



# UNIVERSITÀ DI PISA

Dipartimento di Informatica  
Corso di Laurea Triennale in Informatica

Corso 2° anno - 6 CFU

## Calcolo Numerico

**Professore:**  
Prof. Luca Germignani

**Autore:**  
Matteo Giuntori  
Filippo Ghirardini

---

Anno Accademico 2023/2024

## Contents

<b>1</b>	<b>Aritmetica di Macchina</b>	<b>3</b>
1.1	Teorema di rappresentazione . . . . .	3
1.2	Errore di rappresentazione . . . . .	4
1.3	Operazioni di macchina . . . . .	5
<b>2</b>	<b>Calcolo degli errori</b>	<b>6</b>
<b>3</b>	<b>Matrici</b>	<b>7</b>
3.1	Norme vettoriali . . . . .	7
<b>4</b>	<b>Condizionamento</b>	<b>8</b>
4.1	Metodi diretti . . . . .	8
4.1.1	Matrice diagonale . . . . .	8
4.1.2	Matrice triangolare . . . . .	8
<b>5</b>	<b>Metodi iterativi</b>	<b>10</b>
5.1	Convergenza . . . . .	10
5.2	Metodo di Jacobi . . . . .	12
5.3	Metodo di Gauss-Seidel . . . . .	13
5.4	Criteri di arresto . . . . .	14
<b>6</b>	<b>Equazioni non lineari</b>	<b>16</b>
6.1	Tecnica della separazione . . . . .	17
6.2	Metodo di bisezione . . . . .	20
6.2.1	Approssimazione della radice . . . . .	21
6.2.2	Svantaggi . . . . .	21
6.3	Metodo del punto fisso . . . . .	22
6.4	Metodo delle tangenti . . . . .	25
6.4.1	Convergenza locale . . . . .	26
6.4.2	Convergenza lineare . . . . .	26
6.4.3	Convergenza quadratica . . . . .	26
6.4.4	Convergenza globale . . . . .	28
<b>7</b>	<b>Teoremi</b>	<b>32</b>
7.1	Maggiorazione dell'errore di macchina . . . . .	32
7.1.1	Dimostrazione . . . . .	32
7.2	Errore analitico e totale . . . . .	32
7.3	Localizzazione degli autovalori . . . . .	33
7.3.1	Dimostrazione . . . . .	33
7.4	Esistenza della fattorizzazione LU . . . . .	34
7.4.1	Dimostrazione . . . . .	34
7.5	Convergenza con norma . . . . .	34
7.5.1	Dimostrazione . . . . .	34
7.6	Convergenza con raggio spettrale . . . . .	35
7.6.1	Dimostrazione . . . . .	35
7.7	Predominanza diagonale . . . . .	35
7.7.1	Dimostrazione . . . . .	35
7.8	Convergenza del metodo di bisezione . . . . .	36
7.8.1	Dimostrazione . . . . .	36
7.9	Limite metodo funzionale . . . . .	37
7.9.1	Dimostrazione . . . . .	37
7.10	Convergenza locale metodo funzionale . . . . .	38
7.10.1	Dimostrazione . . . . .	38
7.11	Convergenza locale metodo delle tangenti . . . . .	38
7.11.1	Dimostrazione . . . . .	38

# Calcolo Numerico

Realizzato da: Giuntoni Matteo e Filippo Ghirardini

A.A. 2023-2024

---

# 1 Aritmetica di Macchina

Per una macchina la scrittura  $(x + y) + z$  è diverso da  $x + (y + z)$ . Vediamo dunque che ci sono alcuni punti focali da considerare per far sì che una macchina funzioni correttamente:

- Trovare uno standard per come **memorizzare** i numeri.
- Trovare uno standard per come **manipolare** i numeri.

Da questi due punti possiamo ricondurci ad un solo problema, come andare a **rappresentare** i numeri.

## 1.1 Teorema di rappresentazione

**Teorema 1.1.1.** Dato  $x \in \mathbb{R}, x \neq 0$ <sup>1</sup> e una base di numerazione  $B, B \in \mathbb{N}, B > 1$  esistono e sono univocamente determinati:

1. Un intero  $p \in \mathbb{Z}$  detto **esponente** della rappresentazione
2. Una successione di numeri naturali  $\{d_i\}_{i=1}^{+\infty}$  con  $d_i \neq 0, 0 \leq d_i \leq B - 1$  e  $d_i$  non definitivamente uguali a  $B - 1$ , dette cifre della rappresentazione tali per cui  $x$  si scrive in modo **unico** nella seguente forma:

$$x = \text{sign}(x) B^p \sum_{i=1}^{+\infty} d_i B^{-i}. \quad (1)$$

dove la sommatoria viene chiamata **mantissa**

**Esempio 1.1.1** (Esempio in base 10). Poniamo come numero da rappresentare 0.1 in base 10 è

$$0.1 = +10^0(0.1)$$

Andiamo ora ad analizzare il significato di questo teorema. Esso descrive quella che viene chiamata **rappresentazione in virgola mobile**, in quanto l'esponente  $p$  on è determinato in modo da avere la parte intera nulla. Le cose da considerare in questo teorema sono:

- La condizione  $d_i \neq 0$  e  $d_i$  non definitivamente uguale a  $B - 1$  sono introdotte per garantire l'unicità delle rappresentazioni. Ad esempio:

$$B = 10 \text{ abbiamo } 1 = +10^1(1 \cdot 10^{-1}) = +10^2(0 \cdot 10^{-1} + 1 \cdot 10^{-1})$$

Quindi due rappresentazioni diverse per lo stesso numero, però considerando le condizioni scritte sopra la seconda non risulta accettabile perché la prima cifra è nulla.

Questa clausola ci garantisce anche l'unicità delle rappresentazioni nei numeri **periodici**:

$$0.\bar{9} = 10^0(0.99 \dots 9)$$

Non è ammissibile in quanto è definitivamente uguale a  $B - 1$ .

- Questa rappresentazione si estende anche all'insieme dei numeri complessi del tipo  $z = a + ib$ , utilizzando una rappresentazione come coppie di numeri reali del tipo  $(a, b)$ .

Possiamo dedurre che visto che stiamo lavorare con registri di memoria di un calcolatore con memoria a numero finito, anche la quantità di cifre rappresentabili saranno a numero finito esso viene chiamato **insieme dei numeri di macchina**.

Dal teorema di rappresentazione in base di un numero reale può avvenire assegnando delle posizioni di memoria per il segno, per l'esponente e per le cifre della rappresentazione.

<sup>1</sup>Lo zero viene utilizzato dalla macchina per alcune operazioni come il confronto, quindi deve averlo ma lo rappresenta in un modo particolare

**Definizione 1.1.1** (Insieme dei numeri di macchina). Si definisce l'insieme dei numeri di macchina in rappresentazione floating point con  $t$  cifre, base  $B$  e range  $-m, M$  l'insieme dei numeri reali.

$$\mathbb{F}(B, t, m, M) = \{0\} \cup \{s \in \mathbb{R} : x = \pm B^p \sum_{i=1}^t d_i B^{-i}, d_1 \neq 0, -m \leq p \leq M\}$$

Si osserva in questa definizione che:

- L'insieme  $\mathbb{F}$  ha cardinalità **finita**:  $N = 2B^{t-1}(B-1)(M+m+1) + 1$ .
- L'insieme dei numeri di macchina  $\mathbb{F}(B, t, m, M)$  è **simmetrico** rispetto all'origine.
- Possiamo definire  $\Omega = B^M \sum_{i=1}^t (B-1)B^{-i}$  come il **più grande** numero macchina e  $\omega = +B^{-m}B^{-1}$  come invece il **più piccolo** positivo.
- Posto un  $x = B^p \sum_{i=1}^t d_i B^{-i}$  possiamo definire il suo **successivo** numero di macchina come  $y = B^p (\sum_{i=1}^{t-1} d_i B^{-i} + (d_t + 1)B^{-t})$ .  
Da qui vediamo che la distanza  $y - x = B^p - t$  porta i numeri ad essere **non equidistanti** fra di loro, quindi la distanza aumenta con l'avvicinarsi a  $\Omega$ .  
Questo ci fa comodo perché ci interessa l'**errore relativo**, quindi su numeri piccoli ci serve un errore piccolo mentre su numeri grandi posso fare errori grandi.

**Esempio 1.1.2.** Facciamo ora un esempio in cui andiamo a rappresentare il numero successivo di  $x = B^p \sum_{i=1}^t d_i B^{-i}$ . Esso si può scrivere come  $y = B^p \left( \sum_{i=1}^{t-1} d_i B^{-i} + (d_t + 1)B^{-t} \right)$ .

Mentre si può scrivere la distanza fra questi due valori come  $y - x = B^p - t$ .

E' stato fissato uno standard IEEE 754 negli anni 70/80 che dice che, visto ci sono macchine che hanno metodi di rappresentazione diversi, bisogna fissare un standard, ovvero  $B = 2$  con registri a 32 o 64 bit.

Questa rappresentazione ha uno svantaggio che può sembrare minimo ma non lo è, lo 0 si rappresenta due volte con  $-0, +0$ . Per ovviare a questo problema si è andato ad abbandonare questa rappresentazione in esponenti ma si rappresentato i numeri nel seguente modo:  $p_1 2^0 + p_1 2^1 + \dots + p_1 1 s^1 0$  che rappresentano numeri da 0 a  $2^{11} - 1$  quindi 2047 numeri, mentre lo 0 si può scrivere come:

- 0 tenendo tutti i valori a 0
- Oppure tendendo tutti i valori a 1

In entrambi i casi abbiamo un range di valori che va da  $[-1022, 1024]$ . A questo punto ho  $2^{P-1022}$  numeri che la macchina rappresenta come  $\pm 2^{P-1022} (0.1d_1 \dots d_{52})$ .

Impostando questo standard abbiamo  $\Omega = 2^{1024} (01 \dots 1)_2$  e  $\omega = 2^{-1022} (101)_2$ .

**Osservazione 1.1.1.** Quando  $p = 0$  abbiamo i numeri che si trovano nella porzione della retta dei numeri che è compresa fra  $-\omega$  e  $\omega$  e possiamo qui avere anche tutti 0 e quindi si introduce il caso dei numeri denormalizzati.

Se abbiamo l'esponente uguale a tutti 1, la convenzione è che tutte le cifre della mantissa sono tutti uguali a 0/1 questo numero indica il  $\pm\infty$  altrimenti sta a significare NaN (not a number). Questi valori ci permettono di gestire forme indeterminate.

## 1.2 Errore di rappresentazione

Quando si va a rappresentare un numero reale non nullo  $x \in \mathbb{R}$  e con  $x \neq 0$  si può andare a commettere degli errori di rappresentazione detto anche **errore relativo di approssimazione**, e si definisce come, prendendo un  $\tilde{x} \in \mathbb{F}(B, t, m, M)$

$$\epsilon_x = \frac{\tilde{x} - x}{x} = \frac{\eta x}{x}, x \neq 0$$

Definiamo  $|\epsilon_x| = \left| \frac{\tilde{x} - x}{x} \right| \leq \frac{B^{P-t}}{B^{P-1}} \leq \frac{B^{P-t}}{B^{P-1}} = B^{1-t} = u$  la  $u$  è definita come **precisione di macchina**.

Andiamo inoltre a definire le condizioni di underflow e overflow. Dato un  $x \in \mathbb{R}, x \neq 0$  abbiamo che:

1. Se  $|x| < \omega$  o  $|x| > \Omega$  overflow. In questo caso si va ad associare il  $+\infty$ .
2. Se invece  $\omega \leq |x| \leq \Omega$  abbiamo underflow. In questo caso allora prendiamo una  $x = B^p \sum_{i=1}^{\infty} d_i B^{-1} \rightarrow B^p \sum_{i=1}^t d_i B^{-1} = \tilde{x}$  che è una approssimazione

### 1.3 Operazioni di macchina

Consideriamo ora due numeri  $x, y \in \mathbb{F}$  e chiediamoci perché la macchina non possa fare l'operazione  $x + y$ . La risposta è che i risultati da questa operazione di ritornano fra i numeri di macchina. Per ovviare a questo problema dovremo usare le Operazioni di macchina che si identificano come  $\oplus \ominus \otimes \oslash$ . Nel nostro caso l'addizione di macchina  $x \oplus y = \text{troncamento}(x + y) = (x + y)(1 + \epsilon_1)$  con  $|\epsilon_1| \leq u$  con  $e_1$  detto errore locale dell'operazione.

