

# Bibbia di Calcolo Numerico

Questo documento è stato redatto per fornire una guida completa e corretta, in italiano, agli argomenti principali del corso di Calcolo Numerico del Prof. Federico Piazzon (laurea triennale in Informatica). La struttura segue l'ordine logico di uno studio teorico: in ogni sezione sono presenti le definizioni, i concetti chiave, gli algoritmi più importanti e, dove richiesto, le dimostrazioni essenziali. Si presta particolare attenzione alla chiarezza del linguaggio e alla coerenza formale.

## 1 Rappresentazione dei numeri e aritmetica in virgola mobile

### 1.1 Sistema binario e standard IEEE 754

I calcolatori rappresentano i numeri in forma binaria, cioè come sequenze di 0 e 1. La norma IEEE 754 definisce il formato dei numeri in virgola mobile: un numero reale è rappresentato come

$$(-1)^s \times m \times 2^e$$

dove  $s$  è il bit di segno,  $m$  è la mantissa normalizzata (numero binario compreso tra 1 e 2) e  $e$  è l'esponente con opportuno bias. I formati più usati sono single precision (32 bit) e double precision (64 bit); in entrambi i casi solo un sottoinsieme infinito dei numeri reali è rappresentabile. La precisione macchina  $\varepsilon$  (machine epsilon) misura l'errore relativo massimo introdotto dall'arrotondamento: per il formato double, ad esempio,  $\varepsilon \approx 2^{-53}$  <sup>1</sup>.

Ogni operazione elementare viene eseguita secondo la regola

$$\text{fl}(a \circ b) = (a \circ b)(1 + \varepsilon), \text{ con } |\varepsilon| \leq \varepsilon,$$

dove  $\circ \in \{+, -, \times, /\}$  e fl indica il risultato dell'operazione macchina <sup>1</sup>. Poiché l'errore  $\varepsilon$  si somma ai vari passi di un algoritmo, la propagazione degli errori è un tema centrale del calcolo numerico.

### 1.2 Errori assoluto e relativo

Dato un valore "vero"  $A$  e un'approssimazione  $B$ , l'errore assoluto è  $|A - B|$  e l'errore relativo è  $|A - B| / |A|$  <sup>2</sup>. Nello studio dell'aritmetica in virgola mobile si misura spesso l'errore relativo in termini di  $\varepsilon$ .

### 1.3 Errore di rappresentazione e cancellazione numerica

Un numero reale che non è rappresentabile esattamente viene sostituito con il numero macchina più vicino, generando un errore di rappresentazione. Nei calcoli in virgola mobile la cancellazione numerica si verifica quando si sottraggono due quantità quasi uguali: il risultato può perdere molte cifre significative, amplificando gli errori relativi presenti negli operandi. Per mitigare questo fenomeno occorre evitare sottrazioni tra numeri quasi identici e preferire formule algebricamente equivalenti ma numericamente più stabili.

## 2 Stabilità numerica e condizionamento

### 2.1 Problema matematico ben posto e stabilità degli algoritmi

Un problema matematico (ad esempio la soluzione di  $Ax = b$  o la ricerca di una radice  $f(x)=0$ ) è ben posto se soddisfa tre condizioni (criteri di Hadamard): esistenza della soluzione, unicità della soluzione e dipendenza continua dai dati <sup>3</sup>. Se una di queste proprietà viene meno il problema è mal posto, e anche piccoli errori nei dati possono provocare grandi variazioni nella soluzione.

Un algoritmo si dice numericamente stabile se il risultato che produce può essere interpretato come la soluzione esatta del problema con dati leggermente perturbati: formalmente un algoritmo  $alg$  è backward-stable per la funzione  $f$  se  $alg(x) = f(x + \delta x)$  con  $|\delta x|$  piccolo rispetto a  $|x|$  <sup>4</sup>. Se le operazioni macchina che compongono l'algoritmo sono eseguite in modo che gli errori si comportino come piccole perturbazioni dei dati, l'algoritmo è stabile.

