of numbers.

Problem

Select a finite number of representatives to approximate real numbers.

Naive attempt for numbers: **uniform distribution**. Choose $\varepsilon > 0$, N even and

$$\mathcal{E} = \left\{ \varepsilon k : -\frac{N}{2} < k \leq \frac{N}{2}, k \in \mathbb{Z} \right\}.$$

Naive attempt for numbers: **uniform distribution**. Choose $\varepsilon > 0$, N even and

$$\mathcal{E} = \left\{ \varepsilon k : -\frac{N}{2} < k \leq \frac{N}{2}, k \in \mathbb{Z} \right\}.$$

We can represent real numbers in the set

$$S = \left\{ -\varepsilon \frac{N}{2} < x \leq \varepsilon \frac{N}{2} \right\} \subset \mathbb{R}.$$

Naive attempt for numbers: **uniform distribution**. Choose $\varepsilon > 0$, N even and

$$\mathcal{E} = \left\{ \varepsilon k : -\frac{N}{2} < k \leq \frac{N}{2}, k \in \mathbb{Z} \right\}.$$

We can represent real numbers in the set

$$S = \left\{ -\varepsilon \frac{N}{2} < x \leq \varepsilon \frac{N}{2} \right\} \subset \mathbb{R}.$$

Representation function (from real numbers to the chosen set)

$$e : S \rightarrow \mathcal{E}, \quad x \rightarrow k\varepsilon, \quad (k-1)\varepsilon < x \leq k\varepsilon.$$

Naive attempt for numbers: **uniform distribution**. Choose $\varepsilon > 0$, N even and

$$\mathcal{E} = \left\{ \varepsilon k : -\frac{N}{2} < k \leq \frac{N}{2}, k \in \mathbb{Z} \right\}.$$

We can represent real numbers in the set

$$S = \left\{ -\varepsilon \frac{N}{2} < x \leq \varepsilon \frac{N}{2} \right\} \subset \mathbb{R}.$$

Representation function (from real numbers to the chosen set)

$$e : S \rightarrow \mathcal{E}, \quad x \rightarrow k\varepsilon, \quad (k-1)\varepsilon < x \leq k\varepsilon.$$

Representation absolute error

$$e(x) - x, \quad \max_{x \in S} |e(x) - x| < \varepsilon$$

$$\mathcal{E} = \left\{ \varepsilon k : -\frac{N}{2} < k \leq \frac{N}{2}, k \in \mathbb{Z} \right\}, \quad S = \left\{ -\varepsilon \frac{N}{2} < x \leq \varepsilon \frac{N}{2} \right\}$$

Representation function (from real numbers to the chosen set)

$$e : S \rightarrow \mathcal{E}, \quad x \rightarrow k\varepsilon, \quad (k-1)\varepsilon < x \leq k\varepsilon.$$

Representation error

$$e(x) - x, \quad \max_{x \in S} |e(x) - x| \leq \varepsilon$$

This choice makes the absolute error uniformly small.

Absolute vs. Relative Error

- Absolute error $\tilde{x} - x$.
- Relative error $\frac{\tilde{x} - x}{x}$, for $x \neq 0$.

Absolute vs. Relative Error

- Absolute error $\tilde{x} - x$.
- Relative error $\frac{\tilde{x} - x}{x}$, for $x \neq 0$.

Troubles with absolute error:

- large numbers \rightarrow correct digits
- small numbers \rightarrow wrong digits

Absolute vs. Relative Error

- Absolute error $\tilde{x} - x$.
- Relative error $\frac{\tilde{x} - x}{x}$, for $x \neq 0$.

Troubles with absolute error:

- large numbers \rightarrow correct digits
- small numbers \rightarrow wrong digits

x	\tilde{x}	absolute	relative
1.234	1.235	$1 \cdot 10^{-3}$	$8.1 \cdot 10^{-4}$
1234	1235	1	$8.1 \cdot 10^{-4}$
0.001234	0.001235	$1 \cdot 10^{-6}$	$8.1 \cdot 10^{-4}$

We wish to have a small relative error.

New problem

Select a finite number of representatives to approximate real numbers with a uniformly bounded relative error.

Idea: consider the digit representation of a real number.

Given a numeration basis $\beta \geq 2$, there is a correspondence between real numbers in $[0, 1)$ and sequences of digits

$$\mathbb{R} \iff \{d_i\}_{i=1,2,\dots}, \quad d_i \in \{0, \dots, \beta - 1\}.$$

Given a numeration basis $\beta \geq 2$, there is a correspondence between real numbers in $[0, 1)$ and sequences of digits

$$\mathbb{R} \iff \{d_i\}_{i=1,2,\dots}, \quad d_i \in \{0, \dots, \beta - 1\}.$$

Not a bijection: two sequences may represent the same real number

Given a numeration basis $\beta \geq 2$, there is a correspondence between real numbers in $[0, 1)$ and sequences of digits

$$\mathbb{R} \iff \{d_i\}_{i=1,2,\dots}, \quad d_i \in \{0, \dots, \beta - 1\}.$$

Not a bijection: two sequences may represent the same real number ($0.0\bar{9} = 0.1\bar{0}$).

Given a numeration basis $\beta \geq 2$, there is a correspondence between real numbers in $[0, 1)$ and sequences of digits

$$\mathbb{R} \iff \{d_i\}_{i=1,2,\dots}, d_i \in \{0, \dots, \beta - 1\}.$$

Not a bijection: two sequences may represent the same real number ($0.0\bar{9} = 0.1\bar{0}$).

- fixed point (natural): $x \rightarrow (\{d_i\}_{i=1,2,\dots}, M)$ (d_i is the i th digit and M says where the point is).
- floating point (scientific): $x \rightarrow (0.d_1d_2d_3 \cdots)\beta^p$, with $d_1 \neq 0$.

$$12.38 = 0.1238 \cdot 10^2.$$

Theorem (Basis representation theorem)

Let $x \in \mathbb{R} \setminus \{0\}$ and let $\beta \geq 2$ be a numeration basis, there exist unique $p \in \mathbb{Z}$ and a sequence $\{d_i\}_{i=1,2,\dots}$ such that

Theorem (Basis representation theorem)

Let $x \in \mathbb{R} \setminus \{0\}$ and let $\beta \geq 2$ be a numeration basis, there exist unique $p \in \mathbb{Z}$ and a sequence $\{d_i\}_{i=1,2,\dots}$ such that

- (i) $d_i \in \{0, 1, \dots, \beta - 1\};$

Theorem (Basis representation theorem)

Let $x \in \mathbb{R} \setminus \{0\}$ and let $\beta \geq 2$ be a numeration basis, there exist unique $p \in \mathbb{Z}$ and a sequence $\{d_i\}_{i=1,2,\dots}$ such that

- (i) $d_i \in \{0, 1, \dots, \beta - 1\}$;
- (ii) $d_1 \neq 0$; (To avoid ambiguity, e.g., $0.03 \cdot 10^2$ and $0.3 \cdot 10^1$)

Theorem (Basis representation theorem)

Let $x \in \mathbb{R} \setminus \{0\}$ and let $\beta \geq 2$ be a numeration basis, there exist unique $p \in \mathbb{Z}$ and a sequence $\{d_i\}_{i=1,2,\dots}$ such that

- (i) $d_i \in \{0, 1, \dots, \beta - 1\}$;*
- (ii) $d_1 \neq 0$;*
- (iii) d_i not definitely equal to $\beta - 1$; (There exists an infinite set of indices K such that $d_k \neq \beta - 1$ for every $k \in K$)*

Theorem (Basis representation theorem)

Let $x \in \mathbb{R} \setminus \{0\}$ and let $\beta \geq 2$ be a numeration basis, there exist unique $p \in \mathbb{Z}$ and a sequence $\{d_i\}_{i=1,2,\dots}$ such that

- (i) $d_i \in \{0, 1, \dots, \beta - 1\}$;
- (ii) $d_1 \neq 0$;
- (iii) d_i not definitely equal to $\beta - 1$;

such that

$$x = \text{sign}(x)\beta^p \sum_{i=1}^{\infty} \beta^{-i} d_i.$$

The number $\sum_{i=1}^{\infty} \beta^{-i} d_i$ is said to be **mantissa**.

Theorem (Basis representation theorem)

Let $x \in \mathbb{R} \setminus \{0\}$ and let $\beta \geq 2$ be a numeration basis, there exist unique $p \in \mathbb{Z}$ and a sequence $\{d_i\}_{i=1,2,\dots}$ such that

- (i) $d_i \in \{0, 1, \dots, \beta - 1\}$;
- (ii) $d_1 \neq 0$;
- (iii) d_i not definitely equal to $\beta - 1$;

such that

$$x = \text{sign}(x)\beta^p \sum_{i=1}^{\infty} \beta^{-i} d_i.$$

The number $\sum_{i=1}^{\infty} \beta^{-i} d_i$ is said to be **mantissa**.

Idea: consider numbers with finite mantissa

Given a numeration basis $\beta \geq 2$, the number $t > 0$ of digits of the mantissa, and m, M positive integers, we define a set of **floating point numbers**

$$\mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

These number have zero digits from the $(t + 1)$ -st on.

Given a numeration basis $\beta \geq 2$, the number $t > 0$ of digits of the mantissa, and m, M positive integers, we define a set of **floating point numbers**

$$\mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

These number have zero digits from the $(t + 1)$ -st on.

We will use this set to represent real numbers.

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$?

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$?

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$
- $\pi \in \mathcal{F}$?

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$
- $\pi \in \mathcal{F}$? no $\iff \pi = 0.31415\dots \cdot 10^1$ (more than two digits)

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$
- $\pi \in \mathcal{F}$? no $\iff \pi = 0.31415\dots \cdot 10^1$ (more than two digits)
- $12.345 \in \mathcal{F}$?

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$
- $\pi \in \mathcal{F}$? no $\iff \pi = 0.31415\dots \cdot 10^1$ (more than two digits)
- $12.345 \in \mathcal{F}$? no $\iff 12.345 = 0.12345 \cdot 10^2$ (more than two digits)

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$
- $\pi \in \mathcal{F}$? no $\iff \pi = 0.31415\dots \cdot 10^1$ (more than two digits)
- $12.345 \in \mathcal{F}$? no $\iff 12.345 = 0.12345 \cdot 10^2$ (more than two digits)
- $1000 \in \mathcal{F}$?

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$
- $\pi \in \mathcal{F}$? no $\iff \pi = 0.31415\dots \cdot 10^1$ (more than two digits)
- $12.345 \in \mathcal{F}$? no $\iff 12.345 = 0.12345 \cdot 10^2$ (more than two digits)
- $1000 \in \mathcal{F}$? no $\iff 1000 = 0.10 \cdot 10^4$ (too big)

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$
- $\pi \in \mathcal{F}$? no $\iff \pi = 0.31415\dots \cdot 10^1$ (more than two digits)
- $12.345 \in \mathcal{F}$? no $\iff 12.345 = 0.12345 \cdot 10^2$ (more than two digits)
- $1000 \in \mathcal{F}$? no $\iff 1000 = 0.10 \cdot 10^4$ (too big)
- $1/10000 \in \mathcal{F}$?

