

Calcolo Numerico

Luca De Paulis

10 marzo 2021

INDICE

1	ARITMETICA DI MACCHINA	3
1.1	Rappresentazione dei numeri di macchina	3
1.1.1	Virgola fissa	3
1.1.2	Virgola mobile	3
1.1.3	Insieme dei numeri di macchina	6
1.1.4	Standard IEEE (754)	7
1.2	Errori di approssimazione	9
1.3	Aritmetica di macchina	11
1.4	Teoria degli errori	14
1.4.1	Errore inerente	14
1.4.2	Errore algoritmico	17
1.4.3	Errore totale	18
1.5	Esempi di condizionamento e stabilità	19
1.5.1	Problema del calcolo della somma	20
2	ALGEBRA LINEARE NUMERICA	22
2.1	Norme vettoriali	22
2.1.1	Norme vettoriali principali	23
2.1.2	Equivalenza topologica tra le norme di \mathbb{F}^n	24
2.2	Norme matriciali	24
A	APPENDICE A	26
A.1	"9-periodico"	26
A.2	La precisione di macchina è un limite superiore stretto	26
A.3	Dimostrazione dell'equivalenza topologica tra norme	27
B	PREREQUISITI DI ALGEBRA LINEARE	29
B.1	Spazi vettoriali	29
B.2	Indipendenza e basi	30
B.3	Applicazioni lineari	31

1

ARITMETICA DI MACCHINA

1.1 RAPPRESENTAZIONE DEI NUMERI DI MACCHINA

Il primo problema da risolvere quando si vuole fare analisi numerica è scegliere un metodo per rappresentare i numeri reali su una macchina. Infatti un generico numero reale ha potenzialmente una scrittura decimale infinita, dunque essendo le risorse disponibili in una macchina *finite* dobbiamo trovare un modo di approssimarlo.

1.1.1 Virgola fissa

Il primo metodo è il metodo a **virgola fissa**: in questo metodo si rappresentano tutti i numeri nella loro forma decimale normale e si considerano esattamente k cifre dopo la virgola.

Questo metodo è molto semplice e ci consente di fare operazioni elementari (come le somme o i prodotti) immediatamente ("in colonna"), tuttavia ha anche degli svantaggi evidenti, come

- il range dei numeri rappresentabili su n cifre/bit è piccolo;
- siccome il numero di bit dedicato ai numeri dopo la virgola è basso, la precisione è molto bassa e assolutamente non adeguata ad applicazioni di analisi numerica.

Per questo è stata inventata la rappresentazione in virgola mobile.

1.1.2 Virgola mobile

Innanzitutto dobbiamo trovare un modo standard per rappresentare un qualsiasi numero reale. Per far ciò ci viene in aiuto il seguente teorema.

Teorema Teorema di rappresentazione in base.

1.1.1 Sia $x \in \mathbb{R}$, $x \neq 0$ e sia $\beta \geq 2$ una base di rappresentazione.

Allora esistono e sono univocamente determinati un **esponente** $p \in \mathbb{Z}$ e una successione di **cifre** $(d_i)_{i \in \mathbb{N}}$ (con $d_i \in \mathbb{Z}$ per ogni i) per cui vale che

(i) $d_1 \neq 0$;

(ii) $0 \leq d_i \leq \beta - 1$ per ogni i ;

(iii) la successione d non è definitivamente uguale a $\beta - 1$, ovvero per ogni $k > 0$ esiste un $j \geq k$ tale che $d_j \neq \beta - 1$;

(iv) si ha che

$$x = \text{sgn}(x) \cdot \beta^p \cdot \left(\sum_{i=1}^{\infty} d_i \beta^{-i} \right).$$

La rappresentazione di x nella forma data dal [Teorema 1.1.1](#) si dice **rappresentazione normalizzata in virgola mobile**.

Esempio 1.1.2. Consideriamo ad esempio il numero reale 123. Per il [Teorema 1.1.1](#) possiamo esprimere 123 come

$$123 = +10^3 \cdot (0.123) = +10^3 \cdot (1 \cdot 10^{-1} + 2 \cdot 10^{-2} + 3 \cdot 10^{-3}).$$

Questa rappresentazione è l'unica possibile se imponiamo che d_1 sia diverso da 0 e che le cifre della rappresentazione siano tutte considerate come *cifre decimali* moltiplicate per una certa potenza della base (in questo caso 3).

La rappresentazione data dal [Teorema 1.1.1](#) è essenzialmente il concetto di *notazione scientifica*: in notazione scientifica portiamo il numero reale in modo che abbia una singola cifra diversa da zero prima della virgola, mentre nella notazione data dal Teorema quella cifra diventa la prima cifra dopo la virgola. Inoltre, il Teorema ci garantisce che questa rappresentazione non è solamente valida in base 10, ma lo è in qualsiasi base (noi lo useremo in particolare in base 2).

Il fattore β^p viene detto **esponente**, mentre la parte decimale (corrispondente alla sommatoria) viene detta **mantissa**.

Osservazione 1.1.1. Le ipotesi (i) e (iii) del [Teorema](#) sono fondamentali per l'unicità della rappresentazione. In effetti

- se venisse a mancare la prima condizione avremmo che

$$123 = +10^3 \cdot (0.123) = +10^4 \cdot (0.0123) = +10^5 \cdot (0.00123) = \dots;$$

- se venisse a mancare la terza condizione (che ci dice che la successione d non può terminare con una sequenza infinita di $(\beta - 1)$, ovvero il numero non può finire con $\beta - 1$ periodico) avremmo dei problemi più subdoli, che derivano dalla non esistenza del "9 periodico" (dimostrata in appendice, [Proposizione A.1.1](#)).

Infatti in qualsiasi base β il numero $0.\overline{(\beta - 1)}$ è uguale a 1, quindi se ammettessimo entrambe le rappresentazioni verrebbe meno l'unicità.

Notazione. Useremo la notazione $(n)_\beta$ per riferirci al numero n espresso in base β .

Ad esempio il numero $(1011)_2$ si riferisce al numero

$$1 \cdot 2^3 + 0 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = 8 + 2 + 1 = 11 \text{ (in base 10).}$$

Per rappresentare i numeri reali in macchina dobbiamo quindi risolvere ancora due problemi:

- la macchina usa (nella maggior parte dei casi) la base 2, quindi dobbiamo poter trasformare da base 10 in base 2;
- dobbiamo rappresentare i numeri infiniti in modo approssimato.

Per quanto riguarda il primo, l'algoritmo per trasformare un numero decimale (in base 10) in un numero in base 2 è il seguente:

1. si trasforma la parte intera (le cifre prima della virgola) in base 2 tramite divisioni successive per 2;
2. si trasforma la parte decimale (le cifre dopo la virgola) in base 2 tramite *moltiplicazioni* successive per 2.

Esempio 1.1.3. Trasformiamo 3.15 in base 2: ovviamente $(3)_{10} = (11)_2$ e quindi rimane solo da trasformare la parte decimale.

L'algoritmo consiste nel prendere il numero decimale 0.15 e moltiplicarlo per 2 ripetutamente: se il risultato ha come cifra prima della virgola uno 0 aggiungiamo uno 0 alla rappresentazione in base 2, altrimenti un 1. Nella pratica:

$$\begin{aligned} 0.15 \cdot 2 &= \boxed{0}.30 \\ 0.30 \cdot 2 &= \boxed{0}.60 \\ 0.60 \cdot 2 &= \boxed{1}.20 \\ 0.20 \cdot 2 &= \boxed{0}.40 \\ 0.40 \cdot 2 &= \boxed{0}.80 \\ 0.80 \cdot 2 &= \boxed{1}.60 \\ 0.60 \cdot 2 &= \boxed{1}.20 \\ &\vdots \end{aligned}$$

Osserviamo che ripetendo il procedimento la sequenza (0.60, 0.20, 0.40, 0.80) si ripete all'infinito, dunque il numero in base 2 termina con le cifre 1001 ripetute periodicamente. La parte decimale in base 2 corrisponde quindi a $(0.00\overline{1001})_2$.

Il numero completo è quindi

$$(11.00\overline{1001})_2 = 2^2 \cdot (0.1100\overline{1001})_2.$$

Per rappresentare i numeri infiniti dobbiamo approssimare la parte decimale, considerando solo un numero finito t di cifre dopo la virgola. Per far ciò esistono due metodi (che esamineremo più nel dettaglio nel seguito):

- nel caso del **troncamento** si considerano le prime t cifre dopo la virgola e si scartano tutte le successive;
- nel caso dell'**arrotondamento** analizziamo la cifra decimale di posto $t + 1$:
 - se d_{t+1} appartiene all'intervallo $\left[0, \frac{\beta}{2}\right)$ (dove β è la base) consideriamo semplicemente le prime t cifre;
 - altrimenti (se $d_{t+1} \in \left[\frac{\beta}{2}, \beta\right]$) aggiungiamo 1 all'ultima cifra decimale (come riporto, quindi lo propaghiamo in caso sia necessario).

Esempio 1.1.4. Riprendiamo il numero considerato precedentemente, ovvero

$$(3.15)_{10} = 2^2 \cdot (0.1100\overline{1001})_2$$

e approssimiamolo a $t = 8$ cifre decimali.

Nel caso del troncamento basta considerare le prime 8 cifre decimali, quindi il numero troncato è

$$2^2 \cdot (0.11001001)_2.$$

Nel caso dell'arrotondamento invece dobbiamo considerare la nona cifra decimale: espandendo la rappresentazione decimale osserviamo che la nona cifra è

$$0.11001001\underline{1}00\dots;$$

siccome $1 \in \left[\frac{2}{2}, 2\right] = [1, 2]$ dobbiamo *approssimare per eccesso*, ottenendo

$$2^2 \cdot (0.11001010)_2.$$

1.1.3 Insieme dei numeri di macchina

Siccome un singolo numero può occupare una quantità fissa in memoria (tipicamente 32 o 64 bit) dobbiamo fissare dei limiti per l'esponente e per il numero di cifre decimali che possiamo usare per la rappresentazione. Osserviamo inoltre che, se le nostre risorse sono limitate, aumentando il numero di cifre decimali disponibili dobbiamo necessariamente diminuire lo spazio dedicato a memorizzare l'esponente e viceversa.

Definizione 1.1.5 **Insieme dei numeri di macchina.** Sia $\beta \geq 2$ una base, t il numero di cifre decimali utilizzabili e siano $m, M \in \mathbb{Z}$ gli estremi per l'esponente (ovvero ogni esponente p rappresentabile deve appartenere a $[-m, M]$). Si dice allora **insieme dei numeri di macchina** l'insieme

$$\mathcal{F}(\beta, t, m, M) := \{0\} \cup \left\{ x \in \mathbb{R} : x = \text{sgn}(x) \cdot \beta^p \sum_{i=1}^{\infty} d_i \beta^{-i}, \right. \\ \left. \begin{aligned} & d_1 \neq 0, \\ & 0 \leq d_i \leq \beta - 1, \\ & -m \leq p \leq M \end{aligned} \right\}.$$

Osservazione 1.1.2. Nelle condizioni in cui specifichiamo quali numeri sono rappresentabili compaiono le ipotesi (i) e (ii) del [Teorema di rappresentazione in base](#), ma non l'ipotesi (iii). Quest'ultima deriva dal fatto che i numeri di macchina hanno una precisione finita (ovvero hanno al più t cifre decimali), dunque non è possibile avere un numero che termina in $\beta - 1$ periodico.