**Esempio 1.3.1.** Supponiamo di dover calcolare in macchina la funzione  $f(x) = \frac{x-1}{x}$ . In macchina questa funzione corrisponderebbe a  $g(\tilde{x}) = (\tilde{x} \ominus 1) \oslash \tilde{x}$ . Abbiamo quindi:

$$g(\tilde{x}) = \frac{(x(1 + \epsilon_x) - 1)(1 + \epsilon_1)}{x(1 + \epsilon_x)} \cdot (1 + \epsilon_1) = \frac{(x(1 + \epsilon_x) - 1)(1 + \epsilon_1 + \epsilon_2)}{x(1 + \epsilon_x)} = \frac{(x(1 + \epsilon_x) - 1)(1 + \epsilon_1 + \epsilon_2 - \epsilon_x)}{x}$$

$$g(\tilde{x}) = (\tilde{x} \oplus 1) \oslash \tilde{x} = \frac{(x(1 + \epsilon_x) - 1)(1 + \epsilon_1 + \epsilon_2 - \epsilon_x)}{x}$$

$$\frac{g(\tilde{x}) - f(x)}{f(x)} = \frac{((x - 1)/x) + (\epsilon_1 + \epsilon_2)((x - 1)/x) + \epsilon_2/x - ((x - 1)/x)}{(x - 1)/x} = \epsilon_1 + \epsilon_2 - \frac{\epsilon_x}{x - 1}$$

**Esempio 1.3.2.** Supponiamo ora di calcolare la funzione  $f(x) = \frac{x-1}{x}$  in un altro modo,  $g_2(\tilde{x}) = \frac{g_2(\tilde{x}) - f(x)}{f(x)}$  ed andiamo a fare l'analisi dell'errore.

$$\frac{g_1(\tilde{x}) - f(x)}{f(x)} \doteq \epsilon_1 + \epsilon_2 + \frac{\epsilon_1}{(x - 1)} \quad \text{Questo è il risultato di un analisi al primo ordine}$$

$$\begin{aligned} g_2(\tilde{x}) &= 1 \ominus \frac{1}{\tilde{x}}(1 + \delta_1) = 1 \ominus \frac{1}{x}(1 + \delta_1)(1 - \epsilon_x) = [1 - \frac{1}{x}(1 - \delta_1)(1 - \epsilon_1)](1 + \delta_2) \\ &\doteq (1 + \delta_1) - \frac{1}{x}(1 + \delta_1 + \delta_2 + \epsilon_x) \doteq (1 - \frac{1}{x}) + \delta_2(1 - \frac{1}{x}) - \frac{\delta_1}{x} + \frac{\epsilon_x}{x} \doteq \delta_2 - \frac{\delta_1}{x - 1} + \frac{\epsilon_x}{x - 1} \\ \frac{g_2(\tilde{x}) - f(x)}{f(x)} &= \delta_2 - \frac{\delta_1}{x - 1} + \frac{\epsilon_x}{x - 1} \end{aligned}$$

Questo è il risultato finale dove  $\delta_2 - \frac{\delta_1}{x-1}$  viene definita come parte stabilità mentre  $\delta_2 - \frac{\delta_1}{x-1}$  viene chiamato condizionamento, il risultato finale viene definito invece numero stabile.

## 2 Calcolo degli errori

Supponiamo di avere una funzione  $f : [a, b] \rightarrow \mathbb{R}$  e  $f \neq 0$ , per andare a calcolare questa funzione come già visto usiamo un algoritmo che esprime tale valore come risultato di una sequenza di operazioni aritmetiche. Questa rappresentazione come abbiamo già potuto verificare con esempi produce degli errori di approssimazione. Questi errori possono essere suddivisi in 3 tipologie.

**Definizione 2.0.1** (Errore inerente o inevitabile). *Si dice errore inerente o inevitabile generato nel calcolo di  $f(x) \neq 0$  la quantità:*

$$\epsilon_{in} = \frac{f(\tilde{x}) - f(x)}{f(x)}$$

**Definizione 2.0.2** (Errore algoritmico). *Si dice errore algoritmico generato nel calcolo di  $f(x) \neq 0$  la quantità:*

$$\epsilon_{alg} = \frac{g(\tilde{x}) - f(\tilde{x})}{f(\tilde{x})}$$

**Definizione 2.0.3** (Errore totale). *Si dice errore algoritmico totale nel calcolo di  $f(x) \neq 0$  mediante l'algoritmo specificato da  $g$  la quantità:*

$$\epsilon_{tot} = \frac{g(\tilde{x}) - f(x)}{f(x)}$$

**Osservazione 2.0.1.** Vediamo che se  $|\epsilon_{in}|$  è grande il problema si definisce **problema mal condizionato**. Mentre se  $|\epsilon_{alg}|$  è grande l'algoritmo si dice che **algoritmo è numericamente instabile**.

### 3 Matrici

#### 3.1 Norme vettoriali

Sono uno strumento che ci permette di definire una distanza tra due vettori.

**Definizione 3.1.1** (Norma vettoriale). È una funzione del tipo  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  che deve soddisfare tre proprietà:

1.  $f(x) \geq 0 \wedge f(x) = 0 \Leftrightarrow x = 0$
2.  $f(\alpha x) = |\alpha|f(x) \quad \forall \alpha \in \mathbb{R} \quad \forall x \in \mathbb{R}^n$
3. **Disuguaglianza triangolare:**  $f(x + y) \leq f(x) + f(y) \quad \forall x, y \in \mathbb{R}^n$

e la indichiamo come

$$f(x) = \|x\|$$

Detto questo possiamo definire una distanza come:

$$dist(x, y) = \|x - y\| \quad (2)$$

Le tre proprietà ci danno alcune informazioni sulla distanza:

1. La distanza deve essere non negativa e valere 0 solo se i due vettori coincidono
2. La distanza tra  $x$  e  $y$  deve essere uguale a quella tra  $y$  e  $x$
3.  $\|x - y\| = \|(x - a) + (a - y)\| \leq \|x - a\| + \|a - y\|$

**Definizione 3.1.2** (Distanza euclidea - Norma 2).

$$f(x) = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} = \|x\|_2 \quad (3)$$

**Definizione 3.1.3** (Norma infinito).

$$f(x) = \max |x_i| = \|x\|_\infty \quad (4)$$

**Definizione 3.1.4** (Norma 1).

$$f(x) = \sum_{i=1}^n |x_i| = \|x\|_1 \quad (5)$$

**Esempio 3.1.1.** Prendiamo due vettori

$$\begin{Bmatrix} 1 \\ 1 \end{Bmatrix} \quad \begin{Bmatrix} 2 \\ -2 \end{Bmatrix}$$

e calcoliamo le varie norme:

$$\|x\|_2 = \sqrt{2} \quad \|x\|_\infty = 1 \quad \|x\|_1 = 2$$

**Definizione 3.1.5** (Norma matriciale). È una funzione del tipo  $f : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$  che deve soddisfare tre proprietà:

1.  $f(A) \geq 0 \wedge f(A) = 0 \Leftrightarrow A = 0$
2.  $f(\alpha A) = |\alpha|f(A) \quad \forall \alpha \in \mathbb{R} \quad \forall A \in \mathbb{R}^{n \times n}$
3. **Disuguaglianza triangolare:**  $f(A + B) \leq f(A) + f(B) \quad \forall A, B \in \mathbb{R}^{n \times n}$
4.  $f(A \cdot B) \leq f(A)f(B)$

e la indichiamo come

$$f(A) = \|A\|$$



## 4 Condizionamento

Studiare il condizionamento di un problema in forma  $Ax = b$  significa chiedersi di quanto cambia la soluzione perturbando di poco  $A$  e  $B$ .

### 4.1 Metodi diretti

Dato un sistema lineare  $Ax = b$  le soluzioni possono essere trovate tramite

$$x_i = f_i(a_{11}, \dots, a_{1n}, b_1, \dots, b_n) \quad i : 1 \dots n$$

Si può partire studiando la forma di  $A$  per cercare dei sistemi risolvibili facilmente.

#### 4.1.1 Matrice diagonale

Il primo caso è quando  $A$  è una matrice **diagonale**:

$$a_{ij} = 0 \iff i \neq j$$

Il determinante di una matrice diagonale è  $\det(A) = \prod_{i=1}^n a_i$  e il determinante è diverso da 0 solo se tutti gli elementi della diagonale lo sono. Possiamo riscrivere la matrice in un sistema di equazioni lineari:

$$Ax = b \Leftrightarrow \begin{cases} a_1 x_1 = b_1 \Leftrightarrow x_1 = \frac{b_1}{a_1} \\ \vdots \\ a_n x_n = b_n \Leftrightarrow x_n = \frac{b_n}{a_n} \end{cases}$$

Questo sistema lo risolviamo in  $O(n)$ . Dato però che ridurre una matrice normale in una diagonale è un processo complesso, non conviene farlo.

#### 4.1.2 Matrice triangolare

Una matrice è **triangolare inferiore** quando  $a_{ij} = 0$  per  $j > i$ , mentre è **triangolare superiore** quando  $a_{ij} = 0$  per  $i > j$ .

Per calcolare il determinante di una matrice triangolare superiore con Laplace facciamo:

$$A = \begin{bmatrix} a_{11} & 0 & \cdots & 0 \\ \vdots & \ddots & & \vdots \\ \vdots & \ddots & \ddots & 0 \\ a_{n1} & \cdots & \cdots & a_{nn} \end{bmatrix}$$

Per risolvere  $Ax = b$  con una matrice triangolare, vediamo che il sistema associato ci porta ad avere prima un'equazione con tutte le incognite fino all'ultima con una sola incognita. È intuitivo che per calcolarne la soluzione è sufficiente partire dall'ultima ed eseguire una **sostituzione all'indietro**.

---

```
function [x]=backward_substitution(a,b)
    n=length(b);
    x = zeros(n,1);
    for k=n:-1:1
        s=0;
        for j=k+1:n
            s=s+a(k,j)*x(j);
        end
        x(k)=(b(k)-s)/a(k,k);
    end
```

---

Dato che l'operazione moltiplicativa in macchina richiede leggermente più lunga di quella della somma, è sufficiente considerare le prime per calcolare la complessità di questo programma.

La complessità al caso peggio è quindi  $\frac{n \cdot (n+1)}{2} = O(n^2) = \frac{n^2}{2} + O(n)$ .

**Esempio 4.1.1** (Matrice bidiagonale superiore). Supponiamo di avere la seguente matrice bidiagonale superiore:

$$A = \begin{bmatrix} a_{11} & a_{12} & 0 & 0 \\ 0 & \ddots & \ddots & 0 \\ 0 & 0 & \ddots & a_{n-1n} \\ 0 & 0 & 0 & a_{nn} \end{bmatrix} = \begin{bmatrix} a_1 & b_1 & 0 & 0 \\ 0 & \ddots & \ddots & 0 \\ 0 & 0 & \ddots & b_{n-1} \\ 0 & 0 & 0 & a_n \end{bmatrix}$$

In questo caso il codice si può cambiare sostituendo il secondo ciclo con l'operazione

---

```
s=a(k,k+1)*x(k+1);
```

---

e la complessità diventa  $O(n)$ .

## 5 Metodi iterativi

Una classe di metodi che permette di risolvere sistemi lineari è quella dei metodi iterativi. L'idea è di costruire una successione di vettori partendo dalla matrice che convergano alla soluzione del sistema lineare.

$$\{X_k\}_{k \in \mathbb{N}}$$

**Osservazione 5.0.1.** Non potendo generare infiniti termini, ad un certo punto ci sarà bisogno di fermarsi quando si pensa di essere sufficientemente vicini alla soluzione.

Diciamo che

$$\lim_{k \rightarrow +\infty} x_k \Leftrightarrow \lim_{k \rightarrow +\infty} \|x_k - x\|_\infty = 0 \quad (6)$$

Questo è vero perché:

$$\forall j = 1, \dots, n \quad 0 \leq |x_j^{(k)} - x_j| \leq \|x^{(k)} - x\|_\infty$$

Abbiamo una matrice invertibile  $A$ . Supponiamo di scomporre la matrice come:

$$A = M - N \quad (7)$$

con l'assunzione che anche  $M$  sia invertibile ( $\det(M) \neq 0$ ). A questo punto possiamo dire che:

$$Ax = b \Leftrightarrow (M - N)x = b \Leftrightarrow Mx = Nx + b \Leftrightarrow x = M^{-1}Nx + M^{-1}b \Leftrightarrow x = Px + q$$

Dovendo trovare la soluzione di  $x = Px + q$ , possiamo scegliere un vettore iniziale  $x_0 \in \mathbb{R}^n$  e costruirci la successione di vettori

$$x_{k+1} = Px_k + q \quad (8)$$

Se questa successione converge ho trovato la soluzione del sistema lineare.

