

# Capitolo 2

## Analisi degli errori

### 2.1 Rappresentazione dei numeri

L'utilizzo in modo corretto del calcolatore per fare calcoli di tipo scientifico, richiede la conoscenza di come sono rappresentati i numeri e degli errori che derivano da questa rappresentazione. L'uso dei numeri reali richiede una attenzione particolare, essendo questi infiniti, mentre il calcolatore ci dà la possibilità di rappresentarne solo un numero finito.

La nostra notazione per rappresentare i numeri è una notazione posizionale a base 10. Ciò significa che se scriviamo 123 intendiamo esprimere il numero:

$$1 \cdot 10^2 + 2 \cdot 10^1 + 3 \cdot 10^0.$$

Cioè la cifra che occupa la posizione più a destra deve essere moltiplicata per  $10^0 (= 1)$ , quella che occupa la posizione centrale per 10, ecc. Questa convenzione è ormai universalmente adottata, al contrario di altre “convenzioni” che quotidianamente usiamo per la comunicazione di nostre idee ed esperienze come, ad esempio, le lingue, gli alfabeti, le direzioni di scrittura, la musica, ecc.

Proprio per la sua universale adozione, la notazione posizionale a base 10 può apparire a molti come l'unica possibile. Ciò non è vero. Anzi, dal punto di vista storico risulta che tale rappresentazione si è affermata solo da pochi secoli. In precedenza erano state usate tante altre rappresentazioni, alcune in base 10 ma non posizionali (ad esempio quella romana), altre posizionali

ma non a base 10 (ad esempio la sumerica). Torniamo alla rappresentazione attuale.

Un qualunque numero lo possiamo esprimere nella forma:

$$\pm\alpha_p\alpha_{p-1}\dots\alpha_0.\beta_1\beta_2\dots\beta_q\dots \quad (2.1)$$

sottointendendo:

$$\pm\alpha_p \cdot 10^p + \dots + \alpha_0 \cdot 10^0 + \beta_1 \cdot 10^{-1} + \dots + \beta_q \cdot 10^{-q} \dots$$

I coefficienti  $\alpha_i$  e  $\beta_i$  sono uno dei seguenti simboli: 0, 1, 2, ..., 9. Nulla vieta di usare al posto di 10 un altro numero, diciamo  $N$ . In tale caso la stringa (intesa come successione ordinata di simboli) descritta dalla 2.1 assumerebbe un altro significato. Infatti, supposto che i simboli  $\alpha$  e  $\beta$  siano compresi tra 0 ed  $N - 1$ , essa sarebbe equivalente a:

$$\pm\alpha_p \cdot N^p + \dots + \alpha_0 \cdot N^0 + \beta_1 \cdot N^{-1} + \dots + \beta_q \cdot N^{-q} \dots$$

Si ottiene così la rappresentazione posizionale di un numero in base  $N$ .

Nel campo scientifico si preferisce usare una notazione teoricamente equivalente alla 2.1, cioè<sup>1</sup>:

$$\pm\gamma_0.\gamma_1\gamma_2\dots\gamma_t\dots \cdot N^q \quad (2.2)$$

con  $\gamma_0 \neq 0$ . Quindi il numero  $3.57 \cdot 10^2$  è espresso in notazione esponenziale normalizzata, mentre i numeri  $73.54 \cdot 10^2$  o  $0.57 \cdot 10^{-3}$  non lo sono. Con questa rappresentazione si possono rappresentare tutti i numeri eccetto lo zero. Essa non è univoca. Ad esempio in base 10 i due numeri:

$$\begin{aligned} &1.200\dots \\ &1.199\dots \end{aligned}$$

rappresentano lo stesso numero. La parte  $\gamma_0.\gamma_1\gamma_2\dots\gamma_t\dots$  viene chiamata mantissa e noi la indicheremo con  $m$ . Il numero  $q$  viene chiamato esponente, mentre  $N$  è la base. Ogni numero reale non nullo sarà quindi indicato in questa rappresentazione, detta *rappresentazione normalizzata*, con:

$$\pm m \cdot N^q$$

---

<sup>1</sup>Seguiamo la convenzione più recente in alternativa a quella equivalente  $\pm.\gamma_0\gamma_1\dots\gamma_t\dots \cdot N^{q+1}$  che si trova più frequentemente nei libri.

