

PART-1 → Computer arithmetic and errors

In some situations, when working with a PC, we may face the somehow embarrassing situation where

$$1 + \epsilon = 1$$

and $\epsilon \neq 0$. This is because not all numbers can be represented exactly in a computer, and arithmetic mistakes are a consequence. The information in a **single-precision** PC is built upon vectors composed by “0/1”32 bits (4 bytes). A large number of integers can be thus represented exactly, in base 2.

2^{31}	2^{30}					2^3	2^2	2^1	2^0		
0	0	0	0	\dots	\dots	0	0	1	0	1	0

The situation with real numbers $x \neq 0$ is more complicated. The representation of x is structured as

$$x = \pm n \cdot b^e,$$

where n is the **mantissa**, b is the **base** (always 2) and e is the exponent. This representation is also called **floating point**.



It is important to notice that, no matter the number of bits employed, a limited bits are assigned for e and n , therefore not all real numbers can be represented.

NUMERI MACCHINA := sottoinsieme di \mathbb{R} rappresentabile in aritmetica di macchina := tutti quanti i reali rappresentabili con il th. di rappresentazione in base

$$Y(\beta, t, m, M) = \{0\} \cup \left\{ x \in \mathbb{R} : x = \text{sgn}(x) \cdot \beta^p \sum_{i=1}^t d_i \beta^{-i} \right\} \quad \text{con} \quad \begin{matrix} -m \leq p \leq M \\ 0 \leq d_i \leq \beta^{-1} \end{matrix}, d_1 \neq 0$$

TH. DI RAPPRESENTAZIONE IN BASE

Per ogni $x \in \mathbb{R}$, $x \neq 0$, $\exists!$ $p \in \mathbb{Z}$ e $\{d_i\}$ successione di interi tali che:

- (i) $d_i \neq 0$
(ii) $0 \leq d_i \leq \beta - 1$
(iii) $\forall k > 0 \exists j \geq k$ tale che $d_j \neq \beta - 1$
per cui $x = \text{sgn}(x) \cdot \beta^p \sum_{i=1}^{+\infty} d_i \beta^{-i}$

Prop: $0.\bar{9} = 1$
 Dim: $0.\bar{9} = \sum_{i=1}^{+\infty} 9 \cdot 10^{-i} = 9 \cdot 10^{-1} + 9 \cdot 10^{-2} + \dots =$
 $= \frac{9}{10} \cdot \left(1 + \frac{1}{10} + \frac{1}{10^2} + \dots \right) =$
 questa è la serie geometrica $\sum_{n=0}^{+\infty} q^n$
 di ragione $q = \frac{1}{10}$ che converge a $\frac{1}{1-q}$
 $\sum_{n=0}^{+\infty} \left(\frac{1}{10}\right)^n = \frac{1}{1-\frac{1}{10}} = \frac{1}{\frac{9}{10}} = \frac{10}{9}$
 $= \frac{9}{10} \cdot \frac{10}{9} = 1 \quad \square$

CARDINALITA' INSIEME NUMERI MACCHINA

$$Y(\beta, t, m, M) = \{0\} \cup \left\{ \text{sgn}(x) \cdot \beta^p \sum_{i=1}^t d_i \beta^{-i}, -m \leq p \leq M, d_i \neq 0 \right\}$$

$+1$ $x_2(+,-)$ /

$i \beta^p$ Sono tanti
quanti i p :

da $-m$ a -1 $(=m)$
+ da 1 a M $(=M)$
+ più 60 $(=1)$
 $= m + M + 1$

Conto tutte le configurazioni possibili ($= \beta^t$)
poi ci tolgo quelle che iniziano con 0 ($= \beta^{t-1}$)
quindi $\beta^t - \beta^{t-1} = \beta^t (1 - \beta^{-1})$

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

MASSIMO (Ω) & MINIMO (ω) INSIEME NUMERI MACCHINA

$$\begin{aligned}\Omega &= \beta^H (0, d_1, d_2, \dots, d_t) = \\ &= \beta^H ((\beta-1)\beta^{-1} + (\beta-1)\beta^{-2} + \dots + (\beta-1)\beta^{-t}) = \\ &= \beta^H (\beta-1) (\beta^{-1} + \beta^{-2} + \dots + \beta^{-t}) = \\ &= \beta^H (\beta-1) \left(-1 + \sum_{n=0}^t \left(\frac{1}{\beta} \right)^n \right) =\end{aligned}$$

serie geometrica di ragione $q = \frac{1}{\beta}$
troncata a t :

