

## Exemple d'errors

2

Càlcul aproximat de la massa de la Terra

Usant la llei de la gravitació universal de Newton i la llei de la caiguda lliure dels cossos de Galileu, s'obté la fórmula:

$$M = \frac{gR^2}{G}, \tag{1}$$

on  $g$  és l'acceleració de la gravetat,  $R$  el radi de la Terra, i  $G$  la constant de la gravitació universal. Es disposa dels valors experimentals següents:

$$\overline{g} = 9.80665 \text{ m s}^{-2}, \quad \overline{G} = 6.67428 \cdot 10^{-11} \text{ m}^3\text{kg}^{-1}\text{s}^{-2}, \quad \overline{R} = 6371.0 \text{ km}.$$

Aplicant la fórmula anterior, resulta l'aproximació  $\overline{M} = 5.9639 \cdot 10^{24} \text{ kg}$ .

**Nota**  $M = 5.9736 \cdot 10^{24} \text{ kg}$  (Wikipedia, NASA).

$M = 5.9742 \cdot 10^{24} \text{ kg}$  (J.M.A. Danby, *Fundamentals of Celestial Mechanics*, Willmann-Bell, Inc., 1992).

Caldria estudiar els errors comesos atenent a les aproximacions donades dels valors experimentals.

## Fonts d'error

3

- **Errors de modelització**: els models matemàtics són aproximacions de la realitat.
- **Errors de truncament**: els mètodes numèrics fan aproximacions del model matemàtic.
- **Errors experimentals**: les mesures de les dades del problema no són exactes i porten errors de diversos tipus:
  - **errors aleatoris**: les mesures estan afectades per factors "aleatoris" que no podem controlar;
  - **errors sistemàtics**: les mesures provenen, per exemple, d'una calibració incorrecta de l'aparell de mesura;
  - **errors aberrants**: deguts a errors humans, a canvis sobtats en les condicions de l'experiment, etc.
- **Errors d'arrodoniment**: les operacions es realitzen amb un nombre finit de xifres (amb l'ajuda d'una calculadora o d'un ordinador).

## Error absolut i error relatiu

4

Definicions i notacions

Sigui  $x$  el valor exacte d'una quantitat i  $\bar{x}$  un valor aproximat.

- **Error absolut** en  $x$ :

$$e_a(x) := e_a(\bar{x}, x) = x - \bar{x},$$

- **Error relatiu** en  $x$ :

$$e_r(x) := e_r(\bar{x}, x) = \frac{e_a(x)}{x} \simeq \frac{e_a(x)}{\bar{x}} = \frac{x - \bar{x}}{\bar{x}}.$$

Fites d'error:

- $e_a(x) := e_a(\bar{x}, x)$  és una **fita de l'error absolut** en  $x$  si  $|e_a(x)| \leq \varepsilon_a(x)$ ,
- $e_r(x) := e_r(\bar{x}, x)$  és una **fita de l'error relatiu** en  $x$  si  $|e_r(x)| \leq \varepsilon_r(x)$ .

Notació usual:

- $x = \bar{x} \pm \varepsilon_a(x) \iff x \in [\bar{x} - \varepsilon_a(x), \bar{x} + \varepsilon_a(x)]$ ,
- $x = \bar{x} (1 \pm \varepsilon_r(x)) \iff x \in [\bar{x} - \varepsilon_r(x)|\bar{x}|, \bar{x} + \varepsilon_r(x)|\bar{x}|]$ .

## Error absolut i error relatiu

5

Exemples

Sigui  $a = \sqrt{20000} = 141.4213562\dots$  i  $\bar{a} = 141.4$ .

$$e_a(\bar{a}, a) = e_a(141.4, \sqrt{20000}) = 0.02135\dots < 0.022 \equiv \varepsilon_a(a),$$

$$e_r(\bar{a}, a) = e_r(141.4, \sqrt{20000}) \simeq \frac{0.02135}{141.4} = 0.00015099\dots < 0.00016 \equiv \varepsilon_r(a).$$

La primera fita indica que l'error no afecta el primer dígit fraccionari i la segona, que l'error no afecta el tercer dígit significatiu (tot i que tampoc el quart).

Sigui  $b = \sqrt{800000} = 894.42719\dots$  i  $\bar{b} = 894.4$ .

$$e_a(\bar{b}, b) = e_a(894.4, \sqrt{800000}) = 0.02719\dots < 0.028 \equiv \varepsilon_a(b),$$

$$e_r(\bar{b}, b) = e_r(894.4, \sqrt{800000}) \simeq \frac{0.02719}{894.4} = 0.000030399\dots < 0.000031 \equiv \varepsilon_r(b).$$

La primera fita indica que l'error no afecta el primer dígit fraccionari i la segona fita, que l'error no afecta el quart dígit significatiu.

