

CONSIGLIO NAZIONALE DELLE RICERCHE

QUADERNI DEL GRUPPO NAZIONALE PER L'INFORMATICA MATEMATICA

CLAUDIO ESTATICO

GRADIENTE CONIUGATO

E

REGOLARIZZAZIONE DI PROBLEMI MAL POSTI

FIRENZE , 1996

Indice

Introduzione	1
1 Inversa generalizzata di operatori lineari	5
1.1 Introduzione	5
1.2 Pseudosoluzione e soluzione generalizzata	6
1.3 Definizione operatoriale di inversa generalizzata	8
1.4 Proprietà della inversa generalizzata	14
1.5 Inversa generalizzata di operatori compatti e matrici	17
2 Metodi di regolarizzazione per problemi mal posti	21
2.1 Introduzione	21
2.2 Problemi mal posti e problemi inversi	21
2.3 Metodi di regolarizzazione	24
2.4 Regolarizzazione di Tikhonov	27
2.5 Soluzioni vincolate, approssimate regolarizzate e metodo della discrepanza	30
2.6 Finestre spettrali e regolarizzazione di ordine superiore	39
2.7 Metodi iterativi di regolarizzazione	45
2.8 Trattamento numerico dei metodi di regolarizzazione	50
3 L'algoritmo regolarizzante del gradiente coniugato	63
3.1 Introduzione	63
3.2 Il metodo del gradiente (o della massima discesa)	64
3.3 Il metodo del gradiente coniugato	69
3.4 Regolarizzazione mediante gradiente coniugato	92
A	117
A.1 Sistema singolare di operatori compatti e teorema di Picard	117
B	119
B.1 Risoluzione spettrale di operatori autoaggiunti	119
C	131
C.1 Differenziale di Fréchet	131
Bibliografia	135

Introduzione

La ricerca matematica è spinta in primo luogo dalla necessità di risolvere problemi concreti. Generalmente l'utilizzo di nuovi strumenti matematici non viene immediatamente supportato da un preciso studio teorico, ma spesso l'applicazione di nuove metodologie precede di parecchi anni la loro formulazione teorica completa.

In questo lavoro viene analizzato l'algoritmo del gradiente coniugato per il trattamento di problemi mal posti, che ha avuto uno sviluppo di questo tipo.

In prima analisi un problema matematico viene detto mal posto quando non dipende con continuità dai dati, ossia quando piccole perturbazioni del dato in ingresso danno luogo a soluzioni assai diverse.

Appare evidente che un problema di questo tipo comporta nella pratica grandi difficoltà di gestione, poichè in natura i dati a disposizione sono sempre conosciuti in maniera approssimata entro tolleranze finite.

Il trattamento di problemi mal posti si realizza per mezzo di metodologie che hanno lo scopo di rendere il problema poco sensibile a perturbazioni sui dati. Queste metodologie vengono dette "regolarizzanti".

Scopo del lavoro risulta proprio l'analisi di recenti studi che garantiscono come l'algoritmo del gradiente coniugato sia un metodo di regolarizzazione. In accordo con quanto detto, si sottolinea che il gradiente coniugato veniva già utilizzato numericamente come metodo per il trattamento di problemi mal posti, sebbene non fosse stata ancora dimostrata la sua appartenenza alla classe degli algoritmi regolarizzanti.

La monografia è articolata in modo da proporre i concetti fondamentali riguardanti la regolarizzazione e quelli più approfonditi sul gradiente coniugato.

Il contesto in cui si sviluppa lo studio è quello degli spazi di Hilbert. Più precisamente si considera una generica equazione lineare

$$Tx = y$$

dove $T : X \rightarrow Y$ è un operatore lineare limitato tra gli spazi X, Y di Hilbert, $y \in Y$ è il dato del problema e $x \in X$ la soluzione da determinare.

Poichè come equazioni lineari mal poste si intendono anche equazioni in cui la soluzione non esiste o non è unica, la risoluzione è ricondotta alla determinazione dell'inversa generalizzata, ossia della soluzione ai minimi quadrati di norma minima.

Il primo capitolo ha lo scopo di introdurre la nozione di inversa generalizzata. L'esposizione affianca all'usuale definizione basata sull'introduzione dell'insieme delle soluzioni ai minimi quadrati, una diversa trattazione, che viene detta "operatoriale", basata su proprietà di restrizione e estensione dell'operatore T a opportuni sottospazi.

In particolare si sostituisce l'operatore T non invertibile, con un operatore \tilde{T} invertibile tale che $\tilde{T}|_{N(T)^\perp} = T|_{N(T)^\perp}$ dove $N(T)$ rappresenta il nucleo dell'operatore T . L'invertibilità dell'operatore \tilde{T} consente di determinare l'elemento $\tilde{x} = \tilde{T}^{-1}y \in Y$, la

cui proiezione su $N(T)^\perp$ corrisponde alla soluzione ai minimi quadrati di norma minima. In questo modo si ha una visione differente del concetto di inversa generalizzata che, oltre a chiarirne il significato, rappresenta un primo approccio alle difficoltà create da operatori il cui insieme immagine, o range, non sia chiuso.

Nel secondo capitolo vengono trattati i problemi mal posti e i metodi di regolarizzazione. Le nozioni acquisite nel capitolo precedente trovano adesso una naturale collocazione poiché le equazioni lineari mal poste vengono risolte proprio con la determinazione di soluzioni che approssimano la soluzione generalizzata.

Come già accennato, il capitolo ha lo scopo di fornire una breve introduzione all'argomento, che risulta vastissimo e in continuo sviluppo.

Nella prima parte del capitolo l'analisi è indirizzata all'algoritmo regolarizzante di A. Tikhonov. Questo particolare metodo consente di introdurre in maniera semplice il cosiddetto criterio della discrepanza, di fondamentale importanza per il capitolo successivo; inoltre consente di verificare la stretta relazione tra regolarizzazione e filtraggio in frequenza.