$$S = \sum_{n=0}^t q^n = 1 + q + q^2 + \dots + q^t$$

$$S - qS = S(1-q) = (1 + \cancel{q} + \cancel{q^2} + \dots + \cancel{q^t}) - (\cancel{q} + \cancel{q^2} + \dots + \cancel{q^t} + q^{t+1}) = 1 + q^{t+1}$$

$$S = \frac{1 - q^{t+1}}{1 - q}$$

$$= \beta^H (\beta-1) \left(\frac{1 - \left(\frac{1}{\beta} \right)^{t+1}}{1 - \frac{1}{\beta}} - 1 \right) =$$

$$= \beta^H (\beta-1) \left(\frac{\beta^{t+1} - 1}{\beta^{t+1}} \cdot \frac{\beta}{\beta-1} - 1 \right) =$$

$$= \beta^H (\cancel{\beta-1}) \frac{\beta^{t+1} - 1 - \beta^t(\beta-1)}{\beta^t(\beta-1)} =$$

$$= \beta^H \frac{\beta^t(\cancel{\beta} - \cancel{\beta} + 1) - 1}{\beta^t} =$$

$$= \beta^H \frac{\beta^t - 1}{\beta^t} = \beta^H (1 - \beta^{-t})$$

$$\omega = \beta^{-m} (0, 1) =$$

$$= \beta^{-m} (1 \cdot \beta^{-1}) =$$

$$= \beta^{-(m+1)}$$

TRONCAMENTO & ARROTONDAMENTO

$x \in \mathbb{R} \rightarrow \tilde{x} \in \mathcal{M}$ [ovvero da un numero reale x devo passare ad un numero di macchina \tilde{x}]

$$x = \underbrace{\text{sgn}(x)}_{t, -} \cdot \underbrace{\beta^p}_{10^p} \cdot \underbrace{\sum_{i=1}^{\infty} d_i \beta^{-i}}_{0.d_1 d_2 d_3 \dots}$$

$$\tilde{x} = \text{trunc}(x) = \text{sgn}(x) \cdot \beta^p \cdot \sum_{i=1}^t d_i \beta^{-i} \quad [\text{limite (tronco) la } \Sigma \text{ al } t\text{-esimo}]$$

$$\tilde{x} = \text{arr}(x) = \text{sgn}(x) \cdot \text{trunc} \left(\beta^p \cdot \left(\sum_{i=1}^{t+1} d_i \beta^{-i} \right) + \frac{1}{2} \beta^{-t} \right) =$$

$$= \begin{cases} \text{trunc}(x) & \text{se } d_{t+1} < \frac{1}{2} \beta \\ \text{trunc}(x) + \beta^{p-t} & \text{se } d_{t+1} \geq \frac{1}{2} \beta \end{cases}$$

$$\left[\begin{array}{l} \tilde{x} = 10^3 \cdot 0.\overset{1}{2}\overset{3}{3}\overset{4}{4}\overset{5}{5}\overset{6}{6}29\dots = \\ p=3 \quad t=5 \\ = \text{trunc}(x) + \beta^{p-t} = \\ = 10^3 \cdot 0.37562 + 10^{-2} = \\ = 375.62 + 0.01 \\ = 375.63 \end{array} \right]$$

Nota: mentre con il troncamento sto sempre dentro l'intervallo $[w, \Omega]$, con l'arrotondamento potrei andare oltre $\Omega \Rightarrow \text{OVERFLOW}$

OVERFLOW := $p > M$
UNDERFLOW := $p < -m$

ERRORI DI RAPPRESENTAZIONE

Errore assoluto := $\tilde{x} - x$

Errore relativo := $\frac{\tilde{x} - x}{x} \rightarrow$ chiamato anche Errore di rappresentazione E_x in $\mathcal{M}(\beta, t, m, M)$
Quindi studiare il problema significa andare a vedere come questo errore si propaga all'interno dell'algoritmo.

Nota: se voglio rappresentare l'errore sulle cifre della mantissa ($d_1 d_2 d_3 \dots$), l'errore assoluto mi dipende da β^p (ovvero dalla grandezza del numero), mentre l'errore relativo no, quindi a parità di mantissa ottengo lo stesso errore

ESEMPIO: i numeri di macchina $\mathcal{M}(2, 3, 1, 1)$ sono una buona definizione perché gli errori sono contenuti

base 2
min e max dell'exp.
(in questo caso sempre = 1)
arrivo fino alla 3ª cifra di mantissa

\Rightarrow $2 \cdot 0.d_1 d_2 d_3$ dentro questa \mathcal{M} , tutti i numeri sono rappresentati in questo modo