### 2.2 Condizionamento numerico

Il condizionamento misura quanto la soluzione di un problema ben posto varia al variare dei dati. Sia  $f: D \rightarrow \mathbb{R}^m$  una funzione e sia  $x$  l'input: l'indice di condizionamento relativo in  $x$  è definito come

$$\kappa_f(x) = \lim_{\epsilon \rightarrow 0} \sup_{\|\delta x\| \leq \epsilon \|x\|} \left[ \frac{\|f(x + \delta x) - f(x)\|}{\|f(x)\|} \right] / \left[ \frac{\|\delta x\|}{\|x\|} \right].$$

Per i sistemi lineari  $Ax = b$  con  $A$  invertibile, il numero di condizionamento della matrice in norma compatibile è  $\kappa(A) = \|A^{-1}\| \|A\|$  <sup>5</sup>. Tale valore è il rapporto tra la massima e la minima variazione della soluzione al variare dei dati: se  $\kappa(A)$  è grande il problema è mal condizionato e gli errori nei dati (o nell'aritmetica) possono amplificarsi; se è vicino a 1 il problema è ben condizionato. Nel caso di norma spettrale,  $\kappa_2(A)$  è il rapporto tra i valori singolari massimo e minimo della matrice <sup>6</sup>.

Per il problema della radice di una funzione  $f$ , il numero di condizionamento assoluto per una radice semplice  $r$  è  $1/|f'(r)|$  <sup>7</sup>; se la derivata è piccola la radice è mal condizionata. Se la radice ha molteplicità  $m > 1$  la derivata prima si annulla e il problema diventa estremamente mal condizionato <sup>8</sup>.

### 2.3 Stima dell'errore nella soluzione di sistemi lineari

Dato il sistema  $Ax = b$  ben posto e una soluzione approssimata  $\tilde{x}$ , se si perturba solo il termine noto  $b$  in  $\tilde{b} = b + \delta b$ , la teoria dell'analisi degli errori mostra che l'errore relativo sulla soluzione è limitato da

$$\|x - \tilde{x}\| / \|x\| \leq \kappa(A) (\|\delta b\| / \|b\|) + O(\epsilon),$$

cioè è proporzionale al numero di condizionamento della matrice. Se si perturba anche la matrice  $A$ , la situazione diventa più complessa: la stima contiene termini additivi che dipendono da  $\|A^{-1}\| \delta A$  e possono peggiorare drasticamente la stabilità.

## 3 Soluzione diretta di sistemi lineari

### 3.1 Sostituzione avanti e indietro

Se  $L$  è una matrice triangolare inferiore con diagonale non nulla, il sistema  $Lx = b$  può essere risolto per ricorrenza. L'algoritmo di sostituzione avanti calcola prima  $x_1 = b_1/l_{11}$ , quindi  $x_2 = (b_2 - l_{21}x_1)/l_{22}$  e così via; in generale

$$x_k = [b_k - \sum_{i=1}^{k-1} l_{ki} x_i] / l_{kk} \quad k = 1, \dots, n \quad 9$$

Per una matrice triangolare superiore  $U$  si applica la sostituzione indietro, risolvendo a partire dalla riga finale. Questi algoritmi richiedono che gli elementi diagonali non siano nulli e hanno costo  $O(n^2)$ . Sono numericamente stabili se la matrice è ben condizionata.

### 3.2 Fattorizzazione LU senza pivoting

Per una matrice quadrata  $A$  invertibile si cerca una decomposizione  $A = LU$  con  $L$  triangolare inferiore e diagonale unitaria,  $U$  triangolare superiore. L'algoritmo di eliminazione di Gauss senza pivoting funziona solo se tutti i minori principali (determinanti delle submatrici in alto a sinistra) sono diversi da zero <sup>10</sup>; in caso contrario il processo si arresta. In presenza di elementi diagonali molto piccoli la cancellazione numerica rende questa decomposizione instabile.

### 3.3 Fattorizzazione LU con pivoting parziale

La fattorizzazione con pivoting parziale esegue scambi di righe per assicurare che il pivot (l'elemento utilizzato per eliminare) sia il massimo in valore assoluto nella colonna corrente. In questo modo si ottiene una decomposizione  $PA = LU$  con matrice di permutazione  $P$ , molto più stabile numericamente: in pratica tutte le matrici quadrate ammettono questa decomposizione e il metodo è stabile tranne che per casi patologici <sup>11</sup>.

## 4 Metodi iterativi per sistemi lineari

### 4.1 Punti fissi e lemma delle contrazioni

Un punto fisso di una funzione  $g$  è un numero  $x^*$  tale che  $g(x^*) = x^*$ . Il teorema del punto fisso di Banach afferma che se  $X$  è uno spazio metrico completo e  $g$  è una contrazione (esiste un  $q < 1$  con  $d(g(x), g(y)) \leq q d(x, y)$  per ogni  $x, y \in X$ ), allora esiste un unico punto fisso e l'iterazione  $x_{k+1} = g(x_k)$  converge a tale punto per ogni scelta di  $x_0$  <sup>12</sup>. Inoltre la distanza dal punto fisso si riduce geometricamente:  $d(x_k, x^*) \leq q^k d(x_0, x^*)$ .