## Representació de nombres en una base

6

Definició i exemples

La **representació en base  $b \geq 2$**  d'un nombre real  $x \neq 0$  és

$$\begin{aligned} x &= \pm a_{q-1}a_{q-2}\dots a_0.a_{-1}a_{-2}\dots b_{-q} \quad [q \text{ xifres abans del punt}] \\ &= \pm(a_{q-1}b^{q-1} + a_{q-2}b^{q-2} + \dots + a_0 + a_{-1}b^{-1} + \dots), \quad |x| \geq 1 \quad (q \geq 1) \\ x &= \pm 0.0\dots 0a_{q-1}a_{q-2}\dots b_{-q} \quad [-q \text{ zeros després del punt}] \\ &= \pm(a_{q-1}b^{q-1} + a_{q-2}b^{q-2} + \dots), \quad |x| < 1 \quad (q < 1) \end{aligned}$$

on els  $a_j$  ( $j < q$ ) són **xifres** en la base  $b$ :  $0 \leq a_j < b$  amb la primera xifra significativa  $a_{q-1} \neq 0$ . Quan  $b = 10$ , no cal indicar la base.

$$\begin{aligned} 107.125 &= 1 \cdot 10^2 + 0 \cdot 10^1 + 7 \cdot 10^0 + 1 \cdot 10^{-1} + 2 \cdot 10^{-2} + 5 \cdot 10^{-3} \\ 0.00333\dots &= 3 \cdot 10^{-3} + 3 \cdot 10^{-4} + 3 \cdot 10^{-5} + \dots \\ \pi &= 3 \cdot 10^0 + 1 \cdot 10^{-1} + 4 \cdot 10^{-2} + 1 \cdot 10^{-3} + 5 \cdot 10^{-4} + \dots \\ 0.1_2 &= 0.5 \\ 0.1_{10} &= 0.000\overline{1}_2 \end{aligned}$$

## Representació de nombres en una base

7

Exemple de representació en base 2

Representació de  $125.1_{10}$  en base 2.

$$125.1_{10} = a_{q-1}a_{q-2}\dots a_0.a_{-1}a_{-2}a_{-3}\dots_2.$$

amb els bits  $a_j \in \{0, 1\}$ , ( $j < q = 3$ ).

**Representació de la part entera  $125_{10}$  en base 2.**

$$\begin{aligned} 125/2 &= 62 \quad (\text{prenem residu } 1) \\ 62/2 &= 31 \quad (\text{prenem residu } 0) \\ 31/2 &= 15 \quad (\text{prenem residu } 1) \\ 15/2 &= 7 \quad (\text{prenem residu } 1) \\ 7/2 &= 3 \quad (\text{prenem residu } 1) \\ 3/2 &= 1 \quad (\text{prenem quocient } 1) \quad (\text{prenem residu } 1) \end{aligned}$$

$$125_{10} = 1111101_2.$$

## Representació de nombres8

Exemple de representació en base 2

Representació de la part fraccionària  $0.1_{10}$  en base 2.

$$\begin{aligned} 0.1 \times 2 &= 0.2 && \text{(prenem part entera 0)} \\ 0.2 \times 2 &= 0.4 && \text{(prenem part entera 0)} \\ 0.4 \times 2 &= 0.8 && \text{(prenem part entera 0)} \\ 0.8 \times 2 &= 1.6 && \text{(prenem part entera 1)} \\ 0.6 \times 2 &= 1.2 && \text{(prenem part entera 1)} \\ 0.2 \times 2 &= 0.4 && \text{(prenem part entera 0)} \\ &\dots && \end{aligned}$$

$$0.1_{10} = 0.0001\overline{1}_2$$

Representació binària completa

$$125.1_{10} = 1111101.0001\overline{1}_2 .$$

## Representació decimal en punt flotant9

Definició

Fent flotar el punt decimal  $q$  posicions en la representació decimal d'  $x \neq 0$ :

$$x = \pm 0.\alpha_1\alpha_2 \dots \cdot 10^q = \pm m \cdot 10^q, \quad \alpha_j = a_{q-j} \in \{0, 1, \dots, 9\} \quad (j > 0)$$

Aquesta **representació decimal en punt flotant** d'  $x$  ve donada per:

- el nombre enter  $q$ , anomenat **exponent**, i
- el nombre real  $m$ , tal que  $0.1 \leq m < 1$ , anomenat **mantissa**.

El primer dígit fraccionari  $\alpha_1 \neq 0$  d'  $m$  és el primer dígit significatiu d'  $x$ .

Exemples:

- $g = 9.80665 = 0.980665 \cdot 10^1$ , l'exponent és 1 i la mantissa, 0.980665.
- $G = 6.67428 \cdot 10^{-11} = 0.667428 \cdot 10^{-10}$ , l'exponent és -10 i la mantissa, 0.667428.

## Representació decimal en punt flotant10

Representació aproximada en calculadores

Les calculadores poden emmagatzemar una quantitat finita de dígits. Per tant, no es poden emmagatzemar **totes** les mantisses (ni tots els exponents).