Il numero zero sarà rappresentato con 0. Ovviamente nella rappresentazione di un numero intero tutti i valori  $\gamma_i$  saranno nulli da un certo indice (maggiore di uno) in poi. È ovvio che:

$$1 \leq m < N. \quad (2.3)$$

Una stessa stringa rappresenta numeri diversi a seconda della base usata. Ad esempio la stringa 123 in base quattro rappresenta il numero:

$$1 \cdot 4^2 + 2 \cdot 4 + 3 = \text{ventisette}$$

cioè il numero che in base 10 indichiamo con 27. Nel caso in cui  $N$  fosse più grande di dieci, il numero di simboli che usualmente si usano per indicare le cifre tra 0 ed  $N - 1$ , cioè 0,1,2,...,9 non bastano più. Si fa ricorso alle lettere maiuscole dell'alfabeto per integrare i simboli mancanti. Ad esempio se  $N = 16$ , i simboli usati sono 0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F. Intendendo con questo che il numero dieci viene rappresentato con A, l'undici con B, ecc.. Il numero sedici sarà ovviamente indicato con 10.

Esiste una convenienza nell'uso di una base al posto di un'altra? Diamo uno sguardo alla tabella seguente in cui il numero (non la stringa!) centoventitrè è rappresentato in diverse basi:

base	rappresentazione
16	7B
10	123
5	443
4	1323
2	1111011

È evidente che lo stesso numero richiede un numero di cifre crescente al diminuire della base. Ciò può risultare scomodo. Vi è però certamente un vantaggio nell'usare le basi più piccole. Infatti i calcoli risultano semplificati. La tavola pitagorica, che in base dieci è formata da un quadrato di 100 elementi, nella base 2 risulterebbe formata da 4 elementi: chi potrebbe negare che ciò sarebbe un notevole sollievo per i bambini delle scuole elementari?

La base due è ancora più interessante perché in tale base tutti i numeri sono rappresentati da successioni dei due simboli 0 ed 1 che possono essere

messe in corrispondenza con ogni sistema che abbia solo due stati (ad esempio acceso-spento, magnetizzato-smagnetizzato ecc.). Difatti è per questo motivo che essa risulta indispensabile per la codifica di informazioni nei calcolatori. Ogni informazione, per essere memorizzata nel calcolatore, viene opportunamente codificata nel linguaggio comprensibile alla macchina che è appunto basato sui simboli del sistema di numerazione binario, 0 ed 1. Abbiamo già avuto occasione di dire che questi due simboli rappresentano lo stato fisico dei circuiti di cui è costituita la memoria, che possono essere spenti (=0) o accesi (=1).

Un numero reale rappresentato tramite la 2.2 deve essere prima espresso in cifre binarie 0 ed 1 per poter essere memorizzato in un calcolatore.

Ovviamente la rappresentazione 2.2 necessita un numero infinito di cifre e quindi sarebbe impossibile memorizzarla in un calcolatore. Nei calcolatori ad ogni numero viene riservato uno spazio fisso di memoria, per cui sia il numero di cifre della mantissa che il numero di cifre dell'esponente non possono superare dei limiti prefigurati. I numeri di macchina, (detti anche numeri *floating point*), sono del tipo:

$$\pm \gamma_0 \cdot \gamma_1 \gamma_2 \dots \gamma_t \cdot N^q \quad (2.4)$$

con  $\gamma_0 \neq 0$  e  $q_{min} \leq q \leq q_{max}$ , più, naturalmente, lo zero.