### 4.2 Schema generale dei metodi iterativi lineari

Si voglia risolvere il sistema lineare  $Ax = b$ . Si sceglie uno splitting  $A = Q - R$  con  $Q$  invertibile e si definisce l'iterazione

$$x^{(k+1)} = Q^{-1} R x^{(k)} + Q^{-1} b. \quad (1)$$

In forma compatta  $x^{(k+1)} = B x^{(k)} + c$ , dove  $B = Q^{-1} R$  è la matrice di iterazione. La sequenza converge alla soluzione se e solo se lo spettro di  $B$  è contenuto nel disco unitario ( $\rho(B) < 1$ ), condizione

equivalente a richiedere che  $\|B\| < 1$  per una norma indotta <sup>13</sup>. Quando ciò si verifica, la convergenza è garantita dal lemma delle contrazioni.

Vantaggi dei metodi iterativi: richiedono memoria  $O(n)$  e si adattano bene a matrici sparse; possono essere fermati in anticipo quando l'approssimazione è soddisfacente.

### 4.3 Metodo di Richardson

Il metodo di Richardson si ottiene scegliendo  $Q = (1/\alpha) I$  nello splitting  $A = Q - R$ . L'iterazione diventa

$$x^{(k+1)} = x^{(k)} + \alpha (b - A x^{(k)}), \quad (2)$$

che può essere interpretata come una discesa lungo la direzione del residuo. La convergenza dipende dallo spettro di  $A$ : la scelta ottimale di  $\alpha$  è  $2/(\lambda_{\min} + \lambda_{\max})$  se  $A$  è simmetrica definita positiva. Dal teorema della serie di Neumann si ricava che la matrice  $I - \alpha A$  è invertibile se  $|1 - \alpha \lambda_i| < 1$  per tutti gli autovalori  $\lambda_i$ , e l'inverso può essere sviluppato come serie geometrica <sup>14</sup>. Se  $A$  è mal condizionata conviene preconditionare il sistema scegliendo una matrice  $M \approx A$  facilmente invertibile e applicando il metodo a  $M^{-1} A x = M^{-1} b$  (precondizionamento di Richardson). I metodi di Jacobi e Gauss-Seidel si ottengono scegliendo  $Q$  rispettivamente come la diagonale di  $A$  o la parte triangolare inferiore.

### 4.4 Convergenza dei metodi di Jacobi e Gauss-Seidel

Nel metodo di Jacobi si pone  $Q = \text{diag}(A)$  e  $R = A - \text{diag}(A)$ . La matrice di iterazione è  $B_J = I - \text{diag}(A)^{-1} A$ . Una condizione sufficiente per la convergenza è che  $A$  sia strettamente a diagonale dominante, cioè  $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$ .  $A$  ha spettro entro il disco unitario, garantendo la convergenza del metodo. Nel metodo di Gauss-Seidel si usa  $Q$  uguale alla parte triangolare inferiore di  $A$ : la convergenza è più rapida ma richiede un ordine sequenziale.} per ogni riga. Dal lemma dei cerchi di Gershgorin ogni autovalore di  $A$  appartiene a un disco centrato nel coefficiente diagonale con raggio pari alla somma degli elementi fuori diagonale <sup>15</sup>. Se  $A$  è strettamente diagonale dominante, tutti i dischi stanno lontani da zero e  $\text{diag}(A)^{-1} A$

### 4.5 Criteri di arresto

Nel calcolo pratico non si può attendere la convergenza esatta. Un criterio comune è basato sul residuo relativo: si arresta l'iterazione quando  $\|r^{(k)}\|/\|b\| < \text{tol}$ , dove  $r^{(k)} = b - A x^{(k)}$ . Un altro criterio confronta due iterati successivi: si arresta quando  $\|x^{(k+1)} - x^{(k)}\| / \|x^{(k+1)}\| < \text{tol}$ . È importante scegliere una tolleranza compatibile con la precisione macchina e con l'errore di condizionamento.