TH. Se non si verifica overflow o underflow allora:

$$|\text{trunc}(x) - x| < \beta^{p-t}$$

$$|\text{arr}(x) - x| \leq \frac{1}{2} \beta^{p-t}$$

DIM. Prendo $a, b \in \mathcal{M}$ consecutivi:

$$a = \beta^p \cdot \sum_{i=1}^t d_i \beta^{-i}$$

$$b = \beta^p \cdot \left(\sum_{i=1}^t d_i \beta^{-i} \right) + \beta^{-t} \quad (\text{ho semplicemente sommato 1 alla cifra } t\text{-esima})$$

Quindi se ho $a \leq x < b$, indipendentemente da quanto sono vicino a b , casco sempre in a :

$$(i) \text{ trunc}(x) = a \Rightarrow |\text{trunc}(x) - x| = |a - x| = |a - x| < |a - b| = \beta^{p-t}$$

$$(ii) \text{ arr}(x) = \begin{cases} a & \text{se } x < \frac{a+b}{2} \\ b & \text{se } x \geq \frac{a+b}{2} \end{cases} \Rightarrow |\text{arr}(x) - x| \leq \frac{b-a}{2} = \frac{1}{2} \beta^{p-t} \quad (\text{l'uguaglianza vale solo se } x = \frac{a+b}{2} \Rightarrow \text{Sono a met  intervallo})$$

□

TH. (limitazione dell'errore di rappresentazione)

se non si verifica overflow o underflow, allora la quantit  μ detta precisione di macchina   cos  definita:

$$\left| \frac{\tilde{x} - x}{x} \right| < \mu = \begin{cases} \beta^{1-t} & \text{se } \tilde{x} = \text{trunc}(x) \\ \frac{1}{2} \beta^{1-t} & \text{se } \tilde{x} = \text{arr}(x) \end{cases}$$

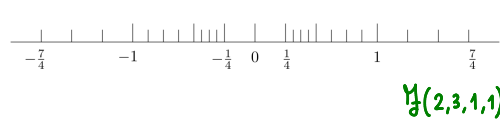
μ non dipende dalla x che sto rappresentando ma   una limitazione dell'errore relativo (x non deve essere in over/under-flow)
Quindi posso dire che μ dipende dalla macchina e non dalla rappresentazione di x !

$$\text{DIM. } d_1 \neq 0 \text{ (definizione di } \mathcal{M}(\beta, t, m, M)) \Rightarrow x \geq \beta^p \cdot d_1 \beta^{-1} \geq \beta^{p-1} \Rightarrow \left| \frac{\tilde{x} - x}{x} \right| \leq \frac{|\tilde{x} - x|}{\beta^{p-1}}$$

$$(i) \tilde{x} = \text{trunc}(x) \Rightarrow \left| \frac{\tilde{x} - x}{x} \right| < \frac{\beta^{p-t}}{\beta^{p-1}} = \beta^{1-t}$$

per maggiore, diminuisco il denominatore con $\beta^p (0.1) = \beta^{p-1}$

NOTA: β^{p-t}   la distanza tra due numeri di macchina consecutivi



$$(ii) \tilde{x} = \text{arr}(x) \Rightarrow \left| \frac{\tilde{x} - x}{x} \right| < \frac{1}{2} \frac{\beta^{p-t}}{\beta^{p-1}} = \frac{1}{2} \beta^{1-t}$$

anche qui   come sopra ma ricorda che:

$|\text{arr}(x) - x| \leq \frac{1}{2} \beta^{p-t}$, ma l'uguaglianza vale solo quando

$$\begin{cases} d_{t+1} = \frac{\beta}{2} \\ d_{t+i} = 0 \quad \forall i \geq 2 \end{cases} ; \text{ quindi } x \geq \beta^p \left(d_1 \beta^{-1} + d_{t+1} \beta^{-(t+1)} \right) \geq \beta^{p-1}$$

$$\text{posso riscrivere tutto: } \left| \frac{\text{arr}(x) - x}{x} \right| \leq \frac{1}{2} \frac{\beta^{p-t}}{|x|} \leq \frac{1}{2} \frac{\beta^{p-t}}{\beta^{p-1}} = \frac{1}{2} \beta^{1-t}$$

□

NOTA: (ϵ_x vs μ)

$$\epsilon_x = \frac{\tilde{x} - x}{x}, \quad |\epsilon_x| < \mu$$

$\Rightarrow \tilde{x} = x(1 + \epsilon_x) \Rightarrow$ Quindi \tilde{x} non è altro che una perturbazione di x !