Si sols disposem de  $t$  dígits per a la mantissa, la representació

$$x = \pm 0.\alpha_1\alpha_2 \dots \alpha_t\alpha_{t+1} \dots \cdot 10^q = \pm m \cdot 10^q, \text{ amb } \alpha_1 \neq 0,$$

pot ser aproximada, arrodonint el darrer dígit  $t$ , per  $\text{fl}_t(x)$ , flotant d'  $x$  amb  $t$  dígits significatius i arrodoniment:

- $\text{fl}_t(x) = \pm 0.\alpha_1\alpha_2 \dots \alpha_t \cdot 10^q$ , si  $\alpha_{t+1} < 5$ ;
- $\text{fl}_t(x) = \pm \text{fl}_t((0.\alpha_1\alpha_2 \dots \alpha_t + 10^{-t}) \cdot 10^q)$ , si  $\alpha_{t+1} \geq 5$ .

**Exemple:** Sigui  $x = 0.999527 \cdot 10^1$ . Llavors:

$$\text{fl}_5(x) = 0.99953 \cdot 10^1, \quad \text{fl}_4(x) = 0.9995 \cdot 10^1, \quad \text{fl}_3(x) = 0.100 \cdot 10^2$$

## Representació decimal en punt flotant11

Error d'arrodoniment

Observem que, atenent a les definicions:

$$|\epsilon_a(\text{fl}_t(x), x)| = |x - \text{fl}_t(x)| \leq \frac{1}{2} 10^{-t} 10^q = \frac{1}{2} 10^{q-t} =: \epsilon_a(\text{fl}_t(x), x) .$$

$\frac{1}{2} 10^{q-t}$  és una fita de l'**error absolut** en la **representació en punt flotant amb  $t$  dígits significatius i arrodoniment** de qualsevol nombre real  $x \neq 0$  amb exponent  $q$ .

Com que  $x \neq 0$  i  $m \geq 0.1$ :

$$|\epsilon_r(\text{fl}_t(x), x)| = \frac{|\epsilon_a(\text{fl}_t(x), x)|}{|x|} \leq \frac{1}{2} \frac{10^{q-t}}{m \cdot 10^q} \leq \frac{1}{2} 10^{1-t} =: \epsilon_r(\text{fl}_t(x), x) .$$

$\frac{1}{2} 10^{1-t}$  és una fita de l'**error relatiu** en la **representació en punt flotant amb  $t$  dígits significatius i arrodoniment** de qualsevol nombre real  $x \neq 0$ .

## Representació decimal en punt flotant12

Exemples d'errors d'arrodoniment

- $\text{fl}_6(g) = 0.980665 \cdot 10^1$ :  $t = 6$ ,  $q = 1$ .

$$\epsilon_a(g) = \frac{1}{2} 10^{1-6} = \frac{1}{2} 10^{-5}, \quad \epsilon_r(g) = \frac{1}{2} 10^{1-6} = \frac{1}{2} 10^{-5} .$$

- $\text{fl}_6(G) = 0.667428 \cdot 10^{-10}$ :  $t = 6$ ,  $q = -10$ .

$$\epsilon_a(G) = \frac{1}{2} 10^{-10-6} = \frac{1}{2} 10^{-16}, \quad \epsilon_r(G) = \frac{1}{2} 10^{1-6} = \frac{1}{2} 10^{-5} .$$

## Representació binària en punt flotant13

Representació aproximada en ordinadors i errors d'arrodoniment

Els ordinadors poden emmagatzemar una quantitat finita de bits per a les mantisses i exponents.

La representació en punt binari flotant d'un nombre qualsevol  $x \neq 0$

$$x = \pm 0.\alpha_1\alpha_2 \dots \alpha_t\alpha_{t+1} \dots \cdot 2^q = \pm m \cdot 2^q, \text{ amb } \alpha_1 = 1,$$

pot ser aproximada, emprant  $t$  bits per a la mantisa arrodonint el darrer bit, per  $\text{fl}_t(x)$  (flotant d'  $x$  amb  $t$  bits significatius i arrodoniment ):

- $\text{fl}_t(x) = \pm 0.\alpha_1\alpha_2 \dots \alpha_t \cdot 2^q$ , si  $\alpha_{t+1} = 0$ ;
- $\text{fl}_t(x) = \pm \text{fl}_t((0.\alpha_1\alpha_2 \dots \alpha_t + 2^{-t}) \cdot 2^q)$ , si  $\alpha_{t+1} = 1$ .

Com que  $x \neq 0$  i  $m < \frac{1}{2}$ ,

$$|\epsilon_r(\text{fl}_t(x), x)| = \frac{|\epsilon_a(\text{fl}_t(x), x)|}{|x|} \leq \frac{1}{2} \frac{2^{q-t}}{m \cdot 2^q} \leq \frac{1}{2} 2^{1-t} =: \epsilon_r(\text{fl}_t(x), x)$$

i, per tant,  $2^{-t}$  és una fita de l'**error relatiu** de la **representació en punt flotant amb  $t$  bits significatius i arrodoniment** de qualsevol nombre real  $x \neq 0$ .