## 5 Sistemi sovradeterminati e minimi quadrati

Un sistema lineare  $Ax = b$  si dice sovradeterminato quando  $A$  è una matrice  $m \times n$  con  $m > n$ . In generale non esiste una soluzione esatta, ma si cerca  $x$  che minimizzi l'errore in norma  $\ell^2$ :

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2.$$

Se le colonne di  $A$  sono linearmente indipendenti, esiste un'unica soluzione ai minimi quadrati  $\hat{x}$  che soddisfa le equazioni normali  $A^T A \hat{x} = A^T b$  <sup>16</sup>. Questa equazione esprime la proiezione ortogonale di  $b$  sul sottospazio  $\text{Range}(A)$ : la differenza  $b - A \hat{x}$  è ortogonale a tutte le colonne di  $A$ , ovvero  $A^T (b - A \hat{x}) = 0$ . La dimostrazione segue dalla minimizzazione del polinomio quadratico  $\|Ax - b\|^2$  rispetto a  $x$ .

## 5.1 Teorema delle proiezioni ortogonali

Sia  $V$  uno spazio vettoriale con prodotto scalare  $\langle \cdot, \cdot \rangle$  e  $S \subset V$  un sottospazio chiuso. Per ogni  $v \in V$  esiste un'unica decomposizione  $v = v_S + v_{S^\perp}$  con  $v_S \in S$  e  $v_{S^\perp} \perp S$ . Il vettore  $v_S$  è la proiezione ortogonale di  $v$  su  $S$  ed è il punto di  $S$  più vicino a  $v$  in norma  $\|\cdot\|$ . Applicando questo teorema a  $S = \text{Range}(A)$  si ricava che  $\hat{x}$  risolve le equazioni normali.

## 5.2 Metodo QR per i minimi quadrati

La risoluzione tramite le equazioni normali comporta la formazione di  $A^T A$ , che può peggiorare il condizionamento. Un'alternativa stabile è la fattorizzazione QR: si scrive  $A = Q R$  con  $Q$  matrice  $m \times n$  a colonne ortonormali e  $R$  triangolare superiore <sup>17</sup>. Allora il problema dei minimi quadrati si riduce a risolvere il sistema triangolare  $R x = Q^T b$ . La fattorizzazione può essere ottenuta tramite Householder o Gram-Schmidt modificato, entrambi algoritmi stabili.

# 6 Ricerca delle radici di equazioni non lineari

## 6.1 Esistenza e unicità

Sia  $f: [a, b] \rightarrow \mathbb{R}$  continua. Se  $f(a)$  e  $f(b)$  hanno segni opposti, il teorema degli zeri garantisce l'esistenza di almeno una radice in  $(a, b)$ . In presenza di monotonia stretta, la radice è unica.

## 6.2 Metodo di bisezione

L'algoritmo di bisezione parte da un intervallo  $[a, b]$  con  $f(a)f(b) < 0$  e ad ogni passo dimezza l'intervallo: si calcola  $c = (a + b)/2$ ; se  $f(a)f(c) \leq 0$  si pone  $b \leftarrow c$ , altrimenti  $a \leftarrow c$ . L'intervallo che contiene il segno opposto viene mantenuto e la lunghezza dell'intervallo si dimezza ad ogni iterazione <sup>18</sup>. Il metodo è globalmente convergente e garantisce la conservazione della radice; tuttavia la convergenza è lenta (lineare).

## 6.3 Condizionamento e molteplicità delle radici

Una radice  $r$  è semplice se  $f(r)=0$  e  $f'(r) \neq 0$ ; in tal caso il numero di condizionamento assoluto del problema è  $1/|f'(r)|$  <sup>7</sup>. Se  $f'(r)=0$  (radice di molteplicità  $m > 1$ ), il problema è mal condizionato e la derivata dell'iterazione di Newton sarà prossima a 1, causando convergenza lenta <sup>8</sup>.

## 6.4 Metodo di Newton

Dato un'approssimazione iniziale  $x_0$ , si definisce l'iterazione di Newton

$$x_{k+1} = x_k - f(x_k)/f'(x_k).$$

L'idea euristica è approssimare  $f$  con la tangente nel punto  $(x_k, f(x_k))$ ; l'intersezione di tale retta con l'asse delle ascisse fornisce il nuovo iterato. Se  $f$  è sufficientemente regolare in un intorno della radice  $r$  e  $f'(r) \neq 0$ , si dimostra che esiste un intorno  $I$  di  $r$  tale che, per ogni  $x_0 \in I$ , la sequenza  $x_k$  converge a  $r$  e l'errore si comporta come  $|x_{k+1} - r| \approx C|x_k - r|^2$  (convergenza quadratica) <sup>19</sup>. Nel caso di radice multipla la convergenza è solo lineare e può essere migliorata modificando l'iterazione (ad esempio moltiplicando per  $m$ ).

