

METODI SPETTRALI PER L'EQ. DEL CALORE

La principale caratteristica dei metodi spettrali, che li rende diversi dagli elementi finiti, è che essi utilizzano polinomi globali sul dominio computazionale anziché polinomi a tratti. In particolare grazie alla scelta di polinomi opportuni ed alla regolarità della funzione incognita possono raggiungere velocità di convergenza superiori a quelle degli usuali approcci agli elementi finiti e di differenze finite.

Allo scopo di illustrare in modo diretto i metodi spettrali consideriamo l'eq. del calore in una variabile di spazio

$$u_t = u_{xx}$$

$$u(x, 0) = u_0(x)$$

e condizioni al bordo periodiche, ossia

$$u(0, t) = u(L, t).$$

Inoltre ha u_0 t.c. $u_0(0) = u_0(L)$.

Possiamo rappresentare la soluzione in serie di Fourier

$$(1) \quad \begin{cases} \hat{u}_k(t) = \frac{1}{2\pi} \int_0^{2\pi} u(x,t) e^{-ikx} dx \\ u(x,t) = \sum_{k=-\infty}^{\infty} \hat{u}_k(t) e^{ikx} \end{cases}$$

dove abbiamo scelto l'intervallo $[0, 2\pi]$ per la rappresentazione ed u periodica sullo stesso intervallo.

Analogamente, basta riscrivere (1) ottenendo la rappresentazione sull'intervallo $[0, L]$ dove al posto di x nei termini esponenziali compare il fattore $\frac{2\pi}{L}x$.

Utilizzando (1) ed il fatto che

$$u_{xx} = \sum_{k=-\infty}^{\infty} \hat{u}_k(t) \cdot \frac{d^2}{dx^2} e^{ikx}$$

$$\frac{d^2}{dx^2} e^{ikx} = (ik)^2 e^{ikx} = -k^2 e^{ikx}$$

otteniamo

$$\sum_{k=-\infty}^{\infty} \frac{d\hat{u}_k}{dt} e^{ikx} = - \sum_{k=-\infty}^{\infty} \hat{u}_k k^2 e^{ikx}$$

Identificando termine a termine delle espressioni appena trovate si ha

$$(*) \begin{cases} \frac{d\hat{u}_k}{dt} = -k^2 \hat{u}_k, & -\infty < k < \infty \\ \hat{u}_k(0) = \hat{u}_{0k} \end{cases}$$

La soluzione si calcola facilmente ed è data da

$$\hat{u}_k(t) = \hat{u}_{0k} e^{-k^2 t}, \quad \forall k$$

che sostituita nell'espressione originale delle serie di Fourier porta alla nota soluzione analitica

$$(2) \quad u(x, t) = \sum_{k=-\infty}^{\infty} \hat{u}_{0k} e^{-k^2 t} e^{ikx}.$$

Numericamente non sempre è possibile utilizzare direttamente l'espressione (2) in quanto bisogna avere i valori di tutte le serie. In generale dunque si approssima (2) con la serie troncata ad N termini:

$$P_N u = \sum_{k=-N/2}^{N/2-1} \hat{u}_{0k} e^{-k^2 t} e^{ikx}.$$

Si noti che l'operatore P_N che opera il troncamento è un'operazione di proiezione dallo spazio infinito dimensionale in cui si trova la funzione u ad uno spazio di dimensione finita generato dalle basi dei

polinomi trigonometrici

$$\left\{ e^{ikx} \right\}, \quad k = -\frac{N}{2}, \dots, \frac{N}{2} - 1.$$

In generale data data una funzione u periodica
su $[0, 2\pi]$ e quadrato integrabile sullo stesso intervallo
ovvero $u \in L^2([0, 2\pi])$ allora indichiamo con

$$u_N = P_N u$$

la sua rappresentazione nello spazio dei polinomi generati

da $\left\{ \phi_k \right\}, \quad k = -\frac{N}{2}, \dots, \frac{N}{2} - 1,$

con $\phi_k = e^{ikx}$. Indichiamo tale spazio con \mathbb{P}^N .