Representació binària en punt flotant14

Formats IEEE de representació en precisió simple i doble

$$x = \pm 1.\alpha_2 \dots \alpha_t \alpha_{t+1} \dots_2 \cdot 2^{q-1} = \pm (1 + f) 2^{q-1}$$

es representa per:

- $\text{fl}_t(x) = \pm 1.\alpha_2 \alpha_3 \dots \alpha_t_2 \cdot 2^{q-1}$ , si  $\alpha_{t+1} = 0$ ;
- $\text{fl}_t(x) = \pm \text{fl}_t((1.\alpha_2 \dots \alpha_t_2 + 2^{-t+1})2^{q-1})$ , si  $\alpha_{t+1} = 1$ .

Format	base (b)	digits (t)	$q_{\min} - 1$	$q_{\max} - 1$	bits
IEEE simple	2	24	-126	128	32
IEEE doble	2	53	-1022	1024	64

Taula: Formats IEEE (simple i doble precisió)

IEEE simple	signe (1)	$e = q - 1 + 127$ (8)	mantissa f (23)
IEEE doble	signe (1)	$e = q - 1 + 1023$ (11)	mantissa f (52)

Taula: Distribució de memòria en el format IEEE (simple i doble).

Representació binària en punt flotant15

Formats IEEE de representació en precisió simple i doble

- Es guarden els  $t - 1$  primers bits de la mantissa  $f$ , ja que el primer bit de la matissa  $m$  és 1.
- Si un nombre real  $x$  es pot escriure **exactament**, es diu que és un **nombre de màquina**. Altament tindrà una representació en punt flotant  $\text{fl}_t(x) \neq x$  amb error.
- En **precisió simple**, la fita de l'error relatiu d'arroundiment és  $2^{-24} \approx 0.6 \cdot 10^{-7}$ , es pot garantir gairebé una **precisió de 6 dígits significatius amb arrodoniment**.
- En **precisió doble**, la fita d'error relatiu d'arrodoniment és  $2^{-53} \approx 1.1 \cdot 10^{-16}$ , es pot garantir gairebé una **precisió de 16 dígits significatius**.
- En IEEE simple, els valors  $e = 0$  ( $q = -126$ ) i  $e = 255$  ( $q = 127$ ) es reserven a **NaN** (Not a Number) i **overflow**, respectivament.
- En IEEE doble, els valors  $e = 0$  ( $q = -1022$ ) i  $e = 1023$  ( $q = 1023$ ) es reserven a **NaN** (Not a Number) i **overflow**, respectivament.

Representació binària en punt flotant16

Formats IEEE de representació en precisió simple i doble

- Com es representa  $x = 125.1$  en format IEEE amb precisió simple?

Es té la representació amb punt binari flotant:

$$x = 125.1 = 1111101.00011_2 = 1.11110100011_2 \cdot 2^6$$

L'exponent desplaçat  $e$  es representa en base 2 per

$$e = 6 + 127 = 133 = 10000101_2.$$

Resulta finalment la representació en memòria usant IEEE simple:

IEEE simple	0	10000101	11110100011001100110011
-------------	---	----------	-------------------------

Representació binària en punt flotant17

Èpsilon de la màquina

- $\epsilon = \frac{1}{2}b^{1-t}$  s'anomena **èpsilon de la màquina**. Coincideix amb el nombre positiu més petit que sumat a 1 dóna diferent de 1, és a dir

$$\epsilon = \min\{\varepsilon : \text{fl}_t(1 + \varepsilon) \neq 1\}.$$

Com més petit és, més precisa és la màquina: la precisió indica el nombre  $t$  de xifres significatives correctament representades amb arrodoniment en base  $b$ .

- Notem que  $\text{fl}_t(1 + \varepsilon) = 1$  no vol dir  $\varepsilon$  sigui igual a 0 sinó que és més petit que l'èpsilon de la màquina.
- Treballant amb  $t = 3$  dígits significatius, si

$$x = 0.1 \cdot 10^1 \quad \text{i} \quad y = 0.456 \cdot 10^{-4}$$

llavors  $x + y = x$ , però  $y \neq 0$ .

Problemes numèrics18

Operacions aritmètiques usant representació flotant

Les calculadores i els ordinadors, **degut a la representació dels nombres en punt flotant**, fan els càlculs de manera aproximada.

Això té implicacions importants: **l'ordre de les operacions pot afectar el resultat final!**.