OMEGA GRANDE Chiamiamo Ω il più grande numero di macchina. Possiamo costruirlo in questo modo:

- scegliamo come segno $+$;
- scegliamo come esponente il massimo esponente possibile, ovvero $p = M$;
- poniamo tutte le t cifre decimali uguali al massimo possibile, che in base β è $\beta - 1$.

Segue quindi che

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

Semplificando l'espressione di Ω tramite la formula per la serie geometrica si ottiene

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

Moltiplicando numeratore e denominatore della frazione per β :

$$\begin{aligned} &= \beta^M(\beta - 1) \left(\frac{\beta - (1/\beta)^t}{\beta - 1} - 1 \right) \\ &= \beta^M(\beta - 1) \cdot \frac{\beta - (1/\beta)^t - \beta + 1}{\beta - 1} \\ &= \beta^M(1 - \beta^{-t}). \end{aligned}$$

Esempio 1.1.6. Consideriamo l'insieme $\mathcal{F}(10, 3, 2, 2)$: in questo sistema possiamo rappresentare numeri in base 10 con al massimo 3 cifre decimali e con un esponente compreso tra -2 e 2 .

Il massimo numero rappresentabile in questo sistema è

$$\Omega = +10^2(0.999) = 10^2(1 - 10^{-3}).$$

OMEGA PICCOLO Chiamiamo ω il più piccolo numero *positivo* rappresentabile in un sistema. Possiamo costruirlo in questo modo:

- siccome è positivo, il segno è necessariamente $+$;
- come esponente scegliamo il più piccolo esponente possibile, ovvero $p = -m$;
- siccome le cifre decimali non possono essere tutte uguali a 0 poniamo $d_1 = 1$ e $d_i = 0$ per ogni $i > 1$.

Segue quindi che

$$\omega = +\beta^{-m}(1 \cdot \beta^{-1}) = \beta^{-m}(0.1)_\beta.$$

CARDINALITÀ DELL'INSIEME DEI NUMERI DI MACCHINA Dati β , t , m , M calcoliamo $\#\mathcal{F}(\beta, t, m, M)$. Osserviamo che

1. dalla definizione di $\mathcal{F}(\beta, t, m, M)$, dobbiamo contare lo 0 separatamente;
2. per quanto riguarda i numeri non-nulli, per ogni possibile combinazione di esponente/mantissa abbiamo 2 scelte per il segno (segno positivo o negativo);
3. le scelte per i possibili esponenti corrispondono a tutti i numeri interi nell'intervallo $[-m, M]$, che sono $M + m + 1$;
4. per quanto riguarda la mantissa dobbiamo scegliere t cifre comprese tra 0 e $\beta - 1$ (e lo possiamo fare in β^t modi), ma dobbiamo escludere tutte quelle che iniziano per 0, che sono β^{t-1} : abbiamo quindi $\beta^t - \beta^{t-1}$ scelte per la mantissa di un numero.

In totale si ha quindi che

$$\#\mathcal{F}(\beta, t, m, M) = 1 + 2(M + m + 1)(\beta^t - \beta^{t-1}).$$

Esempio 1.1.7. Ad esempio se $\beta = 2$, $t = 3$ e $m = M = 1$ si ha che

$$\#\mathcal{F}(2, 3, 1, 1) = 1 + 2(1 + 1 + 1)(2^3 - 2^2) = 25.$$

1.1.4 Standard IEEE (754)

Lo standard usato al giorno d'oggi per implementare i numeri a virgola mobile è lo standard IEEE (754). Esso definisce due rappresentazioni distinte, una *single precision*, dove ogni numero è rappresentato su 32 bit, ovvero su 4 byte, e una *double precision*, dove ogni numero è rappresentato su 64 bit (o equivalentemente 8 byte). In particolare noi analizzeremo solo il secondo.

I 64 bit vengono divisi nel seguente modo:

- un singolo bit di segno (0 per i numeri positivi, 1 per i negativi);
- 11 bit per rappresentare l'esponente;
- i restanti 52 bit vengono usati per rappresentare la mantissa.

Tuttavia possiamo ottimizzare un po' lo spazio usato facendo delle semplici osservazioni:

1. siccome siamo in base 2 la condizione $d_1 \neq 0$ ci dice che necessariamente $d_1 = 1$: ciò significa che rappresentare d_1 nella mantissa è inutile (ogni mantissa inizierebbe per 1) e quindi possiamo rappresentare le cifre da d_2 in poi. Questo significa che abbiamo in realtà $t = 53$ cifre decimali a disposizione;
2. per evitare di usare un bit per il segno dell'esponente lo rappresentiamo in *traslazione* (anche detto *bias* o *eccesso a*): invece di rappresentare il vero esponente p rappresentiamo $p + 1022$. Questo ci consente di indicare esponenti negativi senza dover usare strategie di complemento a 2.

L'insieme dei numeri di macchina definito dallo standard IEEE (754) è quindi

$$\mathcal{F}_{\text{IEEE}} := \mathcal{F}(2, 53, 1021, 1024).$$

Osserviamo che questa scelta per i limiti degli esponenti non copre tutte le combinazioni possibili su 11 bit: in particolare possiamo rappresentare $1024 + 1021 + 1 = 2046$ esponenti diversi, mentre su 11 bit potenzialmente potremmo rappresentarne fino a $2^{11} = 2048$.

Questo ci permette di rappresentare dei numeri particolari.

ZERO Abbiamo visto che (per il [Teorema 1.1.1](#)) lo 0 non ammette una rappresentazione normalizzata. Nello standard IEEE (754) esso viene realizzato ponendo a 0 tutti i campi dell'esponente e della mantissa; tuttavia non viene specificato il segno, per cui esistono due rappresentazioni uguali per lo zero (-0 e $+0$).

Questa rappresentazione non va in conflitto con i numeri normalizzati, in quanto ponendo il campo dell'esponente uguale a undici 0 otterrei un esponente effettivo uguale a -1022 , mentre il minimo rappresentabile è -1021 .

NUMERI DENORMALIZZATI Se il campo dell'esponente contiene solo zeri ma la mantissa ha almeno un bit diverso da 0 stiamo rappresentando numeri *denormalizzati*. Un tale numero viene interpretato come se fosse $d_1 = 0$ e quindi ci consente di ottenere dei numeri in valore assoluto più piccoli di ω , risolvendo parzialmente i problemi di *underflow*.

Tuttavia questi numeri hanno una precisione minore dei numeri normalizzati, in quanto i primi k bit della mantissa possono essere tutti uguali a 0. Questo significa che invece di avere 53 cifre di precisione ne abbiamo $53 - k$ (dove k dipende dal numero denormalizzato scelto), dunque dobbiamo porre ancora più attenzione agli eventuali errori di precisione.

Il più grande numero denormalizzato è della forma

$$2^{-1021} (0.01 \dots 1)_2 = 2^{-1022} \overbrace{(0.\overbrace{1 \dots 1}^{52 \text{ volte}})}_2,$$

mentre il più piccolo numero denormalizzato è il numero

$$2^{-1021} (0.0 \dots 01)_2.$$

INFINITI E NAN Ponendo il campo dell'esponente uguale a $(1111111111)_2$ (ovvero undici cifre 1) ottengo una rappresentazione che non appartiene all'insieme definito prima (in quanto il massimo esponente rappresentabile è 1024, che in eccesso a 1022 diventa $2046 = (1111111110)_2$). Sfruttando questo fatto possiamo rappresentare tre "numeri" diversi:

- se il campo della mantissa non è formato da tutti 0 il risultato è NaN (*Not a Number*): questa configurazione viene usata ad esempio per i risultati delle forme indeterminate, come $0/0$;
- se il campo della mantissa è formato da tutti 0 il numero viene interpretato come $\pm\infty$ (a seconda del bit di segno): ciò viene usato per rappresentare l'overflow, ovvero numeri che sono più grandi di Ω (o più piccoli di $-\Omega$).

1.2 ERRORI DI APPROSSIMAZIONE

Abbiamo iniziato ad introdurre il concetto di approssimazione nella [sezione 1.1](#). Definiamolo ora formalmente.

Definizione 1.2.1 **Troncamento e arrotondamento.** Sia $x \in \mathbb{R}$, $x \neq 0$, tale che $\omega < |x| < \Omega$. Consideriamo inoltre la rappresentazione normalizzata di x :

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

Allora definiamo

- (1) il **troncamento** di x alla t -esima cifra, indicato con $\text{trun}(x)$, dato da

$$\text{trun}(x) := \text{sgn}(x)\beta^p \sum_{i=1}^t d_i \beta^{-i};$$

- (2) l'**arrotondamento** di x , indicato con $\text{arr}(x)$, dato da

$$\text{arr}(x) := \begin{cases} \text{trun}(x), & \text{se } 0 \leq d_{t+1} < \beta/2 \\ \text{trun}(x) + \beta^{p-t}, & \text{se } d_{t+1} \geq \beta/2. \end{cases}$$

Osservazione 1.2.1. Aggiungere β^{p-t} ad un numero significa aggiungere 1 alla sua t -esima cifra. Infatti (considerando lo stesso x della definizione) si ha che

$$\beta^p \sum_{i=1}^{\infty} d_i \beta^{-i} + \beta^{p-t} = \beta^p \left(\sum_{i=1}^{\infty} d_i \beta^{-i} + \beta^{-t} \right).$$

Esempio 1.2.2. Consideriamo il numero reale

$$\pi = 3.141592653 \dots = 10^1 \cdot (0.3141592653 \dots)_{10}.$$

Il suo troncamento alla quinta cifra sarà 3.1415, mentre il suo arrotondamento alla stessa cifra (siccome la sesta cifra è $9 \geq 10/2$) sarà

$$3.1415 + 10^{1-5} = 3.1415 + 0.0001 = 3.1416.$$

D'ora in avanti quando parleremo dell'approssimazione di un numero reale x senza voler necessariamente specificare il metodo (troncamento o arrotondamento) usato scriveremo $\text{fl}(x)$ oppure \tilde{x} .

Approssimando numeri reali a numeri di macchina si genera un errore, detto **errore di rappresentazione**.

Definizione 1.2.3 **Errore di rappresentazione.** Sia $x \in \mathbb{R}$, $x \neq 0$ tale che $\omega \leq |x| \leq \Omega$. Chiameremo **errore assoluto** la quantità

$$\eta_x := \tilde{x} - x,$$

mentre chiameremo **errore relativo** o **errore di rappresentazione** la quantità

$$\varepsilon_x := \frac{\tilde{x} - x}{x}.$$

La qualità più importante dell'errore relativo è il fatto che è sempre **superiormente limitato**, come dimostreremo attraverso un lemma ed un teorema.

Lemma 1.2.4 **Limite superiore all'errore assoluto.** Sia $x \in \mathbb{R}$, $x \neq 0$, $\omega \leq x \leq \Omega$. Sia inoltre $x = \text{sgn}(x)\beta^p \sum_{i=1}^{\infty} d_i \beta^{-i}$ la rappresentazione in base di x . Si ha che

$$(i) \quad |\text{trun}(x) - x| < \beta^{p-t}$$

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

Dimostrazione. Se fosse $x = \Omega$ oppure $x = \omega$ allora $\text{trun}(x) = \text{arr}(x) = x$, da cui l'errore assoluto $\tilde{x} - x$ è nullo.