I criteri di arresto per Newton si basano sul confronto  $|x_{k+1} - x_k|$  o sull'ampiezza del residuo  $|f(x_k)|$ . In pratica si impone un numero massimo di iterazioni e si arresta quando la variazione è inferiore ad una certa tolleranza.

## 6.5 Newton come metodo di punto fisso

La mappa di iterazione di Newton può essere vista come una funzione  $g(x) = x - f(x)/f'(x)$ . In un intorno della radice semplice  $r$  vale  $g'(r)=0$ , quindi la mappa è una contrazione e la convergenza segue dal lemma di Banach (con velocità quadratica). Nel caso di radice multipla,  $g'(r)=1-1/m$  e la contrazione è più debole, portando ad una convergenza lineare.

## 7 Interpolazione polinomiale

### 7.1 Problema generale e matrice di Vandermonde

Dato un insieme di nodi distinti  $x_0, \dots, x_n \in [a, b]$  e valori  $y_i = f(x_i)$ , si cerca un polinomio  $p \in \Pi_n$  (lo spazio dei polinomi di grado  $\leq n$ ) che soddisfi  $p(x_i) = y_i$ . Scegliendo la base monomiale  $1, x, x^2, \dots, x^n$ , si ottiene il sistema lineare  $Va = y$ , dove  $V$  è la matrice di Vandermonde con elementi  $V_{ij} = x_i^j$ . Il determinante di  $V$  vale  $\prod_{0 \leq i < j \leq n} (x_j - x_i)$ ; esso è non nullo se e solo se i nodi sono distinti <sup>20 21</sup>. In tal caso la matrice è invertibile e l'interpolante è unico.

### 7.2 Polinomi di Lagrange

Un modo più flessibile per descrivere l'interpolante è tramite i polinomi base di Lagrange. Per ogni nodo  $x_j$  si definisce

$$l_j(x) = \prod_{m \neq j} (x - x_m) / (x_j - x_m),$$

che soddisfi  $l_j(x_j) = 1$  e  $l_j(x_m) = 0$  per  $m \neq j$  <sup>22</sup>. Il polinomio interpolante è allora

$$p(x) = \sum_{j=0}^n y_j l_j(x).$$

Questa rappresentazione è particolarmente utile per valutare l'interpolante e per dedurre formule di quadratura.

### 7.3 Errore di interpolazione

Sotto le ipotesi che  $f \in C^{n+1}[a, b]$ , l'errore tra la funzione e il polinomio interpolante ammette la formula

$$f(x) - p(x) = f^{(n+1)}(\xi(x)) / (n+1)! \times \prod_{i=0}^n (x - x_i), \quad \xi(x) \in (a, b) \quad 23.$$

Questa espressione mostra che l'errore dipende dal valore della  $(n+1)$ -esima derivata e dalla distanza dell'ascissa  $x$  dai nodi di interpolazione.

### 7.4 Stabilità e costante di Lebesgue

L'interpolazione è un operatore lineare  $X_n(f)$  che associa a  $f$  il polinomio interpolante. La costante di Lebesgue è la norma operatore  $\Lambda_n(T) = \sup_{\|f\|_\infty=1} \|X_n(f)\|_\infty$  e misura l'amplificazione degli errori nei dati. La stima d'errore

$$\|f - X_n(f)\|_\infty \leq (\Lambda_n(T) + 1) \inf_{p \in \Pi_n} \|f - p\|_\infty$$

mostra che una costante di Lebesgue grande può rendere l'interpolazione numericamente instabile <sup>24</sup>. Per nodi equispaziati  $\Lambda_n$  cresce esponenzialmente con  $n$ , causando il fenomeno di Runge. Invece, per i nodi di Chebyshev la crescita è solo logaritmica e  $\Lambda_n \approx (2/\pi) \log(n+1)$  <sup>25</sup>. Pertanto i nodi di Chebyshev, ottenuti proiettando punti equispaziati sulla semicirconferenza, sono preferibili per ridurre l'amplificazione degli errori.

## 7.5 Teoremi di Wierstrass e di Jackson

Il teorema di Weierstrass garantisce che ogni funzione continua su  $[a,b]$  può essere approssimata uniformemente da polinomi: per ogni  $\varepsilon > 0$  esiste un polinomio  $P$  tale che  $|f(x) - P(x)| < \varepsilon$  per tutti  $x$  <sup>26</sup>. Il teorema di Jackson fornisce stime quantitative: se  $f$  ha derivata  $r$ -esima limitata e continua, l'errore di approssimazione migliore  $E_n(f)$  è limitato da una costante moltiplicata per  $(n+1)^{-r}$  <sup>27</sup> <sup>28</sup>. Questi risultati motivano l'uso di polinomi di grado crescente per approssimare funzioni lisce.