**Exemple:** Treballant amb  $t = 4$  dígits, si es calcula  $a + b + c$ , on

$$a = 0.5317 \cdot 10^{-2}, \quad b = 0.3387 \cdot 10^2, \quad c = -0.3381 \cdot 10^2,$$

emprant ordenacions diferents, resulta:

$$\text{fl}_4(a + \text{fl}_4(b + c)) = \text{fl}_4(0.5317 \cdot 10^{-2} + 0.6000 \cdot 10^{-1}) = 0.6532 \cdot 10^{-1}$$

$$\text{fl}_4(\text{fl}_4(a + b) + c) = \text{fl}_4(0.3388 \cdot 10^2 - 0.3381 \cdot 10^2) = 0.7000 \cdot 10^{-1}.$$

El resultat exacte és  $a + b + c = 0.65317 \cdot 10^{-1}$  ens diu que és millor la primera ordenació.

Problemes numèrics19

Exemple de cancel·lació

Les dues solucions de l'equació  $x^2 - 18x + 1 = 0$  són

$$x_{1,2} = 9 \pm \sqrt{80} = \begin{cases} x_1 = 0.1794427190999916 \cdot 10^2 \\ x_2 = 0.5572809000084121 \cdot 10^{-1} \end{cases}$$

Si prenem  $\sqrt{80} = 8.9443$  (és a dir  $t = 5$ ) s'obté

$$x_1 = 9 + 8.9443 = 17.9443 = 0.179443 \cdot 10^2 \quad (6 \text{ xifres}),$$

$$x_2 = 9 - 8.9443 = 0.0557 = 0.557 \cdot 10^{-1} \quad (3 \text{ xifres!}).$$

En calcular  $x_2$  hi ha una **cancel·lació** de dígits, perquè restem dues quantitats que són properes i dóna un resultat **significativament erroni**.

## Propagació d'errors

20

### Causas

Hi ha dues raons (o almenys així es pot pensar) **responsables** de la propagació de l'error en un procés de càlcul:

- **Errors en les dades.** Si les dades tenen error, aquest error es propaga al resultat de les operacions.
- **Errors en les operacions** Encara que se sumin dos nombres de màquina  $x, y$ , el resultat representat  $\text{fl}_t(x + y)$  pot ser diferent de  $x + y$ . Les funcions  $f$  internes que s'apliquen tenen també errors que es poden considerar errors en les operacions.
- **Les dues alhora...** Efectivament, els errors de les dades i de les operacions s'acumulen en els resultats intermedis i en el resultat final.

**Per simplificar** se suposa que les operacions no tenen errors i que, per tant, tots els errors es deuen a la propagació dels errors de les dades.

## Fórmula de propagació d'errors

21

### Funcions d'una variable

**Teorema del valor mig:** Sigui  $f : [a, b] \rightarrow \mathbb{R}$  una funció contínua, derivable a  $]a, b[$ . Aleshores existeix un punt  $\xi \in (a, b)$ , tal que

$$f(b) - f(a) = f'(\xi)(b - a).$$

**Aplicació: Propagació d'errors en funcions d'una variable.**  
Sigui  $x \in \mathbb{R}$  i sigui  $\bar{x} \approx x$ .

Del teorema anterior, es té que

$$e_a(f(\bar{x}), f(x)) := f(x) - f(\bar{x}) = f'(\xi)(x - \bar{x}), \quad \xi \in ]\bar{x}, x[.$$

Usant que la funció és contínua i suposant que els errors són petits, es té una **fórmula aproximada de propagació de l'error**:

$$|e_a(f(\bar{x}), f(x))| \approx |f'(\bar{x})| |e_a(\bar{x}, x)|, \\ \varepsilon_a(f(\bar{x}), f(x)) := M \varepsilon_a(\bar{x}, x), \quad \text{on} \quad M = \max_{\xi \in [\bar{x} - \varepsilon_a, \bar{x} + \varepsilon_a]} |f'(\xi)|.$$

## Fórmula de propagació d'errors

22

**Error relatiu:** Coeficient de propagació  
De les darreres expressions, resulta

$$|e_r(f(\bar{x}), f(x))| \approx |\bar{x}| \frac{|f'(\bar{x})|}{|f(\bar{x})|} |e_r(\bar{x}, x)|, \quad (f(x) \neq 0).$$

De fet, el terme  $\varphi(x) = |x| \frac{|f'(x)|}{|f(x)|}$  s'anomena **coeficient de propagació (de l'error relatiu)** i caldria controlar-lo en un procés de càlcul. Si podem fitar-lo, és a dir, si podem dir que

$$|x| \frac{|f'(x)|}{|f(x)|} \leq M$$

per alguna  $M > 0$ , en un entorn de  $\bar{x}$ , llavors

$$\varepsilon_r(f(\bar{x}), f(x)) := M \varepsilon_r(\bar{x}, x).$$

## Fórmula de propagació d'errors

23

### Exemples d'aplicació

Càlcul de les arrels de l'equació  $x^2 - 18x + 1 = 0$ :  $x_{1,2} = 9 \pm \sqrt{80}$ , ara calculant  $x_1$  amb 4 dígits fraccionaris amb arrodoniment:

$$x_1 = 17.9443 \pm \frac{1}{2} \cdot 10^{-4}$$

i després calculant  $x_2$  així:

$$x_2 = f(x_1) = \frac{1}{x_1} \quad \mapsto \quad \bar{x}_2 = f(\bar{x}_1) = \frac{1}{\bar{x}_1} = 0.05572800...$$