Dunque $P_N : L^2([0, 2\pi]) \rightarrow \mathbb{P}^N$. Introduciamo

inoltre il prodotto interno in $L^2([0, 2\pi])$ definito da

$$\langle u, v \rangle = \frac{1}{2\pi} \int_0^{2\pi} u(x) \overline{v(x)} dx.$$

Abbiamo quindi che (4) è equivalente a

$$\hat{u}_k(t) = \langle u, \phi_k \rangle$$

ed il problema (*) risulta

$$\langle u_t, \phi_k \rangle = \langle u_{xx}, \phi_k \rangle$$

analogo alla forma debole usata per gli elementi finiti.

Il metodo numerico usato è detto di Fourier-Galerkin

e può essere risolto tutti l'operatore di proiezione

come

$$\begin{cases} \frac{\partial u_N}{\partial t} = \frac{\partial^2 u_N}{\partial x^2} \\ u_N(x, 0) = u_{0N}(x) \end{cases}$$

da cui eseguendo il prodotto interno si ottiene

$$\begin{cases} \frac{d(\hat{u}_N)_k}{dt} = -k^2 (\hat{u}_N)_k & -N/2 \leq k \leq N/2 - 1 \\ (\hat{u}_N)_k(0) = (\hat{u}_{0N})_k \end{cases}$$

che risolto esattamente fornisce

$$(3) \quad u_N(x, t) = \sum_{k=-N/2}^{N/2-1} (\hat{u}_{0N})_k e^{-k^2 t} e^{ikx}$$

che coincide esattamente con il troncamento ai primi N termini della soluzione esatta (2) del problema.

Questo risultato è vero per tutte le equazioni lineari a coefficienti costanti.

CENNI SULLA CONVERGENZA SPETTRALE

Indichiamo ora con f una generica funzione in $L^2((0, 2\pi))$

e con $P_N f$ la sua serie di Fourier troncata.

Vogliamo dimostrare che $P_N f \rightarrow f$ per $N \rightarrow \infty$

in una norma opportuna.

Ad esempio in norma L^2 , ossia che

$$\int_0^{2\pi} |f - P_N f|^2 dx \rightarrow 0 \quad \text{per } N \rightarrow \infty.$$

Vali l'identità di Parseval

$$\int_0^{2\pi} |f|^2 dx = \sum_{k=-\infty}^{\infty} |\hat{f}_k|^2$$

Da questa possiamo dimostrare l'errore in norma L^2

$$\int_0^{2\pi} |f - P_N f|^2 dx = \int_0^{2\pi} \left| \sum_{\substack{k < -N/2 \\ k > N/2 - 1}} \hat{f}_k e^{ikx} \right|^2 dx = \sum_{\substack{k < -N/2 \\ k > N/2 - 1}} |\hat{f}_k|^2$$

Per l'errore puntuale invece

$$|f(x_0) - P_N f(x_0)| = \left| \sum_{\substack{k < -N/2 \\ k > N/2 - 1}} \hat{f}_k e^{ikx} \right| \leq \sum_{\substack{k < -N/2 \\ k > N/2 - 1}} |\hat{f}_k|$$

Dunque in entrambi i casi l'errore dipende dal comportamento dei coefficienti di Fourier al crescere di $|K|$.

Osserviamo che se la funzione è regolare allora è possibile ottenere stime molto migliori. Infatti se f è derivabile con continuità abbiamo integrando per parti

$$\begin{aligned}\hat{f}_K &= \frac{1}{2\pi} \int_0^{2\pi} f(x) e^{-iKx} dx = -\frac{1}{2\pi iK} (f(2\pi) - f(0)) + \\ &+ \frac{1}{2\pi iK} \int_0^{2\pi} f'(x) e^{-iKx} dx\end{aligned}$$

Passando ai valori assoluti e ricordando che f è periodica ha

$$|\hat{f}_K| \leq \frac{1}{2\pi |K|} \int_0^{2\pi} |f'(x)| dx = O\left(\frac{1}{|K|}\right)$$

ovvero i coefficienti decadono linearmente rispetto a $|K|$.