\Rightarrow l'errore assoluto dipende dalla grandezza del numero
l'errore relativo dipende dalle specifiche delle macchina

ERRORE NELLE OPERAZIONI MACCHINA

$$x \textcircled{op} y = \begin{cases} \text{trunc}(x \text{ op } y) \\ \text{arr}(x \text{ op } y) \end{cases}$$

operazione in
aritmetica di macchina

$$\Rightarrow x \textcircled{op} y = (x \text{ op } y)(1 + \epsilon)$$

con $|\epsilon| < \mu$

$\epsilon :=$ errore locale dell'operazione
op, che impongo essere più
piccolo della precisione di
macchina μ

$$\Rightarrow \epsilon_{\text{Tot}} = \frac{(\tilde{x} \textcircled{op} \tilde{y}) - (x \text{ op } y)}{x \text{ op } y} = \epsilon + \boxed{c_x} \epsilon_x + \boxed{c_y} \epsilon_y$$

COEFF. DI AMPLIFICAZIONE:

mi danno la misura di quanto l'errore
di rappresentazione di un dato si amplifica
nel calcolo di una funzione

In questo modo se ho Troncamento o arrotondamento di una operazione op,
l'errore commesso è dell'ordine della precisione di macchina μ .

FENOMENO DELLA CANCELLAZIONE NUMERICA := quando rappresento con precisione finita dei numeri reali
in un calcolatore, potrei perdere cifre significative, causa
una operazione di sottrazione tra due numeri "quasi uguali"
(= con le prime t cifre uguali)

ESEMPIO:

$$\begin{aligned} p_1 &= 0.12345789 \\ p_2 &= 0.12345666 \end{aligned}$$

$$\begin{aligned} d &= p_1 - p_2 = 0.00000123 \\ &= 0.123 \cdot 10^{-7} \end{aligned}$$

Supponendo $\mathcal{M}(10, 6, m, M)$ e approssimazione con troncamento:

$$\begin{aligned} \text{trunc}(p_1) &= 0.123457 \\ \text{trunc}(p_2) &= 0.123456 \end{aligned} \Rightarrow \text{trunc}(d) = \text{trunc}(p_1) - \text{trunc}(p_2) = 0.000001 = 10^{-6}$$

Quindi se calcolo l'errore totale ϵ_{Tot} delle differenze, ottengo un valore molto alto:

$$\epsilon_{\text{Tot}} = \left| \frac{(p_1 - p_2) - (\tilde{p}_1 - \tilde{p}_2)}{p_1 - p_2} \right| = \left| \frac{d - \text{trunc}(d)}{d} \right| = \left| \frac{0.123 \cdot 10^{-7} - 10^{-6}}{0.123 \cdot 10^{-7}} \right| = \left| 1 - \frac{10^9}{123} \right| \approx 80.3$$

ERRORE NEL CALCOLO DI UNA FUNZIONE RAZIONALE

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$
$$[x] \mapsto x$$

$$f(x) \rightsquigarrow f(\tilde{x}) \rightsquigarrow \Psi(\tilde{x})$$

errore di
rappresentaz.
($x \rightsquigarrow \tilde{x}$)

errore dovuto all'utilizzo
di \odot invece di \otimes
($f \rightsquigarrow \Psi$)

\Rightarrow

Sono partito per calcolare $f(x)$,
vado a calcolare $\Psi(\tilde{x})$

$$\mathcal{E}_{\text{tot}} = \frac{\Psi(\tilde{x}) - f(x)}{f(x)} \quad := (*)$$

$$\mathcal{E}_{\text{in}} = \frac{f(\tilde{x}) - f(x)}{f(x)} \quad := \text{ERRORE INERENTE causato dalla rappresentazione in macchina}$$

$$\mathcal{E}_{\text{alg}} = \frac{\Psi(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})} \quad := \text{ERRORE ALGORITMICO causato dalla rappresentazione in macchina che non produce un numero macchina \in \mathbb{M}(\beta, t, m, M)}$$

$$(*) = \frac{\Psi(\tilde{x})}{f(x)} - 1 = \frac{\Psi(\tilde{x})}{f(\tilde{x})} \frac{f(\tilde{x})}{f(x)} - 1 =$$

$$= (1 + \mathcal{E}_{\text{alg}})(1 + \mathcal{E}_{\text{in}}) - 1 = \cancel{1} + \mathcal{E}_{\text{alg}} + \mathcal{E}_{\text{in}} + \underbrace{\mathcal{E}_{\text{alg}} \mathcal{E}_{\text{in}}}_{\text{questo pezzo lo posso buttare perché opero sotto l'assunzione che la precisione di macchina } |\mu| \ll 1 \text{ e che gli errori } \mathcal{E}_{\text{alg}} \text{ e } \mathcal{E}_{\text{in}} \text{ siano } O(\mu)} - \cancel{1} \doteq \mathcal{E}_{\text{alg}} + \mathcal{E}_{\text{in}}$$

questo pezzo lo posso buttare perché opero
sotto l'assunzione che la precisione di macchina
 $|\mu| \ll 1$ e che gli errori \mathcal{E}_{alg} e \mathcal{E}_{in} siano $O(\mu)$