## 2.2 Standard IEEE

Descriviamo più in dettaglio come viene rappresentato un numero usando lo standard *floating point IEEE*, che costituisce un insieme di regole definito dall'istituto degli ingegneri elettrici e elettronici per la rappresentazione ed elaborazione dei numeri *floating point* nei computer. Lo standard IEEE specifica esattamente cosa sono i numeri *floating point* e come sono rappresentati nell'hardware e ha 4 scopi principali:

1. rendere l'aritmetica *floating point* il più accurata possibile;
2. produrre risultati sensati in situazioni eccezionali;
3. standardizzare le operazioni *floating point* fra i computer;
4. Dard al programmatore un controllo sulla manipolazione delle eccezioni;

Il punto (1) si ottiene in due modi. Lo standard specifica esattamente come un numero floating point dovrebbe essere rappresentato nell'hardware e richiede che le operazioni (addizione, radice, ...) siano il più accurate possibili. Per il punto (2) lo standard produce Inf (Infinito) per indicare che un risultato è più grande del massimo numero floating point rappresentabile e NaN (Not a Number) per indicare che un risultato non ha è definito. Prima dello standard la maggior parte dei computer avrebbero bloccato un programma in tali circostanze. Il punto (3) ha alcune conseguenze:

- (i) I numeri floating point possono essere trasferiti fra due computer che utilizzano lo standard IEEE in binario senza perdere precisione;
- (ii) I dettagli dell'aritmetica floating point dovrebbero essere conosciuti dal programmatore;
- (iii) Lo stesso programma su computer diversi dovrebbe produrre lo stesso risultato, anche se questo non è vero in pratica.

Ricordiamo che l'unità base di informazione immagazzinata da un computer è il bit, una variabile il cui valore è 0 o 1. I bit sono organizzati in gruppi di 8 chiamati byte. L'unità più comune è la parola di 4 byte (32 bit). Per accuratezze più alte è raddoppiata a 8 byte o 64 bit. Vi sono 2 possibili valori per un bit,  $2^8 = 256$  possibili valori per un byte e  $2^{32}$  differenti parole di 4 byte.

I due tipi di numeri rappresentati sono interi (fixed point) e reali (floating point). Un numero reale ha tipo float in c e real in Fortran e Matlab. Nella maggior parte dei compilatori c e in Matlab un float ha per default 8 byte invece di 4. Il formato IEEE sostituisce la base 10 con la base 2.

L'aritmetica con gli interi è molto semplice. Vi sono  $2^{32} \approx 4 \cdot 10^9$  interi a 32 bit che coprono l'intervallo da  $-2 \cdot 10^9$  a  $2 \cdot 10^9$ . Addizione, sottrazione e moltiplicazione sono fatte esattamente se la risposta è compresa nell'intervallo. La maggior parte dei computer danno risultati imprevedibili se il risultato è fuori dal range (overflow). Lo svantaggio dell'aritmetica con gli interi è che non possono essere rappresentate le frazioni e l'intervallo dei numeri è piccolo.

Quando una parola a 32 bit è interpretata come un numero floating point il primo bit è il bit del segno,  $s = +$  o  $s = -$ . I successivi 8 bit formano l'esponente e i rimanenti 23 bit determinano la mantissa. Vi sono 2 possibili segni, 256 esponenti (che variano da 0 a 255) e  $2^{23} \approx 8.4$  milioni di possibili mantisse. Per avere gli esponenti negativi si usa una convenzione: lo zero e

nella posizione 127, dopo ci sono i numeri positivi e prima i numeri negativi. Quindi in memoria viene rappresentato  $q^* = q + 127$ . Inoltre il primo bit della mantissa, che rappresenta  $\gamma_0$ , è sempre uguale a 1, quindi non vi è necessità di memorizzarlo esplicitamente. Nei bit assegnati alla mantissa viene memorizzato  $m^*$  e  $m = 1.m^*$ . Un numero floating point positivo ha quindi il valore  $x = +(1.m^*)_2 2^{q^*-127}$  e la notazione  $(1.m^*)_2$  indica che  $1.m^*$  è interpretata in base 2.