**Osservazione 5.0.2.** Nell'implementazione pratica non userò mai l'equazione 8 ma invece

$$Mx_{k+1} = Nx_k + b \quad (9)$$

### 5.1 Convergenza

**Esempio 5.1.1.** Prendiamo

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad b = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad Ax = b \Leftrightarrow x = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Partiamo definendo le due matrici per la scomposizione:

$$M = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \quad N = \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix}$$

e calcolando la matrice  $P$

$$P = M^{-1}N = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & -\frac{1}{2} \\ -\frac{1}{2} & 0 \end{bmatrix}$$

e il vettore  $q$  Iniziamo il metodo iterativo:

$$\begin{aligned} x^{k+1} &= Px^{(k)} = \begin{bmatrix} 0 & -\frac{1}{2} \\ -\frac{1}{2} & 0 \end{bmatrix} x^{(k)} \\ x^{(1)} &= \begin{bmatrix} 0 & -\frac{1}{2} \\ -\frac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -\frac{1}{2}y \\ -\frac{1}{2}x \end{bmatrix} \\ x^{(2)} &= \begin{bmatrix} 0 & -\frac{1}{2} \\ -\frac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} -\frac{1}{2}y \\ -\frac{1}{2}x \end{bmatrix} = \begin{bmatrix} \frac{1}{4}x \\ \frac{1}{4}y \end{bmatrix} \end{aligned}$$

e notiamo che la successione tende a

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Se prendiamo invece

$$M = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad N = \begin{bmatrix} -2 & 0 \\ 0 & -2 \end{bmatrix}$$

con

$$P = M^{-1}N = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} -2 & 0 \\ 0 & -2 \end{bmatrix} = \begin{bmatrix} 0 & -2 & -2 & 0 \end{bmatrix}$$

abbiamo che la successione **diverge**.

**Definizione 5.1.1.** *Un metodo iterativo è convergente se*

$$\forall x^{(0)} \in \mathbb{R}^n \quad x_k \rightarrow x \quad (10)$$

*Ovvero se per ogni vettore di partenza scelto, il metodo converge.*

**Teorema 5.1.1.** Dato  $x^{(k+1)} = Px^{(k)} + q$ , se

$$\exists \|\cdot\| \text{ t.c. } \|P\| < 1 \quad (11)$$

allora il metodo è convergente.

**Esempio 5.1.2.** Data una matrice

$$A = \begin{bmatrix} 3 & -1 & & \\ -1 & \ddots & \ddots & \\ & \ddots & \ddots & -1 \\ & & -1 & 3 \end{bmatrix}$$

definiamo le matrici per la scomposizione

$$M = \begin{bmatrix} 3 & & \\ & \ddots & \\ & & 3 \end{bmatrix} \quad N = \begin{bmatrix} 0 & 1 & & \\ 1 & \ddots & \ddots & \\ & \ddots & \ddots & 1 \\ & & 1 & 0 \end{bmatrix}$$

Per capire se il metodo è convergente dobbiamo calcolare la norma infinito di  $P$ :

$$M^{-1}N = \begin{bmatrix} 0 & \frac{1}{3} & & \\ \frac{1}{3} & \ddots & \ddots & \\ & \ddots & \ddots & \frac{1}{3} \\ & & \frac{1}{3} & 0 \end{bmatrix} \quad \|P\|_{\infty} = \frac{1}{3} + \frac{1}{3} = \frac{2}{3} < 1$$

Il problema di questo metodo iterativo è che ha complessità  $O(n)$  per ogni iterazione e non è quindi competitivo con l'eliminazione Gaussiana.

**Esempio 5.1.3.** Data una matrice

$$A = \begin{bmatrix} T & I & & \\ -I & \ddots & \ddots & \\ & \ddots & \ddots & I \\ & & -I & T \end{bmatrix} \quad T = \begin{bmatrix} 5 & -1 & & \\ -1 & \ddots & \ddots & \\ & \ddots & \ddots & -1 \\ & & -1 & 5 \end{bmatrix}$$

In questo caso i metodi iterativi sono vantaggiosi poiché anche se la matrice è predominante diagonale (e quindi permette Gauss senza problemi), a causa dell'effetto del fill-in lo rende sconsigliato.

**Esempio 5.1.4.** Data una matrice

$$A = \begin{bmatrix} n & -1 & \dots & -1 \\ -1 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & -1 \\ -1 & \dots & -1 & n \end{bmatrix}$$

calcoliamo  $P$

$$N = \begin{bmatrix} 0 & 1 & \dots & 1 \\ 1 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 1 \\ 1 & \dots & 1 & 0 \end{bmatrix} \quad P = \begin{bmatrix} 0 & \frac{1}{n} & \dots & \frac{1}{n} \\ \frac{1}{n} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \frac{1}{n} \\ \frac{1}{n} & \dots & \frac{1}{n} & 0 \end{bmatrix}$$

La sua norma vale

$$\|P\|_{\infty} = \frac{n-1}{n}$$

e quindi converge.

*Note 5.1.1.* Se la norma vale 1 non posso dire nulla sulla convergenza.

**Teorema 5.1.2.**

$$x^{(k+1)} = Px^{(k)} + q \text{ è convergente} \Leftrightarrow \phi(P) < 1 \quad (12)$$

**Esempio 5.1.5.** Data una matrice

$$A = \begin{bmatrix} 1 & & & \alpha \\ -1 & \ddots & & \\ & \ddots & \ddots & \\ & & -1 & 1 \end{bmatrix}$$

troviamo le matrici per la scomposizione

$$M = I \quad N = \begin{bmatrix} 0 & & & -\alpha \\ 1 & \ddots & & \\ & \ddots & \ddots & \\ & & 1 & 0 \end{bmatrix}$$

Per quali  $\alpha$  il metodo è convergente?

Troviamo la fattorizzazione LU della matrice per il calcolo del determinante e otteniamo così il polinomio caratteristico:

$$x^n + \alpha$$

Dobbiamo poi trovare il modulo degli autovalori per poterli confrontare:

$$x^n = -\alpha$$

$$|\lambda|^n = |-\alpha| \Rightarrow |\lambda| = \sqrt[n]{|\alpha|}$$

$$\sqrt[n]{|\alpha|} < 1 \Leftrightarrow |\alpha| < 1$$

## 5.2 Metodo di Jacobi

In questa tecnica prendiamo  $M$  come la matrice diagonale principale:

$$M = \text{diag}(A) \quad A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \quad M = \begin{bmatrix} a_{11} & & \\ & \ddots & \\ & & a_{nn} \end{bmatrix}$$

Chiaramente questo metodo è applicabile solo quando la diagonale non contiene zeri, in quanto altrimenti non sarebbe invertibile.

Possiamo ottenere la matrice  $N$  facendo  $M - A$ :

$$N = \begin{bmatrix} 0 & -a_{12} & \dots & \dots & -a_{1n} \\ -a_{21} & \ddots & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & -a_{n-1n} \\ -a_{n1} & & & -a_{nn-1} & 0 \end{bmatrix}$$

e il vettore al tempo  $k$ :

$$x_e^{(k+1)} = \frac{1}{a_{ll}} \left[ b_l - \sum_{j=1, j \neq l}^n a_{lj} x_j^{(k)} \right] \quad l = 1, \dots, n \quad (13)$$

che avrà complessità  $O(nnz(A))$  essendo molto parallelizzabile.

**Esempio 5.2.1.** Data la matrice

$$A = \begin{bmatrix} 1 & & -\alpha \\ & \ddots & \vdots \\ & & \ddots & -\alpha \\ -\beta & & & 1 \end{bmatrix}$$

la sua matrice  $J$  è

$$J = \begin{bmatrix} & & \alpha \\ & & \vdots \\ & & \alpha \\ \beta & \dots & \beta & 0 \end{bmatrix}$$

Utilizziamo la fattorizzazione LU per calcolare il polinomio caratteristico:

$$\det(\lambda I - J) = \det \begin{bmatrix} \lambda & & -\alpha \\ & \ddots & \vdots \\ & & \ddots & -\alpha \\ -\beta & \dots & 0 & \lambda \end{bmatrix} = \begin{bmatrix} I & 0 \\ -\frac{\beta}{\lambda} & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \lambda & & -\alpha \\ & \ddots & \vdots \\ & & \lambda & -\alpha \\ & & & u \end{bmatrix}$$

dove  $u = \lambda - \frac{\alpha\beta}{\lambda}$  e il raggio spettrale vale

$$\phi(J) < 1 \Leftrightarrow \sqrt{|\alpha\beta|} < 1$$

Dimostrare che per i valori per cui Jacobi è convergente, la matrice è invertibile.

Sappiamo che  $A$  non è invertibile se e solo se:

$$\exists x \neq 0 | Ax = 0 \Leftrightarrow (M - N)x = 0 \Leftrightarrow Mx = Nx \Leftrightarrow x = M^{-1}Nx \Leftrightarrow x = Px$$

Questo ci fa capire che possiamo scriverlo come  $Px = \lambda x$  che il metodo quindi non è convergente in quanto ha autovalore  $\lambda$  con autovettore  $x$ . Quindi se è convergente allora è invertibile.

### 5.3 Metodo di Gauss-Seidel

Gauss ragiona sul fatto che ad ogni iterazione avrà calcolato tutte le componenti precedenti fino a  $l-1$  al tempo  $k+1$ . Riscrive quindi l'iterazione di Jacobi nella seguente forma:

$$x_e^{(k+1)} = \frac{1}{a_{ll}} \left[ b_e - \sum_{j=1}^{l-1} a_{lj} x_j^{(k+1)} - \sum_{j=l+1}^n a_{ej} x_j^{(k)} \right] \quad (14)$$

La differenza è che questo metodo, rispetto a quello di Jacobi, perde il parallelismo ma converge più velocemente.

**Teorema 5.3.1.** Se  $A$  è predominante diagonale, allora:

- I metodi di Jacobi e Gauss-Seidel sono applicabili
- I metodi di Jacobi e Gauss-Seidel sono convergenti

**Dimostrazione 5.3.1.** Dimostriamo i due punti del teorema:

- Il primo punto è semplice in quanto se la matrice è predominante diagonale per forza di cose avremo almeno un valore diverso da 0 (perché vale la disuguaglianza stretta)
- Per dimostrare che sono convergenti vogliamo mostrare che:

$$\text{Apredominante diagonale} \Rightarrow \phi(J), \phi(GS) < 1 \Rightarrow \text{Convergente}$$

Stiamo usando il fatto che il raggio spettrale sia **sufficiente** a garantire la convergenza.

$$\begin{aligned} \det(\lambda I - P) &= \det(\lambda I - M^{-1}N) \\ &= \det(\lambda M^{-1}M - M^{-1}N) \\ &= \det(M^{-1}(\lambda M - N)) \\ &= \det(M^{-1}) \cdot \det(\lambda M - N) \\ \det(\lambda I - P) = 0 &\Leftrightarrow \det(\lambda M - N) = 0 \end{aligned}$$

Quindi se  $\lambda$  è autovalore di  $P$  allora

$$\det(\lambda M - N) = 0 \quad (15)$$

Assumiamo che esista un autovalore  $\lambda$  di  $P$  con  $|\lambda| \geq 1$ :

$$\lambda M - N = \lambda \begin{bmatrix} a_{11} & & & \\ & \ddots & & \\ & & a_{nn} & \\ & & & \ddots \end{bmatrix} - \begin{bmatrix} 0 & a_{12} & \dots & \cdot \\ & \ddots & \ddots & \vdots \\ & & \ddots & \vdots \\ & & & 0 \end{bmatrix} = \begin{bmatrix} \lambda a_{11} & a_{12} & \dots & a_{1n} \\ & \ddots & \ddots & \vdots \\ & & \ddots & \vdots \\ & & & \lambda a_{nn} \end{bmatrix}$$

dimostriamo ora che la matrice è predominante diagonale

$$|\lambda_{ll}| > \sum_{j=1, j \neq l}^n |a_{lj}| \quad |\lambda||a_{ll}| \geq |a_{ll}| > \sum_{j=1, j \neq l}^n |a_{lj}|$$

ma se è predominante diagonale allora la matrice è per forza invertibile e il determinante non è 0. Quindi abbiamo dimostrato per assurdo che l'uguaglianza 15 non è vera.