As an example

$$\mathcal{F} := \mathcal{F}(10, 2, 2, 3)$$

contains numbers with 2 digits

- $12 \in \mathcal{F}$? yes $\iff 12 = 0.12 \cdot 10^2$
- $0.03 \in \mathcal{F}$? yes $\iff 0.03 = 0.30 \cdot 10^{-1}$
- $\pi \in \mathcal{F}$? no $\iff \pi = 0.31415\dots \cdot 10^1$ (more than two digits)
- $12.345 \in \mathcal{F}$? no $\iff 12.345 = 0.12345 \cdot 10^2$ (more than two digits)
- $1000 \in \mathcal{F}$? no $\iff 1000 = 0.10 \cdot 10^4$ (too big)
- $1/10000 \in \mathcal{F}$? no $\iff 0.001 = 0.10 \cdot 10^{-3}$ (too small in modulus)

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

$\mathcal{F} \subset \mathbb{R}$ is a **finite set** with cardinality

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

$\mathcal{F} \subset \mathbb{R}$ is a **finite set** with cardinality

$$1 + 2(m + M + 1)(\beta - 1)\beta^{t-1}$$

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

$\mathcal{F} \subset \mathbb{R}$ is a **finite set** with cardinality

$$1 + 2(m + M + 1)(\beta - 1)\beta^{t-1}$$

The **largest** number in \mathcal{F} is

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

$\mathcal{F} \subset \mathbb{R}$ is a **finite set** with cardinality

$$1 + 2(m + M + 1)(\beta - 1)\beta^{t-1}$$

The **largest** number in \mathcal{F} is

$$\Omega = \beta^M \sum_{i=1}^t \beta^{-i} (\beta - 1)$$

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

$\mathcal{F} \subset \mathbb{R}$ is a **finite set** with cardinality

$$1 + 2(m + M + 1)(\beta - 1)\beta^{t-1}$$

The **largest** number in \mathcal{F} is

$$\Omega = \beta^M \sum_{i=1}^t \beta^{-i} (\beta - 1) = \beta^M (1 - \beta^{-t}).$$

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

$\mathcal{F} \subset \mathbb{R}$ is a **finite set** with cardinality

$$1 + 2(m + M + 1)(\beta - 1)\beta^{t-1}$$

The **largest** number in \mathcal{F} is

$$\Omega = \beta^M \sum_{i=1}^t \beta^{-i} (\beta - 1) = \beta^M (1 - \beta^{-t}).$$

The **smallest positive** number in \mathcal{F} is

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

$\mathcal{F} \subset \mathbb{R}$ is a **finite set** with cardinality

$$1 + 2(m + M + 1)(\beta - 1)\beta^{t-1}$$

The **largest** number in \mathcal{F} is

$$\Omega = \beta^M \sum_{i=1}^t \beta^{-i} (\beta - 1) = \beta^M (1 - \beta^{-t}).$$

The **smallest positive** number in \mathcal{F} is

$$\omega = \beta^{-m} \beta^{-1} = \beta^{-m-1}.$$

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

$\mathcal{F} \subset \mathbb{R}$ is a **finite set** with cardinality

$$1 + 2(m + M + 1)(\beta - 1)\beta^{t-1}$$

The **largest** number in \mathcal{F} is

$$\Omega = \beta^M \sum_{i=1}^t \beta^{-i} (\beta - 1) = \beta^M (1 - \beta^{-t}).$$

The **smallest positive** number in \mathcal{F} is

$$\omega = \beta^{-m} \beta^{-1} = \beta^{-m-1}.$$

For $x \geq \Omega$ we cannot represent numbers, for $0 < x < \omega$ we make a large relative error.

$\mathcal{F}(2, 2, 1, 1)$ is made of

$\mathcal{F}(2, 2, 1, 1)$ is made of 13 numbers

0

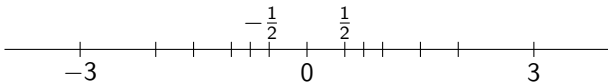
$$\pm 0.10 \cdot 2^{-1} \quad \pm 0.10 \cdot 2^0 \quad \pm 0.10 \cdot 2^1$$

$$\pm 0.11 \cdot 2^{-1} \quad \pm 0.11 \cdot 2^0 \quad \pm 0.11 \cdot 2^1$$

$\mathcal{F}(2, 2, 1, 1)$ is made of 13 numbers

$$\begin{array}{ccccc} 0 & & & & \\ \pm 0.10 \cdot 2^{-1} & \pm 0.10 \cdot 2^0 & \pm 0.10 \cdot 2^1 & & \\ \pm 0.11 \cdot 2^{-1} & \pm 0.11 \cdot 2^0 & \pm 0.11 \cdot 2^1 & & \end{array}$$

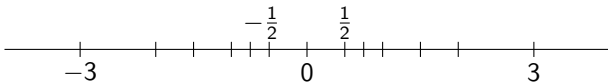
They are **not uniform**: between $1/2$ and $1/4$ and between $1/2$ and 1 there is the same number of elements of \mathcal{F} .



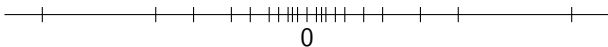
$\mathcal{F}(2, 2, 1, 1)$ is made of 13 numbers

$$\begin{array}{ccccc} 0 & & & & \\ \pm 0.10 \cdot 2^{-1} & \pm 0.10 \cdot 2^0 & \pm 0.10 \cdot 2^1 & & \\ \pm 0.11 \cdot 2^{-1} & \pm 0.11 \cdot 2^0 & \pm 0.11 \cdot 2^1 & & \end{array}$$

They are **not uniform**: between $1/2$ and $1/4$ and between $1/2$ and 1 there is the same number of elements of \mathcal{F} .



$\mathcal{F}(2, 2, 2, 2)$ has 21 numbers



$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

Given $S = \{x \in \mathbb{R} : \omega \leq x \leq \Omega\}$, we construct a **representation function**

$$\text{fl} : \mathbb{R} \longrightarrow \mathcal{F} \cup \{\pm\infty\}$$

with one of the two rules, where $x = \beta^p \sum_{i=1}^{\infty} \beta^{-i} d_i \in S$,

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

Given $S = \{x \in \mathbb{R} : \omega \leq x \leq \Omega\}$, we construct a **representation function**

$$\text{fl} : \mathbb{R} \longrightarrow \mathcal{F} \cup \{\pm\infty\}$$

with one of the two rules, where $x = \beta^p \sum_{i=1}^{\infty} \beta^{-i} d_i \in S$,

- **truncation**: $\tilde{x} = \text{fl}(x) = \beta^p \sum_{i=1}^t \beta^{-i} d_i$.

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

Given $S = \{x \in \mathbb{R} : \omega \leq x \leq \Omega\}$, we construct a **representation function**

$$\text{fl} : \mathbb{R} \longrightarrow \mathcal{F} \cup \{\pm\infty\}$$

with one of the two rules, where $x = \beta^p \sum_{i=1}^{\infty} \beta^{-i} d_i \in S$,

- **truncation**: $\tilde{x} = \text{fl}(x) = \beta^p \sum_{i=1}^t \beta^{-i} d_i$.
- **rounding**: $\tilde{x} = \text{fl}(x)$ the truncation of $x + \beta^{p-t-1}/2$.

$$\mathcal{F} := \mathcal{F}(\beta, t, m, M) = \{0\} \cup \left\{ \pm \beta^p \sum_{i=1}^t \beta^{-i} d_i, \right. \\ \left. -m \leq p \leq M, 0 \leq d_i < \beta \text{ integer for } i = 1, \dots, t, d_1 \neq 0 \right\}.$$

Given $S = \{x \in \mathbb{R} : \omega \leq x \leq \Omega\}$, we construct a **representation function**

$$\text{fl} : \mathbb{R} \longrightarrow \mathcal{F} \cup \{\pm\infty\}$$

with one of the two rules, where $x = \beta^p \sum_{i=1}^{\infty} \beta^{-i} d_i \in S$,

- **truncation**: $\tilde{x} = \text{fl}(x) = \beta^p \sum_{i=1}^t \beta^{-i} d_i$.
- **rounding**: $\tilde{x} = \text{fl}(x)$ the truncation of $x + \beta^{p-t-1}/2$.

If $x > \Omega$ we set $\text{fl}(x) = \infty$ (**overflow**), if $0 \leq x < \omega$ we set $\text{fl}(x) = 0$ (**underflow**) and for negative numbers the extension is straightforward.

We define the **machine precision** as $u = \beta^{-t+1}$ in the case of truncation and $u = \beta^{-t+1}/2$ in the case of rounding.

Theorem

Let $x \in S$ then we have the following bound for the relative error

$$\left| \frac{\text{fl}(x) - x}{x} \right| < u.$$

Proof. For the truncation

We define the **machine precision** as $u = \beta^{-t+1}$ in the case of truncation and $u = \beta^{-t+1}/2$ in the case of rounding.

Theorem

Let $x \in S$ then we have the following bound for the relative error

$$\left| \frac{\text{fl}(x) - x}{x} \right| < u.$$

Proof. For the truncation

$$\frac{|\text{fl}(x) - x|}{|x|} = \frac{x - \text{fl}(x)}{x}$$

Since $x \geq \text{fl}(x)$.

We define the **machine precision** as $u = \beta^{-t+1}$ in the case of truncation and $u = \beta^{-t+1}/2$ in the case of rounding.

Theorem

Let $x \in S$ then we have the following bound for the relative error

$$\left| \frac{\text{fl}(x) - x}{x} \right| < u.$$

Proof. For the truncation

$$\frac{|\text{fl}(x) - x|}{|x|} = \frac{x - \text{fl}(x)}{x} = \frac{\beta^p \sum_{i>t} \beta^{-i} d_i}{\beta^p \sum_{i=1}^{\infty} \beta^{-i} d_i}$$

By definition: $\sum_{i=1}^{\infty} \beta^{-i} d_i - \sum_{i=1}^t \beta^{-i} d_i = \sum_{i=t+1}^{\infty} \beta^{-i} d_i.$

We define the **machine precision** as $u = \beta^{-t+1}$ in the case of truncation and $u = \beta^{-t+1}/2$ in the case of rounding.

Theorem

Let $x \in S$ then we have the following bound for the relative error

$$\left| \frac{\text{fl}(x) - x}{x} \right| < u.$$

Proof. For the truncation

$$\frac{|\text{fl}(x) - x|}{|x|} = \frac{x - \text{fl}(x)}{x} = \frac{\beta^p \sum_{i>t} \beta^{-i} d_i}{\beta^p \sum_{i=1}^{\infty} \beta^{-i} d_i}$$

We define the **machine precision** as $u = \beta^{-t+1}$ in the case of truncation and $u = \beta^{-t+1}/2$ in the case of rounding.

Theorem

Let $x \in S$ then we have the following bound for the relative error

$$\left| \frac{\text{fl}(x) - x}{x} \right| < u.$$

Proof. For the truncation

$$\frac{|\text{fl}(x) - x|}{|x|} = \frac{x - \text{fl}(x)}{x} = \frac{\sum_{i>t} \beta^{-i} d_i}{\sum_{i=1}^{\infty} \beta^{-i} \mathbf{d}_i} \leq \frac{\beta^{-t}}{\beta^{-1}} = u.$$

$$\sum_{i=1}^{\infty} \beta^{-i} d_i \geq \beta^{-1}.$$

We define the **machine precision** as $u = \beta^{-t+1}$ in the case of truncation and $u = \beta^{-t+1}/2$ in the case of rounding.

Theorem

Let $x \in S$ then we have the following bound for the relative error

$$\left| \frac{\text{fl}(x) - x}{x} \right| < u.$$

Proof. For the truncation

$$\frac{|\text{fl}(x) - x|}{|x|} = \frac{x - \text{fl}(x)}{x} = \frac{\sum_{i>t} \beta^{-i} d_i}{\sum_{i=1}^{\infty} \beta^{-i} d_i} \leq \frac{\beta^{-t}}{\beta^{-1}} = u.$$

$$\begin{aligned} \sum_{i=t+1}^{\infty} \beta^{-i} d_i &\leq (\beta - 1) \sum_{i=t+1}^{\infty} \beta^{-i} = (\beta - 1) \left(\sum_{i=0}^{\infty} \beta^{-i} - \sum_{i=0}^t \beta^{-i} \right) \\ &= (\beta - 1) \left(\frac{1}{1 - 1/\beta} - \frac{1 - \beta^{-t-1}}{1 - 1/\beta} \right) = \frac{\beta - 1}{\frac{\beta - 1}{\beta}} \beta^{-t-1} = \beta^{-t}. \end{aligned}$$

We define the **machine precision** as $u = \beta^{-t+1}$ in the case of truncation and $u = \beta^{-t+1}/2$ in the case of rounding.

Theorem

Let $x \in S$ then we have the following bound for the relative error

$$\left| \frac{\text{fl}(x) - x}{x} \right| < u.$$

Proof. For the truncation

$$\frac{|\text{fl}(x) - x|}{|x|} = \frac{x - \text{fl}(x)}{x} = \frac{\sum_{i>t} \beta^{-i} d_i}{\sum_{i=1}^{\infty} \beta^{-i} d_i} \leq \frac{\beta^{-t}}{\beta^{-1}} = u.$$

Floating point on a computer

For the single precision floating point (float 32 bits)

$$\mathcal{F}(2, 24, 126, 127), \quad u = 2^{-24}.$$

Floating point on a computer

For the single precision floating point (float 32 bits)

$$\mathcal{F}(2, 24, 126, 127), \quad u = 2^{-24}.$$

For the double precision floating point (double 64 bits)

$$\mathcal{F}(2, 52, 1022, 1023), \quad u = 2^{-52} \approx 2.2 \cdot 10^{-16}.$$

Bad news: \mathcal{F} has **few algebraic properties**, e.g.,

$$x, y \in \mathcal{F} \not\Rightarrow x + y \in \mathcal{F}$$

We must define floating point operations.

Bad news: \mathcal{F} has **few algebraic properties**, e.g.,

$$x, y \in \mathcal{F} \not\Rightarrow x + y \in \mathcal{F}$$

We must define floating point operations.

We may assume that there exists a floating point sum \oplus such that, if $x, y \in \mathcal{F}$, then (if overflow does not occur)

$$x \oplus y \in \mathcal{F}, \quad x \oplus y = (x + y)(1 + \varepsilon), \quad |\varepsilon| < u$$

Similarly we define \otimes , \ominus , \oslash .