Nella parte finale del capitolo viene brevemente esaminato l'algoritmo regolarizzante di Landweber-Fridman. Nello svilupparsi del percorso che conduce alla regolarizzazione mediante gradiente coniugato, il metodo di Landweber-Fridman assume un'importanza particolare poiché, analogamente al gradiente coniugato, è un metodo iterativo, ossia un metodo che determina in maniera sequenziale approssimazioni successive della soluzione, ognuna calcolata per mezzo delle precedenti.

Il secondo capitolo si chiude con una panoramica su alcuni metodi numerici di regolarizzazione, il cui scopo è dimostrare in che modo le proprietà individuate analiticamente si caratterizzano dal punto di vista numerico.

È importante sottolineare che, come verrà illustrato, l'equazione lineare $Tx = y$ conduce ad un problema mal posto ogniqualvolta l'operatore T non è dotato di range chiuso. Generalmente, a meno di precisazioni particolari, gli operatori considerati nel secondo capitolo soddisfano esclusivamente a questa ipotesi. Le dimostrazioni riportate fanno spesso uso della risoluzione spettrale per operatori autosaggiunti, argomento che per la sua importanza viene trattato in un'apposita appendice.

L'algoritmo del gradiente coniugato viene introdotto nel terzo capitolo.

Il gradiente coniugato è un metodo per la minimizzazione di funzioni; l'introduzione del funzionale $f : X \rightarrow \mathbb{R}$ definito come

$$f(x) = \|Tx - y\|_2^2$$

consente però di utilizzare l'algoritmo come metodo per la risoluzione ai minimi quadrati dell'equazione $Tx = y$.

La prima analisi, di tipo geometrico, viene fatta in uno spazio euclideo n -dimensionale. In questo contesto, che risulta il più semplice, si mette in luce il funzionamento dell'algoritmo. L'algoritmo viene successivamente generalizzato alla minimizzazione in spazi di Hilbert e viene confrontato con il metodo della massima discesa e con il metodo delle direzioni coniugate.

Lo studio procede con l'analisi dei teoremi di convergenza formulati negli anni 70 da Kammerer e Nashed (cfr. [21], [22]), i quali asseriscono che il metodo delle direzioni coniugate converge alla soluzione generalizzata.

Verrà innanzitutto dimostrata la convergenza del metodo nel caso in cui il funzionale sia definito per mezzo di operatori limitati invertibili e per mezzo di operatori limitati non invertibili con range chiuso. Successivamente verrà dimostrata la convergenza per operatori limitati con range non chiuso. In quest'ultimo caso le difficoltà nascono dal fatto che l'operatore T ristretto al sottospazio $N(T)^\perp$ possiede inversa non limitata e vengono superate

con l'introduzione del concetto di operatore aggiunto di operatore lineare non limitato.

La convergenza dimostrata per operatori con range non chiuso ha spinto diversi ricercatori a considerare il metodo come possibile algoritmo regolarizzante.

Nel 1986 A. Nemirovskii ha dimostrato che l'algoritmo è un metodo regolarizzante per una particolare classe di soluzioni (i.e. $x \in \cup_{\mu > 0} R(T^\mu)$). Questo risultato, sebbene di notevole rilevanza, non permetteva però di concludere che l'algoritmo fosse un metodo di regolarizzazione, ossia avesse proprietà regolarizzanti che garantissero la convergenza in presenza di soluzioni appartenenti a tutto l'insieme $N(T)^\perp$.

Nel 1990 R. Plato pubblica un articolo in cui, in maniera indiretta, dimostra che l'algoritmo del gradiente coniugato è effettivamente un algoritmo di regolarizzazione e che quindi consente la risoluzione di problemi mal posti in presenza di qualsiasi soluzione appartenente a $N(T)^\perp$ (cfr. [38]).

La parte principale di questa monografia è l'analisi dei risultati di Nemirovskii e Plato.

La dimostrazione del teorema di Nemirovskii seguirà un percorso proposto nel 1994 da M. Hanke (cfr. [17]). Questa dimostrazione fa uso di strumenti di analisi funzionale che possono essere ricondotti in primo luogo alla risoluzione spettrale degli operatori autosaggiunti T^*T e TT^* . Il teorema di Plato, che come già detto risolve il problema dell'appartenenza del gradiente coniugato alla classe degli algoritmi regolarizzanti, viene invece illustrato nella formulazione originale dell'autore.

Una caratteristica di questo lavoro è di fornire un'analisi su "regolarizzazione e gradiente coniugato" breve ma sufficientemente completa. In quest'ottica parecchi teoremi minori, che nelle varie pubblicazioni vengono solitamente solo citati, sono stati verificati e dimostrati, mentre altri sono stati enunciati e dimostrati con formalismi differenti da quelli adottati dai rispettivi autori. Utilizzando le stesse definizioni nei diversi capitoli, l'esposizione è risultata così più omogenea.

Come già accennato sono state sviluppate inoltre appendici matematiche, le quali forniscono le necessarie nozioni di base sugli strumenti di analisi funzionale più sofisticati (in particolare la risoluzione spettrale di operatori autosaggiunti non necessariamente compatti e la derivata di Fréchet di funzionali in spazi di Hilbert).

È utile precisare che l'algoritmo del gradiente coniugato è oggi applicato alla elaborazione di immagini. In questo campo il metodo, insieme a tecniche numeriche di preconditionamento, ha condotto a buoni risultati (cfr. [7], [33]).

Un doveroso ringraziamento va rivolto alla Prof.ssa Paola Briani per i suggerimenti e le proposte di chiarimento di alcuni punti della trattazione.

Capitolo 1

Inversa generalizzata di operatori lineari

1.1 Introduzione

Siano X e Y spazi di Hilbert e $T : X \rightarrow Y$ un operatore lineare limitato. Indichiamo con $D(T)$, $N(T)$, $R(T)$ rispettivamente il dominio, lo spazio nullo e il range dell'operatore T .

Si consideri l'equazione lineare

$$Tx = y \quad (1.1)$$

dove $y \in Y$ è il dato assegnato.

Evidentemente se l'operatore T è invertibile, l'equazione (1.1) ha sempre una e una sola soluzione.

Si supponga però che $R(T) \neq Y$ e $y \notin R(T)$; in tal caso l'equazione non ammette soluzione, infatti $Tx \neq y \quad \forall x \in X$. Inoltre se l'operatore non è iniettivo ($\Leftrightarrow N(T) \neq 0$) la soluzione anche nel caso in cui esista, può non essere unica.