Supponiamo quindi $\omega < x < \Omega$ (il caso $-\Omega < x < -\omega$ è speculare). Siccome valgono le ipotesi del [Teorema 1.1.1](#) possiamo scrivere

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

Esisteranno dunque due numeri di macchina $a, b \in \mathcal{F}$ tali che $a \leq x < b$ e b è il numero di macchina successivo rispetto ad a , ovvero a è il numero di macchina appena più piccolo di x e b è quello appena più grande.

Allora dovrà essere necessariamente

$$a = \beta^p \sum_{i=1}^t d_i \beta^{-i} = \text{trun}(x),$$

mentre b dovrà essere

$$b = \beta^p \left(\sum_{i=1}^t d_i \beta^{-i} + \beta^{-t} \right) = a + \beta^{p-t},$$

da cui segue anche che $b - a = \beta^{p-t}$.

Si ha quindi

$$|\text{trun}(x) - x| = |a - x| < |a - b| = \beta^{p-t},$$

dove il minore centrale deriva dal fatto che $x < b$.

Per quanto riguarda l'arrotondamento, osserviamo invece che

- se $d_{t+1} < \beta/2$ (ovvero se $x < 1/2(a + b)$) allora $\text{arr}(x) = \text{trun}(x) = a$;
- altrimenti $\text{arr}(x) = \text{trun}(x) + \beta^{p-t} = a + \beta^{p-t} = b$.

Nel primo caso si ha

$$|\text{arr}(x) - x| = |x - a| \leq \left| \frac{a+b}{2} - a \right| = \frac{1}{2}|b - a| = \frac{1}{2}\beta^{p-t},$$

mentre nel secondo si ha

$$|\text{arr}(x) - x| = |b - x| \leq \left| b - \frac{a+b}{2} \right| = \frac{1}{2}|b - a| = \frac{1}{2}\beta^{p-t},$$

da cui la tesi. □

Osservazione 1.2.2. Si ha l'uguaglianza $|\text{arr } x - x| = \frac{1}{2}\beta^{p-t}$ se e solo se $x = \frac{1}{2}(a + b)$, ovvero se e solo se x è esattamente il punto medio dell'intervallo $[a, b]$.

Teorema 1.2.5 **Esistenza della precisione di macchina.** *Esiste un numero $u \in \mathbb{R}$, detto **precisione di macchina**, tale che per ogni $x \in \mathbb{R}$, $x \neq 0$, $\omega \leq |x| \leq \Omega$ vale che*

$$|\varepsilon_x| \leq u.$$

In particolare

- (i) se $\text{fl}(x) = \text{trun}(x)$ si ha $u = \beta^{1-t}$;
- (ii) se $\text{fl}(x) = \text{arr}(x)$ si ha $u = \frac{1}{2}\beta^{1-t}$.

Dimostrazione. Osserviamo che se $x = \pm\beta^p \sum_{i=1}^{\infty} d_i \beta^{-i}$ allora varrà sicuramente che

$$|x| \geq \beta^p (0.10\dots)_\beta = \beta^p \cdot (1\beta^{-1}) = \beta^{p-1}.$$

Infatti scegliendo il numero di macchina con lo stesso esponente di x ma la minima mantissa possibile otterremo necessariamente un numero più piccolo (o uguale) del valore assoluto di quello di partenza.

Consideriamo allora i due casi.

TRONCAMENTO Per il [Lemma 1.2.4](#) sappiamo che $|\text{trun}(x) - x| < \beta^{p-t}$. Si ha quindi che

$$|\varepsilon_x| = \frac{|\text{trun}(x) - x|}{|x|} < \frac{\beta^{p-t}}{\beta^{p-1}} = \beta^{1-p}.$$

ARROTONDAMENTO Per il [Lemma 1.2.4](#) sappiamo che $|\text{arr}(x) - x| \leq \frac{1}{2}\beta^{p-t}$. Si ha quindi che

$$|\varepsilon_x| = \frac{|\text{arr}(x) - x|}{|x|} \leq \frac{1/2\beta^{p-t}}{\beta^{p-1}} = \frac{1}{2}\beta^{1-p}. \quad \square$$

Osserviamo che il [Teorema 1.2.5](#) ci dice che u è un limite superiore *debole* (nel senso che $|\varepsilon_x| \leq u$ e non $|\varepsilon_x| < u$), tuttavia almeno nel caso del troncamento si vede dalla dimostrazione che vale anche la disuguaglianza stretta.

In realtà la disuguaglianza stretta vale anche nel caso dell'arrotondamento: la dimostrazione completa è lasciata all'appendice (in particolare nella [sezione A.2](#)).

1.3 ARITMETICA DI MACCHINA

Cerchiamo ora di estendere l'aritmetica tra numeri reali ad un'aritmetica per i numeri di macchina. Il primo tentativo è quello di usare semplicemente le solite operazioni definite tra numeri reali, tuttavia non è detto che il risultato di un'operazione tra numeri di macchina sia ancora un numero di macchina.

Esempio 1.3.1. Consideriamo il sistema $\mathcal{F}(10, 3, -5, 5)$ e i due numeri di macchina

$$a := 123 = 10^3(0.123)_{10}, \quad b := 0.456 = 10^0(0.456).$$

Si ha che $a, b \in \mathcal{F}(10, 3, -5, 5)$, tuttavia $a + b = 123.456 = 10^3(0.123456)_{10}$ non è un numero di macchina poiché è rappresentato con $t = 6$ cifre di precisione.

Dobbiamo quindi definire delle nuove operazioni che siano *interne* all'insieme dei numeri di macchina.

Definizione 1.3.2 **Aritmetica di macchina.** Si definiscono quattro operazioni

$$\oplus, \otimes, \ominus, \oslash : \mathcal{F} \times \mathcal{F} \rightarrow \mathcal{F}$$

tali che per ogni coppia di numeri di macchina $a, b \in \mathcal{F}$ ($b \neq 0$ nel caso di \oslash) si ha che

$$\begin{aligned} a \oplus b &:= \text{fl}(a + b) & a \ominus b &:= \text{fl}(a - b) \\ a \otimes b &:= \text{fl}(ab) & a \oslash b &:= \text{fl}(a/b) \end{aligned}$$

Siamo ora interessati a studiare *quanto errore* commettiamo nel considerare i numeri di macchina invece dei numeri reali non approssimati.

Definizione 1.3.3 **Errore locale.** Sia $*$ un'operazione tra numeri qualsiasi e sia \otimes la corrispondente operazione su numeri di macchina. Dati $x, y \in \mathcal{F}$ si dice **errore locale** dell'operazione la quantità

$$\varepsilon := \frac{(x \otimes y) - (x * y)}{x * y}.$$

Osserviamo che ogni operazione di macchina può essere scritta in termini della sua corrispondente operazione base e dell'errore locale: in effetti

$$x \otimes y = (x * y)(1 + \varepsilon).$$

Esempio 1.3.4 (Errore nel calcolo della somma). Cerchiamo ora di calcolare l'errore totale commesso nell'approssimazione di un'operazione tra numeri reali (definiremo più avanti precisamente il concetto di errore totale). Prendiamo $a, b \in \mathbb{R}$ che approssimati non vadano in overflow (o underflow) e cerchiamo di calcolare $a + b$ in termini di numeri di macchina.

Siccome i numeri a, b non sono necessariamente già numeri di macchina, l'unico modo per svolgere il calcolo è calcolare la quantità

$$\tilde{a} \oplus \tilde{b} = (\tilde{a} + \tilde{b})(1 + \varepsilon).$$

Siccome sappiamo che, se $\varepsilon_a, \varepsilon_b$ sono gli errori di rappresentazione rispettivamente di a e b

$$\tilde{a} = a(1 + \varepsilon_a), \quad \tilde{b} = b(1 + \varepsilon_b)$$

otteniamo che

$$\begin{aligned} \tilde{a} \oplus \tilde{b} &= (a(1 + \varepsilon_a) + b(1 + \varepsilon_b))(1 + \varepsilon) \\ &= a(1 + \varepsilon_a)(1 + \varepsilon) + b(1 + \varepsilon_b)(1 + \varepsilon). \end{aligned}$$

Nello svolgere il prossimo passaggio trascureremo gli addendi in cui compare il prodotto di due errori: effettueremo quindi un'**analisi al primo ordine**, indicata dal simbolo \doteq .

$$\doteq a + b + a\varepsilon_a + b\varepsilon_b + (a + b)\varepsilon.$$

L'errore complessivo che commettiamo è quindi dato dalla differenza tra il risultato ottenuto ($\tilde{a} \oplus \tilde{b}$) e il risultato non approssimato ($a + b$), dividendo per il risultato non approssimato in modo da ottenere un errore relativo:

$$\begin{aligned} \frac{(\tilde{a} \oplus \tilde{b}) - (a + b)}{a + b} &\doteq \frac{a + b + a\varepsilon_a + b\varepsilon_b + (a + b)\varepsilon - (a + b)}{a + b} \\ &= \frac{a}{a + b}\varepsilon_a + \frac{b}{a + b}\varepsilon_b + \varepsilon. \end{aligned}$$

Notiamo che questo errore può essere diviso in due parti diverse:

- la parte data da $\frac{a}{a+b}\varepsilon_a + \frac{b}{a+b}\varepsilon_b$ non dipende dal tipo di operazione che sto cercando di svolgere, ma dal fatto che devo approssimare gli argomenti dell'operazione: lo chiameremo **errore inerente**;

- la parte data da ε dipende dalla sequenza di passi scelta per calcolare la somma: lo chiameremo perciò **errore algoritmico**.

Esempio 1.3.5 (Errore di una funzione). Consideriamo la funzione $f: \mathbb{R} \rightarrow \mathbb{R}$ definita da

$$f(x) = x^2 - 1.$$

Calcoliamo l'errore totale commesso nel calcolare $f(x)$ tramite operazioni di macchina.