Bad news: \mathcal{F} has **few algebraic properties**, e.g.,

$$x, y \in \mathcal{F} \not\Rightarrow x + y \in \mathcal{F}$$

We must define floating point operations.

We may assume that there exists a floating point sum \oplus such that, if $x, y \in \mathcal{F}$, then (if overflow does not occur)

$$x \oplus y \in \mathcal{F}, \quad x \oplus y = (x + y)(1 + \varepsilon), \quad |\varepsilon| < u$$

Similarly we define \otimes , \ominus , \oslash .

A simple idea is to define, for instance, $x \otimes y = \text{fl}(x + y)$, but details are more complicate.

Floating point operations verify just some properties

- $x \oplus y = y \oplus x$ (commutativity of the sum);
- $x \otimes y = y \otimes x$ (commutativity of the product);
- $x \oslash x = 1$.

Floating point operations verify just some properties

- $x \oplus y = y \oplus x$ (commutativity of the sum);
- $x \otimes y = y \otimes x$ (commutativity of the product);
- $x \oslash x = 1$.

They do not verify others

- associativity of the sum and product;
- distributive law;
- simplification (it may happen $x \otimes (y \oslash x) \neq y$);
- no null factor law (it may happen $x \otimes y = z \otimes y$ with $y \neq 0$ and $x \neq z$)

Well-posed problems

Which problems are we interested in?

Well-posed problems

Which problems are we interested in?

But first...

Well-posed problems

What is a problem?

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 1: the **linear system** problem.

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 1: the **linear system** problem.

Given a matrix A (coefficient) and a vector b (right hand side),

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 1: the **linear system** problem.

Given a matrix A (coefficient) and a vector b (right hand side), find all vectors x (unknown)

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 1: the **linear system** problem.

Given a matrix A (coefficient) and a vector b (right hand side), find all vectors x (unknown) such that $Ax = b$.

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 2: the **polynomial equation** problem.

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 2: the **polynomial equation** problem.

Given a polynomial $p(x)$,

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 2: the **polynomial equation** problem.

Given a polynomial $p(x)$, find all complex numbers x (unknowns)

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 2: the **polynomial equation** problem.

Given a polynomial $p(x)$, find all complex numbers x (unknowns) such that $p(x) = 0$.

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 3: the **integration**.

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 3: the **integration**.

Given a continuous function $f(x)$ (integrand) and two real numbers $a \leq b$ (extremes),

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 3: the **integration**.

Given a continuous function $f(x)$ (integrand) and two real numbers $a \leq b$ (extremes), find a real number I (unknown)

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

Example 3: the **integration**.

Given a continuous function $f(x)$ (integrand) and two real numbers $a \leq b$ (extremes), find a real number I (unknown) such that $I = \int_a^b f(x) dx$.

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

In all cases the unknowns are functions of the data.

Well-posed problems

What is a problem?

A problem is made of three parts

- Data
- Unknowns
- Conditions

In all cases the unknowns are functions of the data.

Solutions of problems \Longleftrightarrow Evaluation of a function

In practical cases, only an **approximation** can be provided.

Well-posed problems

Which problems are we interested in?

Well-posed problems

-
-
-

Well-posed problems

Which problems are we interested in?

Well-posed problems

- They have a solution
-
-

Well-posed problems

Which problems are we interested in?

Well-posed problems

- They have a solution (or else what are we computing?)
-
-

Well-posed problems

Which problems are we interested in?

Well-posed problems

- They have a solution
- The solution is unique
-

Well-posed problems

Which problems are we interested in?

Well-posed problems

- They have a solution
- The solution is unique (or else what is the correct answer?)
-

Well-posed problems

Which problems are we interested in?

Well-posed problems

- They have a solution
- The solution is unique
- The solution depends continuously from the data

Well-posed problems

Which problems are we interested in?

Well-posed problems

- They have a solution
- The solution is unique
- The solution depends continuously from the data

Continuous dependence from data is (less obviously) important

- Data in real problems are affected by errors
- Computation is made on finite arithmetic and there are rounding errors

Continuous dependence from data

Consider a problem with data and unknowns in the Banach spaces V and W , respectively.

The problem **strongly** continuously depends from data at $A \in V$ if there exists a neighborhood $\mathcal{U} \subset V$ of A such that the problem has a unique solution $X(B)$ for $B \in \mathcal{U}$ and $\lim_{B \rightarrow A} X(B) = X(A)$.

The problem **weakly** continuously depends from data at $A \in V$ if for $B \in V$, there exists $t_0 > 0$ such that the problem with data $A + tB$ has a unique solution $X(t)$ for $t \in [0, t_0)$ and $\lim_{t \rightarrow 0^+} X(t) = X(0)$.

Strong dependence implies weak dependence, but not the contrary.

Weak dependence is more frequent and sometimes it is enough in applications.

Implicit functions theorem

A problem is often stated as an implicit equation (e. g. systems of equations, zeros of polynomials, ODEs, PDEs).

Theorem (Implicit functions)

Let $F : \Omega \rightarrow \mathbb{R}^n$, with $F \in C^1(\Omega)$, $\Omega \in \mathbb{R}^n \times \mathbb{R}^m$ such that $F_x(x_0, y_0) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is invertible for $(x_0, y_0) \in \Omega^a$.

There exist neighborhoods $\mathcal{U} \ni (x_0, y_0)$ and $\mathcal{V} \ni y_0$ and $g : \mathcal{V} \rightarrow \mathbb{R}^n$, with $g \in C^1(\mathcal{V})$ such that for $(x, y) \in \mathcal{U}$

$$F(x, y) = F(x_0, y_0) \iff x = g(y).$$

^awe use the notation $F(x, y)$ with $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$

An analogous results holds for $F : \Omega \rightarrow \mathbb{C}^n$ and F analytic in $\Omega \in \mathbb{C}^n \times \mathbb{C}^m$. In this case g is analytic as well.

Example: square linear system

Theorem

*Let $A \in \mathbb{C}^{n \times n}$. The linear system $Ax = b$ is **strongly** well posed for $b \in \mathbb{C}^n$ if and only if A is invertible.*

Proof.



Example: square linear system

Theorem

Let $A \in \mathbb{C}^{n \times n}$. The linear system $Ax = b$ is **strongly** well posed for $b \in \mathbb{C}^n$ if and only if A is invertible.

Proof.

Since $(A, b) \in \mathbb{C}^{n \times n} \times \mathbb{C}^n \cong \mathbb{C}^{n^2+n}$, we can see data as belonging to a vector space.

Let $F : \mathbb{C}^{n^2+n} \times \mathbb{C}^n \rightarrow \mathbb{C}^n$ be such that $F((A, b), x) = Ax - b$.

Writing

$$F((A, b), x) = \begin{bmatrix} F_1 \\ F_2 \\ \vdots \\ F_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + \cdots + a_{1n}x_n - b_1 \\ a_{21}x_1 + \cdots + a_{2n}x_n - b_2 \\ \vdots \\ a_{n1}x_1 + \cdots + a_{nn}x_n - b_n \end{bmatrix},$$

we see that F is differentiable, since its components are polynomials.

Example: square linear system

Theorem

Let $A \in \mathbb{C}^{n \times n}$. The linear system $Ax = b$ is **strongly** well posed for $b \in \mathbb{C}^n$ if and only if A is invertible.

Proof.

The derivative with respect to x_i is $\frac{\partial F_i}{\partial F_j} = a_{ij}$ and thus the Jacobian matrix is A and we can say that $F_x((A, b), x)[h] = Ah$, with $h \in \mathbb{C}^n$.

The function F is differentiable with F_x invertible, the implicit functions theorem implies that there exists a neighborhood of $((A, b), x)$ such that $Ax = b$, where

$$\tilde{A}\tilde{x} - \tilde{b} = F((\tilde{A}, \tilde{b}), \tilde{x}) = F((A, b), x) = Ax - b = 0$$

if and only if $\tilde{x} = g(\tilde{A}, \tilde{b})$ with g differentiable. (Note that $g(\tilde{A}, \tilde{b}) = \tilde{A}^{-1}\tilde{b}$.)



Example: square linear system

Theorem

*Let $A \in \mathbb{C}^{n \times n}$. The linear system $Ax = b$ is **strongly** well posed for $b \in \mathbb{C}^n$ if and only A is invertible.*

Proof.

The converse is left as an exercise.



The theorem can be proved using only linear algebra.

Example: square linear system

The linear system $Ax = b$ with A square and invertible.

Example: square linear system

The linear system $Ax = b$ with A square and invertible.

Is the problem well posed? The solution exists and is unique (Cramer's theorem).

Example: square linear system

The linear system $Ax = b$ with A square and invertible.

Is the problem well posed? The solution exists and is unique (Cramer's theorem).

The solution is $x = A^{-1}b$. Let $(A^{-1})_{ij} = \tilde{a}_{ij}$, for $i, j = 1, \dots, n$, be the elements of the inverse

$$x_i = \sum_{j=1}^n \tilde{a}_{ij} b_j$$

Example: square linear system

The linear system $Ax = b$ with A square and invertible.

Is the problem well posed? The solution exists and is unique (Cramer's theorem).

The solution is $x = A^{-1}b$. Let $(A^{-1})_{ij} = \tilde{a}_{ij}$, for $i, j = 1, \dots, n$, be the elements of the inverse

$$x_i = \sum_{j=1}^n \tilde{a}_{ij} b_j$$

with (adjoint's formula, **inefficient!**)

$$\tilde{a}_{ij} = \frac{1}{\det(A)} (-1)^{i+j} \det(M^{(ij)}),$$

where $M^{(ij)}$ is obtained removing the i -th row and the j -th column from A^T .

Example: square linear system

The linear system $Ax = b$ with A square and invertible.

The solution is $x = A^{-1}b$.

$$x_i = \sum_{k=1}^n \tilde{a}_{ik} b_k, \quad \tilde{a}_{ij} = \frac{1}{\det(A)} (-1)^{i+j} \det(M^{(ij)}),$$

Example: square linear system

The linear system $Ax = b$ with A square and invertible.

The solution is $x = A^{-1}b$.

$$x_i = \sum_{k=1}^n \tilde{a}_{ik} b_k, \quad \tilde{a}_{ij} = \frac{1}{\det(A)} (-1)^{i+j} \det(M^{(ij)}),$$

Note that

- the first formula is a polynomial;
- $\det(X)$ is a polynomial of the entries of X ;
- the second formula is a rational function.

Example: square linear system

The linear system $Ax = b$ with A square and invertible.

The solution is $x = A^{-1}b$.

$$x_i = \sum_{j=1}^n \tilde{a}_{ij} b_j, \quad \tilde{a}_{ij} = \frac{1}{\det(A)} (-1)^{i+j} \det(M^{(ij)}),$$

Note that

- the first formula is a polynomial;
- $\det(X)$ is a polynomial of the entries of X ;
- the second formula is a rational function.

Thus, x_i is a rational function (continuous) of the data. We can write

$$x = f(a_{11}, \dots, a_{nn}, b_1, \dots, b_n), \quad f : \mathbb{K}^{n^2+n} \rightarrow \mathbb{K}^n.$$

We have translated a problem to a function.

Example: roots of polynomials

The roots of the polynomial

$$az^2 + bz + c = 0,$$

can be written as

$$z_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}, \quad z_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a},$$

where the complex square root is the one with positive real part (or positive imaginary part if the real part is 0).

They are analytic functions of the coefficients if $b^2 - 4ac \neq 0$ in a neighborhood of the polynomial, while for $b^2 - 4ac = 0$ they lose analyticity (and continuity in some sense).

Example: roots of polynomials

Theorem

*Let $p(z) \in \mathbb{C}[z]$. The polynomial equations $p(z) = 0$ is a **strongly** well posed problem if p has distinct roots.*

Proof.



Example: roots of polynomials

Theorem

Let $p(z) \in \mathbb{C}[z]$. The polynomial equations $p(z) = 0$ is a **strongly** well posed problem if p has distinct roots.

Proof.

We prove the theorem by induction on the degree n of p .

If $n = 1$ then $p(z) = a_1z + a_0$ has one distinct root α . If $\tilde{p}(z) = \tilde{a}_1z + \tilde{a}_0$ lies in a neighborhood of $p(z)$ with $\tilde{a}_1 \neq 0$ then $\tilde{p}(z)$ has a unique root $\tilde{\alpha} = -\tilde{a}_0/\tilde{a}_1$. Note that $\tilde{\alpha}$ is a continuous function of the coefficients of \tilde{p} and $\lim_{\tilde{p} \rightarrow p} \tilde{\alpha} = -a_0/a_1 = \alpha$.^a



^aWe have that $\tilde{p} \rightarrow p$, if and only if $\tilde{a}_0 \rightarrow a_0, \dots, \tilde{a}_n \rightarrow a_n$.