## 7.6 Nodi di Chebyshev e polinomi di Chebyshev

I nodi di Chebyshev per  $n+1$  punti sono  $x_k = \cos[(2k+1)\pi / 2(n+1)]$ ,  $k=0, \dots, n$ . Tali punti minimizzano la norma  $\ell^\infty$  dell'interpolante e attenuano le oscillazioni. Esistono anche i nodi di Chebyshev-Lobatto (estremi compresi). I polinomi di Chebyshev del primo tipo  $T_n(x)$  sono definiti dalla relazione trigonometrica  $T_n(\cos \theta) = \cos(n \theta)$  <sup>29</sup>; soddisfano la ricorrenza  $T_0=1$ ,  $T_1=x$  e  $T_{n+1}(x) = 2x T_n(x) - T_{n-1}(x)$  <sup>30</sup>. I polinomi di Chebyshev del secondo tipo  $U_n$  sono definiti da  $U_n(\cos \theta) \sin \theta = \sin((n+1)\theta)$  <sup>31</sup>. Entrambe le famiglie sono ortogonali su  $[-1,1]$  con pesi rispettivamente  $1/\sqrt{1-x^2}$  e  $\sqrt{1-x^2}$ :

$$\int_{-1}^1 T_n(x) T_m(x) (1/\sqrt{1-x^2}) dx = 0 \text{ se } n \neq m, \pi \text{ se } n = m = 0, \pi/2 \text{ se } n = m \neq 0 \text{ }^{32}.$$

L'ortogonalità consente di rappresentare funzioni tramite serie di Chebyshev con buone proprietà numeriche. Inoltre, normalizzando i polinomi opportunamente, si ottiene la base di minmax: tra tutti i polinomi monici di grado  $n$ ,  $2^{1-n} T_n$  è quello con massimo valore assoluto minimo su  $[-1,1]$  <sup>33</sup>.

## 7.7 Calcolo della base e algoritmi efficienti

Per valutare i polinomi di Lagrange in maniera stabile si usa spesso la formula baricentrica, che riscrive l'interpolante come

$$p(x) = \left[ \sum_{j=0}^n \left\{ \frac{w_j}{(x-x_j)} \right\} y_j \right] / \left[ \sum_{j=0}^n \left\{ \frac{w_j}{(x-x_j)} \right\} \right],^{34}$$

con pesi  $w_j$  precomputati. Tale espressione evita la cancellazione numerica e permette di modificare facilmente i nodi. I polinomi di Chebyshev si possono invece generare con la ricorrenza precedente e valutare con l'algoritmo di Clenshaw.

# 8 Prodotti scalari, Gramiane e basi ortogonali

## 8.1 Prodotti scalari e matrici Gramiane

Un prodotto scalare in uno spazio vettoriale  $V$  è un'applicazione bilineare e simmetrica  $\langle \cdot, \cdot \rangle$  che soddisfa positività definita. Dati vettori  $v_1, \dots, v_n$  in uno spazio con prodotto scalare, la matrice Gramiana  $G$  ha elementi  $G_{ij} = \langle v_i, v_j \rangle$ . Essa è una matrice hermitiana e positiva semidefinita: per ogni vettore  $x$  si ha

$$x^\dagger G x = \|\sum_i x_i v_i\|^2 \geq 0 \quad 34.$$

Il determinante di Gram è nullo se e solo se i vettori sono linearmente dipendenti; quando sono indipendenti la matrice è positiva definita <sup>34</sup>. Ogni matrice positiva semidefinita può essere scritta come  $B^\dagger B$  (decomposizione di Cholesky) <sup>35</sup>.

## 8.2 Basi ortogonali e identità di Parseval

Dato un insieme di vettori ortonormali  $\{e_i\}$ , ogni vettore  $v$  può essere scritto come combinazione  $v = \sum \langle v, e_i \rangle e_i$ . L'identità di Parseval generalizza il teorema di Pitagora: in uno spazio di Hilbert la somma dei quadrati dei coefficienti di una serie di Fourier è uguale all'integrale del quadrato della funzione <sup>36</sup>; per basi ortonormali finite implica  $\|v\|^2 = \sum |\langle v, e_i \rangle|^2$ . La costruzione di basi ortogonali tramite il processo di Gram-Schmidt permette di ortonormalizzare insiemi di vettori mantenendo lo span invariato.