Osserviamo che potenzialmente abbiamo più *algoritmi* (intesi come sequenze di calcoli) per calcolare f : scegliamo ad esempio gli algoritmi di macchina g_1, g_2 definiti da

$$g_1(\tilde{x}) := (\tilde{x} \otimes \tilde{x}) \ominus 1, \quad g_2(\tilde{x}) := (\tilde{x} \ominus 1) \otimes (\tilde{x} \oplus 1).$$

L'algoritmo g_1 prescrive di calcolare innanzitutto $z := \tilde{x} \otimes \tilde{x}$ e dopo calcolare $z \ominus 1$. Se ε_1 è l'errore locale dell'operazione \otimes si ha che

$$z = \tilde{x}^2(1 + \varepsilon_1);$$

chiamando ε_2 l'errore locale dell'operazione \ominus si ottiene

$$g_1(\tilde{x}) = (z - 1)(1 + \varepsilon_2).$$

A questo punto, ricordando che $\tilde{x} = x(1 + \varepsilon_x)$ dove ε_x è l'errore di rappresentazione, si può esprimere $g_1(\tilde{x})$ in termini di x e dei vari errori. In effetti

$$\begin{aligned} & (z - 1)(1 + \varepsilon_2) \\ &= (\tilde{x}^2(1 + \varepsilon_1) - 1)(1 + \varepsilon_2) \\ &= \left((x(1 + \varepsilon_x))^2(1 + \varepsilon_1) - 1 \right)(1 + \varepsilon_2) \\ &= (x^2(1 + \varepsilon_x)^2(1 + \varepsilon_1) - 1)(1 + \varepsilon_2) \\ &\doteq (x^2(1 + 2\varepsilon_x)(1 + \varepsilon_1) - 1)(1 + \varepsilon_2) \\ &\doteq (x^2(1 + 2\varepsilon_x + \varepsilon_1) - 1)(1 + \varepsilon_2) \\ &= x^2(1 + 2\varepsilon_x + \varepsilon_1)(1 + \varepsilon_2) - (1 + \varepsilon_2) \\ &\doteq x^2(1 + 2\varepsilon_x + \varepsilon_1 + \varepsilon_2) - (1 + \varepsilon_2) \\ &= x^2 - 1 + 2x^2\varepsilon_x + x^2\varepsilon_1 + (x^2 - 1)\varepsilon_2. \end{aligned}$$

L'errore totale commesso nell'approssimare il risultato della funzione $f(x)$ con la funzione di macchina $g_1(\tilde{x})$ è dunque

$$\begin{aligned} \varepsilon_{g_1} &:= \frac{x^2 - 1 + 2x^2\varepsilon_x + x^2\varepsilon_1 + (x^2 - 1)\varepsilon_2 - (x^2 - 1)}{x^2 - 1} \\ &= \frac{2x^2}{x^2 - 1}\varepsilon_x + \frac{x^2}{x^2 - 1}\varepsilon_1 + \varepsilon_2. \end{aligned}$$

Ripetiamo il procedimento per l'algoritmo g_2 : siano δ_1, δ_2 e δ_3 i tre errori locali commessi; ovvero

$$\begin{aligned} z_1 &:= \tilde{x} \ominus 1 = (\tilde{x} - 1)(1 + \delta_1) \\ z_2 &:= \tilde{x} \oplus 1 = (\tilde{x} + 1)(1 + \delta_2) \\ g_2(\tilde{x}) &= z_1 \otimes z_2 = z_1 z_2(1 + \delta_3). \end{aligned}$$

Si ha quindi

$$\begin{aligned}
 & z_1 z_2 (1 + \delta_3) \\
 &= (\tilde{x} - 1)(1 + \delta_1)(\tilde{x} + 1)(1 + \delta_2)(1 + \delta_3) \\
 &= (\tilde{x}^2 - 1)(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) \\
 &= (x^2(1 + \varepsilon_x)^2 - 1)(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) \\
 &\doteq (x^2(1 + 2\varepsilon_x) - 1)(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) \\
 &= (x^2(1 + 2\varepsilon_x)(1 + \delta_1) - (1 + \delta_1))(1 + \delta_2)(1 + \delta_3) \\
 &\doteq (x^2(1 + 2\varepsilon_x + \delta_1) - (1 + \delta_1))(1 + \delta_2)(1 + \delta_3) \\
 &= (x^2(1 + 2\varepsilon_x + \delta_1)(1 + \delta_2) - (1 + \delta_1)(1 + \delta_2))(1 + \delta_3) \\
 &\doteq (x^2(1 + 2\varepsilon_x + \delta_1 + \delta_2) - (1 + \delta_1 + \delta_2))(1 + \delta_3) \\
 &= x^2(1 + 2\varepsilon_x + \delta_1 + \delta_2)(1 + \delta_3) - (1 + \delta_1 + \delta_2)(1 + \delta_3) \\
 &\doteq x^2(1 + 2\varepsilon_x + \delta_1 + \delta_2 + \delta_3) - (1 + \delta_1 + \delta_2 + \delta_3) \\
 &= x^2 + 2x^2\varepsilon_x + x^2(\delta_1 + \delta_2 + \delta_3) - 1 - (\delta_1 + \delta_2 + \delta_3) \\
 &= x^2 - 1 + 2x^2\varepsilon_x + (x^2 - 1)(\delta_1 + \delta_2 + \delta_3).
 \end{aligned}$$

L'errore totale generato dalla funzione g_2 è dunque

$$\begin{aligned}
 \varepsilon_{g_2} &:= \frac{x^2 - 1 + 2x^2\varepsilon_x + (x^2 - 1)(\delta_1 + \delta_2 + \delta_3) - (x^2 - 1)}{x^2 - 1} \\
 &= \frac{2x^2}{x^2 - 1} \varepsilon_x + \delta_1 + \delta_2 + \delta_3.
 \end{aligned}$$

Osserviamo che entrambe le funzioni hanno lo stesso errore inerente (dato dal termine contenente ε_x): in effetti l'errore inerente deriva dall'approssimazione di x in \tilde{x} , e ciò non dipende dall'algoritmo scelto.

Il resto del termine è l'errore algoritmico: siccome abbiamo due algoritmi diversi possiamo chiederci quale dei due sia migliore.

Osserviamo che il secondo non dipende da x , mentre il primo sì: in particolare l'errore algoritmico del primo tende ad infinito quando x tende a ± 1 . Si dice quindi che l'algoritmo non è **stabile** per $x \rightarrow \pm 1$, mentre il secondo algoritmo è stabile per ogni $x \in \mathbb{R}$.

1.4 TEORIA DEGLI ERRORI

Diamo ora una definizione precisa dei vari tipi di errori.

Consideriamo nel resto della sezione

- una funzione razionale $f : A \rightarrow \mathbb{R}$ con $A \subseteq \mathbb{R}$, dove per *razionale* si intende che compaiono solo le operazioni di somma, prodotto, sottrazione e divisione
- un numero reale x qualunque e la sua approssimazione \tilde{x}
- un algoritmo $g : \mathcal{F} \rightarrow \mathcal{F}$ per calcolare f , dove per *algoritmo* si intende una sequenza di operazioni tra numeri di macchina usata per calcolare il valore di f .

1.4.1 Errore inerente

Definizione **Errore inerente.** Si dice **errore inerente** del problema la quantità

1.4.1

$$\varepsilon_{\text{in}} := \frac{f(\tilde{x}) - f(x)}{f(x)}.$$

Osserviamo che l'errore inerente misura quanto è grande l'errore nell'approssimare x con il numero di macchina \tilde{x} : se l'errore inerente è grande in un intorno di x nessun algoritmo ci consentirà di calcolare con una buona precisione il risultato della funzione f . In tal caso diremo che il problema è **mal condizionato**, mentre se l'errore inerente è sempre limitato il problema verrà definito **ben condizionato**.

Per studiare il condizionamento di un problema è quindi necessario studiare il suo errore inerente: vogliamo quindi una tecnica per calcolarlo più efficientemente.

Osserviamo che moltiplicando e dividendo per x e per $\tilde{x} - x$ si ha

$$\varepsilon_{\text{in}} = \frac{f(\tilde{x}) - f(x)}{f(x)} = \frac{f(\tilde{x}) - f(x)}{\tilde{x} - x} \cdot \frac{x}{f(x)} \cdot \underbrace{\frac{\tilde{x} - x}{x}}_{\varepsilon_x}.$$

Il primo termine del prodotto ricorda un rapporto incrementale: se la funzione è derivabile in x possiamo semplificare l'espressione usando il Teorema di Taylor.

Supponiamo quindi che f sia continua e derivabile *due volte* in un intervallo $[a, b]$ contenente x (ovvero $f \in \mathcal{C}^2[a, b]$): per il Teorema di Taylor con il resto di Lagrange si ha che esiste un numero ξ_x compreso tra x e \tilde{x} tale che

$$f(\tilde{x}) = f(x) + f'(x)(\tilde{x} - x) + f''(\xi_x) \frac{(\tilde{x} - x)^2}{2!}$$

con $|\xi_x - x| < |\tilde{x} - x|$.

Sostituendo otteniamo

$$\begin{aligned} \varepsilon_{\text{in}} &= \frac{f(\tilde{x}) - f(x)}{\tilde{x} - x} \cdot \frac{x}{f(x)} \cdot \varepsilon_x \\ &= \frac{f'(x)(\tilde{x} - x) + f''(\xi_x) \frac{(\tilde{x} - x)^2}{2!}}{\tilde{x} - x} \cdot \frac{x}{f(x)} \cdot \varepsilon_x \\ &= \left(f'(x) + f''(\xi_x) \frac{\tilde{x} - x}{2!} \right) \cdot \frac{x}{f(x)} \cdot \varepsilon_x \\ &= f'(x) \frac{x}{f(x)} \varepsilon_x + f''(\xi_x) \frac{\tilde{x} - x}{2!} \frac{x}{f(x)} \cdot \varepsilon_x. \end{aligned}$$

Osserviamo ora che il secondo addendo contiene (nascosto) un errore al quadrato, e pertanto in un'analisi al prim'ordine viene approssimato a 0. In effetti moltiplicando e dividendo per x si ha

$$\begin{aligned} & f''(\xi_x) \frac{\tilde{x} - x}{2!} \frac{x}{f(x)} \cdot \varepsilon_x \\ &= f''(\xi_x) \frac{\tilde{x} - x}{2!} \frac{x}{f(x)} \cdot \frac{x}{x} \cdot \varepsilon_x \\ &= f''(\xi_x) \frac{x^2}{2!f(x)} \cdot \frac{\tilde{x} - x}{x} \cdot \varepsilon_x \\ &= f''(\xi_x) \frac{x^2}{2!f(x)} \cdot \varepsilon_x^2 \\ &\doteq 0. \end{aligned}$$

Abbiamo quindi dimostrato il seguente Teorema.

Teorema 1.4.2 **Caratterizzazione dell'errore inerente.** Sia $f : A \subseteq \mathbb{R} \rightarrow \mathbb{R}$ e sia x il punto di cui vogliamo calcolare l'immagine $f(x)$. Se $f \in \mathcal{C}^2(A)$ si ha che

$$\varepsilon_{in} \doteq f'(x) \frac{x}{f(x)} \cdot \varepsilon_x,$$

dove ε_x è l'errore di rappresentazione di x .

La quantità $c_x := f'(x) \frac{x}{f(x)}$ viene detta **coefficiente di amplificazione**.

Abbiamo quindi scoperto che l'errore inerente è sempre della forma $c_x \varepsilon_x$: possiamo quindi definire formalmente il concetto di condizionamento.

Definizione 1.4.3 **Condizionamento di un problema.** Sia $f : A \subseteq \mathbb{R} \rightarrow \mathbb{R}$ con $f \in \mathcal{C}^2(A)$. Se esiste $x_0 \in A$ tale che

$$\lim_{x \rightarrow x_0} |c_x| = +\infty$$

si dice che il problema è **mal condizionato** in un intorno di x_0 , altrimenti si dice che il problema è **ben condizionato**.

Esempio 1.4.4. Riprendiamo l'[Esempio 1.3.5](#): il coefficiente di amplificazione è

$$c_x = \frac{2x^2}{x^2 - 1}$$

il cui valore assoluto tende a $+\infty$ per $x \rightarrow \pm 1$. Si ha quindi che il problema è mal condizionato in un intorno di 1 e in un intorno di -1 .