Example: roots of polynomials

Theorem

Let $p(z) \in \mathbb{C}[z]$. The polynomial equations $p(z) = 0$ is a **strongly** well posed problem if p has distinct roots.

Proof.

If $p(z) = a_0 + a_1z + \cdots + a_nz^n$, and α is a root of p , then

$$p(z) = (z - \alpha) \underbrace{(b_0 + b_1z + \cdots + b_{n-1}z^{n-1})}_{q(z)},$$

for a unique $q(z)$, from which we get the system of equations

$$\begin{cases} a_0 + \alpha b_0 = 0, \\ a_1 + \alpha b_1 - b_0 = 0, \\ \vdots \\ a_{n-1} + \alpha b_{n-1} - b_{n-2} = 0, \\ a_n - b_{n-1} = 0. \end{cases}$$

Example: roots of polynomials

Theorem

Let $p(z) \in \mathbb{C}[z]$. The polynomial equations $p(z) = 0$ is a **strongly** well posed problem if p has distinct roots.

Proof.

$$\begin{cases} a_0 + \alpha b_0 = 0, \\ a_1 + \alpha b_1 - b_0 = 0, \\ \vdots \\ a_{n-1} + \alpha b_{n-1} - b_{n-2} = 0, \\ a_n - b_{n-1} = 0. \end{cases}$$

can be written as $F(a_0, \dots, a_n, \alpha, b_0, \dots, b_{n-1}) = 0$. We have that $F : \mathbb{C}^{n+1} \times \mathbb{C}^{n+1} \rightarrow \mathbb{C}^{n+1}$ is differentiable (it is a polynomial) and

$$\frac{\partial F}{\partial \alpha} = \begin{bmatrix} b_0 \\ \vdots \\ b_{n-1} \\ 0 \end{bmatrix} \in \mathbb{C}^{n+1}, \quad \frac{\partial F}{\partial b} = \begin{bmatrix} \alpha & 0 & \cdots & 0 \\ -1 & \alpha & \ddots & \vdots \\ 0 & -1 & \ddots & 0 \\ \vdots & \ddots & \ddots & \alpha \\ 0 & \cdots & 0 & -1 \end{bmatrix} \in \mathbb{C}^{(n+1) \times n}.$$

Example: roots of polynomials

Theorem

Let $p(z) \in \mathbb{C}[z]$. The polynomial equations $p(z) = 0$ is a **strongly** well posed problem if p has distinct roots.

Proof.

The partial Fréchet derivative is associated with the matrix

$$F_{\alpha, b_0, \dots, b_{n-1}} = \begin{bmatrix} \frac{\partial F}{\partial \alpha} & \frac{\partial F}{\partial b_{n-1}} \end{bmatrix} \in \mathbb{C}^{(n+1) \times (n+1)}.$$

We claim that $F_{\alpha, b_0, \dots, b_{n-1}}$ is invertible if and only if $q(\alpha) \neq 0$, that is α is not a root of q .

Since the last row of $F_{\alpha, b_0, \dots, b_{n-1}}$ contains just a -1 in the last position, it is sufficient to consider

$$M := \begin{bmatrix} b_0 & \alpha & 0 & \cdots & 0 \\ b_1 & -1 & \alpha & \ddots & \vdots \\ \vdots & 0 & \ddots & \ddots & 0 \\ b_{n-1} & \vdots & \ddots & \ddots & \alpha \\ 0 & 0 & \cdots & 0 & -1 \end{bmatrix}$$

Example: roots of polynomials

Theorem

Let $p(z) \in \mathbb{C}[z]$. The polynomial equations $p(z) = 0$ is a **strongly** well posed problem if p has distinct roots.

Proof.

$$M := \begin{bmatrix} b_0 & \alpha & 0 & \cdots & 0 \\ b_1 & -1 & \alpha & \ddots & \vdots \\ \vdots & 0 & \ddots & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & \alpha \\ b_{n-1} & 0 & \cdots & 0 & -1 \end{bmatrix}$$

We prove that M is singular if and only if $q(\alpha) = 0$. If $q(\alpha) = 0$, then

$$\begin{aligned} & \begin{bmatrix} 1 & \alpha & \alpha^2 & \cdots & \alpha^{n-1} \end{bmatrix} M \\ &= \begin{bmatrix} b_0 + b_1\alpha + b_2\alpha^2 + \cdots + b_{n-1}\alpha^{n-1} & \alpha - \alpha & \cdots & \alpha^{n-1} - \alpha^{n-1} \end{bmatrix} \\ &= \begin{bmatrix} q(\alpha) & 0 & \cdots & 0 \end{bmatrix} = 0. \end{aligned}$$

We have proved that M^T has a nonzero kernel, this implies that M^T is singular and thus M is singular. □

Example: roots of polynomials

Theorem

Let $p(z) \in \mathbb{C}[z]$. The polynomial equations $p(z) = 0$ is a **strongly** well posed problem if p has distinct roots.

Proof.

On the contrary, if M is singular, then there exists a nonzero vector $v = [y_0 \ y_1 \ \cdots \ y_{n-1}]$ such that $vM = 0$.

The equation $vM = 0$ is equivalent to

$$\begin{cases} b_0 y_0 + b_1 y_1 + \cdots + b_{n-1} y_{n-1} = 0 \\ \alpha y_0 - y_1 = 0 \\ \alpha y_1 - y_2 = 0 \\ \vdots \\ \alpha y_{n-2} - y_{n-1} = 0, \end{cases}$$

that gives $y_1 = \alpha y_0, \dots, y_{n-1} = \alpha^{n-1} y_0$ and

$$(b_0 + b_1 \alpha + \cdots + b_{n-1} \alpha^{n-1}) y_0 = q(\alpha) y_0 = 0$$

Since $y_0 \neq 0$ (or else $v = 0$) we obtain $q(\alpha) = 0$.



Example: roots of polynomials

Theorem

Let $p(z) \in \mathbb{C}[z]$. The polynomial equations $p(z) = 0$ is a **strongly well posed problem** if p has distinct roots.

Proof.

The hypotheses of the implicit functions theorem are fulfilled if $q(\alpha) \neq 0$.

For any $\tilde{p}(z)$ in a neighborhood of $p(z)$ there exists $\tilde{\alpha}_1$ and $\tilde{q}(z)$ such that $\tilde{p}(z) = (z - \tilde{\alpha})\tilde{q}(z)$ and $\tilde{\alpha}_1$ and $\tilde{b}_0, \dots, \tilde{b}_{n-1}$ are analytic functions of the coefficients of $\tilde{p}(z)$.

The polynomial $q(z)$ has distinct roots. By induction, there exists $\tilde{\alpha}_2, \dots, \tilde{\alpha}_n$ that are solutions of $\tilde{q}(z)$ (in a smaller neighborhood if necessary) and that are analytic with respect to the coefficients of $\tilde{p}(z)$. Since the composition of analytic functions is analytic the roots are analytic functions of the coefficients of $\tilde{p}(z)$. □

Example: roots of polynomials

The roots of the polynomial

$$x^2 - \alpha = 0, \quad \alpha = \rho e^{i\theta} \in \mathbb{C},$$

are $\varphi_1(\alpha) = \sqrt{\rho}e^{i\theta/2}$ and $\varphi_2(\alpha) = -\sqrt{\rho}e^{i\theta/2}$.

Theorem

There exists no square root function continuous in a neighborhood of $0 \in \mathbb{C}$.

Proof.



Example: roots of polynomials

The roots of the polynomial

$$x^2 - \alpha = 0, \quad \alpha = \rho e^{i\theta} \in \mathbb{C},$$

are $\varphi_1(\alpha) = \sqrt{\rho}e^{i\theta/2}$ and $\varphi_2(\alpha) = -\sqrt{\rho}e^{i\theta/2}$.

Theorem

There exists no square root function continuous in a neighborhood of $0 \in \mathbb{C}$.

Proof.

If there exist a continuous function $\lambda : \mathcal{U} \rightarrow \mathbb{C}$ such that $\lambda(z)^2 = z$, where \mathcal{U} is a neighborhood of 0.

We have $\lambda(z)^2 = z$ on a circle $S^1 = \{z \in \mathbb{C} : |z| = \rho\}$, with $\rho > 0$.

Define the two sets

$$\mathcal{U}_1 := \{z \in S^1 : \lambda(z) = \varphi_1(z)\}, \quad \mathcal{U}_2 := \{z \in S^1 : \lambda(z) = \varphi_2(z)\}.$$

We claim that \mathcal{U}_1 is open in S^1 .



Example: roots of polynomials

Theorem

There exists no square root function continuous in a neighborhood of $0 \in \mathbb{C}$.

Proof.

For $z \in \mathcal{U}_1$ there exists a neighborhood \mathcal{V} of z such that $|\lambda(w) - \lambda(z)| < \sqrt{\rho}$ because λ is continuous, but $|\varphi_1 - \varphi_2| = 2\sqrt{\rho}$ and thus $\lambda(z) \equiv \varphi_1(z)$ on \mathcal{V} .

Analogously \mathcal{U}_2 is open.

The two open subset of S^1 are disjoint and their union is S^1 , but since S^1 is a connected set, we have that $S^1 = \mathcal{U}_1$ or $S^1 = \mathcal{U}_2$.

In both cases $\lim_{\theta \rightarrow 0^+} \lambda(z) \neq \lim_{\theta \rightarrow 2\pi^-} \lambda(z)$, while they should coincide if λ is continuous. (For instance if $S^1 = \mathcal{U}_1$ we get $\sqrt{\rho}$ and $-\sqrt{\rho}$ for the two limits.) □

Example: roots of polynomials

Find all solutions of the polynomial equations

$p(x) = a_0 + a_1x + \cdots + a_nx^n = 0$, where $p \in \mathbb{C}_n[x]$ with $a_n \neq 0$.

Is the problem well posed? The equation has n solutions ξ_1, \dots, ξ_n counted with multiplicity (fundamental theorem of algebra).

Example: roots of polynomials

Find all solutions of the polynomial equations

$p(x) = a_0 + a_1x + \cdots + a_nx^n = 0$, where $p \in \mathbb{C}_n[x]$ with $a_n \neq 0$.

Is the problem well posed? The equation has n solutions ξ_1, \dots, ξ_n counted with multiplicity (fundamental theorem of algebra).

If the roots are distinct, let $p_t(x) = p(x) + tq(x)$, for $t \in [-a, a]$ and $q(x)$ polynomial of the same degree as $p(x)$, with $a > 0$, be a polynomial. We have $p_t \rightarrow_{t \rightarrow 0} p(x)$, then

Example: roots of polynomials

Find all solutions of the polynomial equations

$p(x) = a_0 + a_1x + \cdots + a_nx^n = 0$, where $p \in \mathbb{C}_n[x]$ with $a_n \neq 0$.

Is the problem well posed? The equation has n solutions ξ_1, \dots, ξ_n counted with multiplicity (fundamental theorem of algebra).

If the roots are distinct, let $p(x) + tq(x)$, for $t \in [-a, a]$ and $q(x)$ polynomial of the same degree as $p(x)$, with $a > 0$, be a polynomial. We have $p_t \rightarrow_{t \rightarrow 0} p(x)$, then

for a sufficiently small t the polynomial $p_t(x)$ has distinct roots $\zeta_1(t), \dots, \zeta_n(t)$ which are analytic (continuous) functions of t and $\zeta_i(0) = \xi_i$ (a theorem in complex analysis).

The problem is well-posed: we have a function $f : \Omega \rightarrow \mathbb{C}^n$ ($\Omega \subset \mathbb{C}^n$ open)

Example: integrals

Compute $\int_a^b f(x)dx$, with f continuous on $[a, b]$.

Riemann (Lebesgue) integrability guarantees the existence and uniqueness (fundamental theorem of calculus).

Example: integrals

Compute $\int_a^b f(x)dx$, with f continuous on $[a, b]$.

Riemann (Lebesgue) integrability guarantees the existence and uniqueness (fundamental theorem of calculus).

The problem can be stated as a function evaluation

$$f : C[a, b] \rightarrow \mathbb{R}$$

where $C[a, b]$ are continuous functions on $[a, b]$.

Notice that the space $C[a, b]$ as a vector space has infinite dimension!

Problems to be solved

What do we need to solve?

Problems to be solved

What do we need to solve?

A problem can be stated as a function evaluation.

In the general case we have a function

$$f : V \rightarrow W,$$

where V, W are **infinite-dimensional** vector spaces (Banach spaces).

Problems to be solved

What do we need to solve?

A problem can be stated as a function evaluation.

In the general case we have a function

$$f : V \rightarrow W,$$

where V, W are **infinite-dimensional** vector spaces (Banach spaces).