Se ora ipotizziamo f derivabile con continuità p volte, iterando il procedimento otteniamo

$$|\hat{f}_K| \leq \frac{1}{2\pi |K|^p} \int_0^{2\pi} |f^{(p)}(x)| dx = O\left(\frac{1}{|K|^p}\right)$$

Dato x è funzione e' infinitamente derivabile
 allora i coefficienti di Fourier decadono più rapidamente
 di qualche potenza di $|K|$. Questo comporta che
 l'errore commesso dalla serie di Fourier troncata decada
 più rapidamente di qualche potenza di N . Si parla
 in questo caso di convergenza esponenziale o di
accuratezza spettrale.

OSS: In pratica i coefficienti di Fourier sono
 calcolati tramite la trasformata di Fourier discreta
 (DFT) che corrisponde a

$$\tilde{u}_K = \frac{1}{N} \sum_{j=0}^N u(x_j) e^{-iK \frac{2\pi}{N} j}$$

con $x_j = \frac{2\pi}{N} j$.

Tale trasformata per calcolare N coefficienti \tilde{u}_K
 comporterebbe il costo proibitivo di N^2 . Tramite
 l'algoritmo delle Fast Fourier Transform (FFT) il
 costo può essere ridotto a $N \log_2 N$.

ESTENSIONE A PROBLEMI NON LINEARI

Consideriamo l'equazione di Burgers viscosa

$$(1) \begin{cases} u_t + \left(\frac{u^2}{2}\right)_x = u_{xx} \\ u(0,t) = u(2\pi,t) \\ u(x,0) = u_0(x) \end{cases}$$

Ripercorriamo a particolare i passi che servono per la costruzione di un metodo di tipo Fourier-Galerkin.

Per i coefficienti di Fourier osserviamo innanzitutto che se indichiamo con

$$h = \left(\frac{u^2}{2}\right)_x$$

abbiamo

$$\begin{aligned} h = \left(\frac{u^2}{2}\right)_x &= \frac{d}{dx} \left[\frac{1}{2} \sum_{l=-\infty}^{\infty} \hat{u}_l e^{ilx} \cdot \sum_{m=-\infty}^{\infty} \hat{u}_m e^{imx} \right] \\ &= \frac{d}{dx} \left[\frac{1}{2} \sum_{\substack{k=-\infty \\ k=l+m}}^{\infty} \sum_{m=-\infty}^{\infty} \hat{u}_{k-m} \hat{u}_m e^{ikx} \right] \end{aligned}$$

$$= \frac{1}{2} \sum_{k=-\infty}^{\infty} K \sum_{m=-\infty}^{\infty} \hat{u}_{k-m} \hat{u}_m e^{ikx}$$

$$= \sum_{k=-\infty}^{\infty} \hat{h}_k e^{ikx}$$

$$\text{con } \hat{h}_k = \frac{1}{2} K \sum_{m=-\infty}^{\infty} \hat{u}_{k-m} \hat{u}_m$$

ovvia e' rappresentata da una convoluzione discreta.

In Fourier abbiamo dunque

$$(2) \begin{cases} \frac{d \hat{u}_k}{dt} + \hat{h}_k = -k^2 \hat{u}_k & -\infty < k < \infty \\ \hat{u}_k(0) = \hat{u}_{0k} \end{cases}$$

Al solito per ottenere un metodo numerico dobbiamo considerare le serie troncate ad un numero finito di termini.

In questo caso però non tronciamo solo la rappresentazione di u ma anche la rappresentazione di h in quanto la convoluzione discreta deve anch'essa essere protetta sullo spazio finito dimensionale.