## 8.3 Nuclei di riproduzione e proiezioni ortogonali

Nei problemi di approssimazione su spazi di Hilbert, un nucleo di riproduzione  $K$  è una funzione simmetrica tale che  $f(x) = \langle f, K(\cdot, x) \rangle$ . Le proiezioni ortogonali su sottospazi generati da tali funzioni permettono di costruire approssimazioni minimali. La stima dell'errore in norma suprema può essere legata alla costante di Lebesgue del sistema di funzioni utilizzato.

# 9 Quadratura numerica

## 9.1 Idea generale e formule interpolatorie

L'obiettivo della quadratura numerica è approssimare l'integrale  $I = \int_a^b f(x) dx$  con una combinazione lineare di valori di  $f$  in punti scelti:

$$I \approx \sum_{j=0}^n w_j f(x_j).$$

Le formule di Newton-Cotes sono derivate imponendo l'esattezza per i polinomi di grado  $n$ : si utilizza l'interpolante di Lagrange sui nodi  $x_0, \dots, x_n$  e si integra esattamente il polinomio. I pesi  $w_j$  derivano dall'integrazione dei polinomi di base. Il grado di esattezza di una formula è il massimo grado di polinomio per il quale la formula è esatta.

## 9.2 Formula del trapezio

Per  $n=1$  (due nodi  $a$  e  $b$ ) l'interpolante è una retta; integrando si ottiene

$$\int_a^b f(x) dx \approx (b-a)/2 [f(a)+f(b)] \quad 37.$$

Per una partizione  $a=x_0 < \dots < x_N=b$  si ottiene la formula composta

$$\int_a^b f(x) dx \approx \sum_{k=0}^{N-1} [f(x_{k-1})+f(x_k)]/2 \times (x_k - x_{k-1}) \quad 38.$$

L'errore della formula composta vale

$$E = - (b-a)^3/(12 N^2) f''(\xi) \quad 39,$$



per qualche  $\xi \in (a, b)$ . La presenza del secondo derivato mostra che il metodo è di ordine 2; per aumentare la precisione si suddivide l'intervallo (formula composta).

### 9.3 Formula della parabola (Simpson)

Per  $n=2$  (tre nodi equispaziati) l'interpolante è un polinomio quadratico e la formula di Simpson è

$$\int_a^b f(x) dx \approx (b-a)/6 [ f(a) + 4 f((a+b)/2) + f(b) ] \quad 40.$$

Applicando la formula su subintervalli di ampiezza  $h=(b-a)/n$  con  $n$  pari (formula composta), l'ordine di convergenza diventa 4. L'errore per un singolo intervallo di ampiezza  $h$  vale

$$E = -1/90 h^5 f^{(4)}(\xi) = - (b-a)^5/2880 f^{(4)}(\xi) \quad 41,$$

dove  $\xi \in (a, b)$ . Il quarto derivato mostra che Simpson è più preciso del trapezio.

### 9.4 Stabilità e pesi positivi

L'approssimazione numerica è stabile se piccoli errori nelle valutazioni di  $f$  non si amplificano nella somma pesata. Un requisito sufficiente è che i pesi  $w_j$  siano non negativi: sommando numeri positivi l'errore relativo non si amplifica per cancellazioni. Le formule di Newton-Cotes di ordine alto hanno pesi alternati in segno e possono essere instabili; per questo si preferiscono formule composte (trapezi e Simpson) o formule di quadratura di Gauss con pesi positivi. L'utilizzo di nodi di Chebyshev in combinazione con le formule di quadratura (Clenshaw-Curtis) consente di integrare in modo stabile e accurato.

## 10 Conclusioni

Questa Bibbia di Calcolo Numerico raccoglie definizioni, teoremi, algoritmi e dimostrazioni essenziali per prepararsi all'esame di Calcolo Numerico. Sono stati trattati i concetti fondamentali sull'aritmetica floating-point, l'analisi degli errori, la stabilità e il condizionamento, i metodi diretti e iterativi per la risoluzione di sistemi lineari, i metodi per la ricerca delle radici, l'interpolazione polinomiale, la teoria dei minimi quadrati, le basi ortogonali e la quadratura numerica. Ogni sezione è accompagnata da osservazioni sui limiti di stabilità e sui criteri di arresto. Studiare con ordine questi temi permetterà di affrontare con consapevolezza sia gli esercizi pratici sia le dimostrazioni teoriche proposte nell'esame.