## 5.4 Criteri di arresto

Un criterio tipico di arresto per capire quando fermare le nostre iterazioni è:

$$\|x_{k+1} - x_n\| \leq \text{tolleranza} \quad (16)$$

che in codice diventa:

---

```
while(err > tol AND iterazioni <= it_max)
  calcola x_new partendo da x_old
  voluto err = ||x_new - x_old ||
  x_old = x_new
  iterazioni = iterazioni +1
end
```

---

Per capire quando fermarsi possiamo mettere in relazione la differenza al punto 16 con la soluzione come segue:

$$\begin{aligned}
 x_{k+1} - x_k &= \\
 &= x_{k+1} - x + x - x_k \\
 &= Px_k + 1 - Px - q + x - x_k \\
 &= P(x_k - x) + x - x_k \\
 &= (P - I)(x_k - x)
 \end{aligned}$$

dove la matrice  $P - I$  è invertibile dato che, quando sottraiamo l'identità, gli autovalori sono quelli della prima matrice meno 1 e quindi per avere autovalori uguali a 0 dovremmo avere che un autovalore di  $P$  è 1, ma è impossibile perché dato che è convergente abbiamo che il loro valore assoluto è  $< 1$ . Ottengo quindi:

$$\|x_k - x\| \leq \|(P - I)^{-1}\| \cdot \|x_{k+1} - x_k\| \quad (17)$$

**Esempio 5.4.1.** Data la seguente matrice

$$\begin{bmatrix} 1 & \dots & \dots & x \\ & \ddots & & \vdots \\ & & \ddots & x \\ x & \dots & \dots & 1 \end{bmatrix}$$

possiamo dire che è predominante diagonale se  $|x| < 1$ .

Per quali valori di  $x$  il metodo di Gauss-Seidel converge? Bisogna guardare il raggio spettrale. La matrice di iterazione per il metodo GS sarà fatta:

$$GS = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ x & \dots & x & 1 \end{bmatrix}^{-1} \begin{bmatrix} & -x \\ & \vdots \\ & -x \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & \dots & 0 & y_1 \\ \vdots & & \vdots & \vdots \\ \vdots & & \vdots & \vdots \\ 0 & \dots & 0 & y_n \end{bmatrix}$$

Dobbiamo quindi risolvere:

$$\begin{bmatrix} -x \\ \vdots \\ -x \\ 0 \end{bmatrix} = My \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ x & \dots & x & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} -x \\ \vdots \\ -x \\ 0 \end{bmatrix}$$

$$y_n = -x$$

$$y_2 = -x$$

$$\dots$$

$$y_{n+1} = -x$$

$$x \cdot y_1 + x \cdot y_2 + \dots + x \cdot y_{n-1} + y_n = 0$$

quindi abbiamo che il raggio spettrale vale

$$\phi(GS) = (n-1)x^2$$

che deve essere minore di 1:

$$(n-1)x^2 < 1 \Leftrightarrow x^2 < \frac{1}{n-1} \Leftrightarrow -\frac{1}{\sqrt{n-1}} < x < \frac{1}{\sqrt{n-1}}$$



## 6 Equazioni non lineari

Stiamo considerando equazioni del tipo  $f(x) = 0$  dove la funzione  $f$  non è lineare (quindi non è una retta). Di fronte a questo tipo di equazioni, ci sono due difficoltà:

- Non c'è una teoria generale sul *numero* e sull'*esistenza* delle **soluzioni**
- Non esistono metodi diretti di risoluzione

**Esempio 6.0.1.** Determinare il numero di soluzioni reali dell'equazione

$$f(x) = x \log x - 1 = 0$$

Il primo passo è tracciare un grafico approssimativo di questa funzione:

- **Dominio:**  $x > 0$
- **Limiti:**

$$\lim_{x \rightarrow +\infty} x \log x - 1 = +\infty$$

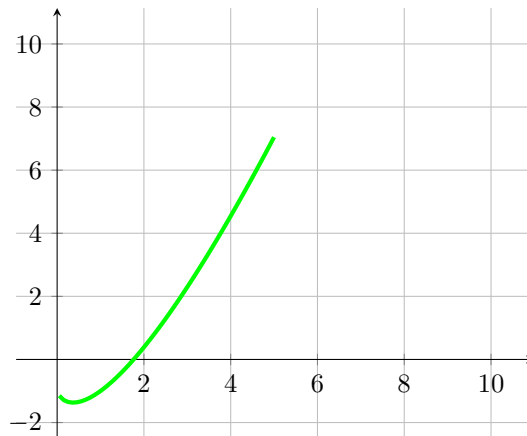
$$\lim_{x \rightarrow 0^+} x \log x = \lim_{x \rightarrow 0^+} \frac{\log x}{\frac{1}{x}} = \lim_{x \rightarrow 0^+} \frac{\frac{1}{x}}{-\frac{1}{x^2}} = \lim_{x \rightarrow 0^+} -\frac{x^2}{x} = 0 \implies \lim_{x \rightarrow 0^+} x \log x - 1 = -1$$

- **Derivata prima:**

$$f'(x) = \log x + x \cdot \frac{1}{x} = \log x + 1$$

$$f'(x) \geq 0 \Leftrightarrow \log x + 1 \geq 0 \Leftrightarrow \log x \geq -1 \Leftrightarrow x \geq \frac{1}{e}$$

- **Derivata seconda:**  $f''(x) = \frac{1}{x} \geq 0 \quad \forall x > 0$



Quindi possiamo dire che

$$\exists! \alpha \in \mathbb{R} \mid f(\alpha) = 0$$

Ci serve dare un **intervallo di localizzazione** della soluzione, ad esempio:

$$\begin{aligned} f(1) &= -1 \\ f(2) &= 2 \log 2 - 1 = \log 4 - 1 \\ \Rightarrow \alpha &\in [1, 2] \end{aligned}$$

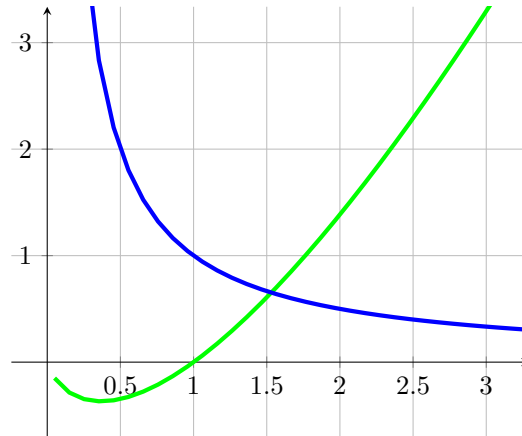
## 6.1 Tecnica della separazione

A partire da una funzione complessa, mi riconduco a funzioni più semplici e vedo dove si intercettano.

**Esempio 6.1.1.** Supponiamo di avere

$$x \log x - 1 = 0 \Leftrightarrow x \log x = 1 \Leftrightarrow \log x = \frac{1}{x}$$

Che sul grafico sono



La studiamo:

- **Dominio:**  $f \in C^\infty(\mathbb{R}^+)$
- **Limiti:**
  - $\lim_{x \rightarrow 0^+} f(x) = -1$
  - $\lim_{x \rightarrow +\infty} f(x) = +\infty$
- **Derivate:**
  - $f'(x) = \log x + 1 \geq 0 \Leftrightarrow \log x \geq -1 \Leftrightarrow x \geq \frac{1}{e}$
  - $f''(x) = \frac{1}{x} > 0 \quad \forall x \in \mathbb{R}^+$

**Esempio 6.1.2.** Data la seguente funzione

$$f(x) = e^x - 2x = 0 \Leftrightarrow e^x = 2x$$

È difficile usare il metodo della separazione perché l'intersezione non è facile da trovare. Usiamo quindi la soluzione grafica:

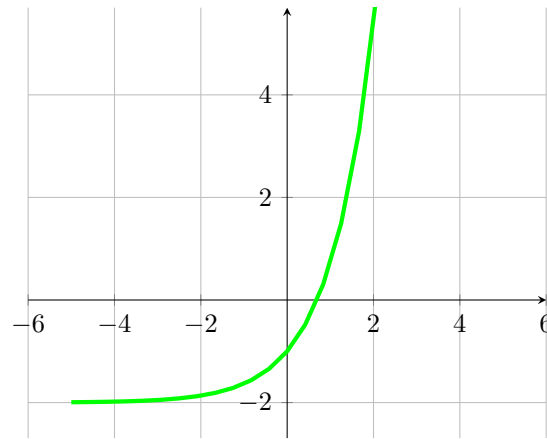
- **Dominio:**  $\forall x \in \mathbb{R}$  che possiamo scrivere anche come  $C^\infty(\mathbb{R})$
- **Limiti:**

$$\begin{aligned} \lim_{x \rightarrow +\infty} e^x - 2x &= \lim_{x \rightarrow +\infty} e^x \left(1 - \frac{2x}{e^x}\right) = 0 \\ \lim_{x \rightarrow -\infty} e^x - 2x &= +\infty \end{aligned}$$

- **Derivata prima:**

$$\begin{aligned} f'(x) &= e^x - 2 \\ f'(x) \geq 0 &\Leftrightarrow e^x \geq 2 \Leftrightarrow x \geq \log 2 \end{aligned}$$

- **Derivata seconda:**  $f''(x) = e^x$



Calcoliamo adesso il valore in  $\log 2$ :

$$f(\log 2) = e^{\log 2} - 2 \log 2 = 2 - 2 \log 2 = 2(1 - \log 2)$$

**Esempio 6.1.3.** Data la funzione:

$$f(x) = x^3 - 6x + 1 = 0 \quad (18)$$

In questo esempio abbiamo un polinomio e abbiamo quindi un'equazione algebrica di terzo grado. Quindi sappiamo il numero di soluzioni *complesse*, nel nostro caso 3. A noi però interessa il numero di soluzioni reali. Sapendo che quelle complesse devono andare sempre in coppia, potremo avere o due soluzioni complesse e una reale oppure tre soluzioni reali. Studiamo la funzione:

- **Dominio:**  $\forall x \in \mathbb{R}$
- **Limiti:**

$$\lim_{x \rightarrow +\infty} x^3 - 6x + 1 = +\infty$$

$$\lim_{x \rightarrow -\infty} x^3 - 6x + 1 = -\infty$$

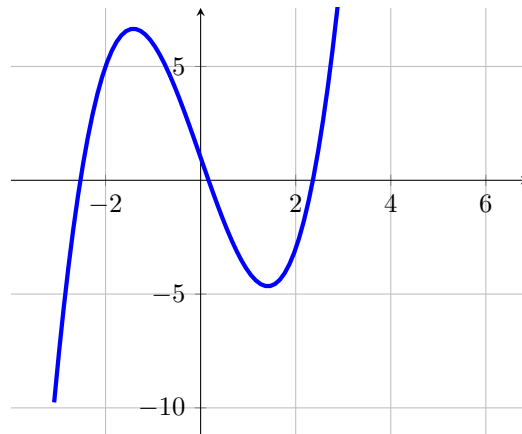
Sappiamo quindi che esiste sicuramente almeno un punto in cui la funzione vale 0 per il teorema dell'esistenza degli zeri.

- **Derivata prima:**

$$f'(x) = 3x^2 - 6$$

$$f'(x) = 0 \Leftrightarrow x^2 = 2 \Leftrightarrow x = \pm\sqrt{2} \rightarrow x < -\sqrt{2} \wedge x > \sqrt{2}$$

- **Derivata seconda:**  $f''(x) = 6x \rightarrow f''(x) > 0 \Leftrightarrow x > 0$



Quindi ci sono tre soluzioni reali, che possiamo localizzare come:

$$\beta \in [0, \sqrt{2}]$$

$$\gamma \in [\sqrt{2}, 3]$$

$$\alpha \in [-3, -2]$$

## 6.2 Metodo di bisezione

Vogliamo risolvere un'equazione  $f(x) = 0$   $f : [a, b] \rightarrow \mathbb{R}$  che assumiamo essere *continua* ( $f \in C^0([a, b])$ ). Supponiamo poi che  $f(a)f(b) < 0$ , quindi la funzione assume valori discordi agli estremi. Questo implica, grazie al teorema di esistenza degli zeri, che  $\exists \alpha \in [a, b] \mid f(\alpha) = 0$ , ovvero che esiste almeno un punto di intersezione con l'asse delle x.

Assumiamo che esista unico il punto di intersezione:

$$\exists! \alpha \in [a, b] \mid f(\alpha) = 0$$

In questa situazione possiamo costruire diverse successioni che convergono ad  $\alpha$  tramite il metodo di bisezione.

---

```

a0 = a; b0=b; ya=f(a); yb=f(b);
for k=1 : inf
    Ck = (a[k-1] + b[k-1])/2; % Calcolo il punto di mezzo dell'intervallo
    y = f(Ck);
    if(y * ya <= 0)
        ak = a[k-1]; bk=Ck; yb=y; % Qui la funzione e' discorde, mi sposto su questo
        intervallo
    else
        ak=Ck; bk=b[k-1]; ya=y; % Qui e' concorde, mi sposto sull'altro intervallo
    end
end
end

```

---

Praticamente prendo il punto medio e controllo quale parte della funzione è concorde o discorde e mi sposto di conseguenza.