Il [Teorema 1.4.2](#) può essere generalizzato per funzioni $\mathbb{R}^n \rightarrow \mathbb{R}$.

Corollario 1.4.5 **Errore inerente su funzioni multivariate.** Sia $f : \Omega \rightarrow \mathbb{R}$ con $\Omega \subseteq \mathbb{R}^n$ aperto e f differenziabile due volte su Ω (ovvero esistono la derivata parziale prima e seconda rispetto ad ogni variabile).

Sia inoltre $\mathbf{x} = (x_1, \dots, x_n)$ il punto di cui vogliamo calcolare $f(\mathbf{x})$. Si ha che

$$\varepsilon_{in} = \sum_{i=1}^n c_{x_i} \varepsilon_{x_i}$$

dove $c_{x_i} := \frac{x_i}{f(\mathbf{x})} \cdot \frac{\partial f}{\partial x_i}$ è il **coefficiente di amplificazione**.

Coefficienti di amplificazione delle operazioni elementari

Il [Corollario 1.4.5](#) può essere molto comodo per calcolare i coefficienti di amplificazione delle operazioni elementari: ciò torna molto utile per il calcolo dell'errore totale.

SOMMA Consideriamo $f(x, y) = x + y$. Si ha che

$$c_x = \frac{x}{x+y} \cdot 1 = \frac{x}{x+y}, \quad c_y = \frac{y}{x+y} \cdot 1 = \frac{y}{x+y}.$$

SOTTRAZIONE Consideriamo $f(x, y) = x - y$. Si ha che

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

PRODOTTO Consideriamo $f(x, y) = xy$. Si ha che

$$c_x = \frac{x}{xy} \cdot y = 1, \quad c_y = \frac{y}{xy} \cdot x = 1.$$

DIVISIONE Consideriamo $f(x, y) = \frac{x}{y}$. Si ha che

$$c_x = \frac{x}{x/y} \cdot \frac{1}{y} = 1, \quad c_y = \frac{y}{x/y} \cdot \left(-\frac{x}{y^2}\right) = -1.$$

1.4.2 Errore algoritmico

Definizione 1.4.6 **Errore algoritmico.** Si dice **errore algoritmico** la quantità

$$\varepsilon_{\text{alg}} := \frac{g(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}.$$

L'errore algoritmico misura quanto è grande calcolare il risultato attraverso l'algoritmo di macchina g invece della funzione non approssimata f *partendo da dati già approssimati*.

Per studiare l'errore algoritmico si usa la tecnica conosciuta come *errore in avanti* (*forward error* in inglese), che può essere resa più compatta attraverso l'uso di grafi.

Formalmente un algoritmo è una sequenza di passi *finita* $z^{(1)}, z^{(2)}, z^{(3)}, \dots, z^{(n)}$, tale che ognuno dei passi calcola una singola operazione sfruttando solo i dati dei passi precedenti e i dati del problema (costanti e variabili) e l'ultimo passo calcola la funzione del problema.

Esempio 1.4.7. L'algoritmo $g(x) = ((x \otimes x) \oplus x) \ominus 3$ può essere espresso dalla sequenza di passi

1. $z^{(1)} = x \otimes x$
2. $z^{(2)} = z^{(1)} \oplus x$
3. $z^{(3)} = z^{(2)} \ominus 3$

Ogni passo che calcola un'operazione può quindi essere scritto nella forma

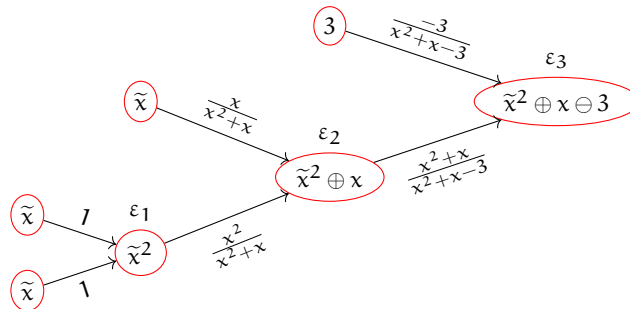
$$z^{(i)} = z^{(h)} \otimes z^{(k)}$$

con $h, k < i$ e \otimes un'operazione qualunque di macchina.

Esempio 1.4.8. Calcoliamo l'errore algoritmico della funzione $f(x) = x^2 + x - 3$ attraverso l'algoritmo di macchina

$$g(\tilde{x}) = ((\tilde{x} \otimes \tilde{x}) \oplus \tilde{x}) \ominus 3.$$

Per far ciò, costruiamo il seguente grafo.



Il grafo è costruito nel seguente modo:

- si considera la sequenza di operazioni da svolgere: in questo caso
 1. $z_1 := \tilde{x} \otimes \tilde{x}$
 2. $z_2 := z_1 \oplus \tilde{x}$
 3. $z_3 := z_2 \ominus 3 = g(\tilde{x})$

- si disegnano i nodi iniziali, ovvero quelli della prima operazione da svolgere, e sopra di essi si mettono gli errori commessi nel rappresentarli: siccome siamo interessati all'errore algoritmico partiamo già con numeri di macchina, quindi l'errore di rappresentazione è 0 e lo omettiamo;
- si crea un nodo corrispondente a z_1 : l'errore commesso nella creazione di questo nodo è l'errore locale del prodotto (qui chiamato ε_1);
- sopra gli archi che portano i nodi \tilde{x} in z_1 indichiamo il coefficiente di amplificazione c_x dei vari nodi: in questo caso sono entrambi uguali a 1 poiché l'operazione è un prodotto;
- ripetiamo il procedimento: creiamo un nuovo nodo con \tilde{x} e "errore di creazione" uguale a 0 e lo usiamo per fare l'operazione di somma, creando un nuovo nodo (che corrisponde a z_2) con errore di creazione uguale all'errore locale della somma ε_2 ;
- per ricavare i coefficienti di amplificazione sfruttiamo le formule ricavate prima, che ci dicono che nella somma $\tilde{x}^2 \oplus \tilde{x}$ i coefficienti di amplificazione sono dati da

$$c_{x^2} = \frac{x^2}{x^2 + x}, \quad c_x = \frac{x}{x^2 + x};$$

- ripetiamo un'ultima volta il ragionamento per quanto riguarda l'ultima operazione.

Possiamo ora calcolare l'errore totale a partire dal grafo: si parte dalla fine del grafo (nodo più a destra) e si calcola la somma tra l'errore del nodo (in questo caso ε_3) e gli errori dei sottoalberi (calcolati ricorsivamente), moltiplicati per il loro coefficiente di amplificazione.

Otteniamo quindi

$$\begin{aligned} \varepsilon_{\text{alg}} &= \varepsilon_3 + \frac{x^2 + x}{x^2 + x - 3}(\dots) + \frac{-3}{x^2 + x - 3} \cdot 0 \\ &= \varepsilon_3 + \frac{x^2 + x}{x^2 + x - 3} \left(\varepsilon_2 + \frac{x^2}{x^2 + x}(\dots) + \frac{x}{x^2 + x} \cdot 0 \right) \\ &= \varepsilon_3 + \frac{x^2 + x}{x^2 + x - 3} \left(\varepsilon_2 + \frac{x^2}{x^2 + x}(\varepsilon_1 + 1 \cdot 0 + 1 \cdot 0) \right) \\ &= \varepsilon_3 + \frac{x^2 + x}{x^2 + x - 3} \left(\varepsilon_2 + \frac{x^2}{x^2 + x} \cdot \varepsilon_1 \right). \end{aligned}$$

1.4.3 Errore totale

Definizione 1.4.9 **Errore totale.** Si dice **errore totale** del problema la quantità

$$\varepsilon_{\text{tot}} := \frac{g(\tilde{x}) - f(x)}{f(x)}.$$

Osserviamo che l'errore totale combina i concetti di errore inerente e algoritmico e ci dice quanto è grande l'errore nel calcolare g con il valore di \tilde{x} rispetto al calcolo esatto di $f(x)$.

Nell'[Esempio 1.3.5](#) abbiamo notato come l'errore totale fosse dato dalla somma dell'errore algoritmico e dell'errore inerente: questa relazione vale sempre, come dimostrato dal prossimo teorema.

Teorema 1.4.10

$$\varepsilon_{\text{tot}} \doteq \varepsilon_{\text{in}} + \varepsilon_{\text{alg}}.$$

Dimostrazione. In effetti si ha

$$\begin{aligned}
 \varepsilon_{\text{tot}} &= \frac{g(\tilde{x}) - f(x)}{f(x)} \\
 &= \frac{g(\tilde{x}) - f(\tilde{x}) + f(\tilde{x}) - f(x)}{f(x)} \\
 &= \frac{g(\tilde{x}) - f(\tilde{x})}{f(x)} + \frac{f(\tilde{x}) - f(x)}{f(x)} \\
 &= \frac{g(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})} \frac{f(\tilde{x})}{f(x)} + \frac{f(\tilde{x}) - f(x)}{f(x)} \\
 &= \varepsilon_{\text{alg}} \cdot \frac{f(\tilde{x})}{f(x)} + \varepsilon_{\text{in}} \\
 &= \varepsilon_{\text{alg}} \cdot \frac{f(\tilde{x}) - f(x) + f(x)}{f(x)} + \varepsilon_{\text{in}} \\
 &= \varepsilon_{\text{alg}} \cdot \left(\frac{f(\tilde{x}) - f(x)}{f(x)} + \frac{f(x)}{f(x)} \right) + \varepsilon_{\text{in}} \\
 &= \varepsilon_{\text{alg}} \cdot (\varepsilon_{\text{in}} + 1) + \varepsilon_{\text{in}} \\
 &\doteq \varepsilon_{\text{alg}} + \varepsilon_{\text{in}}.
 \end{aligned}$$

□

1.5 ESEMPI DI CONDIZIONAMENTO E STABILITÀ

In questa sezione faremo alcuni esempi dello studio del condizionamento di problemi e della stabilità di algoritmi di macchina.

Esempio 1.5.1. Consideriamo la funzione $f(x) = \frac{x-1}{x}$ e studiamone il condizionamento per $x \neq 0$.

Osserviamo che la funzione è derivabile (almeno) due volte nell'intervallo $(0, +\infty)$, dunque possiamo calcolare il coefficiente di amplificazione di un dato $x \in (0, +\infty)$.

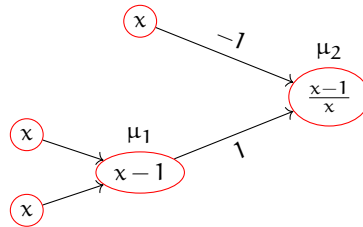
$$c_x = \frac{x}{f(x)} f'(x) = \frac{1}{x-1}.$$

Il problema è ben condizionato se e solo se $|c_x|$ è limitato superiormente per ogni $x \in (0, +\infty)$: tuttavia possiamo osservare che

$$\lim_{x \rightarrow 1} \left| \frac{1}{x-1} \right| = +\infty,$$

dunque il problema è mal condizionato in un intorno di 1.

Esempio 1.5.2. Consideriamo l'algoritmo $g_1(x) = (x \ominus 1) \oslash x$ e studiamone la stabilità. Il grafo di g_1 è



L'errore algoritmico totale è quindi dato da

$$\varepsilon_{\text{alg}} = \mu_2 + 1 \cdot (\mu_1 + 0 + 0) + 0 = \mu_2 + \mu_1.$$

Dimostriamo ora formalmente che l'algoritmo è stabile, ovvero che $|\varepsilon_{\text{alg}}|$ è sempre limitato superiormente.