\Rightarrow Posso trascurare i termini non lineari:
ANALISI DELL'ERRORE AL PRIMO ORDINE (\doteq)

STUDIO ERRORE INERENTE (= condizionamento problema nel calcolo di $f(x)$)

$$\begin{aligned} \mathcal{E}_{in} &= \frac{f(\tilde{x}) - f(x)}{f(x)} = \frac{f(\tilde{x}) - f(x)}{\tilde{x} - x} \frac{\tilde{x} - x}{f(x)} = \\ &= \frac{f(\tilde{x}) - f(x)}{\tilde{x} - x} \frac{x}{f(x)} \frac{\tilde{x} - x}{x} \doteq \boxed{\frac{f'(x)}{f(x)} x} \cdot \mathcal{E}_x \end{aligned}$$

COEFFICIENTE DI AMPLIFICAZIONE C_x
misura come l'errore sul dato viene amplificato
nel calcolo della funzione f

per questa analisi al primo ordine posso pensare allo sviluppo di Taylor
di f in \tilde{x} : $f(\tilde{x}) - f(x) = f'(x)(\tilde{x} - x) + O((\tilde{x} - x)^2)$

Se quindi il C_x è sempre limitato in modulo da una qualche costante
(ovvero posso maggiorare $|C_x|$ con una costante k relativamente piccola alla precisione di macchina),
allora si dice che il problema è BEN CONDIZIONATO

Viceversa, se per $x \rightarrow \tilde{x}$, $|C_x| \rightarrow \infty$, allora il problema è MAL CONDIZIONATO
nel punto \tilde{x}

STUDIO ERRORE ALGORITMICO (= errore nelle operazioni di macchina)

Come visto, è quell'errore che commette l'algoritmo partendo dai valori già trasformati
in numeri macchina:

$$\tilde{x} + \tilde{y} \rightarrow \tilde{x} \oplus \tilde{y} = \begin{matrix} \text{trunc}(\tilde{x} + \tilde{y}) \\ \text{arr}(\tilde{x} + \tilde{y}) \end{matrix}$$

$$\mathcal{E}_{Tot} \doteq \mathcal{E}_{alg} + \mathcal{E}_{in} = \mathcal{E}_{alg} + C_x \mathcal{E}_x + C_y \mathcal{E}_y$$

$$|\mathcal{E}_{Tot}| = \left| \frac{\Psi(\tilde{x}, \tilde{y}) - f(x, y)}{f(x, y)} \right|$$

$$|\mathcal{E}_{alg}| = \left| \frac{\Psi(\tilde{x}, \tilde{y}) - f(\tilde{x}, \tilde{y})}{f(\tilde{x}, \tilde{y})} \right| = \left| \frac{\text{arr}(\tilde{x} \text{ op } \tilde{y}) - (\tilde{x} \text{ op } \tilde{y})}{\tilde{x} \text{ op } \tilde{y}} \right| < \underline{\mu}$$

errore locale limitato dalla
precisione di macchina

Quindi, se non ho avuto overflow
o underflow, quando approssimo un
numero con un numero di macchina,
commetto un errore limitato alla
precisione di macchina μ

COEFFICIENTI DI CONDIZIONAMENTO C_x e C_y PER LE 4 OPERAZIONI

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

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

$$f(x, y) = x + y \Rightarrow C_x = \frac{x}{x+y}, \quad C_y = \frac{y}{x+y}$$

$$f(x, y) = x - y \Rightarrow C_x = \frac{x}{x-y}, \quad C_y = -\frac{y}{x-y}$$

$$f(x, y) = x \cdot y \Rightarrow C_x = 1, \quad C_y = 1$$

$$f(x, y) = x/y \Rightarrow C_x = 1, \quad C_y = -1$$

\Rightarrow in entrambi i casi, per $x \rightarrow y$ il coefficiente scoppia:
ci può essere un caso di cancellazione numerica

ANALISI DELL'ERRORE CON GRAFI

Supponiamo adesso di avere una funzione razionale che non sia una singola operazione aritmetica $\Psi(\tilde{x})$ calcolato al posto di $f(x)$ lo ottengo sostituendo agli x_i i corrispondenti numeri macchina \tilde{x}_i e sostituendo alle operazioni aritmetiche op , le corrispondenti operazioni di macchina \odot nell'ordine in cui vengono eseguite.