Queste difficoltà possono essere superate con l'introduzione della inversa generalizzata; essa consente di estendere il concetto di "unica soluzione".

Nella maggior parte della letteratura, la soluzione generalizzata viene introdotta e definita come l'elemento di norma minima nell'insieme delle soluzioni ai minimi quadrati. Questa definizione, detta variazionale, verrà presentata nel prossimo paragrafo.

Un approccio differente, che poggia sulla determinazione di un operatore invertibile $\tilde{T} : X \rightarrow Y$ tale che $\tilde{T}|_{N(T)^\perp} = T|_{N(T)^\perp}$, conduce nel paragrafo successivo a una nuova definizione di inversa generalizzata che chiameremo operatoriale. Verrà dimostrata l'equivalenza delle due definizioni prima nel caso in cui $R(T)$ è chiuso e successivamente nel caso in cui $R(T)$ è arbitrario, ossia non necessariamente chiuso.

Nel quarto paragrafo la definizione operatoriale consentirà di verificare in maniera semplice diverse proprietà dell'inversa generalizzata; sulla base di queste proprietà verranno poi enunciate interessanti definizioni alternative.

Il capitolo si chiude con una breve analisi della inversa generalizzata per operatori compatti e con la determinazione dell'inversa generalizzata per matrici; quest'ultima risulta la versione numerica della definizione operatoriale per operatori dotati di range chiuso.

1.2 Pseudosoluzione e soluzione generalizzata

Indichiamo con P_R l'operatore di proiezione ortogonale su $R(T)$. Se l'operatore ha range chiuso l'equazione $Tx = Py$ ammette sempre un insieme di soluzioni; infatti, poiché $P_{RY} \in R(T) = R(T)$, allora $\{x \in X : Tx = Py\} \neq \emptyset$. Inoltre, osservando che P_{RY} è l'elemento in $R(T)$ più vicino a $y \in Y$, tali soluzioni sembrano essere le "migliori" ottenibili, ossia quelle la cui immagine dista meno possibile dal dato.

Generalizzando questa osservazione ad operatori a range arbitrario, introduciamo il teorema seguente che permette di definire la pseudosoluzione.

In seguito indicheremo con $\mathcal{L}(X, Y)$ l'insieme degli operatori lineari da X in Y , e con $\mathcal{B}(X, Y)$ l'insieme degli operatori lineari e limitati; l'aggiungimento di un operatore $B \in \mathcal{B}(X, Y)$ verrà indicato con B^* , la chiusura di un insieme S verrà indicata con \bar{S} e lo spazio ortogonale con S^\perp .

Teorema 1.1 Consideriamo $T \in \mathcal{B}(X, Y)$, $u \in X$ e P_R l'operatore di proiezione ortogonale di Y su $R(T)$.

Le seguenti condizioni sono equivalenti:

- (i) $Tu = P_{RY}$
- (ii) $\|Tu - y\| \leq \|Tx - y\| \quad \forall x \in X$
- (iii) $T^*Tu = T^*y$

Dimostrazione.

(i) \Rightarrow (ii) Si osservi innanzitutto che $y \in R(T) \oplus R(T)^\perp$ poichè, dall'equazione (i), $P_{RY} \in R(T)$.

Supposto $y = y_1 + y_2$ con $y_1 \in R(T)$, $y_2 \in R(T)^\perp$, si ha $(P_{RY} - y) = (y_1 - (y_1 + y_2)) = y_2 \in R(T)^\perp$. Sia $x \in X$, allora

$$\begin{aligned} \|Tx - y\|^2 &= \|Tx - P_{RY}\|^2 + \|P_{RY} - y\|^2 \\ &= \|Tx - P_{RY}\|^2 + \|Tu - y\|^2 \geq \|Tu - y\|^2 \end{aligned}$$

(ii) \Rightarrow (iii)

Poichè $\|Tu - y\| \leq \|Tx - y\| \quad \forall x \in X$ possiamo considerare Tu come quell'elemento $\hat{w} \in R(T)$ t.c.

$$\|\hat{w} - y\| \leq \|w - y\| \quad \forall w \in R(T)$$

Per il teorema della proiezione si ha che $\hat{w} - y \in R(T)^\perp = N(T^*)$, quindi $T^*(\hat{w} - y) = 0 \Leftrightarrow T^*\hat{w} = T^*y$ cioè

$$T^*Tu = T^*y$$

(iii) \Rightarrow (i)

$Tu - y \in N(T^*) = R(T)^\perp$ quindi

$$P_R(Tu - y) = 0 \Leftrightarrow P_R Tu = P_{RY} \Leftrightarrow Tu = P_{RY}$$

Definizione 1.1 Ogni elemento $u \in X$ che verifica le condizioni (i) - (iii) del teorema 1.1 è detto pseudosoluzione o soluzione nel senso dei minimi quadrati.

È facile osservare da (ii) che l'immagine di ogni pseudosoluzione approssima al meglio il dato dell'equazione (1.1), infatti minimizza la distanza $\|Tx - y\|$ su tutto lo spazio X . Il teorema 1.1 considera un generico operatore $T \in \mathcal{B}(X, Y)$ a range arbitrario. Se $R(T)$ non è chiuso la proiezione P_{RY} può non appartenere a $R(T)$. In tal caso, osservando (i), risulta evidente che non esistono pseudosoluzioni del problema (1.1).

Proposizione 1.1 Il problema (1.1) ammette pseudosoluzioni se e solo se $y \in R(T) \oplus R(T)^\perp$.

Dimostrazione. Se $y \in R(T) \oplus R(T)^\perp$ allora $y = y' + y^\perp$, con $y' \in R(T)$ e $y^\perp \in R(T)^\perp$. L'insieme S formato dalle pseudosoluzioni del problema (1.1) è la controimmagine secondo T dell'elemento y' cioè $S = \{u \in X, Tu = y'\} \neq \emptyset$. Viceversa se $y \in R(T) \setminus R(T)$, allora $P_{RY} \notin R(T)$ e quindi l'equazione (i) del teorema 1.1 non ha soluzioni.