Per esempio il numero  $2.752 \cdot 10^3 = 2572$  può essere scritto:

$$\begin{aligned} 2752 &= 2^{11} + 2^9 + 2^7 + 2^6 = \\ &= 2^{11}(1 + 2^{-2} + 2^{-4} + 2^{-5}) = \\ &= 2^{11}(1 + (0.01)_2 + (0.0001)_2 + (0.00001)_2) = \\ &= 2^{11}(1.01011)_2 \end{aligned}$$

allora la rappresentazione di questo numero avrebbe segno + esponente  $q^* = q + 127 = 138 = (10001010)_2$  e  $m^* = (010110\dots 0)_2$ .

Il caso  $e = 0$  ( che corrisponde a  $2^{-127}$ ) e il caso  $e = 255$  (che corrisponde a  $2^{128}$ ) hanno una interpretazione differente e complessa che rende la IEE diversa dagli altri standard.

Se  $q^* = 0$  il valore del numero memorizzato è  $x = \pm(0.m^*)_2 2^{-126}$ . Questo è chiamato underflow graduale (l'underflow è la situazione in cui il risultato di una operazione è diversa da zero ma è più vicina a zero di qualsiasi numero floating point). I numeri corrispondenti vengono chiamati denormalizzati. L'underflow graduale ha la conseguenza che due numeri floating point sono uguali se e solo se sottraendone uno dall'altro si ha esattamente zero.

Se escludiamo i numeri denormali allora il più piccolo numero positivo floating point è  $a = realmin = q^{-126}$ , il successivo numero positivo  $b$  ha  $q^* = 1$  e  $m^* = 0.0\dots 01$ , la distanza fra  $a$  e  $b$  è  $2^{23}$  volte più piccola, della distanza fra  $a$  e zero. Senza l'underflow graduale vi è un gap fra zero ed il numero floating point più vicino.

L'altro caso è  $q^* = 255$  che ha due sottocasi, *Inf* se  $m^* = 0$  e *NaN* se  $m^* \neq 0$ . Sia il  $c$  che il Matlab stampano *Inf* o *NaN* quando si stampa una variabile floating point che contiene questo risultato. Il computer produce *Inf* se il risultato di una operazione è più grande del più grande numero rappresentabile.

Operazioni invalide che producono come risultato *NaN* sono:  $Inf/Inf$ ,  $0/0$ ,  $Inf - Inf$ . Operazioni con *Inf* hanno il significato usuale:  $Inf + finito = Inf$ ,  $Inf/Inf = NaN$ ,  $Finito/Inf = 0$ ,  $Inf - Inf = NaN$ .

È Chiaro che l'accuratezza delle operazioni Floating Point è determinata dalla grandezza degli errori di arrotondamento. Questo errore di arrotondamento è determinato dalla distanza di un numero dal numero floating point più vicino. Eccetto che per i numeri denormalizzati, i numeri floating point differiscono di un bit, l'ultimo bit di  $m^*$ . Cioè  $2^{-23} \approx 10^{-7}$  in singola precisione. Questo è l'errore relativo e non l'errore assoluto.