### Estimació de les fites dels errors

$$\varepsilon_a(\bar{x}_2, x_2) = \varepsilon_a\left(\frac{1}{\bar{x}_1}, \frac{1}{x_1}\right) \simeq \left| \frac{-1}{x_1^2} \right| \varepsilon_a(\bar{x}_1, x_1) \simeq \frac{1}{17.9443^2} \frac{1}{2} \cdot 10^{-4} \simeq 0.16 \cdot 10^{-6}$$

$$\varepsilon_r(\bar{x}_2, x_2) = \frac{\varepsilon_a(\bar{x}_2, x_2)}{|\bar{x}_2|} \simeq 17.9443 \cdot 0.16 \cdot 10^{-6} \simeq 0.29 \cdot 10^{-5}$$

Així,  $\varepsilon_a(\bar{x}_2, x_2) \leq 0.5 \cdot 10^{-6}$  indica que  $\bar{x}_2$  té 6 xifres decimals correctes i  $\varepsilon_r(\bar{x}_2, x_2) \leq 0.5 \cdot 10^{-5}$  indica que  $\bar{x}_2$  té 5 xifres significatives correctes.  
Hi ha una millora respecte al càlcul anterior de  $x_2$ : s'han guanyat dues xifres correctes en el resultat evitant els efectes de cancel·lació.

## Fórmula de propagació d'errors

24

### Exemples d'aplicació

Fita de l'error comès en avaluar  $f(x) = \ln \cos^2(x)$  en un punt  $x$  del qual només **conexim tres dígits correctes**  $\bar{x} = 0.735$ .

- **Valor aproximat:**  $\ln \cos^2(\bar{x}) = -0.5972683...$
- **Fita de l'error en  $x$ :**  $\varepsilon_a(x) = \frac{1}{2} 10^{-3}$ .
- **Derivada de la funció:**  $f'(x) = -2 \tan(x)$ .
- **Fita del valor absolut de la derivada:**  $|f'(\xi)|$  per a  $\xi \in [0.7345, 0.7355]$  (i.e.,  $[\bar{x} - \varepsilon_a, \bar{x} + \varepsilon_a]$ ): com que la funció tangent és creixent i positiva a tot l'interval  $]0, \pi/2[$ , llavors  $\tan(\xi) \leq \tan(0.7355) \lesssim 0.905$ , i

$$|f'(\xi)| \leq 2 \cdot 0.905 = 1.810.$$

- **Fita de l'error absolut en  $f(x)$ :** aplicant la fórmula de propagació de l'error:

$$\varepsilon_a(f(\bar{x} = 0.735), f(x)) = 1.810 \frac{1}{2} 10^{-3} = 0.905 \cdot 10^{-3}.$$

$$f(x) = -0.5972683 \pm 0.000905.$$

## Fórmula de propagació d'errors

25

### Exemple 2

El **coeficient de propagació de l'error relatiu**  $\varphi$  resulta ser

$$\varphi(x) = x \frac{f'(x)}{f(x)} = x \frac{-2 \tan(x)}{\log \cos^2(x)}.$$

Es té que  $|\varphi(x)| \leq 2.2$  si  $x \approx 0.735$  i per tant

$$\varepsilon_r(f(0.735)) \leq 2.2 \varepsilon_r(\bar{x}, x) \leq 2.2 \frac{0.5 \cdot 10^{-3}}{0.735} \approx 1.5 \cdot 10^{-3}. \\ f(x) = -0.5972683(1 \pm 0.0015).$$

Fórmula de propagació d'errors26

Propagació d'errors en diverses variables27

Funcions de diverses variable

**Teorema del valor mig en diverses variables.**

Sigui  $G$  un obert de  $\mathbb{R}^n$ , i  $f : G \rightarrow \mathbb{R}$  una funció diferenciable sobre  $G$ . Siguin  $x = (x_1, \dots, x_n)$ ,  $y = (y_1, \dots, y_n)$  dos punts de  $G \subset \mathbb{R}^n$  tals que el segment que els uneix està contingut a  $G$ . Aleshores existeix un punt  $\xi$  d'aquest segment tal que

$$f(y) - f(x) = \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\xi)(y_i - x_i).$$

**Aplicació: Propagació d'errors en funcions de diverses variables**

Sigui  $x \in \mathbb{G} \subset \mathbb{R}^n$ , sigui  $\bar{x} \approx x$ .