Si noti che l'insieme $R(T) \oplus R(T)^\perp$ è denso in Y e coincide con esso se $R(T)$ è chiuso. Quindi se l'operatore T possiede range chiuso la pseudosoluzione esiste sempre. Il passo successivo nella definizione di inversa generalizzata è la scelta di un elemento particolare nell'insieme delle pseudosoluzioni.

Proposizione 1.2 Se $y \in R(T) \oplus R(T)^\perp$, l'insieme $S = \{x \in X, Tx = P_{RY}\}$ delle pseudosoluzioni è non vuoto e convesso.

Dimostrazione. Se $y \in R(T) \oplus R(T)^\perp$ allora $P_{RY} \in R(T)$, quindi, come già osservato nella proposizione precedente, l'insieme S è non vuoto, poichè esiste necessariamente almeno un $x \in X$ tale che $Tx = P_{RY}$.

Se T è iniettivo, S è costituito da un solo elemento. In questo caso generalizzando la definizione di insieme convesso, consideriamo verificata la tesi.

Supponiamo allora che S contenga almeno 2 elementi, $x', x'' \in S$. Dimostriamo che $x''' = \alpha x' + (1 - \alpha)x'' \in S$ con $\alpha \in (0, 1)$.

$$\begin{aligned} \|Tx''' - y\| &= \|T(\alpha x' + (1 - \alpha)x'') - y\| = \|\alpha Tx' + (1 - \alpha)Tx'' - y\| \\ &= \|\alpha T(x' - x'') + T(x'' - y)\| = \|T(x'' - y)\| \end{aligned}$$

Pr l'ultima uguaglianza si osservi che, da $T^*Tx' = T^*y$ e $T^*Tx'' = T^*y$, si ha $T^*T(x' - x'') = 0$ quindi

$$T(x' - x'') \in N(T^*) = R(T)^\perp \Rightarrow T(x' - x'') = 0$$

Considerando nuovamente (ii) del teorema (1.1) si conclude che $x''' \in S$.

Siccome l'insieme S delle pseudosoluzioni è chiuso (infatti è la controimmagine di un elemento rispetto a un operatore continuo) e convesso, esso ammette un elemento di norma

minima. In accordo con quanto detto nell'introduzione, definiamo l'inversa generalizzata come la pseudosoluzione di norma minima.

Definizione 1.2 VARIAZIONALE
 L'inversa generalizzata T^\dagger di T è una applicazione $T^\dagger : D(T^\dagger) = R(T) \oplus R(T)^\perp \rightarrow X$ definita come

$$T^\dagger y = x^\dagger$$

$$\text{dove}$$

$$(i) \quad x^\dagger \in S = \{ u \in X : \|Tu - y\| \leq \|Tx - y\| \quad \forall x \in X \} \subset X$$

$$(ii) \quad \|x^\dagger\| \leq \|u\| \quad \forall u \in S.$$

La pseudosoluzione x^\dagger viene detta soluzione generalizzata del problema (1.1).

Si noti che, se l'operatore è invertibile, $S = \{y^{-1}\}$, $x^\dagger = y^{-1}$, $T^\dagger = T^{-1}$.

1.3 Definizione operatoriale di inversa generalizzata

La definizione di inversa generalizzata adottata nel precedente paragrafo, che poggia sul teorema 1.1, viene detta variazionale. Al fine di analizzarne le proprietà, introduciamo lo stesso concetto di inversa generalizzata in un altro modo.

Consideriamo innanzitutto il caso in cui l'operatore $T \in B(X, Y)$ ha range chiuso. Quello che ci apprestiamo a fare è definire condizioni necessarie e sufficienti affinché tale operatore possa essere scritto come $T = P_R \tilde{T}$ con P_R operatore di proiezione su $R(T)$ e $\tilde{T} \in B(X, Y)$ invertibile. Definiremo così l'inversa generalizzata come $T^\dagger = P_{N(T)^\perp} \tilde{T}^{-1}$. Vediamo la cosa in dettaglio.

Teorema 1.2 Consideriamo $T \in B(X, Y)$. Supponiamo che

$$T = P \tilde{T} \tag{1.2}$$

dove P è un operatore di proiezione, $\tilde{T} \in B(X, Y)$ invertibile con inversa anch'essa limitata. Allora:

- (i) $R(T)$ è chiuso
- (ii) P è la proiezione su $R(T)$, ossia $P = P_R$
- (iii) $\{y \in Y : y = \tilde{T}x, x \in N(T)\} = R(T)^\perp$
 cioè $N(T)$ è in corrispondenza biunivoca, rispetto a \tilde{T} , con $R(T)^\perp$

Viceversa se $R(T)$ è chiuso e $\dim(N(T)) = \dim(R(T)^\perp)$, allora T ammette la rappresentazione (1.2). L'operatore \tilde{T} può essere scelto in modo tale che $\|\tilde{T}\| = \|\tilde{T}\|$

Il teorema precedente, con dimostrazione, enunciato nel caso più generale in cui X, Y sono spazi di Banach, si trova in [4].

L'operatore \tilde{T} appare qui come una sorta di estensione invertibile dell'operatore T . Si osservi infatti che:

$$\tilde{T}x \neq 0 \quad \forall x \in N(T) \setminus \{0\}$$

$$\tilde{T}x = Tx \quad \forall x \in N(T)^\perp$$

Si osservi che se vale la rappresentazione (1.2), X e Y hanno necessariamente la stessa dimensione (\tilde{T} infatti è invertibile). Il teorema può comunque essere generalizzato al caso in cui $\dim(X) \neq \dim(Y)$; si tratterà di immergere l'operatore T in modo opportuno, a seconda che $\dim(X) > \dim(Y)$ o $\dim(X) < \dim(Y)$. Il prossimo teorema asserisce che l'estensione di cui abbiamo bisogno esiste sempre.