L'aritmetica IEEE in doppia precisione usa 8 byte (64 bit). Vi è un bit del segno, 11 bit per l'esponente e 52 bit per la mantissa. Quindi in doppia precisione  $\epsilon = 2^{-52} \approx 4 \cdot 10^{-16}$ . Cioè l'aritmetica in doppia precisione da circa 15 digit decimali di accuratezza, invece di 7. Vi sono  $2^{11}$  possibili esponenti che variano da 1023 a -1022. Il più grande numero è circa  $10^{307} \approx 2^{1023}$

## 2.3 Errore assoluto e relativo

Consideriamo il numero reale positivo:

$$x = \gamma_0 \cdot \gamma_1 \gamma_2 \dots \gamma_t \dots \cdot N^q$$

Se consideriamo i numeri di macchina con  $t$  cifre per la mantissa,  $x$  sarà compreso fra:

$$x_1 = \gamma_0 \cdot \gamma_1 \gamma_2 \dots \gamma_t \cdot N^q$$

e

$$x_2 = x_1 + 0.0 \dots 01 \cdot N^q = x_1 + N^{-t} \cdot N^q$$

graficamente si ha:

$$\begin{array}{c} \gamma_0 \cdot \gamma_1 \gamma_2 \dots \gamma_t \dots \cdot N^q \\ \hline \begin{array}{ccc} | & & | \\ x_1 & & x_2 \end{array} \end{array}$$

Vi sono due modi diversi per approssimare questo numero mediante i numeri floating point, in cui solo  $t$  cifre della mantissa sono rappresentabili:

1. il troncamento;
2. l'arrotondamento.

Nel primo caso il numero  $x$  viene sempre approssimato con  $x_1$  trascurando tutte le cifre successive a  $t$ ; nel secondo caso il numero viene approssimato con  $x_1$  se  $x < (x_1 + x_2)/2$ , cioè se  $x$  si trova nella prima metà dell'intervallo  $[x_1x_2]$ , altrimenti viene approssimato con  $x_2$ .

Denotiamo con  $fl(x)$  l'approssimazione di  $x$ . Se viene usato il troncamento, qualunque sia la posizione di  $x$  nell'intervallo,  $x$  viene approssimato con il numero di macchina alla sua sinistra e quindi l'errore assoluto, cioè la distanza in valore assoluto di  $fl(x)$  da  $x$ , sarà sempre minore dell'ampiezza dell'intervallo  $[x_1x_2]$ :

$$|fl(x) - x| < N^{-t} \cdot N^q$$

Nel secondo caso invece l'errore assoluto metà dell'ampiezza dell'intervallo:

$$|fl(x) - x| < \frac{1}{2}N^{-t} \cdot N^q$$

L'errore relativo di un numero  $x \neq 0$ , approssimato con  $fl(x)$  è definito da:

$$\left| \frac{fl(x) - x}{x} \right|$$

**Teorema 2.3.1** *Se  $x \neq 0$  e usiamo  $t$  cifre per rappresentare la mantissa, allora:*

$$\left| \frac{fl(x) - x}{x} \right| \leq u \tag{2.5}$$

dove:

$$u = \begin{cases} N^{-t} & \text{in caso di troncamento} \\ \frac{1}{2}N^{-t} & \text{in caso di arrotondamento} \end{cases}$$

**Dimostrazione.** Sia  $x = \gamma_0.\gamma_1\gamma_2 \dots \gamma_t \dots \cdot N^q$ , allora  $|x| \geq N^q$  e, nel caso di troncamento, si ha:

$$\frac{|fl(x) - x|}{|x|} \leq \frac{N^{q-t}}{N^q} \leq N^{-t},$$

da cui segue la tesi.

Si procede analogamente nel caso di arrotondamento e nel caso di numeri negativi.



Il numero  $u$  è detto *unità di arrotondamento* (o *precisione di macchina*). Secondo lo standard IEEE un numero “reale” in doppia precisione occupa 64 bits (= 8 bytes) di memoria di cui 52 sono per la mantissa. Se la base usata è  $N = 2$ , la precisione di macchina risulta essere  $u = 2^{-52} \approx 2.22044 \cdot 10^{-16}$ .