What can we solve?

Problems to be solved

What do we need to solve?

A problem can be stated as a function evaluation.

In the general case we have a function

$$f : V \rightarrow W,$$

where V, W are **infinite-dimensional** vector spaces (Banach spaces).

What can we solve?

By hand or with a computer, we can evaluate only **rational functions**

$$f : \mathbb{R}^n \rightarrow \mathbb{R}, \quad f = \frac{p(x_1, \dots, x_n)}{q(x_1, \dots, x_n)}$$

with p, q polynomials (in n variables).

Moreover we cannot use all real numbers but just a fistful of them.

Problems to be solved

What do we need to solve?

A problem can be stated as a function evaluation.

In the general case we have a function

$$f : V \rightarrow W,$$

where V, W are **infinite-dimensional** vector spaces (Banach spaces).

What can we solve?

By hand or with a computer, we can evaluate only **rational functions**

$$f : \mathbb{R}^n \rightarrow \mathbb{R}, \quad f = \frac{p(x_1, \dots, x_n)}{q(x_1, \dots, x_n)}$$

with p, q polynomials (in n variables).

Moreover we cannot use all real numbers but just a fistful of them.

This is awkward!

Approximation

We accept an approximated solution.

It is perhaps surprising that most of continuous problems can be well approximated using only elementary operations.

Approximation

We accept an approximated solution.

It is perhaps surprising that most of continuous problems can be well approximated using only elementary operations.

First steps of the approximation

- Go from infinite dimensional spaces to finite dimensional (discretization);

Approximation

We accept an approximated solution.

It is perhaps surprising that most of continuous problems can be well approximated using only elementary operations.

First steps of the approximation

- Go from infinite dimensional spaces to finite dimensional (discretization);
- Go from a generic function to a rational function.

We make the following approximation

$$\hat{f} : V \rightarrow W \implies \varphi : \mathbb{R}^n \rightarrow \mathbb{R} \implies f : \mathbb{R}^n \rightarrow \mathbb{R}$$

with f rational.

Approximation

We accept an approximated solution.

It is perhaps surprising that most of continuous problems can be well approximated using only elementary operations.

First steps of the approximation

- Go from infinite dimensional spaces to finite dimensional (discretization);
- Go from a generic function to a rational function.

We make the following approximation

$$\hat{f} : V \rightarrow W \implies \varphi : \mathbb{R}^n \rightarrow \mathbb{R} \implies f : \mathbb{R}^n \rightarrow \mathbb{R}$$

with f rational. The function f evaluated on finite arithmetic is said to be a **numerical algorithm**.

We can define the relative analytic error, when $\varphi(x) \neq 0$, as

$$\varepsilon_{AN} = \frac{f(x) - \varphi(x)}{\varphi(x)}$$

Analytic error

Let $\text{conv}(a, b) = \text{conv}(\{a, b\})$ be the interval $[a, b]$ if $a < b$; the interval $[b, a]$ if $a > b$; or the point a if $a = b$.

Analytic error

Let $\text{conv}(a, b) = \text{conv}(\{a, b\})$ be the interval $[a, b]$ if $a < b$; the interval $[b, a]$ if $a > b$; or the point a if $a = b$.

Theorem (Taylor's formula)

Let $f \in C^{r+1}(\mathbb{R})$ and let $x_0 \in \mathbb{R}$. Then for $x \in \mathbb{R}$, there exists $\xi(x) \in \text{conv}(x_0, x)$, such that

$$f(x) = p(x) + r(x)$$

where

$$p(x) = f(x_0) + f'(x_0)(x - x_0) + \cdots + \frac{f^{(r)}(x_0)}{r!}(x - x_0)^r$$

and

$$r(x) = \frac{f^{(r+1)}(\xi(x))}{(r+1)!}(x - x_0)^{r+1}$$

Analytic error

Theorem (Taylor's formula)

Let $f \in C^{r+1}(\mathbb{R})$ and let $x_0 \in \mathbb{R}$. Then for $x \in \mathbb{R}$, there exists $\xi(x) \in \text{conv}(x_0, x)$, such that

$$f(x) = p(x) + r(x)$$

where

$$p(x) = f(x_0) + f'(x_0)(x - x_0) + \cdots + \frac{f^{(r)}(x_0)}{r!}(x - x_0)^r$$

and

$$r(x) = \frac{f^{(r+1)}(\xi(x))}{(r+1)!}(x - x_0)^{r+1}$$

Analytic error

Theorem (Taylor's formula)

Let $f \in C^{r+1}(\mathbb{R})$ and let $x_0 \in \mathbb{R}$. Then for $x \in \mathbb{R}$, there exists $\xi(x) \in \text{conv}(x_0, x)$, such that

$$f(x) = p(x) + r(x)$$

where

$$p(x) = f(x_0) + f'(x_0)(x - x_0) + \cdots + \frac{f^{(r)}(x_0)}{r!}(x - x_0)^r$$

and

$$r(x) = \frac{f^{(r+1)}(\xi(x))}{(r+1)!}(x - x_0)^{r+1}$$

$p(x)$ is the Taylor polynomial of degree r at x_0 .

Analytic error

Theorem (Taylor's formula)

Let $f \in C^{r+1}(\mathbb{R})$ and let $x_0 \in \mathbb{R}$. Then for $x \in \mathbb{R}$, there exists $\xi(x) \in \text{conv}(x_0, x)$, such that

$$f(x) = p(x) + r(x)$$

where

$$p(x) = f(x_0) + f'(x_0)(x - x_0) + \cdots + \frac{f^{(r)}(x_0)}{r!}(x - x_0)^r$$

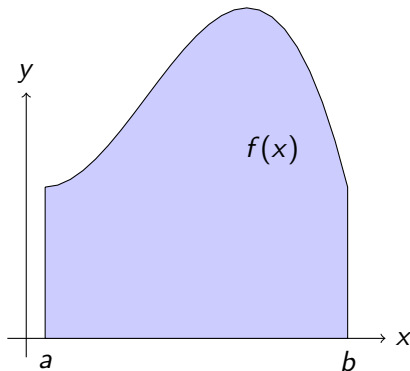
and

$$r(x) = \frac{f^{(r+1)}(\xi(x))}{(r+1)!}(x - x_0)^{r+1}$$

$r(x)$ is the remainder (Lagrange's remainder) and it is an **absolute** analytic error (that is small for x near to x_0).

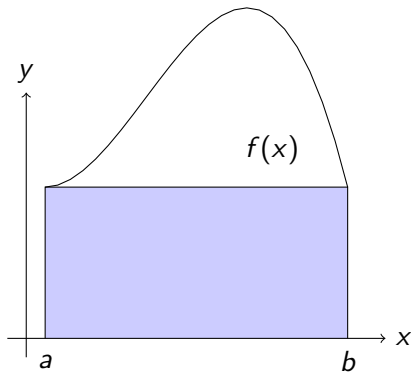
Analytic error

The integral of $f \in C[a, b]$, with $f > 0$ can be approximated by a trapezoid (trapezoidal rule).



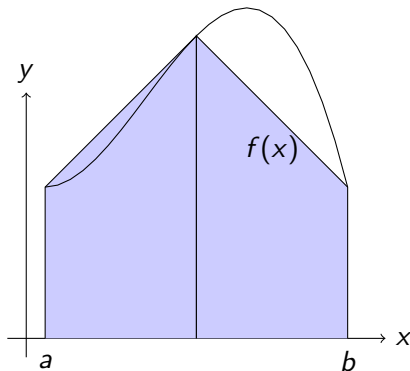
Analytic error

The integral of $f \in C[a, b]$, with $f > 0$ can be approximated by a trapezoid (trapezoidal rule).



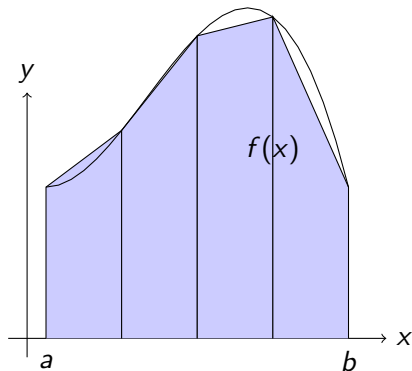
Analytic error

The integral of $f \in C[a, b]$, with $f > 0$ can be approximated by a trapezoid (trapezoidal rule).

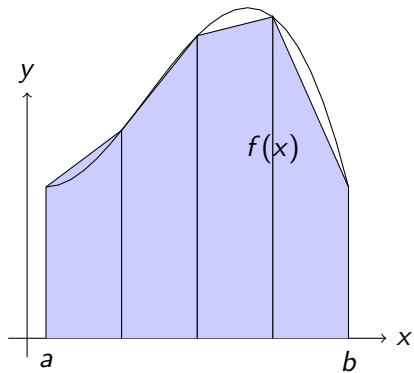


Analytic error

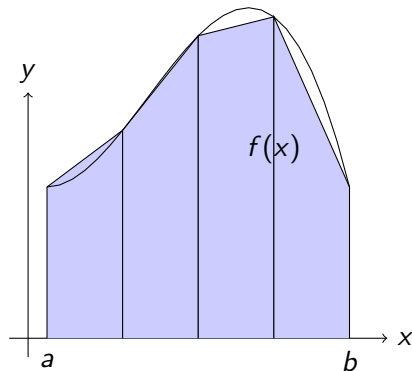
The integral of $f \in C[a, b]$, with $f > 0$ can be approximated by a trapezoid (trapezoidal rule).



Analytic error



Analytic error



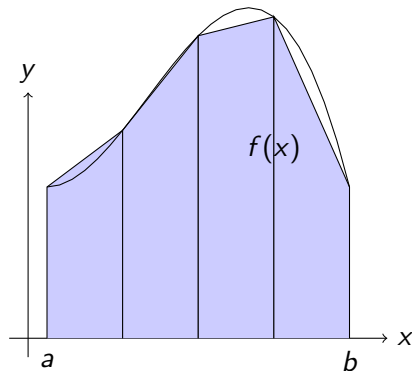
Let $h = (b - a)/n$, and let $x_i = a + \frac{b-a}{n}i$, then

$$\int_a^b f(x)dx = h\left(\frac{1}{2}f(x_0) + f(x_1) + \cdots + f(x_{n-1}) + \frac{1}{2}f(x_n)\right) + r(x)$$

It can be proved that, if $f \in C^2[a, b]$, there exists $\xi \in (a, b)$, such that

$$r(x) = \frac{(b-a)^3}{12n^2}f''(\xi).$$

Analytic error



Let $h = (b - a)/n$, and let $x_i = a + \frac{b-a}{n}i$, then

$$\int_a^b f(x)dx = h\left(\frac{1}{2}f(x_0) + f(x_1) + \cdots + f(x_{n-1}) + \frac{1}{2}f(x_n)\right) + r(x)$$

It can be proved that, if $f \in C^2[a, b]$, there exists $\xi \in (a, b)$, such that

$$r(x) = \frac{(b-a)^3}{12n^2}f''(\xi).$$

Disclaimer

We will not discuss anymore of the analytic error, and we will only consider rational functions (or elementary functions).

Error for rational functions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational, that is $f = p/q$ with p, q polynomials.

From mathematical analysis we know that f is defined and differentiable for $q \neq 0$. (We may assume that p and q are prime).

Examples of rational functions:

$$f(x) = \frac{x^2 + 2x}{x + 1}, \quad f(x) = \frac{xy + x^2 + 5x^2yz}{z + \pi y^2}$$

But also

Error for rational functions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational, that is $f = p/q$ with p, q polynomials.

From mathematical analysis we know that f is defined and differentiable for $q \neq 0$. (We may assume that p and q are prime).

Examples of rational functions:

$$f(x) = \frac{x^2 + 2x}{x + 1}, \quad f(x) = \frac{xy + x^2 + 5x^2yz}{z + \pi y^2}$$

But also

- the Taylor polynomial (one variable);
- the trapezoidal rule, if the integrand function is a rational function;
- the determinant of A .

Error for rational functions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational, that is $f = p/q$ with p, q polynomials.

From mathematical analysis we know that f is defined and differentiable for $q \neq 0$. (We may assume that p and q are prime).

Examples of rational functions:

$$f(x) = \frac{x^2 + 2x}{x + 1}, \quad f(x) = \frac{xy + x^2 + 5x^2yz}{z + \pi y^2}$$

But also

- the Taylor polynomial (one variable);
- the trapezoidal rule, if the integrand function is a rational function;
- the determinant of A .

There are two types of error in the evaluation of f at x on floating point numbers.

Error for rational functions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational, that is $f = p/q$ with p, q polynomials.