Tali ulteriori approssimazioni, che non analizzano i dettagli, non alterano comunque le proprietà di accuratezza del metodo.

In breve il metodo di Fourier-Galerkin diventa

$$(3) \quad \begin{cases} \frac{\partial \bar{u}_N}{\partial t} + P_N \left(\frac{\partial \bar{u}_N^2}{\partial x} \right) = \frac{\partial^2 \bar{u}_N}{\partial x^2} \\ \bar{u}_N(0, t) = \bar{u}_N(2\pi, t) \\ \bar{u}_N(x, 0) = P_N u_0(x) \end{cases}$$

Si noti che dovendo per l'appunto proiettare $\frac{\partial}{\partial x} \left(\frac{u_N^2}{2} \right)$ sullo spazio dei polinomi trigonometrici di grado N tramite P_N la soluzione $\bar{u}_N \neq u_N = P_N u$, ossia non coincide con la serie troncata delle sol. esatte. Il metodo viene quindi

$$\frac{d(\hat{\bar{u}}_N)_k}{dt} + (\hat{\bar{h}}_N)_k = -k^2 (\hat{\bar{u}}_N)_k, \quad -\frac{N}{2} \leq k \leq \frac{N}{2} - 1$$

$$\text{con } (\hat{\bar{h}}_N)_k = \frac{1}{2} k \sum_{m=-\frac{N}{2}}^{\frac{N}{2}-1} \hat{u}_{k-m} \hat{u}_m \quad (*)$$

dove si è supposto $\hat{u}_e = 0$ per $e > \frac{N}{2} - 1$, $e < -\frac{N}{2}$.

Si noti che $(\hat{h}_n)_n$ è una sommatoria di
convoluzione che può essere risolta efficientemente
tramite l'uso algoritmico della FFT.

In sostanza l'uso della FFT in questo caso
equivale ad operare nello spazio fisico con il calcolo
di $(\hat{h}_n)_n$ originando un metodo pseudo-spettrale.

Osserviamo inoltre che se avessimo calcolato

$$h = u \cdot u_x = \sum_{l=-\infty}^{\infty} \hat{u}_l e^{i(l+n)x} \sum_{m=-\infty}^{\infty} m \hat{u}_m e^{i(l+m)x}$$

$$= \sum_{k=-\infty}^{\infty} \hat{u}_{k-m} m \hat{u}_m e^{ikx} = \sum_{k=-\infty}^{\infty} \hat{h}_k e^{ikx}$$

$$\text{con } \hat{h}_k = \sum_{m=-\infty}^{\infty} \hat{u}_{k-m} m \hat{u}_m \quad \text{diverso da } \hat{h}_k.$$

Tale scelta origina comunque un metodo equivalente sia
in termini di costo computazionale che di accuratezza,

CENNI SUI METODI AI VOLUMI FINITI

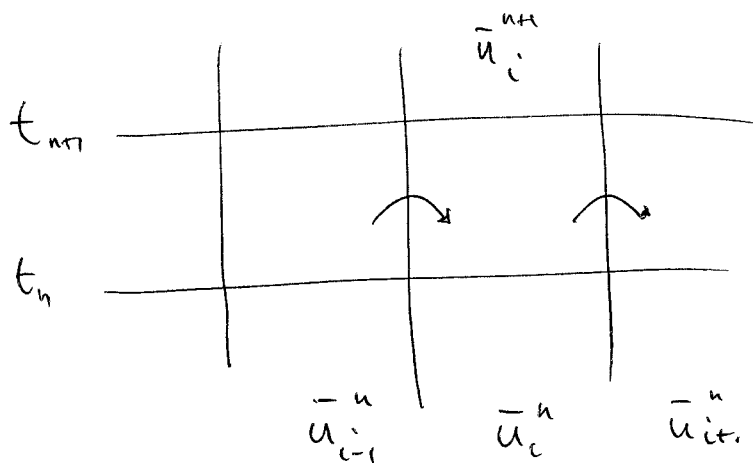
Definiamo innanzitutto il concetto di "media di cella".