Teorema 1.3

• Sia $T' \in B(X', Y)$, $R(T')$ chiuso.
 Se $\dim(X') < \dim(Y)$ allora esiste uno spazio di Hilbert X , $X \supset X'$ e un'operatore $T \in B(X, Y)$ t.c.

$$(i) \quad Tx = T'x \quad \forall x \in X'$$

$$Tx = 0 \quad \forall x \in X \ominus X'$$

$$(ii) \quad R(T') = R(T)$$

(iii) T ammette rappresentazione (1.2)

• Sia $T' \in B(X, Y')$, $R(T')$ chiuso.

Se $\dim(X) > \dim(Y')$ allora esiste uno spazio di Hilbert Y , $Y \supset Y'$ e un'operatore $T \in B(X, Y)$ t.c.

$$(i') \quad Tx = T'x \quad \forall x \in X$$

$$(ii') \quad R(T') = R(T)$$

(iii') T ammette rappresentazione (1.2)

Per la dimostrazione si consulti [4].

Definizione 1.3 OPERATORIALE

Definiamo inversa generalizzata per l'operatore $T \in B(X, Y)$ dotato di range chiuso, e la indichiamo con T^\dagger , l'applicazione lineare $T^\dagger : Y \rightarrow X$ definita come

$$T^\dagger = P_{N(T)^\perp} \tilde{T}^{-1} \tag{1.3}$$

Riassumendo, l'inversa generalizzata T^\dagger viene determinata nel modo seguente. L'operatore T , con range chiuso, viene "esteso" con invertibilità a \tilde{T} ponendo in corrispondenza biunivoca $N(T) \subset X$ con $R(T)^\perp \subset Y$. Invertendo \tilde{T} e considerando come immagini solo le componenti in $N(T)^\perp$, le uniche significative rispetto a T , si ottiene

$$T : X \rightarrow Y$$

$$N(T)^\perp \xleftrightarrow{T^\dagger} R(T)$$

$$N(T) \rightarrow 0 \tag{1.4}$$

(1.5)

$$\begin{aligned} \tilde{T} : X &\longrightarrow Y \\ N(T)^\perp &\longleftarrow R(T) \\ N(T) &\longleftarrow R(T)^\perp \end{aligned}$$

(1.6)

$$\begin{aligned} T^\dagger : Y &\longrightarrow X \\ R(T) &\longleftarrow N(T)^\perp \\ R(T)^\perp &\longrightarrow 0 \end{aligned}$$

Si osservi che la trattazione adottata in questo paragrafo è indipendente da quella del paragrafo precedente, ossia la definizione 1.3 non fa uso di strumenti definiti nel paragrafo 1.2. Occorre comunque dimostrare che le due definizioni sono equivalenti e non creano alcuna inconsistenza. Questo problema è risolto dal seguente teorema.

Teorema 1.4 Se T ammette rappresentazione (1.2), allora $\forall y \in Y$ considerando la definizione 1.3, si ha

- (a) $\inf_{x \in X} \|Tx - y\| = \|Tx_0 - y\|$ con $x_0 = \tilde{T}^{-1}y$
 L'estremo inferiore è raggiunto da tutti e soli gli $x \in \Delta_0 = \{x : Tx = Tx_0\} \subset X$
- (b) (i) $\inf_{x \in \Delta_0} \|x\| = \|x^\dagger\|$, dove $x^\dagger = P_{N(T)^\perp} \tilde{T}^{-1}y = T^\dagger y$
 (ii) $\|x\| > \|x^\dagger\| \quad \forall x \in \Delta_0, \quad x \neq x^\dagger$
- (c) L'operatore T^\dagger che soddisfa le proprietà (i) - (ii) è unico.

Dimostrazione.

Si osservi che, dal teorema 1.2 (i), $R(T)$ risulta necessariamente chiuso. Osservando che $Tx_0 - y$ e $T(x - x_0)$ sono ortogonali, infatti $T(x - x_0) \in R(T)$ e $Tx_0 - y = T\tilde{T}^{-1}y - y = P_{R\tilde{T}^{-1}}y - y = (P_R - I)y = P_{R(T)^\perp}y \in R(T)^\perp$, si ha

$$\|Tx - y\|^2 = \|Tx_0 - y\|^2 + \|T(x - x_0)\|^2$$

Quindi l'estremo inferiore in (a) assume il valore $\|Tx_0 - y\|$, che si ottiene se e solo se $Tx = Tx_0$. (a) risulta così dimostrato.

Considerando lo spazio quoziente $X \setminus N$, che si ottiene quotizzando rispetto alla relazione di equivalenza

$$x' \sim x'' \Leftrightarrow x' - x'' \in N(T) \Leftrightarrow T(x' - x'') = 0$$

l'insieme Δ_0 può essere visto come la classe d'equivalenza contenente x_0 . Siccome $TP_{N(T)^\perp} = T$ (vedi (1.4)), allora $x^\dagger \in [\Delta_0]$, infatti dalla definizione 1.3

$Tx^\dagger = TP_{N(T)^\perp} \tilde{T}^{-1}y = T\tilde{T}^{-1}y = Tx_0$. L'elemento x^\dagger può essere preso come rappresentante della classe $[\Delta_0]$; quindi

$$\|x\|^2 = \|x^\dagger\|^2 + \|z\|^2 \quad \forall x \in [\Delta_0]$$

Si ottiene $\|x\| > \|x^\dagger\| \quad \forall x \neq x^\dagger$, che dimostra (b).

Per dimostrare l'unicità dell'operatore T^\dagger , consideriamo un operatore lineare T' che soddisfi le proprietà (b). Evidentemente $T'y = x^\dagger$, poiché x^\dagger è l'unico elemento che ha la proprietà di minimo (b) - (ii).

Allora

$$T^\dagger y = T'y \quad \forall y \in Y$$

da cui $(T^\dagger - T')y = 0 \quad \forall y \in Y$, cioè $T^\dagger = T'$ che conclude la dimostrazione.

Dal teorema 1.4 (a) si osserva immediatamente che l'insieme $\Delta_0 \subset X$, i cui elementi minimizzano la distanza $\|Tx - y\|$, coincide con l'insieme delle soluzioni ai minimi quadrati del paragrafo 1.1. Come dimostrato al punto (b), l'inversa generalizzata 1.3 è l'unico elemento di norma minima in Δ_0 . L'equivalenza tra le due definizioni (variazionale e operatoriale) nel caso di operatori con range chiuso è così dimostrata.