**Osservazione 6.2.1.** Possiamo fare le seguenti osservazioni:

1.  $a_k \leq b_k \quad \forall k$
2.  $\alpha \in [a_k, b_k] \quad \forall k$
3.  $b_k - a_k = \frac{b_{k-1} - a_{k-1}}{2} = \dots = \frac{b_0 - a_0}{2^k}$

che implicano:

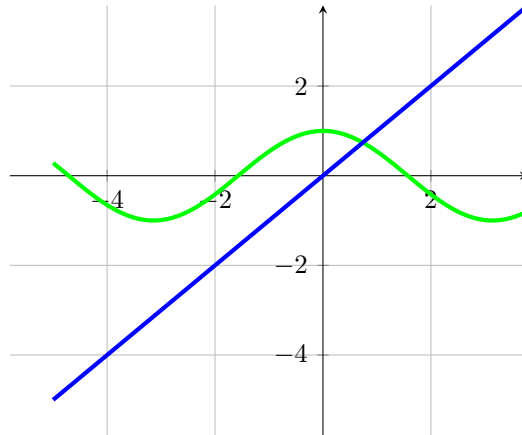
$$\begin{aligned}
 0 &\leq |\alpha - a_k| \leq b_k - a_k = \frac{b_0 - a_0}{2^k} \\
 0 &\leq |\alpha - b_k| \leq b_k - a_k = \frac{b_0 - a_0}{2^k} \\
 0 &\leq |C_k - \alpha| \leq b_{k-1} - a_{k-1} \leq \frac{b_0 - a_0}{2^k}
 \end{aligned}$$

e portando  $k$  all'infinito,  $a_k$ ,  $b_k$  e  $C_k$  tendono tutti ad  $\alpha$ .

**Esempio 6.2.1.** Supponiamo di voler risolvere

$$f(x) = x - \cos x = 0$$

Proviamo prima con la **tecnica della separazione**:



Osservando attentamente notiamo che

$$\exists! \alpha \in [0, 1]$$

poiché fuori da questo intervallo la funzione  $x$  assume valori maggiori di 1 o minori di  $-1$  che non sono valori che può assumere la funzione  $\cos(x)$ .

Per poter usare il **metodo della bisezione** devo verificare che  $f(x)$  sia continua e la sua monotonia. Come estremi per il metodo possiamo porre  $[a, b] = [0, 1]$ .

$$|C_k - \alpha| \leq \frac{b_0 - a_0}{2^k} \leq \epsilon$$

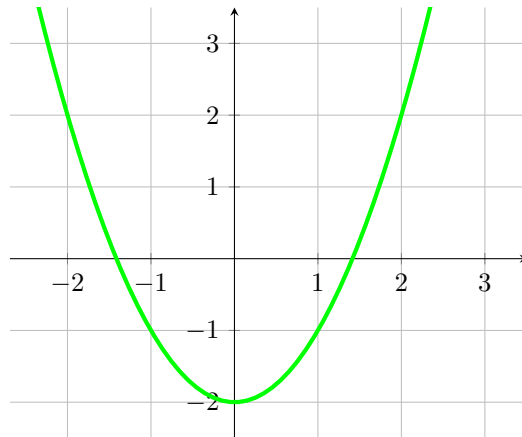
Posso a priori decidere la distanza dal mio numero finale facendo  $\log_2 \frac{1}{\epsilon}$  passi.

### 6.2.1 Approssimazione della radice

Supponiamo di voler risolvere

$$x^2 - 2 = 0$$

Sappiamo che le soluzioni sono  $x = \pm\sqrt{2}$  e che il grafico è:



Il calcolatore non sa come calcolare  $\sqrt{2}$ , il modo migliore per farglielo fare è approssimare questa funzione con il metodo di bisezione nell'intervallo  $[0, 2]$ .

### 6.2.2 Svantaggi

Questo metodo ha fondamentalmente due svantaggi:

- Ha difficoltà ad essere esteso ai numeri **complessi**, ma per il nostro corso non ci riguarda
- Il **costo computazionale** è molto influenzato dal numero di valutazioni della funzioni che faccio e, dato che non conosco precisamente quali calcoli dovrà fare, potrebbe essere un problema. Nel metodo di bisezioni faccio una valutazione per passo e il numero di passi potrebbe essere elevato.

### 6.3 Metodo del punto fisso

Questo metodo prevede di riscrivere il nostro problema come:

$$f(x) \iff x = g(x)$$

**Esempio 6.3.1.** Alcuni esempi:

$$\begin{aligned} f(x) = 0 &\iff x = x - f(x) & g(x) &= x - f(x) \\ f(x) = 0 &\iff x = x - \frac{f(x)}{h(x)} & g(x) &= x - \frac{f(x)}{h(x)} \end{aligned}$$

Poi si genera una successione da un punto iniziale:

$$\begin{cases} x_0 \in \mathbb{R} \\ x_{\alpha+1} = g(x_\alpha) \end{cases}$$

**Esempio 6.3.2.** Partiamo dalla seguente equazione:

$$x - \sqrt{x} = 0 \quad x \geq 0$$

Che posso riscrivere come:

$$x = \sqrt{x} \iff x^2 = x \iff x = 0 \vee x = 1$$

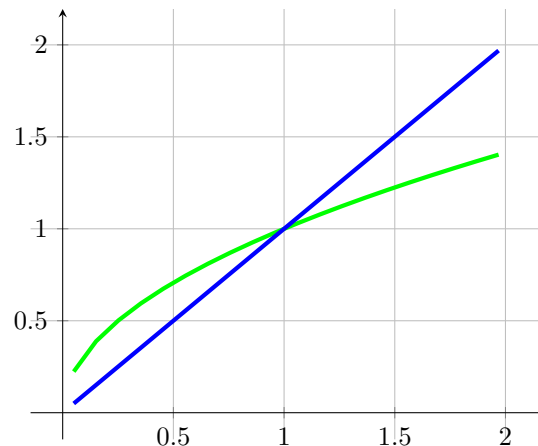
Notiamo come abbiamo ad esempio due possibili modi per scrivere l'equazione:

$$\begin{aligned} x - \sqrt{x} = 0 &\iff x = \sqrt{x} & g_1(x) &= \sqrt{x} \\ x - \sqrt{x} = 0 &\iff x^2 = x & g_2(x) &= x^2 \end{aligned}$$

Vediamo la successione generata da  $g_1(x)$ :

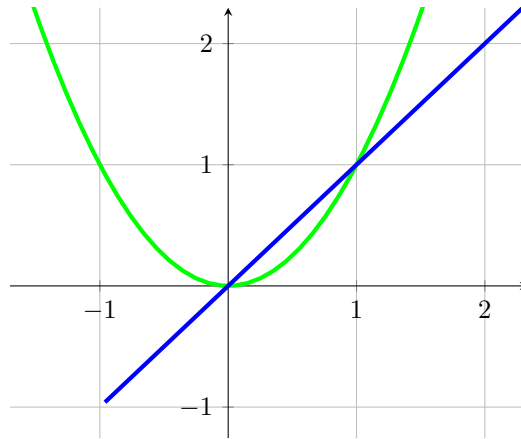
$$\begin{cases} x_0 \in \mathbb{R}^+ \\ x_{\alpha+1} = g_1(x_\alpha) = \sqrt{x_\alpha} \end{cases}$$

Che graficamente è



Con questa iterazione, se faccio i vari calcoli, da qualunque punto iniziale io scelga finirò alla soluzione 1 (intersezione di  $x$  e  $\sqrt{x}$ ).

Se scelgo invece la successione generata da  $g_2(x)$ , ottengo l'effetto contrario e trovo solo la soluzione 0 quando parto da  $0 < x_0 < 1 \vee -1 < x_0 < 0$ , altrimenti la successione diverge a  $+\infty$  per  $x_0 > 1 \vee x_0 < -1$ .



**Teorema 6.3.1** (Teorema del punto fisso). Sia  $g \in C^1([a, b])$ ,  $\alpha \in (a, b)$  e  $g(\alpha) = \alpha$  un **punto fisso**.

*Ipotesi:* se  $\exists \rho > 0 \mid \forall x \in [\alpha - \rho, \alpha + \rho] = I_\alpha \geq [a, b]$  si ha che  $|g'(x)| < 1$ .

*Tesi:* allora  $\forall x_0 \in I_\alpha$  la successione generata dal metodo  $x_{k+1} = g(x_k)$  soddisfa le seguenti proprietà:

- $x_k \in I_\alpha \quad \forall k$
- $x_k \rightarrow \alpha$

**Osservazione 6.3.1.** È importante nel precedente teorema che l'intervallo  $[a, b]$  sia centrato per garantire che da qualunque punto si inizi si riduca la distanza dal punto fisso senza però uscire dall'intervallo.

**Dimostrazione 6.3.1.** Sappiamo che se  $g(x)$  è continua allora anche  $|g'(x)|$  è continua. Per il **teorema di Weirstrass** sappiamo che una funzione continua su un intervallo limitato ammette *massimo*, nel nostro caso  $\lambda < 1$ .

Dimostriamo per induzione che  $x_0 \in I_\alpha \Rightarrow |x_k - \alpha| \leq \lambda^k \rho \quad \forall k \geq 0$ :

- *Passo base:* per  $k = 0$  abbiamo che:

$$|x_0 - \alpha| \leq \lambda^0 \rho = \rho$$

che è vero per le ipotesi:  $x_0 \in I_\alpha$ .

- *Passo induttivo:* posso scrivere:

$$|x_{k+1} - \alpha| = |g(x_k) - g(\alpha)| = |g'(\xi_k)(x_k - \alpha)| = |g'(\xi_k)| |x_k - \alpha| \leq \lambda |x_k - \alpha| = \lambda^{k+1} \rho$$

**Esempio 6.3.3.** Partiamo dalla funzione dell'esempio 6.3.2, in particolare da  $g_1$ :

$$g_1'(x) = \frac{1}{2\sqrt{x}} \Rightarrow |g_1'(x)| = \left| \frac{1}{2\sqrt{x}} \right| = \frac{1}{2\sqrt{x}}$$

Vediamo dove il suo modulo è minore di 1:

$$\frac{1}{2\sqrt{x}} < 1 \Leftrightarrow 2\sqrt{x} > 1 \Leftrightarrow \sqrt{x} > \frac{1}{2} \Leftrightarrow x > \frac{1}{4}$$

Quindi comunque si scelga un intervallo centrato in 1 tale che  $x > \frac{1}{4}$  (al massimo potrò avere  $[\frac{1}{4}, \frac{5}{4}]$ ) allora la successione converge.

Vediamo la funzione  $g_2$ :

$$g_2'(x) = 2x \Rightarrow |g_2'(x)| = |2x| = 2|x|$$

Vediamo dove il suo modulo è minore di 1:

$$2|x| < 1 \Leftrightarrow |x| < \frac{1}{2}$$

Quindi comunque si scelga un intervallo centrato in 0 tale che  $x$  è compreso tra  $-\frac{1}{2}$  e  $\frac{1}{2}$ , allora la successione converge.



**Teorema 6.3.2.** Prendiamo  $g \in C^1([a, b])$   $g(\alpha) = \alpha$   $\alpha \in (a, b)$ .  
Se  $|g'(\alpha)| < 1$ , prendiamo  $h(x) = |g'(x)| - 1$  quindi:

$$h(\alpha) = |g'(\alpha)| - 1 < 0 \Rightarrow \exists [\alpha - \rho, \alpha + \rho] = I_\alpha \mid \forall x \in I_\alpha \quad h(x) < 0 \Leftrightarrow |g'(x)| < 1$$

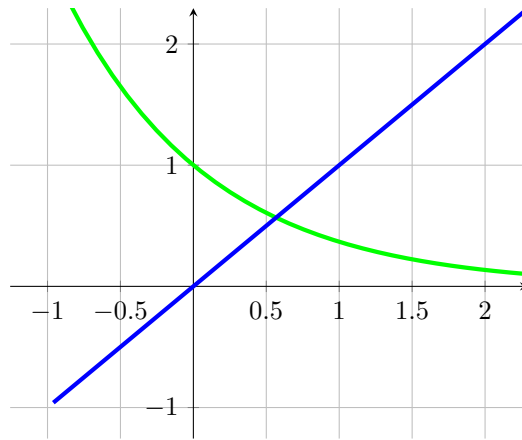
Quindi posso dire che  $\alpha$  è **punto attrattivo** (se parto vicino ad  $\alpha$  convergo su di esso, esiste un intorno e c'è convergenza locale). Altrimenti è **punto repulsivo**. Se è uguale, allora è **neutrale**.