In effetti si ha

$$|\varepsilon_{\text{alg}}| = |\mu_1 + \mu_2| \stackrel{(1)}{\leq} |\mu_1| + |\mu_2| \stackrel{(2)}{<} 2u$$

dove la disuguaglianza (1) è la disuguaglianza triangolare dei moduli, mentre la (2) viene dal fatto che gli errori locali sono sempre minori della precisione di macchina u .

1.5.1 Problema del calcolo della somma

Studiamo la funzione $f : \mathbb{R}_{>0}^n \rightarrow \mathbb{R}$ data da

$$f(x_1, \dots, x_n) := \sum_{i=1}^n x_i.$$

Osserviamo che la funzione ha come dominio $\mathbb{R}_{>0}^n$, dunque tutti gli x_i devono essere strettamente positivi. Inoltre per semplicità scriveremo $\mathbf{x} = (x_1, \dots, x_n)$.

CONDIZIONAMENTO Studiamo inizialmente il condizionamento del problema studiando l'errore inerente.

$$\begin{aligned} \varepsilon_{\text{in}} &\doteq \sum_{i=1}^n c_{x_i} \varepsilon_{x_i} \\ &= \sum_{i=1}^n \frac{x_i}{f(\mathbf{x})} \frac{\partial f}{\partial x_i} \cdot \varepsilon_{x_i} \\ &= \sum_{i=1}^n \frac{x_i}{f(\mathbf{x})} \cdot 1 \cdot \varepsilon_{x_i} \\ &= \sum_{i=1}^n \frac{x_i}{f(\mathbf{x})} \varepsilon_{x_i}. \end{aligned}$$

Per mostrare che il problema è ben condizionato dimostriamo che il valore assoluto dell'errore inerente è limitato superiormente. In effetti si ha

$$\begin{aligned} |\varepsilon_{\text{in}}| &\doteq \left| \sum_{i=1}^n \frac{x_i}{f(\mathbf{x})} \varepsilon_{x_i} \right| \\ &\stackrel{(1)}{\leq} \sum_{i=1}^n \frac{|x_i|}{|f(\mathbf{x})|} |\varepsilon_{x_i}| \\ &\stackrel{(2)}{<} \frac{1}{|f(\mathbf{x})|} \sum_{i=1}^n |x_i| \cdot u \\ &= \frac{u}{|f(\mathbf{x})|} \sum_{i=1}^n |x_i| \\ &\stackrel{(3)}{=} u, \end{aligned}$$

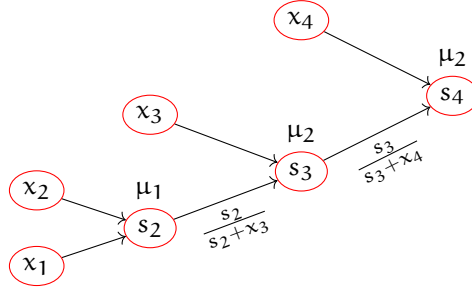
dove (1) è la disuguaglianza triangolare, (2) deriva dal fatto che $|\varepsilon_{x_i}| < u$ per ogni x_i e (3) segue dall'ipotesi che gli x_i siano strettamente positivi: in questo caso infatti il modulo di $f(\mathbf{x}) = x_1 + \dots + x_n$ è uguale alla somma dei moduli degli x_i .

Si ha quindi che il problema è ben condizionato.

STABILITÀ Consideriamo l'algoritmo che consiste nell'effettuare le somme da sinistra a destra, ovvero

$$g_1(\mathbf{x}) = (((x_1 + x_2) + x_3) + \cdots + x_{n-1}) + x_n).$$

Nel caso (ad esempio) di $n = 4$ l'algoritmo può essere rappresentato col seguente grafo:



dove $s_i = s_{i-1} + x_i = x_1 + \cdots + x_i$.

Calcoliamo ad esempio l'errore accumulatosi sul nodo s_3 : si ha

$$\varepsilon_3 = \mu_3 + \frac{s_2}{s_2 + x_3} \underbrace{(\mu_2 + 0 + 0)}_{\varepsilon_2}.$$

Stessa cosa sul nodo s_4 : infatti

$$\varepsilon_4 = \mu_4 + \frac{s_3}{s_3 + x_4} \varepsilon_3.$$

Più in generale si avrà quindi

$$\varepsilon_k = \mu_k + \frac{s_{k-1}}{s_{k-1} + x_k} \varepsilon_{k-1}.$$

Osserviamo che per definizione di s_i possiamo riscrivere $\frac{s_{k-1}}{s_{k-1} + x_k}$ come $\frac{s_{k-1}}{s_k}$ e questa frazione è sempre minore di 1: in questo modo riusciamo a limitare superiormente il modulo dell'errore algoritmo. Infatti

$$\begin{aligned} |\varepsilon_{\text{alg}}| &= |\varepsilon_n| = \left| \mu_n + \frac{s_{n-1}}{s_n} \varepsilon_{n-1} \right| \\ &\leq |\mu_n| + \left| \frac{s_{n-1}}{s_n} \right| |\varepsilon_{n-1}| \\ &< u + 1 \cdot |\varepsilon_{n-1}| \\ &< u + u + |\varepsilon_{n-2}| \\ &< \dots \\ &< (n-1)u. \end{aligned}$$

L'algoritmo è dunque stabile, ma solo se n non è troppo grande.

2 | ALGEBRA LINEARE NUMERICA

2.1 NORME VETTORIALI

Iniziamo ora lo studio dell'algebra lineare numerica: studieremo algoritmi e problemi numerici sugli spazi vettoriali \mathbb{F}^n e $\mathbb{F}^{n \times n}$, dove \mathbb{F} è il campo dei numeri reali (\mathbb{R}) o quello dei numeri complessi (\mathbb{C}).

Uno dei nostri obiettivi sarà quello di risolvere sistemi lineari della forma $Ax = b$ (con A matrice e x, b vettori) sfruttando algoritmi di macchina sufficientemente efficienti e con errori piccoli, quindi vogliamo avere delle quantità che descrivano l'errore commesso.

Osserviamo innanzitutto che, dato che lavoriamo con l'aritmetica di macchina, non riusciremo ad ottenere la soluzione x esatta ma solamente una soluzione approssimata, che indicheremo come al solito con \tilde{x} . L'errore commesso può quindi essere espresso come la differenza $\tilde{x} - x$, tuttavia questa quantità è ancora un vettore e pertanto è difficile confrontarlo con altri errori per vedere quanto è *grande*.

Vogliamo quindi associare ad ogni vettore un *numero reale* non negativo in modo da poter esprimere il concetto di grandezza e di distanza tra vettori.

Definizione 2.1.1 **Norma vettoriale.** Si dice **norma vettoriale** su \mathbb{F}^n una funzione $f : \mathbb{F}^n \rightarrow \mathbb{R}$ che soddisfi le seguenti proprietà:

- (1) $f(v) \geq 0$ per ogni $v \in \mathbb{F}^n$. Inoltre $f(v) = 0$ se e solo se $v = 0$.
- (2) Per ogni $\alpha \in \mathbb{F}$, $v \in \mathbb{F}^n$ si ha $f(\alpha v) = |\alpha|f(v)$.
- (3) Per ogni $v, w \in \mathbb{F}^n$ si ha $f(v + w) \leq f(v) + f(w)$ (**disuguaglianza triangolare**).

Una norma si indica spesso anche con il simbolo $\|\cdot\|$.

Un esempio banale di norma è il valore assoluto su \mathbb{R} , considerando \mathbb{R} come uno spazio vettoriale su se stesso. In effetti:

- (1) il valore assoluto di $x \in \mathbb{R}$ è sempre un numero non negativo, ed in particolare è 0 se e solo se $x = 0$;
- (2) per ogni coppia di valori $x, y \in \mathbb{R}$ si ha che $|xy| = |x| \cdot |y|$;
- (3) vale la disuguaglianza triangolare, e quindi per ogni coppia di valori $x, y \in \mathbb{R}$ si ha $|x + y| = |x| + |y|$.

Definizione 2.1.2 **Distanza.** Sia X un insieme qualsiasi. Si dice **distanza** su X una funzione

$$d : X \times X \rightarrow \mathbb{R}$$

che soddisfi le seguenti proprietà:

- (1) $d(x, y) \geq 0$ per ogni $x, y \in X$. Inoltre $d(x, y) = 0$ se e solo se $x = y$.
- (2) $d(x, y) = d(y, x)$ per ogni $x, y \in X$.
- (3) $d(x, z) \leq d(x, y) + d(y, z)$ per ogni $x, y, z \in X$.

Queste proprietà sono abbastanza semplici se pensiamo al concetto di distanza nella vita di tutti i giorni. Infatti:

- (1) la distanza tra due punti è sempre un numero positivo, ed è 0 se e solo se i due punti sono effettivamente lo stesso punto;
- (2) è indifferente andare dal primo punto al secondo o viceversa;
- (3) la distanza tra i punti x e z è necessariamente minore o uguale alla distanza che percorreremmo se andassimo da x ad un terzo punto y e poi da y a z .

Osserviamo inoltre che nel caso di \mathbb{F}^n esiste una distanza particolarmente facile da ricavare, poiché è direttamente **indotta da una norma**. Infatti se $\|\cdot\|$ è una norma qualsiasi, la funzione $d : \mathbb{F}^n \times \mathbb{F}^n \rightarrow \mathbb{R}$ definita da

$$d(\mathbf{v}, \mathbf{w}) := \|\mathbf{v} - \mathbf{w}\|$$

è una distanza.

Dimostrazione. Dobbiamo mostrare che la funzione d definita sopra sia effettivamente una distanza.

- (1) Per il primo punto della definizione di norma si ha che $d(\mathbf{v}, \mathbf{w}) = \|\mathbf{v} - \mathbf{w}\| \geq 0$ per ogni $\mathbf{v}, \mathbf{w} \in \mathbb{F}^n$. Inoltre questa quantità è 0 se e solo se $\mathbf{v} - \mathbf{w} = \mathbf{0}$, ovvero se e solo se $\mathbf{v} = \mathbf{w}$.
- (2) Si ha che

$$d(\mathbf{v}, \mathbf{w}) = \|\mathbf{v} - \mathbf{w}\| = |-1| \|\mathbf{v} - \mathbf{w}\| \stackrel{(!)}{=} \|\mathbf{w} - \mathbf{v}\| = d(\mathbf{w}, \mathbf{v}),$$

dove il passaggio contrassegnato con (!) deriva dalla seconda proprietà delle norme.

- (3) Per la disuguaglianza triangolare delle norme si ha che

$$d(\mathbf{v}, \mathbf{u}) = \|\mathbf{v} - \mathbf{u}\| = \|(\mathbf{v} - \mathbf{w}) + (\mathbf{w} - \mathbf{u})\| \leq \|\mathbf{v} - \mathbf{w}\| + \|\mathbf{w} - \mathbf{u}\| = d(\mathbf{v}, \mathbf{w}) + d(\mathbf{w}, \mathbf{u}).$$

□

2.1.1 Norme vettoriali principali

Introduciamo in questa sezione le principali norme vettoriali.

