

# Calcolo Numerico

## Foglio di esercizi sull'analisi dell'errore

Gianna Del Corso <gianna.delcorso@unipi.it>

Febbraio 2025

*Esercizio 1.* Si consideri l'insieme dei numeri di macchina  $F(10, 3, 7, 6)$  in cui si assume di operare con troncamento.

- (a) Assegnato  $x = 2022$ , se ne calcoli la rappresentazione  $\tilde{x} \in F$ .
- (b) Si determinino tutti i numeri reali tali che la loro rappresentazione in  $F$  coincida con  $\tilde{x}$ ,
- (c) Si calcoli la precisione di macchina  $u$  e si determinino due numeri  $z \in F$  tali che  $\frac{|z-x|}{|x|} < u$ .

*Esercizio 2.* Si considerino i due insiemi  $F(10, t, m, M)$  con  $t = 3, m = 4, M = 5$  e  $G$  definito come l'unione di  $F$  con l'insieme dei numeri non nulli della forma  $x = \pm 10^{-m}(0.0d_2, \dots, d_t)$ , con  $0 \leq d_i \leq 9$  per  $i=2, \dots, t$  (quindi  $G$  contiene anche numeri piccoli non normalizzati).

- (a) Si calcolino i minimi positivi non nulli  $\omega$  e i massimi  $\Omega$  degli insiemi  $F$  e  $G$ .
- (b) Si determinino le cardinalità degli insiemi  $F$  e  $G$ .
- (c) Quale errore relativo di rappresentazione si commette volendo rappresentare  $1.4 \cdot 10^{-(t+m)}$  in  $F$  e in  $G$ .

*Esercizio 3.* (\*) Sia  $fl(x)$  la rappresentazione, ottenuta per troncamento, nell'insieme  $F(10, 3, m, M)$  di un numero  $x \in \mathbb{R}$ .

- (a) Dimostrare che  $x_1, x_2 \in \mathbb{R}$  e  $0 < x_1 \leq x_2$  implica  $fl(x_1) \leq fl(x_2)$ .
- (b) Dimostrare con un opportuno esempio che  $fl(x_1) \leq fl(x_2)$  non implica che  $x_1 \leq x_2$ .

*Esercizio 4.* Sia  $\mathcal{F} = \mathcal{F}(\beta, t, m, M)$  con  $m = M$ .

- Si dica quali sono il minimo numero positivo  $\omega$  ed il massimo numero positivo  $\Omega$  di  $\mathcal{F}$ .
- Si dica, giustificando la risposta, se i numeri  $b = 1/\omega$  e  $B = 1/\Omega$  appartengono a  $\mathcal{F}$ .
- Si esamini in particolare il caso  $\beta = 2, t = 8$  ed  $m = M = 6$ .

*Esercizio 5.* (\*) Si dica quale è il più grande intero  $N$  tale che tutti gli interi nell'intervallo  $[-N, N]$  sono rappresentabili esattamente in  $\mathcal{F}(2, t, m, M)$ . Quale è il valore di  $N$  corrispondente per l'aritmetica IEEE in doppia precisione?

*Esercizio 6.* Sia  $f(x_1, \dots, x_n) = \sum_{i=1}^n x_i^2$ .

- Si studi il condizionamento del calcolo di  $f(x)$ .
- Si fornisca un algoritmo per il calcolo di  $f(x)$  e se ne studi la stabilità.
- È noto che

$$\sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6}$$

Per  $n = 10^8$  ed  $n = 10^9$  si consideri la somma parziale  $S_n = \sum_{k=1}^n \frac{1}{k^2}$  e si confrontino i valori ottenuti utilizzando con Matlab implementando due differenti algoritmi di somma (dal termine più piccolo a quello più grande e viceversa) e si confrontino gli errori relativi ottenuti con i due programmi.

*Esercizio 7.* Quale delle due espressioni equivalenti (dal punto di vista matematico)

$$\frac{1 - x^4}{1 - x} = 1 + x + x^2 + x^3$$

è più stabile? Si studi anche il condizionamento del calcolo dell'espressione precedente.

*Esercizio 8.* (\*) Per un qualsiasi  $x > 0$  si dica quanto dovrebbe valere  $x$  dopo aver eseguito il seguente frammento di programma

```
for k=1:60
    x=sqrt(x);
end
for k=1:60
    x=x*x;
end
```

Si esegua questo script in Matlab e si veda quale valore assume  $x$  alla fine del programma e si cerchi di dare una motivazione del risultato ottenuto.

*Esercizio 9.* Si vuole calcolare il valore del polinomio

$$p(x) = x^n + 2x^{n-1} + 2^2x^{n-2} + \dots + 2^{n-1}x + 2^n, \quad n > 1,$$

con il metodo di Ruffini-Horner, ovvero utilizzando l'algoritmo individuato nel seguente modo:

$$p(x) = x(x(x(\dots(x+2)+4)\dots) + 2^{n-1}) + 2^n,$$

per  $0 < x < 1$ .

- a) Si studi l'errore algoritmico per il caso  $n = 3$ , in cui  $p(x) = x(x(x+2)+4)+8$ .
- b\*) Si studi l'errore algoritmico per  $n$  generico.