---



---

BASE.M

---



---

Per determinare la base usata da un calcolatore è utile considerare che  $fl(n + u) > n$  per  $1 \leq n < N$  e  $fl(N + u) = N$ . Possiamo quindi utilizzare la seguente procedura

```
base = 1
while (base + eps) > base
    base = base + 1;
end
```

La costante *eps* contiene il valore della precisione di macchina.

---



---

Come conseguenza del Teorema 2.3.1 si ha che l'errore relativo, piuttosto che l'errore assoluto, è legato alle cifre esatte di un numero  $x$ . Se  $|x| < realmin$ , allora  $fl(x) = 0$  e quindi:

$$\left| \frac{x - fl(x)}{x} \right| = 1,$$

cioè l'errore relativo è molto più grande di  $u$ . Lo stesso avviene se  $|x| > realmax$ .

Poniamo:

$$\epsilon = \frac{fl(x) - x}{x}$$

Sappiamo che  $|\epsilon| \leq u$ . Dalla precedente si ricava:

$$fl(x) = x(1 + \epsilon), \tag{2.6}$$

cioè per ogni  $x$  tale che  $realmin \leq |x| \leq realmax$  esiste un numero  $\epsilon$  in valore assoluto minore od uguale ad  $u$  tale che la sua approssimazione mediante un numero di macchina può essere espressa dalla (2.6). Se vale la (2.5) possiamo dire che il numero è rappresentato con  $t$  cifre esatte.

**Definizione 2.3.1** *La precisione di macchina è il più grande numero floating point  $x$  tale che  $x + 1$  non può essere distinto da  $x$ :*

$$u = \max\{x \mid 1 + x = 1 \text{ nell'aritmetica floating point}\}$$

## 2.4 Operazioni con i numeri di macchina

Indichiamo con  $op$  una qualunque operazione aritmetica. Non vogliamo entrare nei dettagli sul modo in cui vengono eseguite le operazioni nel calcolatore. Accenneremo solo brevemente al modo in cui si eseguono le addizioni. Supponiamo di voler eseguire l'addizione tra due numeri  $x = \alpha_0.\alpha_1\alpha_2 \dots \alpha_t \cdot N^p$  e  $y = \beta_0.\beta_1\beta_2 \dots \beta_t \cdot N^q$ , in cui  $p > q$ . Dobbiamo trasformare i due numeri in modo da avere i due esponenti uguali a  $p$ . Per fare ciò bisogna che l'area di memoria in cui si eseguono le operazioni abbia uno spazio riservato alle mantisse superiore a quello usuale. Nel nostro caso si pone  $y = (\beta_0.\beta_1\beta_2 \dots \beta_t \cdot N^{q-p}) \cdot N^p = .0 \dots 0\beta_0\beta_1 \dots \beta_t \cdot N^p$  e poi si sommano le mantisse. Si procede quindi al troncamento o all'arrotondamento. In generale quindi il risultato di una operazione può considerarsi come il risultato esatto dell'operazione tra  $x$  ed  $y$  seguito dall'arrotondamento o dal troncamento, in modo tale da soddisfare in genere la seguente relazione:

$$fl(x \text{ op } y) = (x \text{ op } y)(1 + \epsilon) \quad |\epsilon| < u. \quad (2.7)$$

Le operazioni con i numeri di macchina (aritmetica *floating point*), non godono di tutte le proprietà di cui godono le corrispondenti operazioni con i numeri reali. In generale non vale più la proprietà associativa. Consideriamo ad esempio il seguente caso in cui  $N = 10$  e  $t = 2$ . Eseguiamo la somma  $5.24 \cdot 10^{-2} + 4.04 \cdot 10^{-2} + 1.21 \cdot 10^{-1}$  in due modi diversi:

- 1)  $5.24 \cdot 10^{-2} + (4.04 \cdot 10^{-2} + 1.21 \cdot 10^{-1}) = 5.24 \cdot 10^{-2} + (0.404 + 1.21)10^{-1} = 5.24 \cdot 10^{-2} + 1.61 \cdot 10^{-1} = (0.524 + 1.61)10^{-1} = 2.13 \cdot 10^{-1}$ ;
- 2)  $(5.24 \cdot 10^{-2} + 4.04 \cdot 10^{-2}) + 1.21 \cdot 10^{-1} = 9.28 \cdot 10^{-2} + 1.21 \cdot 10^{-1} = (0.928 + 1.21)10^{-1} = 2.14 \cdot 10^{-1}$ .