Lo stesso percorso che conduce alla definizione 1.3 può essere seguito per operatori a range arbitrario.

Consideriamo $T \in \mathcal{B}(X, Y)$ con range non chiuso.

Si tratta di definire sotto quali condizioni una rappresentazione (1.2) può essere fatta. In questo caso vedremo che l'operatore \tilde{T} ammette inversa non limitata definita su un sottospazio denso in Y . L'inversa generalizzata verrà costruita in maniera analoga, ma il dominio sarà nuovamente un insieme denso in Y . Si osservi che già nel paragrafo precedente l'inversa generalizzata è definita su $R(T) \oplus R(T)^\perp$, che, nel caso in cui $R(T)$ non è chiuso, risulta essere denso in Y .

Alcuni dei teoremi seguenti vengono solo enunciati. Per le dimostrazioni si può consultare [4].

Teorema 1.5 Sia $T \in \mathcal{B}(X, Y)$, con range non chiuso.

Si supponga inoltre che $\dim(N(T)) = \dim(R(T)^\perp)$. Allora

- (i) T ammette la rappresentazione $T = P\tilde{T}$

con P operatore di proiezione su Y , e \tilde{T} dotato di inversa definita su un sottospazio denso in Y .

- (ii) $P = P_{R(T)}$

- (iii) \tilde{T}^{-1} è non limitato

Definizione 1.4 Si definisce inversa generalizzata T^\dagger per l'operatore $T \in \mathcal{B}(X, Y)$ l'applicazione lineare $T^\dagger : D(T^\dagger) \rightarrow X$ che soddisfa:

- (i) $D(T^\dagger)$ è denso in Y
- (ii) $\forall y \in D(T^\dagger) \quad \inf_{x \in X} \|Tx - y\| = \|T^\dagger y - y\| \quad \text{dove} \quad x^\dagger = T^\dagger y$
- (iii) $\|x\| < \|x^\dagger\| \quad \forall x' \in X, \quad x' \neq x^\dagger \quad \text{tale che} \quad Tx' = Tx^\dagger$

Prima di procedere si osservi che nel caso in cui $R(T)$ è chiuso, l'inversa generalizzata viene definita come $T^\dagger = P_{N(T)^\perp} T^{-1}$ (eq. (1.3)). Il teorema 1.4 dimostra che tale operatore soddisfa le proprietà di minimo enunciate.

La definizione 1.4, che contempla anche il caso in cui $R(T)$ non è chiuso, è stata data invece basandosi sulle proprietà di minimo. Verrà dimostrato in seguito che tale operatore T^\dagger può essere scritto in una forma corrispondente alla (1.3).

Si osservi inoltre che l'unicità dell'operatore, garantita nel caso in cui $R(T)$ è chiuso, qui non è menzionata. È possibile infatti definire diverse inverse generalizzate, considerando diversi domini densi in Y . In realtà due generiche inverse generalizzate sono coincidenti sul sottoinsieme dove sono entrambe definite.

Teorema 1.6 Siano $T^1, T^{1''}$ due inverse generalizzate soddisfacenti la definizione (1.4)
 Sia $D = D(T^1) \cap D(T^{1''})$.

Allora
$$T^1 y = T^{1''} y \quad \forall y \in D$$

Dimostrazione. Supponiamo per assurdo che esista $\bar{y} \in D$ t.c.

$$x^1 = T^1 \bar{y} \neq x^{1''} = T^{1''} \bar{y}$$

Per definizione x^1 e $x^{1''}$ soddisfano entrambi la proprietà di minimo 1.4 - (ii), quindi $Tx^1 = Tx^{1''}$.

Ma allora per la proprietà (iii) si dovrebbe avere $\|x^1\| < \|x^{1''}\|$ e $\|x^{1''}\| < \|x^1\|$, il che è assurdo. ■

Corollario 1.1 Se T^1 è definito su tutto Y , qualsiasi altra inversa generalizzata deve necessariamente essere una restrizione di T^1 .

Come già detto vediamo come l'inversa generalizzata, nel caso in cui $R(T)$ non sia chiuso, ammetta comunque la rappresentazione operatoriale del tipo (1.3).

Teorema 1.7 Sia $T \in B(X, Y)$, $R(T)$ non chiuso e $\dim(N(T)) = \dim(R(T)^\perp)$. Allora, considerando \bar{T} introdotta nel teorema 1.5, l'operatore

$$T^1 = P_{N(T)^\perp} \bar{T}^{-1}$$

è un'inversa generalizzata nel senso della definizione 1.4. Inoltre considerando $x^1 = T^1 y$, $y \in D(T^1)$ si ha

$$\{x^1 \in X : \|Tx^1 - y\| \leq \|Tx - y\| \quad \forall x \in X\} = \{x^1 \in X : Tx^1 = Tx\}$$

T^1 è non limitata.

Dimostrazione. Utilizzando la decomposizione del teorema 1.5, dimostriamo che T^1 è effettivamente una inversa generalizzata come definita in 1.4.

$D(\bar{T}^{-1})$ è denso in Y e $D(\bar{T}^{-1}) = D(T^1)$, quindi T^1 soddisfa (i) della definizione 1.4. Dimostriamo ora che T^1 soddisfa (ii) della definizione 1.4.

Poiché $T = P_{R(T)} \bar{T}$, allora $R(\bar{T}) \subseteq R(T) \oplus R(T)^\perp$; quindi $D(T^1) = D(\bar{T}^{-1}) = R(\bar{T}) \subseteq R(T) \oplus R(T)^\perp$.

Verifichiamo innanzitutto che, $\forall y \in D(T^1)$, $TT^1 y = P_{R(T)} y$. Si ha

$$TT^1 y = TP_{N(T)^\perp} \bar{T}^{-1} y = P_{R(T)} \bar{T} \bar{T}^{-1} y = P_{R(T)} y \quad \forall y \in D(T^1).$$

Sia $x^1 \in D(T)$ e $y \in D(T^1) = D(\bar{T}^{-1})$; allora

$$\|Tx^1 - y\|^2 = \|Tx^1 - y\|^2 + \|Tx^1 - Tx^1\|^2$$

infatti $Tx^1 - y = TT^1 y - y = (TT^1 - I)y = (P_{R(T)} - I)y = -P_{R(T)^\perp} y \in R(T)^\perp$

e $Tx^1 - Tx^1 = T(x^1 - x^1) \in R(T)$.