- **Inherent error.** We do not evaluate $f(x)$ but rather $f(\tilde{x})$ where $\tilde{x} = \text{fl}(x)$

$$\varepsilon_{IN} = \frac{f(\tilde{x}) - f(x)}{f(x)}, \quad f(x) \neq 0.$$

Error for rational functions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational, that is $f = p/q$ with p, q polynomials.

- **Inherent error.** We do not evaluate $f(x)$ but rather $f(\tilde{x})$ where $\tilde{x} = \text{fl}(x)$

$$\varepsilon_{IN} = \frac{f(\tilde{x}) - f(x)}{f(x)}, \quad f(x) \neq 0.$$

- **Algorithmic error.** We do not evaluate $f(\tilde{x})$ but another function $\tilde{f}(\tilde{x})$, since the arithmetic operations of f are computed as machine operations.

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}, \quad f(\tilde{x}) \neq 0.$$

Error for rational functions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational, that is $f = p/q$ with p, q polynomials.

- **Inherent error.** We do not evaluate $f(x)$ but rather $f(\tilde{x})$ where $\tilde{x} = \text{fl}(x)$

$$\varepsilon_{IN} = \frac{f(\tilde{x}) - f(x)}{f(x)}, \quad f(x) \neq 0.$$

- **Algorithmic error.** We do not evaluate $f(\tilde{x})$ but another function $\tilde{f}(\tilde{x})$, since the arithmetic operations of f are computed as machine operations.

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}, \quad f(\tilde{x}) \neq 0.$$

Inherent error can be defined also for non-rational functions.

Algorithmic error can be defined also for elementary functions, which are treated as operations.

Error for rational functions

- **Inherent error.** We do not evaluate $f(x)$ but rather $f(\tilde{x})$ where $\tilde{x} = \text{fl}(x)$

$$\varepsilon_{IN} = \frac{f(\tilde{x}) - f(x)}{f(x)}, \quad f(x) \neq 0.$$

If the inherent error is relatively small we say that the problem is **well-conditioned** or else **ill-conditioned**.

- **Algorithmic error.** We do not evaluate $f(\tilde{x})$ but another function $\tilde{f}(\tilde{x})$, since the arithmetic operations of f are computed as machine operations.

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}, \quad f(\tilde{x}) \neq 0.$$

If the algorithmic error is relatively small we say that algorithm f is **numerically stable** or else **numerically unstable**.

A bit imprecise. Can be made more precise assuming $u \rightarrow 0$ and asking $\lim_{u \rightarrow 0} \varepsilon_{IN/ALG} = 0$.

Error for rational functions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational, that is $f = p/q$ with p, q polynomials.

The total error (or forward error) is

$$\varepsilon_{TOT} = \frac{\tilde{f}(\tilde{x}) - f(x)}{f(x)}, \quad f(x) \neq 0,$$

which gives a genuine measure of the error in the evaluation.

The ideal situation is $|\varepsilon_{TOT}| < u$, where u is the machine precision.

But, in practice, it is sufficient $|\varepsilon_{TOT}| \leq Mu$, with M constant.

Error for rational functions

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational, that is $f = p/q$ with p, q polynomials.

The total error (or forward error) is

$$\varepsilon_{TOT} = \frac{\tilde{f}(\tilde{x}) - f(x)}{f(x)}, \quad f(x) \neq 0,$$

which gives a genuine measure of the error in the evaluation.

The ideal situation is $|\varepsilon_{TOT}| < u$, where u is the machine precision.
But, in practice, it is sufficient $|\varepsilon_{TOT}| \leq Mu$, with M constant.

What is the relationship between these errors?

Error for rational functions

Theorem

Let $x \in \mathbb{R}^n \setminus \{0\}$, and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational with $f(x) \neq 0$ and $f(\tilde{x}) \neq 0$, where $\tilde{x} = \text{fl}(x)$, then

$$\varepsilon_{TOT} = \varepsilon_{IN} + \varepsilon_{ALG} + \varepsilon_{IN}\varepsilon_{ALG}.$$

Error for rational functions

Theorem

Let $x \in \mathbb{R}^n \setminus \{0\}$, and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational with $f(x) \neq 0$ and $f(\tilde{x}) \neq 0$, where $\tilde{x} = \text{fl}(x)$, then

$$\varepsilon_{TOT} = \varepsilon_{IN} + \varepsilon_{ALG} + \varepsilon_{IN}\varepsilon_{ALG}.$$

Proof.

$$\begin{aligned}\varepsilon_{TOT} &= \frac{\tilde{f}(\tilde{x}) - f(\tilde{x}) + f(\tilde{x}) - f(x)}{f(x)} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(x)} + \frac{f(\tilde{x}) - f(x)}{f(x)} \\ &= \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})} \frac{f(\tilde{x})}{f(x)} + \varepsilon_{IN} = \varepsilon_{ALG} \left(\frac{f(\tilde{x})}{f(x)} - 1 + 1 \right) + \varepsilon_{IN} \\ &= \varepsilon_{ALG} \left(\frac{f(\tilde{x}) - f(x)}{f(x)} + 1 \right) + \varepsilon_{IN} = \varepsilon_{ALG}(\varepsilon_{IN} + 1) + \varepsilon_{IN} \\ &= \varepsilon_{IN} + \varepsilon_{ALG} + \varepsilon_{IN}\varepsilon_{ALG}.\end{aligned}$$



Error for rational functions

Theorem

Let $x \in \mathbb{R}^n \setminus \{0\}$, and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ rational with $f(x) \neq 0$ and $f(\tilde{x}) \neq 0$, where $\tilde{x} = \text{fl}(x)$, then

$$\varepsilon_{TOT} = \varepsilon_{IN} + \varepsilon_{ALG} + \varepsilon_{IN}\varepsilon_{ALG}.$$

If ε_{IN} and ε_{ALG} are small (tends to zero as $u \rightarrow 0$, we have)

$$\varepsilon_{TOT} = \varepsilon_{IN} + \varepsilon_{ALG} + o(u) \doteq \varepsilon_{IN} + \varepsilon_{ALG}$$

Error analysis

Error analysis consists in relating the forward error with the error in the representation (and thus on the machine precision)

Formulae cannot be used directly

- Inherent error, we use the derivative (if the function is differentiable).
- Algorithmic error, we use a diagram analysis

Inherent error

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$, rational,

$$\varepsilon_{IN} = \frac{f(\tilde{x}) - f(x)}{f(x)}, \quad f(x) \neq 0.$$

where $\tilde{x} = \text{fl}(x)$.

If $x_i \neq 0$, for $i = 1, \dots, n$, we want to relate ε_{IN} with the representation errors

$$\varepsilon_i = \frac{\tilde{x}_i - x_i}{x_i},$$

and we make the usual assumption

$$|\varepsilon_i| < u,$$

where u is the machine precision. Ignoring overflow and underflow for simplicity.

Inherent error

We try to use the definition.

Example 1.

$$f(x) = x^2,$$

Let $\tilde{x} = \text{fl}(x) = x(1 + \varepsilon_1)$, with $|\varepsilon_1| < u$, obtained as

$$\frac{\tilde{x} - x}{x} = \varepsilon_1 \iff \tilde{x} - x = x\varepsilon_1 \iff \tilde{x} = x(1 + \varepsilon_1).$$

We have, for $x \neq 0$,

$$\varepsilon_{IN} = \frac{\tilde{x}^2 - x^2}{x^2} = \frac{x^2(1 + \varepsilon_1)^2 - x^2}{x^2} = \frac{x^2(1 + 2\varepsilon_1 + \varepsilon_1^2 - 1)}{x^2} = 2\varepsilon_1 + \varepsilon_1^2.$$

Since we are mostly interested on what happens as $u \rightarrow 0$, we can consider only first order terms

$$\varepsilon_{IN} \doteq 2\varepsilon_1$$

or $\varepsilon_{IN} = 2\varepsilon_1 + o(u)$. **Well-conditioned!**

Inherent error

In general using the definition is complicated. A simpler technique is obtained using the following.

Theorem

Let $x, \tilde{x} \in \mathbb{R} \setminus \{0\}$ and let $f \in C^1(\text{conv}(x, \tilde{x}))$, with $f(x) \neq 0$. There exists $\xi \in \text{conv}(x, \tilde{x})$, such that

$$\varepsilon_{IN} = \frac{x}{f(x)} f'(\xi) \varepsilon_x,$$

where $\varepsilon_x = \varepsilon_1$ is the representation error on x . If, moreover $f \in C^2(\text{conv}(x, \tilde{x}))$, then

$$\varepsilon_{IN} = \frac{x}{f(x)} f'(x) \varepsilon_x + o(u),$$

that is $\varepsilon_{IN} \doteq \frac{x}{f(x)} f'(x) \varepsilon_x$.

The formula above is very useful.

Inherent error

Proof.

$$\frac{f(\tilde{x}) - f(x)}{f(x)} = \frac{f'(\xi)(\tilde{x} - x)}{f(x)} \frac{x}{x} = \frac{xf'(\xi)}{f(x)} \frac{\tilde{x} - x}{x} = \frac{x}{f(x)} f'(\xi) \varepsilon_1.$$

The mean value theorem has been used.



Inherent error

$$\varepsilon_{IN} = \frac{f(\tilde{x}) - f(x)}{f(x)} \doteq \frac{x}{f(x)} f'(x) \varepsilon_x.$$

The term $c_x = \frac{x}{f(x)} f'(x)$ is said to be **amplification factor** and measure the amplification of the error.

We use the formula with the function $f(x) = x^2$.

$$\varepsilon_{IN} \doteq \frac{x}{x^2} 2x \varepsilon_1 = 2\varepsilon_1$$

Much easier!

Inherent error

$$\varepsilon_{IN} = \frac{f(\tilde{x}) - f(x)}{f(x)} \doteq \frac{x}{f(x)} f'(x) \varepsilon_x.$$

The term $c_x = \frac{x}{f(x)} f'(x)$ is said to be **amplification factor** and measure the amplification of the error.

We use the formula with the function $f(x) = x^2$.

$$\varepsilon_{IN} \doteq \frac{x}{x^2} 2x \varepsilon_1 = 2\varepsilon_1$$

Much easier!

What happens for more than one variable?

Inherent error

If $x = (x_1, \dots, x_n)$, and consider $f(x)$, we have the formula (with $x_i \neq 0$, $f(x) \neq 0$)

$$\varepsilon_{IN} \doteq \frac{x_1}{f(x)} \frac{\partial f}{\partial x_1}(x) \varepsilon_1 + \frac{x_2}{f(x)} \frac{\partial f}{\partial x_2}(x) \varepsilon_2 + \dots + \frac{x_n}{f(x)} \frac{\partial f}{\partial x_n}(x) \varepsilon_n$$

where

$$\varepsilon_i = \frac{\tilde{x}_i - x_i}{x_i}, \quad c_i = \frac{x_i}{f(x)} \frac{\partial f}{\partial x_i}(x),$$

are the representation error and the amplification coefficients, respectively.

Inherent error on operations

$$\varepsilon_{IN} \doteq c_1 \varepsilon_1 + \cdots c_n \varepsilon_n, \quad c_i = \frac{x}{f(x)} \frac{\partial f}{\partial x_i}(x).$$

Let $f(x, y) = xy$, the multiplication. Then $\varepsilon_{IN} \doteq c_x \varepsilon_x + c_y \varepsilon_y$, where

Inherent error on operations

$$\varepsilon_{IN} \doteq c_1 \varepsilon_1 + \cdots c_n \varepsilon_n, \quad c_i = \frac{x}{f(x)} \frac{\partial f}{\partial x_i}(x).$$

Let $f(x, y) = xy$, the multiplication. Then $\varepsilon_{IN} \doteq c_x \varepsilon_x + c_y \varepsilon_y$, where

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{xy} y = 1,$$

Inherent error on operations

$$\varepsilon_{IN} \doteq c_1 \varepsilon_1 + \cdots c_n \varepsilon_n, \quad c_i = \frac{x}{f(x)} \frac{\partial f}{\partial x_i}(x).$$

Let $f(x, y) = xy$, the multiplication. Then $\varepsilon_{IN} \doteq c_x \varepsilon_x + c_y \varepsilon_y$, where

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{xy} y = 1, \quad c_y = \frac{y}{f(x, y)} \frac{\partial f}{\partial y}(x, y) = \frac{y}{xy} x = 1.$$

Inherent error on operations

$$\varepsilon_{IN} \doteq c_1 \varepsilon_1 + \cdots c_n \varepsilon_n, \quad c_i = \frac{x}{f(x)} \frac{\partial f}{\partial x_i}(x).$$

Let $f(x, y) = xy$, the multiplication. Then $\varepsilon_{IN} \doteq c_x \varepsilon_x + c_y \varepsilon_y$, where

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{xy} y = 1, \quad c_y = \frac{y}{f(x, y)} \frac{\partial f}{\partial y}(x, y) = \frac{y}{xy} x = 1.$$

Thus we have

$$c_x = 1, \quad c_y = 1,$$

and the problem is well-conditioned. Obvious, really!