Schematicamente Ψ è ottenuta implementando un algoritmo costituito da p passi del tipo:

$$z^{(i)} = y_1^{(i)} \odot y_2^{(i)} \quad \text{per } i=1,2,\dots,p$$

Sono i dati iniziali o i risultati delle operazioni ai passi precedenti:

$$y_1^{(i)}, y_2^{(i)} \in \{x_1, x_2, \dots, x_n, z^{(1)}, z^{(2)}, \dots, z^{(i-1)}\}$$

$$f(x) = z^p$$

ESEMPIO: $f(x_1, x_2) = x_1^2 - x_2^2$

la posso calcolare con il seguente algoritmo:

$$\begin{aligned} z^{(1)} &= x_1 \cdot x_1 \\ z^{(2)} &= x_2 \cdot x_2 \\ z^{(3)} &= z^{(1)} - z^{(2)} \end{aligned}$$

Adesso, per calcolare l'errore $\varepsilon_{\text{Tot}}^{(i)}$ del risultato all' i -esimo passo:

$$\varepsilon_{\text{Tot}}^{(i)} = \varepsilon^{(i)} + c_1 \varepsilon_1^{(i)} + c_2 \varepsilon_2^{(i)} \quad \text{per } i=1,2,\dots,p$$

errore generato dalla i -esima op.

errori presenti negli operandi della i -esima operazione

\Rightarrow possono essere causati dai dati iniziali o accumulati nei risultati intermedi del calcolo

Quello che ancora non sappiamo calcolare è l'errore algoritmico ε_{Alg} .
Procediamo con un grafo che descrive la sequenza delle operazioni dell'algoritmo:

(i) i nodi corrispondono ai $\left\{ \begin{array}{l} \text{dati } x_i, \text{ con } i=1,2,\dots,n \\ \text{risultati intermedi } z_i^{(i)}, \text{ con } i=1,2,\dots,p \end{array} \right.$

(ii) in corrispondenza dei nodi x_i sono riportati gli errori di rappresentazione dei dati ε_i

(iii) in corrispondenza dei nodi $z^{(i)}$ sono riportati gli errori locali delle operazioni aritmetiche $\varepsilon^{(i)}$

(iv) ai nodi $z^{(i)}$ arrivano gli archi orientati provenienti dai nodi corrispondenti agli operandi della i -esima op.

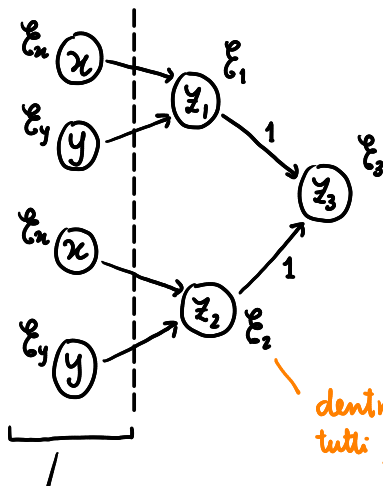
(v) in corrispondenza a ciascun arco è riportato il relativo coefficiente di amplificazione

\Rightarrow L'errore relativo totale si ottiene percorrendo il grafo dall'ultimo nodo verso i nodi iniziali: ad ogni nodo $\varepsilon_{\text{Tot}}^{(i)}$ si calcola sommando l'errore $\varepsilon^{(i)}$ che corrisponde al nodo e gli errori precedenti accumulati nei nodi ad esso collegati, moltiplicati per i coefficienti di amplificazione corrispondenti agli archi percorsi

ESEMPIO:

$$f(x, y) = \underbrace{(x-y)}_{z_1} \underbrace{(x+y)}_{z_2}$$

$$\begin{aligned} z_1 &= x \ominus y \\ z_2 &= x \oplus y \\ z_3 &= z_1 \otimes z_2 \end{aligned}$$



$$\epsilon_{alg} = \epsilon_3 + 1(\epsilon_1) + 1(\epsilon_2)$$

$$|\epsilon_{alg}| = |\epsilon_3 + \epsilon_1 + \epsilon_2| < \underbrace{|\epsilon_3|}_{< \mu} + \underbrace{|\epsilon_1|}_{< \mu} + \underbrace{|\epsilon_2|}_{< \mu} < 3\mu$$

$\Rightarrow |\epsilon_{alg}| < 3\mu$ e posso quindi affermare che l'errore algoritmico sia più piccolo della precisione di macchina 2^{-p} !

dentro questi ϵ_i mi porto dietro tutti gli errori calcolati prima

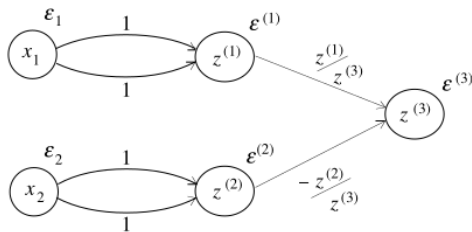
Tutto questo pezzo del grafo lo tratto con la formula di ϵ_{in} . Quindi in questo caso l'algoritmo risulta malcondizionato per $|x| \rightarrow |y|$:

$x \approx y \Rightarrow z_1 = x - y$ è piccolo e sensibile a piccoli errori nei valori di x e y (cancellazione numerica)

Esempio. L'errore inerente della funzione $f(x_1, x_2) = x_1^2 - x_2^2$ è dato da

$$\epsilon_{in} \doteq \frac{2x_1^2}{x_1^2 - x_2^2} \epsilon_1 - \frac{2x_2^2}{x_1^2 - x_2^2} \epsilon_2.$$

Quindi il problema del calcolo di $f(x_1, x_2)$ è mal condizionato quando il modulo di $x_1^2 - x_2^2$ è molto più piccolo di x_1^2 e di x_2^2 . Si vuole ora determinare l'errore totale che si produce nel calcolo di $f(x_1, x_2)$, con l'algoritmo (1.15). Dal grafo riportato nella prossima figura si ricava l'errore totale:



$$\begin{aligned} \epsilon_{tot1} &\doteq \epsilon^{(3)} + \frac{z^{(1)}}{z^{(3)}} (\epsilon^{(1)} + 2\epsilon_1) - \frac{z^{(2)}}{z^{(3)}} (\epsilon^{(2)} + 2\epsilon_2) \\ &\doteq \epsilon^{(3)} + \frac{x_1^2}{x_1^2 - x_2^2} \epsilon^{(1)} - \frac{x_2^2}{x_1^2 - x_2^2} \epsilon^{(2)} + \frac{2x_1^2}{x_1^2 - x_2^2} \epsilon_1 - \frac{2x_2^2}{x_1^2 - x_2^2} \epsilon_2, \end{aligned}$$

dove $\epsilon^{(1)}$, $\epsilon^{(2)}$, $\epsilon^{(3)}$ sono gli errori locali delle operazioni. Quindi l'errore algoritmico è

$$\epsilon_{alg1} \doteq \epsilon^{(3)} + \frac{x_1^2}{x_1^2 - x_2^2} \epsilon^{(1)} - \frac{x_2^2}{x_1^2 - x_2^2} \epsilon^{(2)}.$$

La stessa funzione può essere calcolata anche con il seguente algoritmo:

$$\begin{aligned} v^{(1)} &= x_1 + x_2 \\ v^{(2)} &= x_1 - x_2 \\ v^{(3)} &= v^{(1)} \cdot v^{(2)}. \end{aligned}$$

Dal grafo corrispondente si ricava l'errore totale:

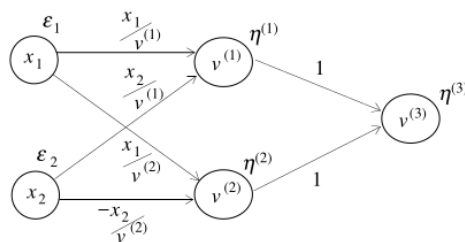
$$\begin{aligned} \epsilon_{tot2} &\doteq \eta^{(3)} + \left[\eta^{(1)} + \frac{x_1}{v^{(1)}} \epsilon_1 + \frac{x_2}{v^{(1)}} \epsilon_2 \right] + \left[\eta^{(2)} + \frac{x_1}{v^{(2)}} \epsilon_1 - \frac{x_2}{v^{(2)}} \epsilon_2 \right] \\ &\doteq \eta^{(1)} + \eta^{(2)} + \eta^{(3)} + \frac{2x_1^2}{x_1^2 - x_2^2} \epsilon_1 - \frac{2x_2^2}{x_1^2 - x_2^2} \epsilon_2, \end{aligned}$$

dove $\eta^{(1)}$, $\eta^{(2)}$, $\eta^{(3)}$ sono gli errori locali delle operazioni. Quindi l'errore algoritmico è

$$\epsilon_{alg2} \doteq \eta^{(1)} + \eta^{(2)} + \eta^{(3)}.$$

Si confrontano i due errori algoritmici. Risulta

$$|\epsilon_{alg1}| < \left(1 + \frac{x_1^2 + x_2^2}{|x_1^2 - x_2^2|} \right) u, \quad |\epsilon_{alg2}| < 3u,$$



perciò il secondo algoritmo è più stabile del primo se $x_1^2 + x_2^2 > |x_1^2 - x_2^2|$. \square

L'analisi dell'errore condotta con il grafo permette di valutare sia l'errore algoritmico che l'errore inerente; poiché però per l'errore inerente conviene utilizzare la (1.13), si può semplificare l'analisi utilizzando il grafo per il calcolo del solo errore algoritmico, assumendo nulli gli errori di rappresentazione dei dati, come sarà fatto negli esempi che seguono.