Quindi $\|Tx^1 - y\| \leq \|Tx^1 - y\|$, che dimostra (ii).

Per dimostrare infine (iii), si osservi che l'insieme $\Delta = \{x \in D(T) : Tx = Tx\}$ contiene tutti e soli gli elementi che minimizzano la relazione al punto (ii).

Ogni elemento $\bar{x} \in \Delta$ può essere scritto nella forma $\bar{x} = x^1 + z$, con $z \in N(T)$, infatti

$$Tx = Tx^1 \Leftrightarrow T(x - x^1) = 0 \Leftrightarrow x - x^1 \in N(T).$$

Inoltre, per come è definito,

$$x^1 = T^1 y = P_{N(T)^\perp} (\bar{T}^{-1} y) \in N(T)^\perp$$

Quindi $\|x^1\|^2 = \|x^1\|^2 + \|z\|^2$ che verifica (iii).

Per completare, verifichiamo che T^1 risulta necessariamente non limitata. Accenniamo brevemente come si svilupperà la dimostrazione.

Se per assurdo T^1 fosse limitata, vedremmo che l'operatore T ammetterebbe una rappresentazione $T = P_{R(T)} \bar{T}_0$, con \bar{T}_0 operatore limitato invertibile con inversa limitata. In questo caso il teorema 1.2 garantirebbe che T possiede range chiuso, il che contraddice l'ipotesi su T . Vediamo la cosa in dettaglio.

Supponiamo per assurdo che $T^1 = P_{N(T)^\perp} \bar{T}^{-1}$ sia limitata. Definiamo $T_0 : X \rightarrow Y$

$$T_0 = T + UP_{N(T)}$$

dove U è un'isometria che immerge $N(T) \subset X$ in $R(T)^\perp \subset Y$, che sono della stessa dimensione per ipotesi.

Poiché $P_{R(T)} UP_{N(T)} = 0$, si ha che

$$T = P_{R(T)} T_0$$

con T_0 limitato poichè somma di T e $UP_{N(T)}$ entrambi limitati. Dimostriamo che l'inverso di T_0 è l'operatore $B : D(B) \rightarrow X$

$$B = T^1 + U^{-1} P_{R(T)^\perp}$$

Osserviamo innanzitutto che, essendo T^1 limitato, anche B è un operatore limitato e che $D(B) = D(T^1)$, infatti $D(B) = D(T^1) \cap D(U^{-1} P_{R(T)^\perp}) = D(T^1) \cap Y = D(T^1)$.

Inoltre, poichè $TU^{-1} P_{R(T)^\perp} = 0$, $UP_{N(T)} T^1 = 0$ e $UP_{N(T)} U^{-1} P_{R(T)^\perp} = P_{R(T)^\perp}$ (per l'ultima si osservi che $UP_{N(T)} U^{-1} |_{R(T)^\perp} = I_{R(T)^\perp}$), si ha che $(T_0 B)y = (TT^1 + P_{R(T)^\perp})y \quad \forall y \in D(B)$.

Siccome, come già visto in precedenza, $TT^1 = P_{R(T)}$ su $D(B) = D(T^1)$, allora

$$T_0 B y = TT^1 y + P_{R(T)^\perp} y = (P_{R(T)} + P_{R(T)^\perp})y = y \quad \forall y \in D(B),$$

quindi B è inversa destra di T_0 .

Vediamo che B è anche inversa sinistra di T_0 e che quindi corrisponde all'inversa di T_0 .

Si ha $BT_0 B y = B y \quad \forall y \in D(B)$.

$R(B)$ è lo spazio generato da $R(T^1) = N(T)^\perp \cap D(T)$ e $N(T)$, che equivale esattamente a $D(T) = D(T_0)$. Quindi $\forall x \in D(T_0) \exists y \in D(B)$ tale che $x = B y$, da cui $BT_0 x = x' \quad \forall x' \in D(T_0)$.

B è anche inversa sinistra di T_0 , quindi è l'unica inversa di T_0 , ossia $B = T_0^{-1}$.

Quindi T ammette una rappresentazione del tipo (1.2) con T_0 invertibile con inversa limitata. Per il teorema 1.2, $R(T)$ risulta essere necessariamente chiuso, in contraddizione con l'ipotesi su $R(T)$. Ne segue che T^1 non può essere limitata. ■

Si osservi che nel caso in cui l'operatore è invertibile la rappresentazione (1.2) si caratterizza con $P_{R(T)} = I$, $\bar{T} = T$. Si ottiene $T^1 = P_{N(T)^\perp} \bar{T}^{-1} = IT^{-1} = T^{-1}$, come ovviamente ci si aspettava.

A conclusione di questo paragrafo è utile sottolineare che tutti i risultati ottenuti valgono nel contesto degli operatori lineari tra spazi di Hilbert. Per un'analisi dettagliata dell'inversa generalizzata in spazi di Banach si può consultare [3] e [4]; mentre nel caso in cui gli spazi siano solo topologici si può consultare [49].

1.4 Proprietà della inversa generalizzata

Enunciamo e verifichiamo alcune proprietà dell'operatore T^\dagger .

Teorema 1.8 Se $T \in \mathcal{B}(X, Y)$ è dotato di range chiuso, allora

$$R(T^\dagger) = R(T^*) = R(T^\dagger T)$$

Dimostrazione. La dimostrazione poggia sulla nota relazione $R(T^*) = N(T)^\perp$, valida poiché $R(T)$ è chiuso.

- $R(T^\dagger) \subseteq R(T^*)$

La verifica è ovvia considerando la definizione operatoriale 1.3.

Vediamo una dimostrazione equivalente che utilizza le relazioni esposte nella definizione variazionale 1.2.

Dimostriamo che, fissato $y \in Y$, si ha $T^\dagger y \in N(T)^\perp = R(T^*)$.