Supponiamo ora di avere due numeri reali  $x$  ed  $y$  e di volere fare la somma  $x + y$ . Nel calcolatore  $x$  sarà rappresentato da  $\bar{x} = x(1 + \epsilon_1)$  e  $y$  da  $\bar{y} = y(1 + \epsilon_2)$ . Inoltre  $fl(\bar{x} + \bar{y}) = (\bar{x} + \bar{y})(1 + \epsilon)$ . L'errore relativo sarà quindi, trascurando i termini contenenti  $\epsilon\epsilon_1$  ed  $\epsilon\epsilon_2$ ,

$$\begin{aligned} \frac{(x+y) - (\bar{x} + \bar{y})(1+\epsilon)}{x+y} &= \frac{x+y - (x+x\epsilon_1 + y+y\epsilon_2)(1+\epsilon)}{x+y} \approx \\ &\approx -\epsilon - \epsilon_1 \frac{x}{x+y} - \epsilon_2 \frac{y}{x+y}. \end{aligned}$$

Dal precedente risultato risulta evidente che se il denominatore non è piccolo, come avviene certamente nel caso in cui i due numeri siano dello stesso segno, l'errore relativo è dello stesso ordine degli errori relativi  $\epsilon$ ,  $\epsilon_1$  ed  $\epsilon_2$  e quindi in valore assoluto minore di  $u$ . Ciò non avviene nel caso in cui il denominatore risulti molto piccolo. Ciò risulterà evidente nel seguente esempio.

Supponiamo infatti di voler effettuare la differenza tra i due numeri reali  $x_1 = 0.147554326$  ed  $x_2 = 0.147251742$  ed operare con  $t = 5$  ed  $N = 10$ . Usando l'arrotondamento sarà  $fl(x_1) = 1.47554 \cdot 10^{-1}$  e  $fl(x_2) = 1.47252 \cdot 10^{-1}$ . Si ha  $fl(fl(x_1) - fl(x_2)) = 3.02000 \cdot 10^{-4}$ , mentre la vera differenza è  $3.02584 \cdot 10^{-4}$ . Le ultime tre cifre della mantissa risultano errate. Ciò è dovuto al fatto che dopo aver eseguito la differenza  $fl(x_1) - fl(x_2) = 0.000302$ , la rappresentazione normalizzata  $3.02000 \cdot 10^{-4}$  ha introdotto tre zeri alla fine della mantissa. Ovviamente questi tre zeri non sono esatti. L'errore relativo risultante è:

$$\frac{3.02584 - 3.02000}{3.02584} \approx 10^{-3}$$

che è piuttosto elevato. Quindi bisogna fare particolare attenzione a questo fenomeno, che viene chiamato *cancellazione numerica*. Essa si presenta quando si esegue la sottrazione tra numeri molto vicini.

Spesso è comodo vedere la (2.7) da un punto di vista leggermente diverso. Consideriamo per esempio il caso in cui  $op$  sia l'addizione. Sarà:

$$fl(x+y) = (x + \epsilon x) + (y + \epsilon y),$$

cioè il risultato dell'operazione di macchina può essere visto come l'operazione esatta su dati perturbati  $(x + \delta x)$  e  $(y + \delta y)$ . Ovviamente nel caso considerato sarà  $\delta x = \epsilon x$  e  $\delta y = \epsilon y$ . Nel caso in cui  $op$  sia la moltiplicazione si ha:

$$fl(xy) = xy(1 + \epsilon) = (x + \delta x)(y + \delta y) \text{ con } \delta x = 0, \delta y = \epsilon y.$$

La tecnica di considerare la perturbazione nei dati che porterebbe allo stesso risultato finale se le operazioni fossero eseguite con precisione infinita è molto utile in Analisi Numerica e viene chiamata *analisi degli errori all'indietro* (*backward*).