---

<sup>1</sup> <sup>4</sup> CS267: Supplementary Notes on Floating Point

<https://people.eecs.berkeley.edu/~demmel/cs267-1995/lecture21/lecture21.html>

<sup>2</sup> Absolute and Relative Error | Calculus II

<https://courses.lumenlearning.com/calculus2/chapter/absolute-and-relative-error/>

<sup>3</sup> terminology - What is a well-posed problem? - Computational Science Stack Exchange

<https://scicomp.stackexchange.com/questions/40966/what-is-a-well-posed-problem>

<sup>5</sup> <sup>6</sup> Condition number - Wikipedia

[https://en.wikipedia.org/wiki/Condition\\_number](https://en.wikipedia.org/wiki/Condition_number)

<sup>7</sup> <sup>8</sup> The rootfinding problem — Fundamentals of Numerical Computation

<https://fncbook.github.io/v1.0/nonlineqn/rootproblem.html>

- 9 **Triangular matrix - Wikipedia**  
[https://en.wikipedia.org/wiki/Triangular\\_matrix](https://en.wikipedia.org/wiki/Triangular_matrix)
- 10 11 **LU decomposition - Wikipedia**  
[https://en.wikipedia.org/wiki/LU\\_decomposition](https://en.wikipedia.org/wiki/LU_decomposition)
- 12 **Banach fixed-point theorem - Wikipedia**  
[https://en.wikipedia.org/wiki/Banach\\_fixed-point\\_theorem](https://en.wikipedia.org/wiki/Banach_fixed-point_theorem)
- 13 14 **www.math.csi.cuny.edu**  
<https://www.math.csi.cuny.edu/~verzani/Courses/Old/F2015/MTH335/Notes/Chapter4/series-iteration.html>
- 15 **Gershgorin circle theorem in nLab**  
<https://ncatlab.org/nlab/show/Gershgorin%20circle%20theorem>
- 16 **Orthogonal least squares**  
<https://understandinglinearalgebra.org/sec-least-squares.html>
- 17 **QR decomposition - Wikipedia**  
[https://en.wikipedia.org/wiki/QR\\_decomposition](https://en.wikipedia.org/wiki/QR_decomposition)
- 18 **Bisection method - Wikipedia**  
[https://en.wikipedia.org/wiki/Bisection\\_method](https://en.wikipedia.org/wiki/Bisection_method)
- 19 **Newton's Method - Department of Mathematics at UTSA**  
<https://mathresearch.utsa.edu/wiki/index.php>
- 20 21 **Vandermonde matrix - Wikipedia**  
[https://en.wikipedia.org/wiki/Vandermonde\\_matrix](https://en.wikipedia.org/wiki/Vandermonde_matrix)
- 22 **Lagrange polynomial - Wikipedia**  
[https://en.wikipedia.org/wiki/Lagrange\\_polynomial](https://en.wikipedia.org/wiki/Lagrange_polynomial)
- 23 **Interpolation - The Beginning**  
<https://www.cs.odu.edu/~tkennedy/cs417/f24/Public/interpolationBeginning/index.html>
- 24 25 **Lebesgue constant - Wikipedia**  
[https://en.wikipedia.org/wiki/Lebesgue\\_constant](https://en.wikipedia.org/wiki/Lebesgue_constant)
- 26 **Weierstrass Approximation Theorem -- from Wolfram MathWorld**  
<https://mathworld.wolfram.com/WeierstrassApproximationTheorem.html>
- 27 28 **Jackson's Theorem -- from Wolfram MathWorld**  
<https://mathworld.wolfram.com/JacksonsTheorem.html>
- 29 30 31 32 33 **Chebyshev polynomials - Wikipedia**  
[https://en.wikipedia.org/wiki/Chebyshev\\_polynomials](https://en.wikipedia.org/wiki/Chebyshev_polynomials)
- 34 35 **Gram matrix - Wikipedia**  
[https://en.wikipedia.org/wiki/Gram\\_matrix](https://en.wikipedia.org/wiki/Gram_matrix)
- 36 **Parseval's theorem - Wikipedia**  
[https://en.wikipedia.org/wiki/Parseval's\\_theorem](https://en.wikipedia.org/wiki/Parseval's_theorem)
- 37 38 39 **Trapezoidal rule - Wikipedia**  
[https://en.wikipedia.org/wiki/Trapezoidal\\_rule](https://en.wikipedia.org/wiki/Trapezoidal_rule)
- 40 41 **Simpson's rule - Wikipedia**  
[https://en.wikipedia.org/wiki/Simpson's\\_rule](https://en.wikipedia.org/wiki/Simpson's_rule)