**Esempio 6.3.4.** Prendiamo la funzione

$$f(x) = e^{-x} - x = 0$$

1. Utilizziamo il metodo di separazione grafica per determinare che esiste unica la soluzione:

$$e^{-x} - x = 0 \Leftrightarrow e^{-x} = x$$



E ho trovato un intervallo di separazione  $\alpha \in [0, 1]$ .

2. Consideriamo il metodo iterativo

$$x_{k+1} = e^{-x_k} \quad k \geq 0$$

- Dire se l'intervallo è localmente convergente. Sappiamo che

$$g'(x) = -e^{-x} \Rightarrow |g'(x)| = e^{-x}$$

quindi

$$e^{-x} < 1 \Leftrightarrow -x < 0 \Leftrightarrow x > 0$$

e posso dire che  $|g'(\alpha)| < 1$  e quindi il metodo è **localmente convergente**.

- Si dica come scegliere un punto iniziale a partire dal quale il metodo converge. Per localizzare  $\alpha$ , utilizzo il metodo di bisezione:

$$x = \frac{1}{2} \Rightarrow g\left(\frac{1}{2}\right) = \frac{1}{\sqrt{e}} > \frac{1}{2}$$

Quindi posso scegliere 1 perché ho abbastanza spazio dall'altro lato per centrare l'intervallo in  $\alpha$ .

**Esempio 6.3.5.** Prendiamo la funzione dell'esempio 6.1.1 e vediamo che

$$\alpha \in [1, 2]$$

Consideriamo due metodi iterativi:

•

$$x_{k+1} = \frac{1}{\ln x_k} \quad k \geq 0$$

•

$$x_{k+1} = e^{\frac{1}{x_k}} \quad k \geq 0$$

Verifichiamo le derivate:

•

$$g'(x) = \frac{-\frac{1}{x}}{\ln^2 x} \Rightarrow \left| \frac{-\frac{1}{x}}{\ln^2 x} \right| = \frac{1}{x \ln^2 x}$$

La valuto in  $\alpha$  tenendo a mente che è punto fisso ( $\alpha = \frac{1}{\alpha}$ ):

$$|g'(\alpha)| = \frac{1}{\alpha \ln^2 \alpha} = \frac{\alpha^2}{\alpha} = \alpha > 1$$

E quindi per questo metodo è un **punto repulsivo**.

•

$$g'(x) = e^{\frac{1}{x}} \left( -\frac{1}{x^2} \right)$$

La valuto in  $\alpha$ :

$$|g'(\alpha)| = e^{\frac{1}{\alpha}} \left( -\frac{1}{\alpha^2} \right) = \frac{1}{\alpha} < 1$$

E quindi per questo metodo è un **punto attrattivo**.

## 6.4 Metodo delle tangenti

In questo metodo (anche chiamato *di Newton*) si utilizza la **tangente** per approssimare il punto, in quanto è una buona approssimazione locale, e poi si interseca questa con l'asse delle  $x$ .

Partiamo da un punto:

$$P = (x_0, f(x_0))$$

di cui posso scrivere la retta tangente come:

$$y - f(x_0) = f'(x_0)(x - x_0)$$

*Note 6.4.1.* La funzione per questo metodo deve essere non solo continua ma anche **derivabile**. Per ottenere quindi maggiore precisione dobbiamo porre più restrizioni sulla funzione:

$$f \in C^1([a, b])$$

Intersechiamo la tangente con l'asse delle  $x$ :

$$\begin{cases} y - f(x_0) = f'(x_0)(x - x_0) \\ y = 0 \end{cases} \rightarrow x_1 - x_0 = -\frac{f(x_0)}{f'(x_0)} \Rightarrow x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

*Note 6.4.2.* Diventa quindi necessario aggiungere la seguente restrizione:

$$f'(\alpha) \neq 0$$

### 6.4.1 Convergenza locale

Il metodo delle tangenti può essere visto come un metodo di iterazione funzionale in cui:

$$x_{k+1} = g(x_k) \quad g(x) = x - \frac{f(x)}{f'(x)}$$

Verifichiamo che sia **localmente convergente**:

$$g'(x) = 1 - \frac{f'(x) \cdot f'(x) - f''(x) \cdot f(x)}{(f'(x))^2} = \frac{f''(x) \cdot f(x)}{(f'(x))^2}$$

$$|g'(\alpha)| = \left| \frac{f''(\alpha) \cdot f(\alpha)}{(f'(\alpha))^2} \right| = 0 < 1$$

*Note 6.4.3.* Questo ci porta ad aggiungere un'altra restrizione alla funzione, ovvero che sia derivabile due volte.

$$f \in C^2([a, b])$$

**Esempio 6.4.1.** Data la funzione

$$f(x) = x^2$$

che si noti non rispettare  $f'(\alpha) \neq 0$ .

Scriviamo il metodo delle tangenti:

$$x_{k+1} = x_k - \frac{x_k^2}{2x_k} = x_k - \frac{1}{2}x_k = \frac{1}{2}x_k$$

e quindi

$$x_k = \left(\frac{1}{2}\right)^k x_0$$

Vediamo quindi che può convergere anche se non viene rispettata quella clausola. Noi comunque la manteniamo per semplificarci i calcoli.

### 6.4.2 Convergenza lineare

Supponiamo di avere una successione che converge ad  $\alpha$ :

$$x_k \rightarrow \alpha \quad x_n \neq \alpha \forall n$$

Vediamo che converge linearmente se:

$$\lim_{n \rightarrow \infty} \left| \frac{x_{k+1} - \alpha}{x_n - \alpha} \right| = l \quad 0 < l < 1$$

Questo significa che ad ogni passo l'errore si riduce di un fattore  $l$ .

Se assumiamo  $l = \frac{1}{2}$  vediamo che, data una tolleranza  $\epsilon$ , il numero di passi è:

$$\left(\frac{1}{2}\right)^n \leq \epsilon \Leftrightarrow n \log \frac{1}{2} \leq \log \epsilon \Leftrightarrow -n \leq \log \epsilon \Leftrightarrow n \geq \log \frac{1}{\epsilon}$$

### 6.4.3 Convergenza quadratica

Una successione converge almeno quadraticamente se

$$\lim_{n \rightarrow \infty} \frac{|x_{k+1} - \alpha|}{|x_n - \alpha|^2} = l \in \mathbb{R}$$

Se assumiamo  $l = 1$ , data una tolleranza  $\epsilon$ , il numero di passi da fare è:

$$\left(\frac{1}{2}\right)^{2^n} \leq \epsilon \Leftrightarrow -2^n \leq \log \epsilon \Leftrightarrow 2^n \geq \log \frac{1}{\epsilon} \Leftrightarrow n \geq \log \log \frac{1}{\epsilon}$$

Quindi con un metodo che converge quadraticamente, anche nel caso in cui abbiamo un fattore di riduzione peggiore rispetto che ad un metodo lineare, abbiamo un numero comunque ridotto di passi da eseguire.

**Dimostrazione 6.4.1** (Convergenza quadratica del metodo delle tangenti).

$$\begin{aligned}
 0 &= f(\alpha) = f(x_n) + f'(x_n)(\alpha - x_n) + f''(\psi_n)(x_n - \alpha)^2 \quad \psi_n \in [x_n, \alpha] \\
 -\frac{f(x_n)}{f'(x_n)} &= (\alpha - x_n) + \frac{f''(\psi_n)(x_n - \alpha)^2}{2f'(x_n)} \\
 &= \alpha + \frac{f''(\psi_n)(x_n - \alpha)^2}{2f'(x_n)} \\
 \left| \frac{x_{n+1} - \alpha}{(x_n - \alpha)^2} \right| &= \left| \frac{f''(\psi_n)}{2f'(\psi_n)} \right| \rightarrow \frac{|f'(\alpha)|}{|2f'(\alpha)|} = l \in \mathbb{R}
 \end{aligned}$$

**Esempio 6.4.2.** Data la funzione:

$$e^x + x + 2 = 0$$

Se studiamo la funzione vediamo:

- *Dominio:*  $f \in C^\infty(\mathbb{R})$
- *Limiti:*
  - $\lim_{x \rightarrow +\infty} f(x) = +\infty$
  - $\lim_{x \rightarrow -\infty} f(x) = -\infty$
- *Derivate:*
  - $f'(x) = e^x + 1 \geq 0 \quad \forall x \in \mathbb{R}$
  - $f''(x) = e^x \geq 0 \quad \forall x \in \mathbb{R}$



Possiamo dire che:

$$\exists! \alpha \in \mathbb{R} \mid f(\alpha) = 0 \quad f \in C^2(\mathbb{R}) \quad f'(\alpha) \neq 0$$

e che quindi questo metodo è **localmente convergente**.

**Esempio 6.4.3.** Data la funzione:

$$xa - 1 = 0 \quad a \in \mathbb{R} \quad a \neq 0$$

Possiamo facilmente vedere che

$$x = \frac{1}{a}$$

Con un metodo lineare converge in un passo. Proviamo però a vederla come:

$$a = \frac{1}{x} \Leftrightarrow \frac{1}{x} - a = 0$$

e applichiamo il metodo delle tangenti:

$$x_{k+1} = x_k + \frac{\left(\frac{1}{x_k} - a\right)}{\frac{1}{x_k^2}} = x_k + x_k^2 \cdot \left(\frac{1}{x_k} - a\right) = 2x_k - ax_k^2$$

Che ha bisogno di un numero di passi maggiore. Però non fa divisioni ma solamente prodotti e addizioni. Questo nei primi calcolatori che avevano solo addizioni, sottrazioni e moltiplicazioni, era il metodo che veniva utilizzato per approssimare le divisioni.

**Esempio 6.4.4.** Data la funzione dell'esercizio 6.1.1. Vediamo che se partiamo con il metodo delle tangenti da  $\frac{1}{e}$ , la successione diverge. Se invece si sceglie un valore di partenza tra 0 e  $\frac{1}{e}$ , la successione passa dall'aritmetica reale a quella complessa dato che la macchina gestisce anche i logaritmi di numeri negativi.

#### 6.4.4 Convergenza globale

**Teorema 6.4.1** (Teorema di convergenza in largo). Sia:

$$f \in C^2([a, b]) \quad f(\alpha) = 0, \alpha \in (a, b)$$

Se  $\exists \delta > 0$  tale che  $\forall x \in (\alpha, \alpha + \delta] \subseteq (a, b)$  si ha:

1.  $f'(x) \neq 0$
2.  $f(x)f''(x) > 0$

Allora  $\forall x_0 \in (\alpha, \alpha + \delta]$  la successione generata dal metodo della tangente con punto iniziale  $x_0$  converge ad  $\alpha$ .