NORMA 2 O EUCLIDEA La norma più importante è chiamata **norma 2** o **norma euclidea** poiché dà origine alla distanza euclidea. È definita in questo modo:

$$\|\mathbf{v}\|_2 := \sqrt{\sum_{i=1}^n |v_i|^2}.$$

Nel caso in cui lo spazio vettoriale sia \mathbb{R}^2 la norma euclidea ci dà la lunghezza euclidea del vettore \mathbf{v} , o equivalentemente la distanza del punto \mathbf{v} dall'origine degli assi:

$$\|\mathbf{v}\|_2 = \sqrt{v_1^2 + v_2^2}.$$

Il valore assoluto presente nella formula è inutile se lo spazio vettoriale è \mathbb{R}^n ma è necessario se lavoriamo sui complessi: ricordiamo infatti che se $z = a + ib$ è un numero complesso il suo modulo è il numero reale

$$|a + ib| = \sqrt{a^2 + b^2}$$

ed è necessario in quanto il risultato della norma deve essere un numero reale.

Osserviamo inoltre che in \mathbb{R}^n si ha $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}}$, dove \mathbf{x}^T è il trasposto del vettore colonna \mathbf{x} , mentre in \mathbb{C}^n si ha $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^H \mathbf{x}}$, dove \mathbf{x}^H è il trasposto coniugato del vettore \mathbf{x} .

NORMA 1 La **norma 1** è la norma definita da

$$\|\mathbf{v}\|_1 = \sum_{i=1}^n |v_i|.$$

NORMA INFINITO la **norma infinito** è la norma definita da

$$\|\mathbf{v}\|_\infty = \max_{i=1,\dots,n} |v_i|.$$

2.1.2 Equivalenza topologica tra le norme di \mathbb{F}^n

Come abbiamo detto all'inizio del capitolo, useremo le norme vettoriali per determinare *quanto* errore abbiamo fatto nel fare un determinato calcolo. Tuttavia non esiste un'unica norma su \mathbb{F}^n e ovviamente norme diverse danno risultati diversi.

Esempio 2.1.3. Consideriamo il vettore $\mathbf{x} = (1, -1, 1)^T \in \mathbb{R}^3$. Allora

$$\|\mathbf{x}\|_2 = \sqrt{3}, \quad \|\mathbf{x}\|_1 = 3, \quad \|\mathbf{x}\|_\infty = 1.$$

In realtà le norme si comportano tutte *più o meno* allo stesso modo, come ci viene garantito dal seguente Teorema.

Teorema 2.1.4 **Equivalenza topologica tra le norme.** Siano f, g due norme su \mathbb{F}^n . Allora esistono due costanti $\alpha, \beta \in [0, +\infty)$ tali che per ogni $\mathbf{v} \in \mathbb{F}^n$ si ha

$$\alpha g(\mathbf{v}) \leq f(\mathbf{v}) \leq \beta g(\mathbf{v}).$$

Dimostrazione. Data nella ?? dell'Appendice. □

Questo Teorema è molto utile quando vogliamo ad esempio mostrare che un errore tende a 0. Infatti se ad esempio ε è un vettore che misura l'errore commesso e so che $\|\varepsilon\|_2 \rightarrow 0$, allora per qualsiasi altra norma p dovranno esistere due costanti $\alpha, \beta \in \mathbb{R}^+$ tali che

$$\alpha \|\varepsilon\|_2 \leq p(\varepsilon) \leq \beta \|\varepsilon\|_2.$$

Ma il membro sinistro e il membro destro di questa disequazione tendono a 0, dunque per il Teorema dei Carabinieri anche il membro centrale dovrà tendere a 0, ovvero la norma dell'errore tende a 0 *a prescindere dalla norma scelta*.

2.2 NORME MATRICIALI

Come per i vettori (ovvero per lo spazio vettoriale \mathbb{F}^n) possiamo definire un concetto di norma anche per lo spazio di matrici quadrate $\mathbb{F}^{n \times n}$.

Definizione 2.2.1 **Norma matriciale.** Si dice **norma matriciale** su $\mathbb{F}^{n \times n}$ una funzione $f : \mathbb{F}^{n \times n} \rightarrow \mathbb{R}$ che soddisfi le seguenti proprietà:

- (1) $f(A) \geq 0$ per ogni $A \in \mathbb{F}^{n \times n}$. Inoltre $f(A) = 0$ se e solo se $A = 0$.
- (2) Per ogni $\alpha \in \mathbb{F}$, $A \in \mathbb{F}^{n \times n}$ si ha $f(\alpha A) = |\alpha|f(A)$.
- (3) Per ogni $A, B \in \mathbb{F}^{n \times n}$ si ha $f(A + B) \leq f(A) + f(B)$ (**disuguaglianza triangolare**).
- (4) Per ogni $A, B \in \mathbb{F}^{n \times n}$ si ha $f(AB) \leq f(A) \cdot f(B)$.

Una proprietà immediata delle norme matriciali è che per qualsiasi norma $\|\cdot\|$ si ha che $\|I_n\| \geq 1$ (dove $I_n \in \mathbb{F}^{n \times n}$ è la matrice identità di taglia $n \times n$). Infatti

$$\|I_n\| = \|I_n \cdot I_n\| \leq \|I_n\| \cdot \|I_n\| \implies 1 \leq \|I_n\|,$$

dove l'implicazione si ottiene dividendo entrambi i membri per $\|I_n\|$, che è sicuramente non nullo grazie alla proprietà (1) delle norme.

Osserviamo inoltre che, come nel caso delle norme vettoriali, ogni norma matriciale induce una distanza tra matrici $d : \mathbb{F}^{n \times n} \rightarrow \mathbb{F}^{n \times n} \rightarrow \mathbb{R}$, definita da

$$d(A, B) := \|A - B\|.$$

A

APPENDICE A

A.1 "9-PERIODICO"

Proposizione A.1.1 **9-periodico.** In base 10 i numeri $0.\bar{9}$ e 1 sono uguali.

Forniamo due diverse dimostrazioni di questa proposizione.

Prima dimostrazione. Dalle formule per trasformare i numeri periodici in frazioni sappiamo che $0.\bar{9} = 9/9 = 1$. \square

Seconda dimostrazione. Espandendo la definizione di numero periodico otteniamo che

$$0.\bar{9} = 0.999\ldots = 9 \cdot 10^{-1} + 9 \cdot 10^{-2} + 9 \cdot 10^{-3} + \ldots = \sum_{i=1}^{\infty} 9 \cdot 10^{-i}.$$

Sfruttando la formula della serie geometrica si ottiene che

$$\begin{aligned} \sum_{i=1}^{\infty} 9 \cdot 10^{-i} &= 9 \cdot \sum_{i=1}^{\infty} \left(\frac{1}{10}\right)^i \\ &= 9 \cdot \left(\left(\sum_{i=0}^{\infty} \left(\frac{1}{10}\right)^i \right) - \left(\frac{1}{10}\right)^0 \right) \\ &= 9 \cdot \left(\frac{1}{1 - 1/10} - 1 \right) \\ &= 9 \cdot \left(\frac{10}{9} - 1 \right) \\ &= 9 \cdot \frac{1}{9} \\ &= 1. \end{aligned} \quad \square$$

La proposizione vale in generale in una base β qualsiasi ($\beta \geq 2$).

Proposizione A.1.2 In base β ($\beta \geq 2$) vale che $0.\overline{(\beta-1)} = 1$.

Dimostrazione. La dimostrazione è uguale alla seconda dimostrazione della [Proposizione A.1.1](#). \square

A.2 LA PRECISIONE DI MACCHINA È UN LIMITE SUPERIORE STRETTO

In questa sezione mostriamo che la precisione di macchina u è sempre strettamente maggiore all'errore relativo ε_x per ogni $x \in \mathbb{R}$ tale che $x \neq 0$ e $\omega \leq |x| \leq \Omega$, anche usando l'arrotondamento come metodo di approssimazione.

Dalla dimostrazione del [Teorema 1.2.5](#) si ha che

$$|\varepsilon_x| = \frac{1}{2}\beta^{1-p} = u$$

se e solo se la disuguaglianza

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

è in realtà un'uguaglianza, ovvero se e solo se valgono le seguenti due condizioni:

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

Tuttavia, come abbiamo notato nella [Osservazione 1.2.2](#), la prima condizione vale se e solo se

$$|x| = \frac{1}{2}(a + b) = \frac{1}{2}(a + a + \beta^{p-t}) = a + \frac{1}{2}\beta^{p-t}.$$

È quindi necessario che

$$\begin{aligned} \beta^{p-1} &= a + \frac{1}{2}\beta^{p-t} \\ &= \beta^p \left(\sum_{i=1}^t d_i \beta^{-i} \right) + \frac{1}{2}\beta^{p-t} \\ &= \beta^p \left(\sum_{i=1}^t d_i \beta^{-i} + \frac{1}{2}\beta^{-t} \right) \\ \Leftrightarrow \beta^{-1} &= \sum_{i=1}^t d_i \beta^{-i} + \frac{1}{2}\beta^{-t}. \end{aligned}$$

Tuttavia ciò è assurdo: infatti il termine al primo membro ha una singola cifra decimale diversa da zero, ovvero quella di β^{-1} , mentre il termine del secondo membro ha anche delle cifre non nulle nelle posizioni β^{-t} e/o β^{-t-1} .

Segue quindi che $|\varepsilon_x| < u$ in qualsiasi caso.

A.3 DIMOSTRAZIONE DELL'EQUIVALENZA TOPOLOGICA TRA NORME

In questa sezione dimostreremo il [Teorema 2.1.4](#), che per comodità enunciamo nuovamente.

Teorema 2.1.4 **Equivalenza topologica tra le norme.** *Siano f, g due norme su \mathbb{F}^n . Allora esistono due costanti $\alpha, \beta \in [0, +\infty)$ tali che per ogni $v \in \mathbb{F}^n$ si ha*

$$\alpha g(v) \leq f(v) \leq \beta g(v).$$

Per dimostrare il Teorema sfruttiamo il seguente lemma.

Lemma A.3.1 *Sia $\|\cdot\|$ una norma su \mathbb{F}^n . Allora la funzione*

$$x \mapsto \|x\|$$

è una funzione continua.

Dimostrazione. Sia $\mathbf{x}_0 \in \mathbb{F}^n$ un punto qualsiasi e dimostriamo che la norma è continua in \mathbf{x}_0 .

Sia quindi $\varepsilon > 0$ qualsiasi: vogliamo determinare un $\delta > 0$ tale che per ogni $\mathbf{x} \in \mathbb{F}^n$ con $\|\mathbf{x} - \mathbf{x}_0\| < \delta$ si ha $|\|\mathbf{x}\| - \|\mathbf{x}_0\|| < \varepsilon$.