Inherent error on operations

$$\varepsilon_{IN} \doteq c_1 \varepsilon_1 + \cdots c_n \varepsilon_n, \quad c_i = \frac{x}{f(x)} \frac{\partial f}{\partial x_i}(x).$$

Let $f(x, y) = x + y$, the sum. Then $\varepsilon_{IN} \doteq c_x \varepsilon_x + c_y \varepsilon_y$, where

Inherent error on operations

$$\varepsilon_{IN} \doteq c_1\varepsilon_1 + \cdots c_n\varepsilon_n, \quad c_i = \frac{x}{f(x)} \frac{\partial f}{\partial x_i}(x).$$

Let $f(x, y) = x + y$, the sum. Then $\varepsilon_{IN} \doteq c_x\varepsilon_x + c_y\varepsilon_y$, where

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{x + y},$$

Inherent error on operations

$$\varepsilon_{IN} \doteq c_1 \varepsilon_1 + \cdots c_n \varepsilon_n, \quad c_i = \frac{x}{f(x)} \frac{\partial f}{\partial x_i}(x).$$

Let $f(x, y) = x + y$, the sum. Then $\varepsilon_{IN} \doteq c_x \varepsilon_x + c_y \varepsilon_y$, where

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{x + y}, \quad c_y = \frac{y}{f(x, y)} \frac{\partial f}{\partial y}(x, y) = \frac{y}{x + y}.$$

Inherent error on operations

$$\varepsilon_{IN} \doteq c_1 \varepsilon_1 + \cdots c_n \varepsilon_n, \quad c_i = \frac{x}{f(x)} \frac{\partial f}{\partial x_i}(x).$$

Let $f(x, y) = x + y$, the sum. Then $\varepsilon_{IN} \doteq c_x \varepsilon_x + c_y \varepsilon_y$, where

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{x + y}, \quad c_y = \frac{y}{f(x, y)} \frac{\partial f}{\partial y}(x, y) = \frac{y}{x + y}.$$

Thus we have

$$c_x = x/(x + y), \quad c_y = y/(x + y),$$

and the problem becomes ill-conditioned when $x \approx -y$. Perhaps surprising!

Inherent error on operations

Let $f(x, y) = x + y$, the sum. Then $\varepsilon_{IN} \doteq c_x \varepsilon_x + c_y \varepsilon_y$, where

$$c_x = x/(x + y), \quad c_y = y/(x + y),$$

and the problem becomes ill-conditioned when $x \approx -y$. Perhaps surprising!

This phenomenon is said to be **numerical cancellation** and can be understood intuitively.

If we subtract two very near numbers we lose significant digits.

Example. In \mathbb{R} , $\pi - 3.1 \approx 0.0415926$. In $\mathcal{F}(10, 3, *, *)$, $\text{fl}(\pi) = 3.14$ and $\text{fl}(3.1) = 3.1$, and thus $\text{fl}(\pi) - \text{fl}(3.1) = 0.04$ (while $\text{fl}(\pi - 3.1) = 0.0416$) and we have lost two significant digits.

Inherent error for some functions

operation	c_x	c_y
$x + y$	$\frac{x}{x+y}$	$\frac{y}{x+y}$

Inherent error for some functions

operation	c_x	c_y
$x + y$	$\frac{x}{x+y}$	$\frac{y}{x+y}$
$x - y$	$\frac{x}{x-y}$	$\frac{-y}{x-y}$

$$c_x = \frac{x}{x-y} \frac{\partial f}{\partial x} = \frac{x}{x-y}, \quad c_y = \frac{y}{x-y} \frac{\partial f}{\partial y} = \frac{y}{x-y} \cdot (-1)$$

Inherent error for some functions

operation	c_x	c_y
$x + y$	$\frac{x}{x+y}$	$\frac{y}{x+y}$
$x - y$	$\frac{x}{x-y}$	$\frac{-y}{x-y}$
xy	1	1

Inherent error for some functions

operation	c_x	c_y
$x + y$	$\frac{x}{x+y}$	$\frac{y}{x+y}$
$x - y$	$\frac{x}{x-y}$	$\frac{-y}{x-y}$
xy	1	1
x/y	1	-1

$$c_x = \frac{x}{f} \frac{\partial f}{\partial x} = \frac{x}{x/y} \frac{1}{y} = x \frac{y}{x} \frac{1}{y} = 1, \quad c_y = \frac{y}{x/y} \frac{-x}{y^2} = -1,$$

Inherent error for some functions

operation	c_x	c_y
$x + y$	$\frac{x}{x+y}$	$\frac{y}{x+y}$
$x - y$	$\frac{x}{x-y}$	$\frac{-y}{x-y}$
xy	1	1
x/y	1	-1
$\exp(x)$	x	

$$c_x = \frac{x}{\exp(x)} \exp(x) = x,$$

Inherent error for some functions

operation	c_x	c_y
$x + y$	$\frac{x}{x+y}$	$\frac{y}{x+y}$
$x - y$	$\frac{x}{x-y}$	$\frac{-y}{x-y}$
xy	1	1
x/y	1	-1
$\exp(x)$	x	
\sqrt{x}	$\frac{1}{2}$	

$$c_x = \frac{x}{\sqrt{x}} \frac{1}{2\sqrt{x}} = \frac{1}{2},$$

Inherent error for some functions

operation	c_x	c_y
$x + y$	$\frac{x}{x+y}$	$\frac{y}{x+y}$
$x - y$	$\frac{x}{x-y}$	$\frac{-y}{x-y}$
xy	1	1
x/y	1	-1
$\exp(x)$	x	
\sqrt{x}	$\frac{1}{2}$	
x^α	α	

$$c_x = \frac{x}{x^\alpha} \alpha x^{\alpha-1} = \alpha.$$

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the **inherent error** or the **amplification coefficients** and find out when these quantities go to infinity.

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the **inherent error** or the **amplification coefficients** and find out when these quantities go to infinity.

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) =$$

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the **inherent error** or the **amplification coefficients** and find out when these quantities go to infinity.

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{x^2 - y^2} 2x = \frac{2x^2}{x^2 - y^2},$$

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the **inherent error** or the **amplification coefficients** and find out when these quantities go to infinity.

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{x^2 - y^2} 2x = \frac{2x^2}{x^2 - y^2},$$

$$c_y = \frac{y}{f(x, y)} \frac{\partial f}{\partial y}(x, y) =$$

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the **inherent error** or the **amplification coefficients** and find out when these quantities go to infinity.

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{x^2 - y^2} 2x = \frac{2x^2}{x^2 - y^2},$$

$$c_y = \frac{y}{f(x, y)} \frac{\partial f}{\partial y}(x, y) = \frac{y}{x^2 - y^2} (-2y) = \frac{-2y^2}{x^2 - y^2},$$

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the **inherent error** or the **amplification coefficients** and find out when these quantities go to infinity.

$$c_x = \frac{x}{f(x, y)} \frac{\partial f}{\partial x}(x, y) = \frac{x}{x^2 - y^2} 2x = \frac{2x^2}{x^2 - y^2},$$

$$c_y = \frac{y}{f(x, y)} \frac{\partial f}{\partial y}(x, y) = \frac{y}{x^2 - y^2} (-2y) = \frac{-2y^2}{x^2 - y^2},$$

$$\varepsilon_{IN} = \frac{2x^2}{x^2 - y^2} \varepsilon_x + \frac{-2y^2}{x^2 - y^2} \varepsilon_y,$$

where ε_x and ε_y are the representation errors for x and y , respectively.

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the inherent error or the amplification coefficients and **find out when these quantities go to infinity**.

$$\varepsilon_{IN} = \frac{2x^2}{x^2 - y^2} \varepsilon_x + \frac{-2y^2}{x^2 - y^2} \varepsilon_y.$$

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the inherent error or the amplification coefficients and **find out when these quantities go to infinity**.

$$\varepsilon_{IN} = \frac{2x^2}{x^2 - y^2} \varepsilon_x + \frac{-2y^2}{x^2 - y^2} \varepsilon_y.$$

It is apparent that the inherent error is large for $x^2 \approx y^2$, that is

Exercise

Study the inherent error (conditioning) of $f(x, y) = x^2 - y^2$.

Rule of thumb: one must compute the inherent error or the amplification coefficients and **find out when these quantities go to infinity**.

$$\varepsilon_{IN} = \frac{2x^2}{x^2 - y^2} \varepsilon_x + \frac{-2y^2}{x^2 - y^2} \varepsilon_y.$$

It is apparent that the inherent error is large for $x^2 \approx y^2$, that is

$$x \approx \pm y.$$

We expect a large error when x is near to y or $-y$.

Algorithmic error

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}$$

Algorithmic error is very hard to be computed using the definition.

We drop the tilde for simplicity.

Algorithmic error

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}$$

Algorithmic error is very hard to be computed using the definition.

We drop the tilde for simplicity.

Example. $f(x) = x^2$, $\tilde{f}(x) = x \otimes x = x^2(1 + \eta)$,

Algorithmic error

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}$$

Algorithmic error is very hard to be computed using the definition.

We drop the tilde for simplicity.

Example. $f(x) = x^2$, $\tilde{f}(x) = x \otimes x = x^2(1 + \eta)$, where $|\eta| \leq u$ is the error in the operation.

Algorithmic error

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}$$

Algorithmic error is very hard to be computed using the definition.

We drop the tilde for simplicity.

Example. $f(x) = x^2$, $\tilde{f}(x) = x \otimes x = x^2(1 + \eta)$, where $|\eta| \leq u$ is the error in the operation.

$$\varepsilon_{ALG} = \frac{x^2(1 + \eta) - x^2}{x^2} =$$

Algorithmic error

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}$$

Algorithmic error is very hard to be computed using the definition.

We drop the tilde for simplicity.

Example. $f(x) = x^2$, $\tilde{f}(x) = x \otimes x = x^2(1 + \eta)$, where $|\eta| \leq u$ is the error in the operation.

$$\varepsilon_{ALG} = \frac{x^2(1 + \eta) - x^2}{x^2} = \frac{x^2(1 + \eta - 1)}{x^2} =$$

Algorithmic error

$$\varepsilon_{ALG} = \frac{\tilde{f}(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}$$

Algorithmic error is very hard to be computed using the definition.

We drop the tilde for simplicity.

Example. $f(x) = x^2$, $\tilde{f}(x) = x \otimes x = x^2(1 + \eta)$, where $|\eta| \leq u$ is the error in the operation.

$$\varepsilon_{ALG} = \frac{x^2(1 + \eta) - x^2}{x^2} = \frac{x^2(1 + \eta - 1)}{x^2} = \eta.$$

This is a special case: a single operation.

Algorithmic error

Consider the function $f(x, y) = x^2 - y^2$ and the two algorithms

$$\begin{cases} z^{(1)} = x^2, \\ z^{(2)} = y^2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}. \end{cases}$$

Algorithmic error

Consider the function $f(x, y) = x^2 - y^2$ and the two algorithms

$$\left\{ \begin{array}{l} z^{(1)} = x^2, \\ z^{(2)} = y^2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}. \end{array} \right. \quad \left\{ \begin{array}{l} v^{(1)} = x - y, \\ v^{(2)} = x + y, \\ f = v^{(3)} = v^{(1)}v^{(2)}. \end{array} \right.$$

Algorithmic error

Consider the function $f(x, y) = x^2 - y^2$ and the two algorithms

$$\left\{ \begin{array}{l} z^{(1)} = x^2, \\ z^{(2)} = y^2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}. \end{array} \right. \quad \left\{ \begin{array}{l} v^{(1)} = x - y, \\ v^{(2)} = x + y, \\ f = v^{(3)} = v^{(1)}v^{(2)}. \end{array} \right.$$

The two algorithms are different due to the finite arithmetic.

Algorithmic error

Consider the function $f(x, y) = x^2 - y^2$ and the two algorithms

$$\left\{ \begin{array}{l} z^{(1)} = x^2, \\ z^{(2)} = y^2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}. \end{array} \right. \quad \left\{ \begin{array}{l} v^{(1)} = x - y, \\ v^{(2)} = x + y, \\ f = v^{(3)} = v^{(1)}v^{(2)}. \end{array} \right.$$

The two algorithms are different due to the finite arithmetic.

We may assume that there is a local error for any operation

$$\left\{ \begin{array}{ll} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{array} \right.$$

Algorithmic error

Consider the function $f(x, y) = x^2 - y^2$ and the two algorithms

$$\left\{ \begin{array}{l} z^{(1)} = x^2, \\ z^{(2)} = y^2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}. \end{array} \right. \quad \left\{ \begin{array}{l} v^{(1)} = x - y, \\ v^{(2)} = x + y, \\ f = v^{(3)} = v^{(1)}v^{(2)}. \end{array} \right.$$

The two algorithms are different due to the finite arithmetic.