$$e_a(f(\bar{x}), f(x)) \approx \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\bar{x})e_a(\bar{x}_i, x_i)$$

$$\varepsilon_a(f(\bar{x}), f(x)) := \sum_{i=1}^n M_i \varepsilon_a(\bar{x}_i, x_i), \quad \text{on} \quad M_i = \max_{\xi \in [\bar{x} - \varepsilon_a, \bar{x} + \varepsilon_a]^n} \left| \frac{\partial f}{\partial x_i}(\xi) \right|$$

Casos especials

$$f(x, y) = x + y$$

$$\varepsilon_a(\bar{x} + \bar{y}, x + y) = \varepsilon_a(\bar{x}, x) + \varepsilon_a(\bar{y}, y)$$

$$\varepsilon_r(\bar{x} + \bar{y}, x + y) = \left| \frac{x}{x + y} \right| \varepsilon_r(\bar{x}, x) + \left| \frac{y}{x + y} \right| \varepsilon_r(\bar{y}, y)$$

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

$$\varepsilon_a(\bar{x}\bar{y}, xy) \approx |y| \varepsilon_a(\bar{x}, x) + |x| \varepsilon_a(\bar{y}, y)$$

$$\varepsilon_r(\bar{x}\bar{y}, xy) \approx \varepsilon_r(\bar{x}, x) + \varepsilon_r(\bar{y}, y)$$

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

$$\varepsilon_a(\bar{x}/\bar{y}, x/y) \approx 1/|y| \varepsilon_a(\bar{x}, x) + \left| \frac{x}{y^2} \right| \varepsilon_a(\bar{y}, y)$$

$$\varepsilon_r(\bar{x}/\bar{y}, x/y) \approx \varepsilon_r(\bar{x}, x) + \varepsilon_r(\bar{y}, y)$$

Propagació d'errors en diverses variables28

Propagació d'errors en diverses variables29

Exemple de càlcul de la massa de la Terra

Error propagat en  $M(g, G, R) = \frac{gR^2}{G}$ , a partir de les aproximacions de les magnituds:

$$\bar{g} = 9.80665, \quad \bar{G} = 6.67428 \cdot 10^{-11}, \quad \bar{R} = 6371.0 \cdot 10^3.$$

- Errors en les magnituds suposant que són correctes fins a l'última xifra amb arrodoniment:

$$\varepsilon_a(g) = \frac{1}{2} 10^{-5}, \quad \varepsilon_a(G) = \frac{1}{2} 10^{-16}, \quad \varepsilon_a(R) = \frac{1}{2} 10^2.$$

- Fórmula de propagació d'errors en diverses variables:

$$\varepsilon_a(M) \approx \left| \frac{\partial M}{\partial g} \right|(\bar{g}, \bar{G}, \bar{R})\varepsilon_a(g) + \left| \frac{\partial M}{\partial G} \right|(\bar{g}, \bar{G}, \bar{R})\varepsilon_a(G) + \left| \frac{\partial M}{\partial R} \right|(\bar{g}, \bar{G}, \bar{R})\varepsilon_a(R),$$

$$\frac{\partial M}{\partial g} = \frac{R^2}{G}, \quad \frac{\partial M}{\partial G} = -\frac{gR^2}{G^2}, \quad \frac{\partial M}{\partial R} = \frac{2gR}{G}.$$

Exemple de càlcul de la massa de la Terra

- Derivades parcials de  $M(g, G, R)$  en les aproximacions:

$$\frac{\partial M}{\partial g}(\bar{g}, \bar{G}, \bar{R}) \approx 6.0815 \cdot 10^{23},$$

$$\frac{\partial M}{\partial G}(\bar{g}, \bar{G}, \bar{R}) \approx -8.9357 \cdot 10^{34},$$

$$\frac{\partial M}{\partial R}(\bar{g}, \bar{G}, \bar{R}) \approx 1.8722 \cdot 10^{18}.$$

- Fita aproximada de l'error en  $M$ :

$$\varepsilon_a(M) \approx 1.0112 \cdot 10^{20}.$$

$$M \approx 5.96391525 \cdot 10^{24} \pm 1.0112 \cdot 10^{20}$$

$$\iff M \in [5.96381413 \cdot 10^{24}, 5.96401637 \cdot 10^{24}].$$

Propagació dels errors amb aritmètica intervalar30

Mètodes estables i inestables31

Exemple de càlcul de la massa de la Terra

**L'aritmètica intervalar** té una visió diferent a la del teorema del valor mig.

$$\bullet \quad g = 9.80665, \varepsilon_a(g) = \frac{1}{2} 10^{-5} : g \in [9.806645, 9.806655] = [g_m, g_M];$$

$$\bullet \quad G = 6.67428 \cdot 10^{-11}, \\ \varepsilon_a(G) = \frac{1}{2} 10^{-16} : G \in [6.674275 \cdot 10^{-11}, 6.674285 \cdot 10^{-11}] = [G_m, G_M];$$