**Esercizio 2.4.1** Sia  $t = 3$  e supponiamo che le operazioni siano eseguite con  $\bar{t} = 6$  cifre decimali, eseguire le seguenti operazioni aritmetiche dopo aver trasformato i numeri nella rappresentazione floating point:

$$\begin{aligned} &1.235 + 1.527 \cdot 10^{-4}, \\ &.5 \cdot 10^3 - 2.5, \\ &1.7435 \cdot 10^{-5} - 1.7448 \cdot 10^{-5}. \end{aligned}$$

Calcolare l'errore relativo e l'errore assoluto commesso nell'eseguire i calcoli.

## 2.5 Propagazione degli errori

Consideriamo una funzione  $y = f(x)$  e supponiamo di volerla calcolare in un punto particolare  $\bar{x}$  il quale è stato misurato con un errore  $\Delta x$ . Questo errore si ripercuote sulla variabile  $y$  con un errore  $\Delta y$  che, nel caso in cui  $f$  ammetta derivata prima in un intorno di  $\bar{x}$  si può stimare facilmente. Si ha:

$$\bar{y} + \Delta y = f(\bar{x} + \Delta x) \simeq f(\bar{x}) + f'(\bar{x})\Delta x$$

da cui:

$$\frac{\Delta y}{\bar{y}} \simeq \bar{x} \frac{f'(\bar{x})}{f(\bar{x})} \frac{\Delta x}{\bar{x}}.$$

Si deduce che un errore relativo  $\Delta x/\bar{x}$  sulla variabile indipendente produce un errore relativo sulla variabile dipendente che risulterà notevolmente amplificato se la quantità:

$$K(\bar{x}, f) = \left| \bar{x} \frac{f'(\bar{x})}{f(\bar{x})} \right|$$

è molto grande. Questa quantità è detta indice di condizionamento del problema.

**Esempio 2.5.1** Sia  $f(x) = \log x$ . Si ha  $K(x, f) = 1/\log(x)$  da cui risulta che il calcolo del logaritmo è un problema mal condizionato se  $x$  è vicino ad 1.

**Esempio 2.5.2** Consideriamo il problema di calcolare le radici di un polinomio. In questo caso i dati del problema sono i coefficienti del polinomio e il risultato finale le sue radici. Se consideriamo il seguente polinomio:

$$(x - 1)^4 = x^4 - 4x^3 + 6x^2 + 4x + 1$$

sappiamo che la radice è 1 con molteplicità 4. Perturbiamo il termine noto, ottenendo il polinomio:

$$(x - 1)^4 = x^4 - 4x^3 + 6x^2 + 4x + 1 - 10^{-8}$$

il quale ha le radici:

$$\begin{aligned}x_1 &= 1.01 \\x_2 &= 0.99 \\x_3 &= 1 + i0.01 \\x_4 &= 1 - i0.01\end{aligned}$$

quindi una perturbazione relativa di  $1e-8$  sul termine noto produce una perturbazione relativa di  $1e-2$  sul risultato. Possiamo concludere che il problema è mal condizionato.

Quando il risultato di un problema non si ottiene soltanto calcolando il valore di una funzione regolare in un punto, lo studio della propagazione degli errori e del condizionamento del problema può essere molto laboriosa. Studiare il condizionamento del problema è molto importante nel caso in cui si voglia risolvere il problema usando il calcolatore. In questo caso infatti una perturbazione sui dati viene sempre provocata quando gli stessi vengono memorizzati.