Osserviamo che

$$\|\mathbf{x}\| = \|\mathbf{x} - \mathbf{x}_0 + \mathbf{x}_0\| \leq \|\mathbf{x}_0\| + \|\mathbf{x} - \mathbf{x}_0\|,$$

da cui $\|\mathbf{x}\| - \|\mathbf{x}_0\| \leq \|\mathbf{x} - \mathbf{x}_0\|$. Scambiando i ruoli di \mathbf{x} e \mathbf{x}_0 otteniamo che $\|\mathbf{x}_0\| - \|\mathbf{x}\| \leq \|\mathbf{x}_0 - \mathbf{x}\| = \|\mathbf{x} - \mathbf{x}_0\|$, quindi combinando le due disuguaglianze si ha

$$|\|\mathbf{x}\| - \|\mathbf{x}_0\|| \leq \|\mathbf{x} - \mathbf{x}_0\|.$$

Poniamo dunque $\delta = \varepsilon$. Allora

$$|\|\mathbf{x}\| - \|\mathbf{x}_0\|| \leq \|\mathbf{x} - \mathbf{x}_0\| < \delta = \varepsilon,$$

cioè la norma è continua in \mathbf{x}_0 , da cui segue (per generalità di \mathbf{x}_0) che la norma è una funzione continua. \square

Possiamo quindi dimostrare il [Teorema 2.1.4](#).

B | PREREQUISITI DI ALGEBRA LINEARE

B.1 SPAZI VETTORIALI

Iniziamo rivedendo i concetti fondamentali sugli spazi vettoriali.

Definizione B.1.1 **Campo.** Si dice **campo** un insieme \mathbb{F} dotato di due operazioni, chiamate solitamente *somma* e *prodotto* e indicate con i simboli $+$ e \cdot , tali che

(S1) la somma è associativa;

(S2) la somma è commutativa;

(S3) esiste un elemento $0_{\mathbb{F}} \in \mathbb{F}$ (indicato anche 0 se il campo è sottointeso) che è elemento neutro per la somma;

(S4) ogni elemento ammette un opposto rispetto alla somma, ovvero per ogni $\alpha \in \mathbb{F}$ esiste $-\alpha \in \mathbb{F}$ tale che

$$\alpha + (-\alpha) = 0;$$

(P1) il prodotto è associativo;

(P2) il prodotto è commutativo;

(P3) esiste un elemento $1_{\mathbb{F}} \in \mathbb{F}$ (indicato anche con 1 se il campo è sottointeso) che è elemento neutro per il prodotto;

(P4) ogni elemento diverso da 0 ammette un inverso rispetto al prodotto, ovvero per ogni $\alpha \in \mathbb{F} \setminus \{0\}$ esiste $\alpha^{-1} \in \mathbb{F}$ tale che

$$\alpha \cdot \alpha^{-1} = 1;$$

(D) vale la proprietà distributiva del prodotto rispetto alla somma.

Esempi di campo sono i numeri reali \mathbb{R} , i numeri razionali \mathbb{Q} , i numeri complessi \mathbb{C} . D'ora in avanti indicheremo con \mathbb{F} un generico campo (in particolare \mathbb{F} nella pratica sarà uno tra \mathbb{R} e \mathbb{C}).

Definizione B.1.2 **Spazio vettoriale.** Si dice **spazio vettoriale** sul campo \mathbb{F} un insieme V dotato di due operazioni, una somma $+: V \times V \rightarrow V$ e un *prodotto per scalare* $\cdot: \mathbb{F} \times V \rightarrow V$, che soddisfano le seguenti proprietà:

(S1) la somma è associativa;

(S2) la somma è commutativa;

(S3) la somma ammette un elemento neutro, indicato con 0 o 0_V ;

(S4) ogni elemento di V ammette un opposto rispetto alla somma, ovvero per ogni $v \in V$ esiste $-v \in V$ tale che

$$v + (-v) = 0;$$

(P1) il prodotto per scalare distribuisce sulla somma, ovvero per ogni $\mathbf{v}, \mathbf{w} \in V$ e per ogni $\alpha \in \mathbb{F}$ si ha

$$\alpha \cdot (\mathbf{v} + \mathbf{w}) = \alpha \mathbf{v} + \alpha \mathbf{w};$$

(P2) la somma di scalari distribuisce sul prodotto per scalare, ovvero per ogni $\alpha, \beta \in \mathbb{F}$, $\mathbf{v} \in V$ si ha

$$(\alpha + \beta)\mathbf{v} = \alpha \mathbf{v} + \beta \mathbf{v};$$

(P3) il prodotto per scalari è associativo, ovvero per ogni $\alpha, \beta \in \mathbb{F}$, $\mathbf{v} \in V$ si ha

$$\alpha(\beta \mathbf{v}) = (\alpha\beta)\mathbf{v};$$

(P4) l'elemento neutro del prodotto (in \mathbb{F}) è neutrale rispetto al prodotto per scalare, ovvero per ogni $\mathbf{v} \in V$ si ha

$$1_{\mathbb{F}} \cdot \mathbf{v} = \mathbf{v}.$$

Il campo \mathbb{F} viene spesso detto **campo degli scalari** di V .

Indicheremo solitamente gli elementi del campo degli scalari con lettere greche minuscole (α, β, \dots) e gli elementi dello spazio vettoriale con lettere in grassetto ($\mathbf{v}, \mathbf{w}, \dots$).

Esempi di spazi vettoriali sono

- l'insieme dei vettori colonna a componenti in \mathbb{F} con n elementi, indicato con \mathbb{F}^n ;
- l'insieme dei polinomi a coefficienti in \mathbb{F} , indicato con $\mathbb{F}[x]$;
- l'insieme delle matrici $m \times n$ a coefficienti in \mathbb{F} , indicate con $\mathbb{F}^{m \times n}$.

Definizione B.1.3 **Sottospazio.** Sia V uno spazio vettoriale su \mathbb{F} . Un **sottospazio** di V è un insieme $U \subseteq V$ che sia anch'esso uno spazio vettoriale su \mathbb{F} .

Definizione B.1.4 **Span di un sottoinsieme.** Sia V uno spazio vettoriale su \mathbb{F} e sia $S \subseteq V$. Allora si dice **span** di S l'insieme

$$\text{span}(S) := \{ \alpha_1 \mathbf{v}_1 + \dots + \alpha_k \mathbf{v}_k : \alpha_i \in \mathbb{F}, \mathbf{v}_i \in S \}.$$

Proposizione B.1.5 *Per ogni $S \subseteq V$ si ha che $\text{span}(S)$ è un sottospazio vettoriale di V .*

B.2 INDIPENDENZA E BASI

Assumeremo implicitamente che V sia uno spazio vettoriale su \mathbb{F} .

Definizione B.2.1 **Indipendenza lineare.** I vettori $\mathbf{v}_1, \dots, \mathbf{v}_n \in V$ formano un insieme di **vettori linearmente indipendenti** se si ha che

$$\alpha_1 \mathbf{v}_1 + \dots + \alpha_n \mathbf{v}_n = \mathbf{0} \implies \alpha_1 = \dots = \alpha_n = 0_{\mathbb{F}}.$$

Definizione B.2.2 **Insieme di generatori.** I vettori $\mathbf{v}_1, \dots, \mathbf{v}_n \in V$ formano un insieme di **generatori per V** se per ogni $\mathbf{v} \in V$ esistono $\alpha_1, \dots, \alpha_n \in \mathbb{F}$ tali che

$$\mathbf{v} = \alpha_1 \mathbf{v}_1 + \dots + \alpha_n \mathbf{v}_n.$$

Equivalentemente $\mathbf{v}_1, \dots, \mathbf{v}_n$ sono generatori se

$$\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_n\} = V.$$

Definizione B.2.3 **Base.** L'insieme formato dai vettori v_1, \dots, v_n forma una **base** per V se valgono entrambe le seguenti condizioni:

1. $\{v_1, \dots, v_n\}$ è un insieme di vettori linearmente indipendenti;
2. $\{v_1, \dots, v_n\}$ è un insieme di generatori per V .

Chiameremo **base canonica** di \mathbb{F}^n l'insieme formato dai vettori

$$e_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad e_2 = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \quad \dots, \quad e_n = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}.$$

Non tutti gli spazi vettoriali ammettono una base *finita*: ad esempio nessun sottoinsieme finito di $\mathbb{F}[x]$ lo genera, e quindi non può esistere una base.

Definizione B.2.4 **Dimensione di uno spazio vettoriale.** Sia V uno spazio vettoriale su \mathbb{F} e sia \mathcal{B} una base di V . Si dice **dimensione di V su \mathbb{F}** la quantità

$$\dim_{\mathbb{F}} V := \#\mathcal{B}.$$

Quando il campo degli scalari è sottointeso scriveremo semplicemente $\dim V$.

Ad esempio:

- la dimensione di \mathbb{F}^n è $\dim \mathbb{F}^n = n$ (basta considerare la base canonica);
- la dimensione di $\mathbb{F}^{m \times n}$ è $\dim \mathbb{F}^{m \times n} = mn$.

Teorema B.2.5 **Teorema della dimesnione.** Sia V uno spazio vettoriale su \mathbb{F} di dimensione n . Allora ogni base di V ha n elementi.

B.3 APPLICAZIONI LINEARI

Definizione B.3.1 **Applicazione lineare.** Siano U, V due spazi vettoriali su un campo \mathbb{F} . Si dice **applicazione lineare** una funzione $f: U \rightarrow V$ tale che

1. $f(x + y) = f(x) + f(y)$ per ogni $x, y \in U$;
2. $f(\alpha x) = \alpha f(x)$ per ogni $x \in U, \alpha \in \mathbb{F}$.

Una volta fissate due basi \mathcal{B}_U e \mathcal{B}_V rispettivamente del dominio e del codominio, ad esempio

$$\mathcal{B}_U = \{u_1, \dots, u_n\}, \quad \mathcal{B}_V = \{v_1, \dots, v_m\},$$

riusciamo ad esprimere l'applicazione f in modo univoco tramite una matrice $A_f \in \mathbb{F}^{m \times n}$.

In effetti per ogni vettore u_j nella base \mathcal{B}_U l'immagine del vettore mediante f , ovvero $f(u_j)$, è un elemento di V e può pertanto essere espresso in termini degli elementi della base \mathcal{B}_V di V : esisteranno cioè degli scalari $\alpha_{1j}, \dots, \alpha_{mj}$ tali che

$$f(u_j) = \alpha_{1j}v_1 + \dots + \alpha_{mj}v_m.$$

Questi scalari formeranno la j -esima colonna della matrice associata, che avrà quindi la forma

$$A_f = \begin{pmatrix} \alpha_{11} & \alpha_{12} & \dots & \alpha_{1n} \\ \alpha_{21} & \alpha_{22} & \dots & \alpha_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{m1} & \alpha_{m2} & \dots & \alpha_{mn} \end{pmatrix}$$

Definizione B.3.2 **Immagine e nucleo di un'applicazione lineare.** Sia $f : U \rightarrow V$ un'applicazione lineare.

(1) Si dice **kernel** (o **nucleo**) di f l'insieme

$$\ker f := \{ \mathbf{u} \in U : f(\mathbf{u}) = \mathbf{0}_V \}.$$

(2) Si dice **immagine** di f l'insieme

$$\operatorname{Im} f := \{ f(\mathbf{u}) \in V : \mathbf{u} \in U \}.$$

Si ha che $\ker f$ è un sottospazio di U e $\operatorname{Im} f$ è un sottospazio di V . Inoltre vale il seguente teorema.

Teorema B.3.3 **Teorema della dimensione.** $\dim V = \dim \operatorname{Im} f + \dim \ker f$.