**Dimostrazione 6.4.2.** Partiamo dalla funzione dell'esempio 6.4.2. Possiamo dire che

$$\exists! \alpha \in \mathbb{R} \text{ t.c. } f(\alpha) = 0$$

Se prendiamo  $x_0 > \alpha$ , il teorema può essere applicato senza problemi, e abbiamo quindi la convergenza. Se invece  $x_0 < \alpha$ , il secondo punto del teorema non è rispettato. Possiamo però dimostrare che è sufficiente un passo del metodo delle tangenti per garantirla:

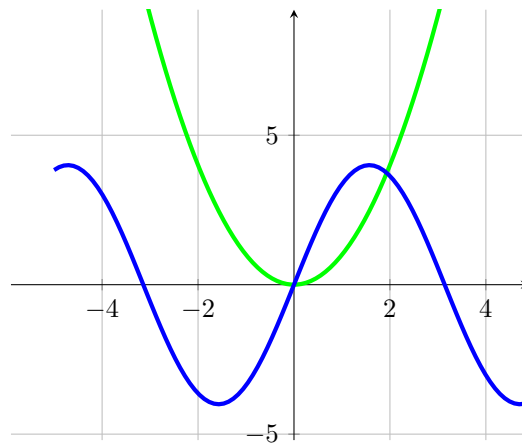
$$\begin{aligned} x_1 - \alpha &> 0 \\ x_1 &> \alpha \\ x_1 - \alpha &= g(x_0) - g(\alpha) = g'(\psi)(x_0 - \alpha) = \\ &= \frac{f(\psi) - f''(\psi)}{(f'(\psi))^2} \cdot (x_0 - \alpha) \\ \frac{f(\psi) - f''(\psi)}{(f'(\psi))^2} &< 0 \quad (x_0 - \alpha) < 0 \end{aligned}$$

**Esempio 6.4.5.** Partiamo dalla funzione:

$$f(x) = x^2 - 4 \sin x = 0 \Leftrightarrow x^2 = 4 \sin x$$

Dal grafico possiamo dedurre che:

- 0 è soluzione
- $\exists! \alpha > 0$  tale che  $f(\alpha) = 0 \wedge \alpha \in [\frac{\pi}{2}, \pi]$



Applichiamo ora il teorema per la convergenza globale per cercare un intervallo da cui partire per il metodo della tangente.

$$f'(x) = 2x - 4 \cos x$$

$$f''(x) = 2 + 4 \sin x$$

L'intervallo che verifica le due ipotesi del teorema è:

$$[0, \pi]$$

**Esempio 6.4.6.** Partiamo dalla seguente funzione:

$$f(x) = \log x + 3x - 5 = 0$$

• **Condizioni di esistenza:**  $x > 0$   $f \in C^\infty(\mathbb{R}^+)$

• **Limiti:**

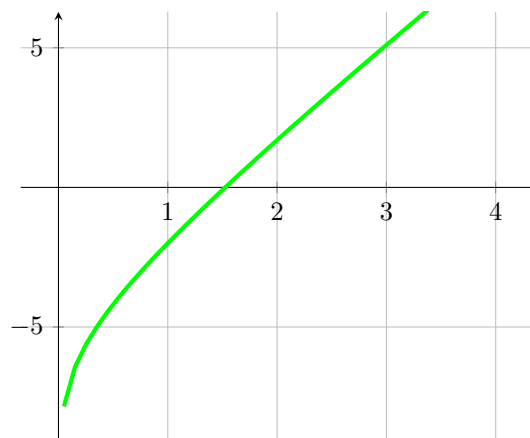
$$- \lim_{x \rightarrow 0^+} f(x) = -\infty$$

$$- \lim_{x \rightarrow +\infty} f(x) = +\infty$$

• **Derivate:**

$$- f'(x) = \frac{1}{x} + 3 > 0 \quad \forall x \in \mathbb{R}^+$$

$$- f''(x) = -\frac{1}{x} < 0 \quad \forall x \in \mathbb{R}^+$$



1. Dire se il metodo delle tangenti è localmente convergente.

Dal grafico vediamo che

$$\exists! \alpha \in \mathbb{R}^+ \text{ tale che } f(\alpha) = 0 \quad \alpha \in [1, 2]$$

e quindi:

$$f \in C^2(\mathbb{R}^+) \wedge f'(\alpha) \neq 0 \implies \text{Il metodo converge}$$

2. Il metodo delle tangenti converge per i seguenti valori?

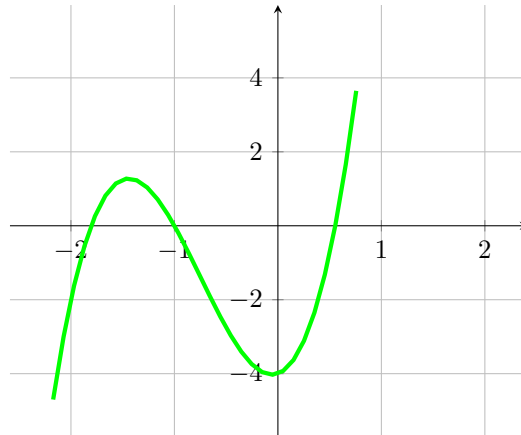
- $x_0 = 1$ : sì perché è nell'intervallo  $[1, \alpha]$
- $x_0 = 2$ : applichiamo un passo del metodo

$$x_1 = 2 - \frac{\log 2 + 1}{\frac{7}{2}} = 2 - \frac{2(\log 2 + 1)^2}{7} = \frac{2(7 - \log 2 - 1)}{7} = \frac{2}{7}(6 - \log 2) > 0$$

e vediamo che anche qui converge

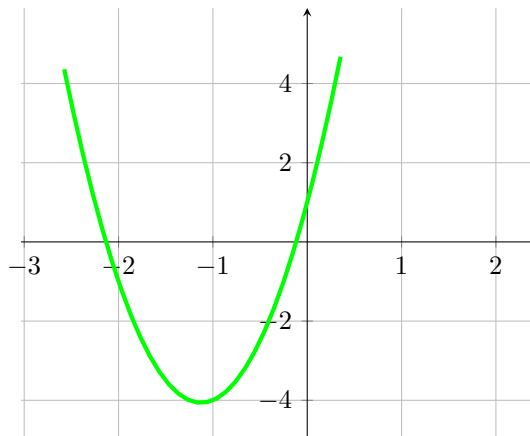
**Esempio 6.4.7.** Partiamo dalla funzione

$$f(x) = a_0 + a_1x + a_2x^2 + a_3x^3 = 0 \quad a_3 > 0$$



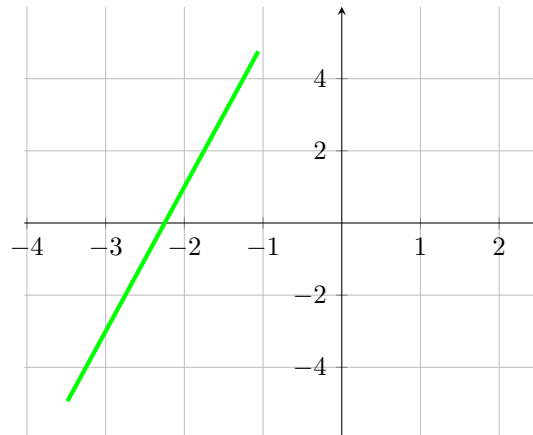
La derivata prima sarà:

$$f'(x) = a_1 + 2a_2x + 3a_3x^2$$



Mentre la derivata seconda:

$$f''(x) = 2a_2 + 6a_3x$$



I tre punti soluzione della funzione avranno tutti convergenza locale. Quella del punto centrale sarà però più complessa da trovare perché si rischia uno stallo e si cerca quindi di evitare di calcolarle con il metodo delle tangenti.

**Esempio 6.4.8.** Partiamo dall'equazione secolare:

$$f(x) = 1 - \sum_{k=1}^{1000} \frac{1}{x - k}$$

1. Determinare il numero di soluzioni reali di questa equazione
2. Pensare ad un metodo per l'approssimazione di queste funzioni



## 7 Teoremi

### 7.1 Maggiorazione dell'errore di macchina

**Teorema 7.1.1.** Sia

$$x \in \mathbb{R}$$

con

$$\omega \leq |x| \leq \Omega$$

vale

$$|\epsilon_x| = \left| \frac{\text{trn}(x) - x}{x} \right| \leq u = B^{1-t} \quad (19)$$

#### 7.1.1 Dimostrazione

Definiamo  $x$  come

$$x = (-1)^s B^p \alpha$$

L'errore assoluto  $|\text{trn}(x) - x|$  è maggiorato dalla distanza tra due numeri di macchina consecutivi e quindi:

$$|\text{trn}(x) - x| \leq B^{p-t}$$

Inoltre vale

$$|x| \geq B^{p-1}$$

Quindi:

$$|\epsilon_x| = \left| \frac{\text{trn}(x) - x}{x} \right| \leq \frac{B^{p-t}}{B^{p-1}} = B^{1-t} = u$$

### 7.2 Errore analitico e totale

**Definizione 7.2.1** (Errore analitico). *Nel calcolo di  $h(x) \neq 0$  tramite la sua approssimazione  $f(x)$ , l'errore analitico è*

$$\epsilon_{an} = \frac{f(x) - h(x)}{h(x)} \quad (20)$$

**Definizione 7.2.2** (Errore totale). *Combinando l'errore analitico con quello **algoritmo e inerente** otteniamo:*

$$\epsilon_{tot} = \frac{g(\tilde{x}) - h(x)}{h(x)} \doteq \epsilon_{an} + (\epsilon_{in} + \epsilon_{alg}) \quad (21)$$

### 7.3 Localizzazione degli autovalori

**Definizione 7.3.1** (Cerchi di Gershgorin). *Sia*

$$A = (a_{i,j}) \in \mathbb{C}^{n \times n}$$

*Definiamo i cerchi di Gershgorin  $K_i$  con  $1 \leq i \leq n$  come*

$$K_i = \{z \in \mathbb{C} : |z - a_{i,i}| \leq \sum_{j=1, j \neq i}^n |a_{i,j}|\} \quad (22)$$

**Teorema 7.3.1.** Vale

$$\lambda \text{ autovalore di } A \implies \lambda \in \bigcup_{i=1}^n K_i \quad (23)$$

#### 7.3.1 Dimostrazione

Sia  $\lambda$  **autovalore** di  $A$  con corrispondente **autovettore destro**  $\mathbf{x}$ .

La relazione  $A\mathbf{x} = \lambda\mathbf{x}$  implica che

$$\sum_{j=1}^n a_{i,j}x_j = \lambda x_i \quad 1 \leq i \leq n$$

da cui

$$(\lambda - a_{i,i})x_i = \sum_{j=1, j \neq i}^n a_{i,j}x_j \quad 1 \leq i \leq n \quad (24)$$

Sia  $p$  l'indice di una componente di modulo massimo di  $\mathbf{x}$ , ad esempio  $|x_p| = \|\mathbf{x}\|_\infty$ . Dato che  $\mathbf{x} \neq 0$ , allora  $|x| > 0$ .

Poniamo  $i = p$  nell'equazione precedente e otteniamo

$$(\lambda - a_{p,p})x_p = \sum_{j=1, j \neq p}^n a_{p,j}x_j$$

che, con i valori assoluti, diventa

$$|(\lambda - a_{p,p})x_p| = |\lambda - a_{p,p}||x_p| = \left| \sum_{j=1, j \neq p}^n a_{p,j}x_j \right| \leq \sum_{j=1, j \neq p}^n |a_{p,j}||x_j|$$

Se dividiamo entrambi i membri per  $|x_p|$  otteniamo

$$|\lambda - a_{p,p}| \leq \sum_{j=1, j \neq p}^n |a_{p,j}| \frac{|x_j|}{|x_p|} \leq \sum_{j=1, j \neq p}^n |a_{p,j}|$$

E questo implica che  $\lambda \in K_p$ .

## 7.4 Esistenza della fattorizzazione LU

**Teorema 7.4.1.** Sia

$$A \in \mathbb{R}^{n \times n}$$

Se  $A(1:k, 1:k)$  è invertibile per  $k = 1, 2, \dots, n-1$ , allora esiste unica la fattorizzazione LU di  $A$ .

### 7.4.1 Dimostrazione

Per **induzione** sulla dimensione  $n$ .

- **Caso base:** per  $n = 1$  vale che  $A = [a] = [1][a]$  che è l'unica fattorizzazione LU
- **Ipotesi induttiva:** supponiamo che il teorema sia vero per matrici di ordine  $m \leq n-1$
- **Passo induttivo:** la fattorizzazione può essere scritta come

$$\begin{bmatrix} A(1:n-1, 1:n-1) & \mathbf{z} \\ \mathbf{v}^T & \alpha \end{bmatrix} = \begin{bmatrix} L(1:n-1, 1:n-1) & 0 \\ \omega^T & 1 \end{bmatrix} \cdot \begin{bmatrix} U(1:n-1, 1:n-1) & \mathbf{y} \\ 0^T & \beta \end{bmatrix}$$

Questo è equivalente al sistema

$$\begin{cases} A(1:n-1, 1:n-1) = L(1:n-1, 1:n-1) \cdot U(1:n-1, 1:n-1) \\ \mathbf{z} = L(1:n-1, 1:n-1)\mathbf{y} \\ \mathbf{v}^T = \omega^T \cdot U(1:n-1, 1:n-1) \\ \alpha = \omega^T \mathbf{y} + \beta \end{cases}$$

1. Per ipotesi induttiva  $A(1:n-1, 1:n-1)$  è invertibile e di conseguenza  $L(1:n-1, 1:n-1)$  e  $U(1:n-1, 1:n-1)$  sono i suoi fattori triangolari
2. Per ipotesi induttiva  $A(1:n-1, 1:n-1)$  è invertibile, quindi lo è anche  $U(1:n-1, 1:n-1)$  lo è, di conseguenza le equazioni 2 e 3 ammettono soluzione unica
3. Dati  $\omega$  e  $\mathbf{y}$  l'equazione 4 permette di identificare univocamente  $\alpha$  e  $\beta$

## 7.5 Convergenza con norma

**Teorema 7.5.1.** Un metodo iterativo è **convergente** se esiste una norma matriciale indotta da una norma vettoriale  $\|\cdot\|$  su  $\mathbb{R}^n$  tale per cui

$$\|P\| < 1$$

### 7.5.1 Dimostrazione

Dalle relazioni

$$\mathbf{x}^{(k+1)} = P\mathbf{x}^{(k)} + \mathbf{q} \quad \mathbf{x} = P\mathbf{x} + \mathbf{q}$$

segue

$$\mathbf{e}^{(k+1)} = \mathbf{x}^{(k+1)} - \mathbf{x} = P(\mathbf{x}^{(k)} - \mathbf{x}) = P\mathbf{e}^{(k)} \quad k \geq 0$$

e quindi

$$\mathbf{e}^{(k+1)} = P^{k+1}\mathbf{e}^{(0)} \quad k \geq 0$$

Passando alla norma vettoriale

$$\|\mathbf{e}^{(k+1)}\| = \|P^{k+1}\mathbf{e}^{(0)}\| \leq \|P^{k+1}\| \|\mathbf{e}^{(0)}\|$$

da cui

$$0 \leq \|\mathbf{e}^{(k+1)}\| \leq \|P\|^{k+1} \|\mathbf{e}^{(0)}\|$$

da cui per il **teorema del confronto**

$$\lim_{k \rightarrow +\infty} \|\mathbf{e}^{(k+1)}\| = \lim_{k \rightarrow +\infty} \|\mathbf{x}^{(k+1)} - \mathbf{x}\| = 0 \quad \forall \mathbf{e}^{(0)}, \mathbf{x}^{(0)}$$

## 7.6 Convergenza con raggio spettrale

**Teorema 7.6.1.** Se il metodo iterativo è convergente allora  $\rho(P) < 1$ .

### 7.6.1 Dimostrazione

Sia  $\lambda$  tale che  $|\lambda| = \rho(P)$  e  $\mathbf{v}$  il corrispondente **autovettore**.