$$C_i = [x_{i-1/2}, x_{i+1/2}] \quad \text{cella } i$$

$$\bar{u}_i^n \approx \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} u(x, t_n) dx = \frac{1}{\Delta x} \int_{C_i} u(x, t_n) dx \quad (\text{media di cella})$$

Consideriamo ora l'equazione

$$u_t + f(u)_x = 0 \quad (\text{legge cons. scalare})$$



Integrando la legge di conservazione sulla cella C_i si ha

$$\frac{d}{dt} \int_{C_i} u(x, t) dx = f(u(x_{i-1/2}, t)) - f(u(x_{i+1/2}, t))$$

Vogliamo approssimare le medie di cella al tempo $n+1$ in

funzione delle medie di cella al tempo n . Possiamo

approssimare le derivate temporali con un metodo Runge-Kutta

per cui dovremo definire il flusso numerico $f(u(x_{i+1/2}, t))$

in funzione delle medie di cella $\bar{u}_i, \bar{u}_{i+1}, \dots$

Alternativamente si integra anche la variabile temporale. 126

Integriamo il tempo da t_n a t_{n+1} e otteniamo

$$\int_{C_i}^{t_{n+1}} u(x, t_{n+1}) dx - \int_{C_i}^{t_n} u(x, t_n) dx = \int_{t_n}^{t_{n+1}} f(u(x_{i+\frac{1}{2}}, t)) dt - \int_{t_n}^{t_{n+1}} f(u(x_{i-\frac{1}{2}}, t)) dt$$

Dividendo per Δx si ha

$$\frac{1}{\Delta x} \int_{C_i}^{t_{n+1}} u(x, t_{n+1}) dx = \frac{1}{\Delta x} \int_{C_i}^{t_n} u(x, t_n) dx - \frac{1}{\Delta x} \left[\int_{t_n}^{t_{n+1}} f(u(x_{i+\frac{1}{2}}, t)) dt - \int_{t_n}^{t_{n+1}} f(u(x_{i-\frac{1}{2}}, t)) dt \right]$$

$$\boxed{\bar{u}_i^{n+1} = \bar{u}_i^n - \frac{\Delta t}{\Delta x} [F_{i+\frac{1}{2}}^n - F_{i-\frac{1}{2}}^n]}$$

$$F_{i+\frac{1}{2}}^n = \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} f(u(x_{i+\frac{1}{2}}, t)) dt$$

In generale calcoliamo $F_{i+\frac{1}{2}}^n$ a' modo approssimato, detto flusso numerico, a partire dalle medie di cella.

Dato che per un sistema iperbolico le soluzioni si propagano con velocità finite è ragionevole assumere che

$$F_{i+1/2}^n = F(\bar{u}_i^n, \bar{u}_{i+1}^n)$$

dove la funzione F caratterizza il flusso numerico.

La forma di F caratterizza il metodo numerico, in questo caso basato su uno "stencil" a tre punti.

Osserviamo che

$$\Delta x \sum_{i=1}^N \bar{u}_i^{n+1} = \Delta x \sum_{i=1}^N \bar{u}_i^n - \frac{\Delta t}{\Delta x} (F_{N+1/2}^n - F_{1/2}^n)$$

la somma delle differenze dei flussi si cancella ad eccezione dei valori al bordo, dove dovranno applicare opportuni confini.

Il metodo può inoltre essere rappresentato come metodo alle differenze finite

$$\frac{\bar{u}_i^{n+1} - \bar{u}_i^n}{\Delta t} + \frac{F_{i+1/2}^n - F_{i-1/2}^n}{\Delta x} = 0$$

ESEMPIO: Otteniamo dunque i corrispondenti flussi utilizzando i metodi qui visti con operatori $L \times F$, $L \times W$
(esercizio scrivere per l'eq. $u_t + au_x = 0$)