ERRORE NEL CALCOLO DI UNA FUNZIONE NON-RAZIONALE (Errore analitico \mathcal{E}_{an})

Se la funzione f non è razionale, è necessario prima approssimarla con una funzione razionale g : questa approssimazione introduce un ERRORE ANALITICO:

$$\mathcal{E}_{an}(x) = \frac{g(x) - f(x)}{f(x)}$$

Se ψ è la funzione effettivamente calcolata al posto della g , posso scrivere:

$$\mathcal{E}_{Tot} = \frac{\psi(\tilde{x}) - f(x)}{f(x)} \doteq \underbrace{\mathcal{E}_{in}}_{\sum_{i=1}^n c_i \varepsilon_i} + \underbrace{\mathcal{E}_{alg}}_{\frac{\psi(\tilde{x}) - g(\tilde{x})}{g(\tilde{x})}} + \underbrace{\mathcal{E}_{an}(\tilde{x})}_{\frac{g(x) - f(x)}{f(x)}}$$

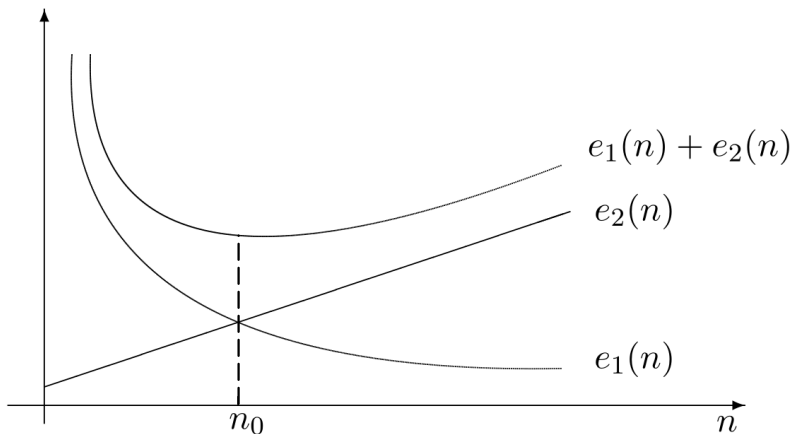
D'ora in avanti si supponerà che $\mathcal{E}_{an}(\tilde{x}) \doteq \mathcal{E}_{an}(x)$ e, come già fatto per \mathcal{E}_{in} e \mathcal{E}_{alg} , ne ometteremo l'argomento:

$$\mathcal{E}_{Tot} \doteq \mathcal{E}_{in} + \mathcal{E}_{alg} + \mathcal{E}_{an}$$

Adesso, come trovo g che approssima f ?

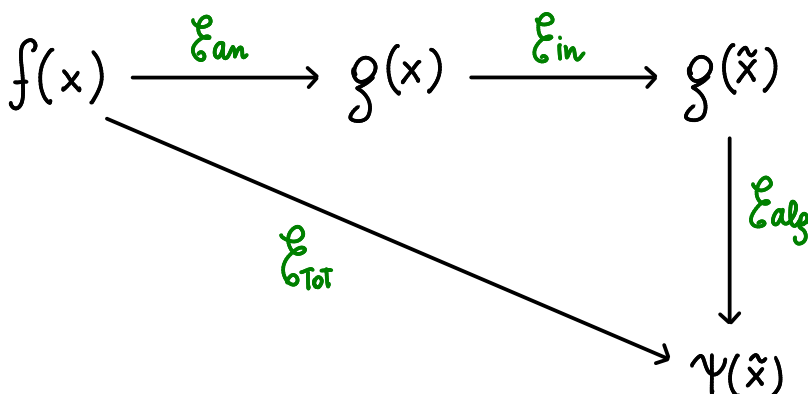
Un modo è quello di usare la formula di Taylor, approssimando $f(x)$ (per semplicità qui monovariata) con un polinomio di grado n ottenuto troncando all' $(n+1)$ -esimo termine la formula di Taylor di $f(x)$.

In generale l'errore analitico \mathcal{E}_{an} tende a diminuire quanto più elevato è n , mentre l'errore algoritmico \mathcal{E}_{alg} tende ad aumentare con n :



$e_1(n)$, $e_2(n)$ indicano le maggiorazioni di $|\mathcal{E}_{an}|$ e $|\mathcal{E}_{alg}|$ al variare di n .

Quindi conviene scegliere un valore di n vicino a n_0 , perché per n molto più grande di n_0 , ad un maggior volume di calcolo, non corrisponde una diminuzione dell'errore effettivamente prodotto.



Quando approssimo f irrazionale con la serie di Taylor troncata, l'errore commesso è proprio il resto di Laplace!