We may assume that there is a local error for any operation

$$\left\{ \begin{array}{ll} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{array} \right. \quad \left\{ \begin{array}{ll} v^{(1)} = x - y, & \eta_1, \\ v^{(2)} = x + y, & \eta_2, \\ f = v^{(3)} = v^{(1)}v^{(2)} & \eta_3. \end{array} \right.$$

where $|\varepsilon_i| \leq u$, $|\eta_i| \leq u$.

Algorithmic error

Studying the algorithmic error using the definition is horrible

$$\begin{cases} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{cases}$$

$$\varepsilon_{ALG} = \frac{(x \otimes x) \ominus (y \otimes y) - (x^2 - y^2)}{x^2 - y^2}$$

Algorithmic error

Studying the algorithmic error using the definition is horrible

$$\begin{cases} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{cases}$$

$$\begin{aligned} \varepsilon_{ALG} &= \frac{(x \otimes x) \ominus (y \otimes y) - (x^2 - y^2)}{x^2 - y^2} \\ &= \frac{(x^2(1 + \varepsilon_1) - y^2(1 + \varepsilon_2))(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \end{aligned}$$

Algorithmic error

Studying the algorithmic error using the definition is horrible

$$\begin{cases} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{cases}$$

$$\begin{aligned} \varepsilon_{ALG} &= \frac{(x \otimes x) \ominus (y \otimes y) - (x^2 - y^2)}{x^2 - y^2} \\ &= \frac{(x^2(1 + \varepsilon_1) - y^2(1 + \varepsilon_2))(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\ &= \frac{(x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2)(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \end{aligned}$$

Algorithmic error

Studying the algorithmic error using the definition is horrible

$$\begin{cases} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{cases}$$

$$\begin{aligned} \varepsilon_{ALG} &= \frac{(x \otimes x) \ominus (y \otimes y) - (x^2 - y^2)}{x^2 - y^2} \\ &= \frac{(x^2(1 + \varepsilon_1) - y^2(1 + \varepsilon_2))(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\ &= \frac{(x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2)(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\ &= \frac{x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2 + x^2\varepsilon_3 - x^2\varepsilon_1\varepsilon_3 - y^2\varepsilon_3 - y^2\varepsilon_2\varepsilon_3 - x^2 - y^2}{x^2 - y^2} \end{aligned}$$

Algorithmic error

$$\begin{aligned}\varepsilon_{ALG} &= \frac{(x \otimes x) \ominus (y \otimes y) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{(x^2(1 + \varepsilon_1) - y^2(1 + \varepsilon_2))(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{(x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2)(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2 + x^2\varepsilon_3 - x^2\varepsilon_1\varepsilon_3 - y^2\varepsilon_3 - y^2\varepsilon_2\varepsilon_3 - x^2 + y^2}{x^2 - y^2}\end{aligned}$$

Algorithmic error

$$\begin{aligned}\varepsilon_{ALG} &= \frac{(x \otimes x) \ominus (y \otimes y) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{(x^2(1 + \varepsilon_1) - y^2(1 + \varepsilon_2))(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{(x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2)(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2 + x^2\varepsilon_3 - x^2\varepsilon_1\varepsilon_3 - y^2\varepsilon_3 - y^2\varepsilon_2\varepsilon_3 - x^2 + y^2}{x^2 - y^2} \\&= \frac{x^2\varepsilon_1 - y^2\varepsilon_2 + x^2\varepsilon_3 - y^2\varepsilon_3 - x^2\varepsilon_1\varepsilon_3 - y^2\varepsilon_2\varepsilon_3}{x^2 - y^2}\end{aligned}$$

Algorithmic error

$$\begin{aligned}\varepsilon_{ALG} &= \frac{(x \otimes x) \ominus (y \otimes y) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{(x^2(1 + \varepsilon_1) - y^2(1 + \varepsilon_2))(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{(x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2)(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{\cancel{x^2} - x^2\varepsilon_1 - \cancel{y^2} - y^2\varepsilon_2 + x^2\varepsilon_3 - x^2\varepsilon_1\varepsilon_3 - y^2\varepsilon_3 - y^2\varepsilon_2\varepsilon_3 - \cancel{x^2} + \cancel{y^2}}{x^2 - y^2} \\&= \frac{x^2\varepsilon_1 - y^2\varepsilon_2 + x^2\varepsilon_3 - y^2\varepsilon_3 - \cancel{x^2\varepsilon_1\varepsilon_3} - \cancel{y^2\varepsilon_2\varepsilon_3}}{x^2 - y^2} \\&\stackrel{.}{=} \frac{x^2\varepsilon_1 - y^2\varepsilon_2 + (x^2 - y^2)\varepsilon_3}{x^2 - y^2} = \frac{(x^2 - y^2)\varepsilon_3}{x^2 - y^2} + \frac{x^2\varepsilon_1 - y^2\varepsilon_2}{x^2 - y^2}\end{aligned}$$

Algorithmic error

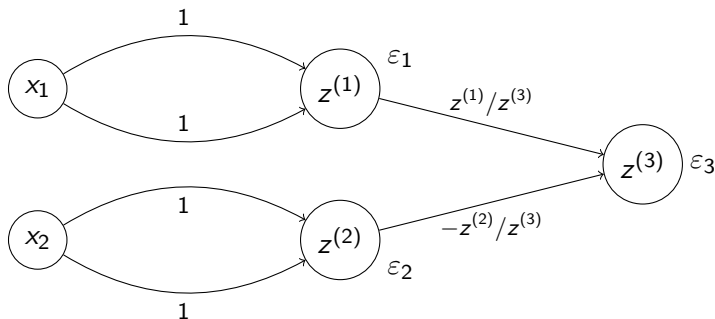
$$\begin{aligned}\varepsilon_{\text{ALG}} &= \frac{(x \otimes x) \ominus (y \otimes y) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{(x^2(1 + \varepsilon_1) - y^2(1 + \varepsilon_2))(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{(x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2)(1 + \varepsilon_3) - (x^2 - y^2)}{x^2 - y^2} \\&= \frac{x^2 - x^2\varepsilon_1 - y^2 - y^2\varepsilon_2 + x^2\varepsilon_3 - x^2\varepsilon_1\varepsilon_3 - y^2\varepsilon_3 - y^2\varepsilon_2\varepsilon_3 - x^2 + y^2}{x^2 - y^2} \\&= \frac{x^2\varepsilon_1 - y^2\varepsilon_2 + x^2\varepsilon_3 - y^2\varepsilon_3 - x^2\varepsilon_1\varepsilon_3 - y^2\varepsilon_2\varepsilon_3}{x^2 - y^2} \\&\stackrel{\cdot}{=} \frac{x^2\varepsilon_1 - y^2\varepsilon_2 + (x^2 - y^2)\varepsilon_3}{x^2 - y^2} = \frac{(x^2 - y^2)\varepsilon_3}{x^2 - y^2} + \frac{x^2\varepsilon_1 - y^2\varepsilon_2}{x^2 - y^2} \\&= \varepsilon_3 + \frac{x^2\varepsilon_1 - y^2\varepsilon_2}{x^2 - y^2}\end{aligned}$$

Algorithmic error

$$\begin{cases} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{cases}$$

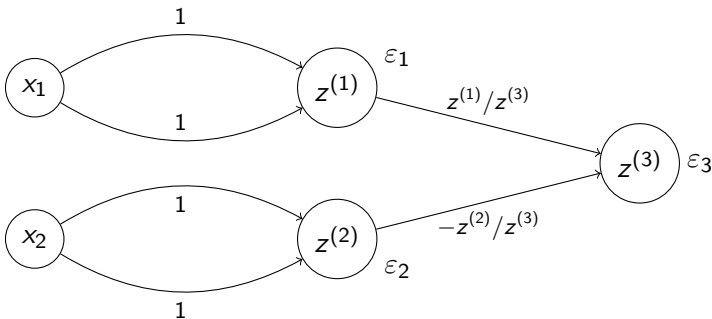
Algorithmic error

$$\begin{cases} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{cases}$$



Algorithmic error

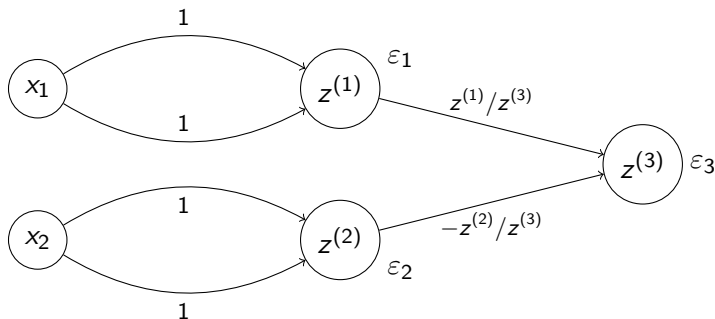
$$\begin{cases} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{cases}$$



$$\varepsilon_{ALG} \doteq \varepsilon_3 + \frac{z^{(1)}}{z^{(3)}} \varepsilon_1 - \frac{z^{(2)}}{z^{(3)}} \varepsilon_2 =$$

Algorithmic error

$$\begin{cases} z^{(1)} = x^2, & \varepsilon_1, \\ z^{(2)} = y^2, & \varepsilon_2, \\ f = z^{(3)} = z^{(1)} - z^{(2)}, & \varepsilon_3. \end{cases}$$



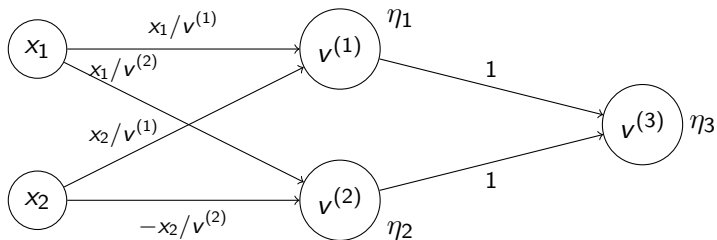
$$\varepsilon_{ALG} \doteq \varepsilon_3 + \frac{z^{(1)}}{z^{(3)}} \varepsilon_1 - \frac{z^{(2)}}{z^{(3)}} \varepsilon_2 = \varepsilon_3 + \frac{x^2 \varepsilon_1 - y^2 \varepsilon_2}{x^2 - y^2}$$

Algorithmic error

$$\left\{ \begin{array}{ll} v^{(1)} = x + y, & \eta_1, \\ v^{(2)} = x - y, & \eta_2, \\ f = v^{(3)} = v^{(1)}v^{(2)}, & \eta_3. \end{array} \right.$$

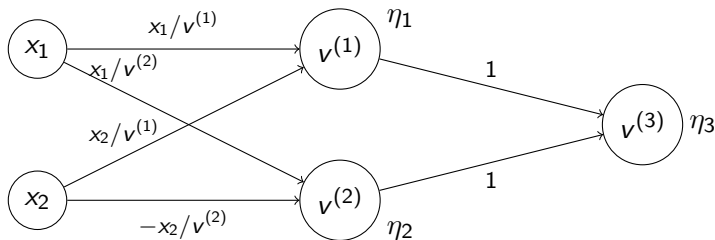
Algorithmic error

$$\begin{cases} v^{(1)} = x + y, & \eta_1, \\ v^{(2)} = x - y, & \eta_2, \\ f = v^{(3)} = v^{(1)}v^{(2)}, & \eta_3. \end{cases}$$



Algorithmic error

$$\begin{cases} v^{(1)} = x + y, & \eta_1, \\ v^{(2)} = x - y, & \eta_2, \\ f = v^{(3)} = v^{(1)}v^{(2)}, & \eta_3. \end{cases}$$



$$\varepsilon_{ALG} \doteq \eta_3 + \eta_1 + \eta_2$$

Algorithmic error

Which is the best?

$$\varepsilon_{ALG}^{(1)} \varepsilon_3 + \frac{x^2 \varepsilon_1 - y^2 \varepsilon_2}{x^2 - y^2} \quad \varepsilon_{ALG}^{(2)} \doteq \eta_3 + \eta_1 + \eta_2.$$

Algorithmic error

Which is the best?

$$\varepsilon_{ALG}^{(1)} \varepsilon_3 + \frac{x^2 \varepsilon_1 - y^2 \varepsilon_2}{x^2 - y^2} \quad \varepsilon_{ALG}^{(2)} \doteq \eta_3 + \eta_1 + \eta_2.$$

We can bound

$$|\varepsilon_{ALG}^{(1)}| \leq \left(1 + \frac{x^2 + y^2}{|x^2 - y^2|}\right) u, \quad |\varepsilon_{ALG}^{(2)}| \leq 3u.$$

The second is the best in most cases.