$$\bullet \quad R = 6371.0 \cdot 10^3, \\ \varepsilon_a(R) = \frac{1}{2} 10^2 : R \in [6.37095 \cdot 10^6, 6.37105 \cdot 10^6] = [R_m, R_M].$$

$$M \in \left[ \frac{g_m R_m^2}{G_M}, \frac{g_M R_M^2}{G_m} \right] \rightarrow M \in [5.9638141 \cdot 10^{24}, 5.9640164 \cdot 10^{24}]$$

coincideix pràcticament amb l'anterior ja que els errors són petits

$$M \in [5.9638143 \cdot 10^{24}, 5.96401637 \cdot 10^{24}].$$

Exemple de recurrència inestable

Es volen calcular les integrals

$$R_n = \int_0^1 x^n e^{x-1} \, dx.$$

S'observa que:

- $R_0 = 1 - \exp(-1)$ .
- $0 < R_n < \frac{1}{n+1}$  per a tot  $n \geq 0$ .
- $R_n = 1 - nR_{n-1}$  (integració per parts).

Es proposa el **mètode de càlcul recurrent** següent per al càlcul d'un  $R_N$ :

- $R_0 = 1 - \exp(-1)$ .
- $R_n = 1 - nR_{n-1}$ ,  $n = 0, \dots, N$ .

## Mètodes estables i inestables

32

Exemple de recurrència inestable

$n$	$R_n$
1	3.678794411714423340e-01
2	2.642411176571153320e-01
3	2.072766470286540041e-01
4	1.708934118853839834e-01
5	1.455329405730800829e-01
10	8.387707005829270202e-02
15	5.903379364190186607e-02
18	-2.945367075153626502e-02

**Taula:**  $R_{18} < 0$  no té cap xifra significativa correcta!. Tots els càlculs s'han fet usant format de dades double en llenguatge C.

## Mètodes estables i inestables

33

Exemple de recurrència inestable

### Anàlisi dels errors

Si  $e_0 = R_0 - \bar{R}_0$  l'error absolut inicial en la dada  $R_0$ ,

$$e_n = R_n \bar{R}_n - R_n = 1 - nR_{n-1} - 1 + n\bar{R}_{n-1} = -n(R_{n-1} - \bar{R}_{n-1}) = -ne_{n-1},$$

L'error de  $R_N$ ,

$$e_N = (-1)^N N! e_0,$$

es fa molt gran quan  $N$  augmenta, **independentment** d' $e_0$ .  
El mètode recurrent és inestable i no permet el càlcul de les integrals per a  $N$  gran.

## Algorismes estables i inestables

34

Exemple de recurrència estable

Capgirant la recurrència, es té la recurrència inversa:

$$R_n = 1 - nR_{n-1} \rightarrow R_{n-1} = \frac{1 - R_n}{n}$$

Per calcular un  $R_N$ , es proposa utilitzar la recurrència inversa a partir d'un  $R_M$  apropiat:

- $R_M = 0$ .
- $R_{n-1} = \frac{1 - R_n}{n}$ ,  $n = M, \dots, N + 1$ .

Anàlisi de l'error propagat des de  $R_M$  fins a  $R_N$ .

$$e_{n-1} = -\frac{1}{n}e_n \Rightarrow e_N = (-1)^{M-N} \frac{1}{M(M-1) \cdots (N+1)} e_M.$$

Com que l'error inicial és  $e_M < \frac{1}{M+1}$ , la recurrència inversa troba  $R_N$  amb un error de propagació, que es fa més petit a cada pas, i que es pot fitar per:

$$|e_N| < \frac{1}{(M+1)M(M-1) \cdots (N+1)}.$$

## Mètodes estables i inestables

35

Exemple de recurrència estable

$n$	$R_n$	$ e_n $
40	0	2.3e-2
35	2.704628971076339372e-02	2.9e-10
30	3.127967393216807279e-02	7.4e-18
25	3.708621442373923743e-02	4.3e-25
20	4.554488407581805398e-02	6.8e-32
18	5.011985495809425512e-02	1.83-34
15	5.901754087929777376e-02	3.6e-38
10	8.387707010339416625e-02	1.02e-43

**Taula:** Els càlculs fets amb un mètode estable (recurrència inversa).

## Problemes mal condicionats

36

Són problemes on la solució depèn de manera **molt sensible de les dades**.

### Exemple:

El sistema d'equacions

$$\begin{aligned} 2.0000x + 0.6667y &= 2.6667, \\ 1.0000x + 0.3333y &= 1.3333, \end{aligned}$$

té solució  $x = 1.0000$ ,  $y = 1.0000$ , mentre que el sistema

$$\begin{aligned} 2.0000x + 0.6665y &= 2.6667, \\ 1.0000x + 0.3333y &= 1.3333, \end{aligned}$$

té solució  $x = 1.6666$ ,  $y = -1.0000$ .