Sia  $\mathbf{x}^{(0)} = \mathbf{x} + \mathbf{v}$  con  $\mathbf{x} = A^{-1}\mathbf{b}$  soluzione del sistema  $A\mathbf{x} = \mathbf{b}$ .

La successione con punto iniziale  $\mathbf{x}^{(0)}$  converge a  $\mathbf{x}$ .

Si ha

$$\mathbf{e}^{(k+1)} = P^{k+1}\mathbf{e}^{(0)} = P^{k+1}\mathbf{v} = \lambda^{k+1}\mathbf{v}$$

da cui

$$\|\mathbf{e}^{(k+1)}\| = \|\lambda^{k+1}\mathbf{v}\| = |\lambda|^{k+1}\|\mathbf{v}\|$$

e quindi

$$\lim_{k \rightarrow +\infty} |\lambda|^k = 0$$

che implica

$$|\lambda| < 1$$

## 7.7 Predominanza diagonale

**Teorema 7.7.1.** Se  $A = (a_{i,j}) \in \mathbb{R}^{n \times n}$  è **predominante diagonale** allora valgono:

1.  $A$  è invertibile
2. I metodi di Jacobi e Gauss-Seidel sono applicabili
3. I metodi di Jacobi e Gauss-Seidel sono convergenti

### 7.7.1 Dimostrazione

1. Dal teorema di Gerschgorin vale

$$|0 - a_{i,i}| = |a_{i,i}| > \sum_{j=1, j \neq i}^n |a_{i,j}| \quad 1 \leq i \leq n$$

e dunque

$$0 \notin \bigcup_{i=1}^n K_i$$

- 2.

$$|a_{i,i}| > \sum_{j=1, j \neq i}^n |a_{i,j}| \geq 0 \implies |a_{i,i}| \neq 0 \quad 1 \leq i \leq n$$

3. Dalla relazione

$$\det(P - \lambda I_n) = \det(M^{-1}N - \lambda I_n) = \det(N - \lambda M) = (-1)^n \det(\lambda M - N)$$

Quindi  $\lambda$  è autovalore di  $P$  se e solo se  $\det(\lambda M - N) = 0$ .

Assumiamo che  $|\lambda| \geq 1$ . Vediamo che la matrice  $\lambda M - N$  è predominante diagonale:

$$|a_{i,i}| > \sum_{j=1, j \neq i}^n |a_{i,j}| = \sum_{j=1}^{i-1} |a_{i,j}| + \sum_{j=i+1}^n |a_{i,j}| \quad 1 \leq i \leq n$$

$$|\lambda| |a_{i,i}| > |\lambda| \sum_{j=1}^{i-1} |a_{i,j}| + |\lambda| \sum_{j=i+1}^n |a_{i,j}| \geq |\lambda| \sum_{j=1}^{i-1} |a_{i,j}| + \sum_{j=i+1}^n |a_{i,j}|$$

e quindi vale la predominanza diagonale

$$|\lambda a_{i,i}| > \sum_{j=1}^{i-1} |\lambda a_{i,j}| + \sum_{j=i+1}^n |a_{i,j}|$$

$$|\lambda a_{i,i}| > \sum_{j=1}^{i-1} |a_{i,j}| + \sum_{j=i+1}^n |a_{i,j}|$$

Per il punto 1 sappiamo che allora la matrice è anche invertibile:

$$\forall \lambda \in \mathbb{C}, |\lambda| \geq 1 \implies \det(\lambda M - N) \neq 0$$

Quindi per gli autovalori delle matrici di iterazione deve valere

$$|\lambda| < 1$$

e dunque

$$\rho(P) < 1$$

## 7.8 Convergenza del metodo di bisezione

**Teorema 7.8.1.** Sia  $f : [a, b] \rightarrow \mathbb{R}$  con:

- $f \in C^0([a, b])$
- $f(a)f(b) < 0$

allora si ha, con  $f(\epsilon) = 0$ ,

$$\lim_{k \rightarrow +\infty} a_k = \lim_{k \rightarrow +\infty} b_k = \lim_{k \rightarrow +\infty} c_k = \epsilon \in [a, b] \quad (25)$$

### 7.8.1 Dimostrazione

Per costruzione valgono:

- $a_{k+1} \geq a_k$
- $b_{k+1} \leq b_k$
- $c_k \in [a_k, b_k] \subset [a, b]$
- $0 \leq b_k - a_k \leq \frac{b-a}{2^{k-1}}$
- $f(a_k)f(b_k) \leq 0$
- $k \geq 1$

Quindi esistono  $\epsilon, \eta \in [a, b]$  tali che

$$\lim_{k \rightarrow +\infty} a_k = \epsilon \quad \lim_{k \rightarrow +\infty} b_k = \eta$$

Per il **teorema del confronto** segue che

$$\epsilon = \lim_{k \rightarrow +\infty} a_k = \lim_{k \rightarrow +\infty} b_k = \eta = \lim_{k \rightarrow +\infty} c_k$$

Per la continuità di  $f$  si ha

$$\lim_{k \rightarrow +\infty} f(a_k)f(b_k) = f(\epsilon)^2 \leq 0$$

che implica

$$f(\epsilon) = 0$$

## 7.9 Limite metodo funzionale

**Teorema 7.9.1.** Sia  $g : [a, b] \rightarrow \mathbb{R}$  con:

- $g \in C^1([a, b])$
- $g(\epsilon) = \epsilon$
- $\epsilon \in (a, b)$

Se esiste  $\rho > 0$  tale che:

$$|g'(x)| < 1 \quad \forall x \in [\epsilon - \rho, \epsilon + \rho] = I_\epsilon \subset [a, b]$$

allora per ogni  $x_0 \in I_\epsilon$  la successione generata soddisfa

1.  $x_k \in I_\epsilon \quad \forall k \geq 0$
2.  $\lim_{k \rightarrow +\infty} x_k = \epsilon$

### 7.9.1 Dimostrazione

Dal **teorema di Weirstrass**, dato che  $g'$  è continua e  $I_\epsilon$  chiuso e limitato, vale

$$\lambda = \max_{x \in I_\epsilon} |g'(x)| < 1$$

Si dimostra che la successione generata soddisfa:

$$|x_k - \epsilon| \leq \lambda^k \rho \quad k \geq 0$$

da cui seguono i due punti:

1.  $|x_k - \epsilon| \leq \lambda^k \rho \leq \rho \implies x_k \in I_\epsilon$
2. per il **teorema del confronto**  $0 \leq |x_k - \epsilon| \leq \lambda^k \rho \implies \lim_{k \rightarrow +\infty} |x_k - \epsilon| = 0$

Procediamo per **induzione** su  $k$ :

- **Passo base:** per  $k = 0$  si ha  $|x_0 - \epsilon| \leq \lambda^0 \rho = \rho$
- **Ipotesi induttiva:** assumiamo che valga la disequazione precedente fino a  $k$
- **Passo induttivo:** per il **teorema di Lagrange** vale

$$|x_{k+1} - \epsilon| = |g(x_k) - g(\epsilon)| = |g'(\eta_k)(x_k - \epsilon)| = |g'(\eta_k)| |x_k - \epsilon| \quad |\eta_k - \epsilon| \leq |x_k - \epsilon|$$

Per ipotesi induttiva segue che  $\eta_k \in I_\epsilon$  e quindi

$$|x_{k+1} - \epsilon| = |g'(\eta_k)| |x_k - \epsilon| \leq \lambda |x_k - \epsilon| \leq \lambda \lambda^k \rho = \lambda^{k+1} \rho$$

## 7.10 Convergenza locale metodo funzionale

**Teorema 7.10.1.** Sia  $g : [a, b] \rightarrow \mathbb{R}$  con:

- $g \in C^1([a, b])$
- $g(\epsilon) = \epsilon$
- $\epsilon \in (a, b)$

Se esiste  $|g'(\epsilon)| < 1$  allora il metodo è localmente convergente in  $\epsilon$ .

### 7.10.1 Dimostrazione

Sia  $h : [a, b] \rightarrow \mathbb{R}$  tale che  $h(x) = |g'(x)| - 1$ . Si ha che  $h(x) \in C^0([a, b])$  e  $h(\epsilon) < 0$ . Dal **teorema della permanenza del segno** abbiamo che esiste

$$I_\epsilon = [\epsilon - \rho, \epsilon + \rho] \subset [a, b]$$

tale che

$$h(x) = |g'(x)| - 1 < 0 \quad \forall x \in I_\epsilon$$

La tesi segue dal teorema precedente.

## 7.11 Convergenza locale metodo delle tangenti

**Teorema 7.11.1.** Sia  $f : [a, b] \rightarrow \mathbb{R}$  con:

- $f \in C^2([a, b])$
- $f(\epsilon) = \epsilon$
- $\epsilon \in (a, b)$
- $f'(\epsilon) \neq 0$

allora il metodo è localmente convergente, ovvero

se esiste  $\rho > 0$  tale che per ogni  $x_0 \in [\epsilon - \rho, \epsilon + \rho] = I_\epsilon \subset [a, b]$  la successione soddisfa:

1.  $x_k \in I_\epsilon \quad \forall k \geq 0$
2.  $\lim_{k \rightarrow +\infty} x_k = \epsilon$

Se la successione soddisfa  $x_k \neq \epsilon \quad k \geq 0$  allora la convergenza è almeno **quadratica**.

### 7.11.1 Dimostrazione

Da  $f'(\epsilon) \neq 0$  per il **teorema della permanenza del segno** segue che esiste

$$I'_\epsilon = [\epsilon - \rho', \epsilon + \rho'] \subset [a, b]$$

tale che

$$f'(x) \neq 0 \quad \forall x \in I'_\epsilon$$

Si verifica che la funzione di iterazione  $g : I'_\epsilon \rightarrow \mathbb{R}$  definita come

$$g(x) = x - \frac{f(x)}{f'(x)}$$

soddisfa

1.  $g \in C^1(I'_\epsilon)$
2.  $g'(\epsilon) = \frac{f(\epsilon)f''(\epsilon)}{(f'(\epsilon))^2} = 0$

e quindi la tesi segue dal teorema precedente.

Per la stima della velocità di convergenza, dallo **sviluppo di Taylor** al secondo ordine otteniamo

$$0 = f(\epsilon) = f(x_k) + f'(x_k)(\epsilon - x_k) + \frac{f''(\eta_k)(\epsilon - x_k)^2}{2} \quad |\eta_k - \epsilon| \leq |x_k - \epsilon|$$

da cui

$$x_{k+1} - \epsilon = \frac{f''(\eta_k)(\epsilon - x_k)^2}{2f'(x_k)}$$

e per continuità delle due derivate si ottiene

$$\lim_{k \rightarrow +\infty} \frac{|x_{k+1} - \epsilon|}{|x_k - \epsilon|^2} = \left| \frac{f''(\epsilon)}{2f'(\epsilon)} \right|$$