Per il teorema di decomposizione, siccome $X = N(T) \oplus N(T)^\perp$, consideriamo

$$x^\dagger = T^\dagger y = x_1 + x_2, \text{ con } x_1 \in N(T)^\perp, x_2 \in N(T). \text{ Allora}$$

$$Tx_1 = Tx_1 + Tx_2 = TT^\dagger y = PRy.$$

Per il teorema 1.1, x_1 è una soluzione ai minimi quadrati. Allora si deve avere necessariamente $x_2 = 0$, poiché altrimenti $\|x_1\| < \|x^\dagger\|$.

Si ottiene $x^\dagger \in N(T)^\perp = R(T^*)$.

- $R(T^\dagger) \supseteq R(T^*)$

Anche in questo caso, per mezzo della definizione 1.3, la dimostrazione è immediata (si osservi che $R(\tilde{T}^{-1}) = X$).

Si può comunque verificare il risultato in maniera equivalente utilizzando la definizione 1.2. Supponiamo che $x \in R(T^*)$. Poniamo $y = Tx$.

$Tx = PRTx = PRy$; x è quindi una soluzione ai minimi quadrati. Dimostriamo che è proprio quella di norma minima.

Qualsiasi altra pseudosoluzione x' verifica $Tx' = PRy = Tx$.

Quindi $x' = x + (x' - x)$ con $x \in N(T)^\perp$ e $(x' - x) \in N(T)$ (infatti $T(x' - x) = Tx' - Tx = PRy - PRy = 0$).

Allora $\|x'\|^2 = \|x\|^2 + \|x' - x\|^2 \geq \|x\|^2$.

x è la pseudosoluzione di norma minima, quindi $x = T^\dagger y \in R(T^\dagger)$.

- Per la seconda uguaglianza, ovviamente $R(T^\dagger) \supseteq R(T^\dagger T)$.

Inoltre, se $x \in N(T)^\perp$, cioè $x = T^\dagger y$ con $y \in Y$, si ha $x = T^\dagger PRy \in R(T^\dagger T)$, infatti considerando il teorema 1.1, $Tx = PRy \Leftrightarrow Tx = PR(PRy)$.

Quindi $R(T^\dagger) \subseteq R(T^\dagger T)$.

Teorema 1.9 Sia $T \in \mathcal{B}(X, Y)$. Allora:

- (i) T^\dagger è lineare

- (ii) T^\dagger è limitato $\Leftrightarrow R(T)$ è chiuso

Dimostrazione.

Verifichiamo (i).

Siano $y_1, y_2 \in D(T^\dagger)$ con $D(T^\dagger)$ sottospazio denso in Y .

$$TT^\dagger y_1 + TT^\dagger y_2 = PRy_1 + PRy_2 = PR(y_1 + y_2) = TT^\dagger(y_1 + y_2). \text{ Allora}$$

$$T(T^\dagger y_1 + T^\dagger y_2 - T^\dagger(y_1 + y_2)) = 0.$$

Quindi

- (a) $T^\dagger y_1 + T^\dagger y_2 - T^\dagger(y_1 + y_2) \in N(T) \cap N(T)^\perp = \{0\}$

Inoltre $TT^\dagger(\alpha y_1) = P_R(\alpha y_1) = \alpha P_R y_1 = \alpha TT^\dagger y_1 = T(\alpha T^\dagger y_1)$, con $\alpha \in \mathbb{R}$. Allora, analogamente ad (a), si ottiene

- (b) $T^\dagger(\alpha y_1) = \alpha T^\dagger y_1$.

(a) e (b) verificano (i).

Per (ii) abbiamo già visto e dimostrato nel teorema 1.7 che se $R(T)$ non è chiuso allora T^\dagger risulta essere non limitato.

Dimostriamo ora che se $R(T)$ è chiuso, allora T^\dagger è limitato.

Essendo $R(T)$ chiuso, allora

$$\exists m > 0 \text{ tale che } \|Tx\| \geq m\|x\| \quad \forall x \in N(T)^\perp$$

(cfr. [13], teorema 1.2.1), che implica $\|TT^\dagger y\| \geq m\|T^\dagger y\| \quad \forall y \in Y$.

Si ottiene $\|y\| \geq \|PRy\| = \|TT^\dagger y\| \geq m\|T^\dagger y\|$, quindi

$$\|T^\dagger y\| \leq m^{-1}\|y\|$$

Proposizione 1.8 Se $T \in \mathcal{B}(X, Y)$, $R(T)$ chiuso, allora

- (i) $TT^\dagger = P_{R(T)}$

- (ii) $T^\dagger T = P_{N(T)^\perp}$

- (iii) $T^\dagger TT^\dagger = T^\dagger$

- (iv) $TT^\dagger T = T$

- (v) $(TT^\dagger)^* = TT^\dagger$

- (vi) $(T^\dagger T)^* = T^\dagger T$

- (vii) $(T^\dagger)^\dagger = T$

Dimostrazione. Useremo sistematicamente la definizione operatoriale 1.3. Le stesse dimostrazioni si possono fare utilizzando la definizione variazionale, che risulta però meno efficace.

- (i) $TT^\dagger = T(P_{N(T)^\perp} \tilde{T}^{-1}) = (TP_{N(T)^\perp}) \tilde{T}^{-1} = TT^\dagger = P_R \tilde{T} \tilde{T}^{-1} = P_R$. Si osservi che questa relazione è stata verificata nel caso in cui $R(T)$ non è chiuso nella dimostrazione del teorema 1.7.

- (ii) $T^\dagger T = P_{N(T)^\perp} \tilde{T}^{-1} P_R \tilde{T} = P_{N(T)^\perp}$

infatti, considerando

$$x \in X, x = x' + x'', \quad x' \in N(T)^\perp, x'' \in N(T),$$

$$y' = \tilde{T}x', y'' = \tilde{T}x'', \quad y' \in R(T), y'' \in R(T)^\perp,$$

si ottiene

$$P_{N(T)^\perp} \tilde{T}^{-1} P_R \tilde{T} x = P_{N(T)^\perp} \tilde{T}^{-1} P_R \tilde{T} (x' + x'')$$

$$= P_{N(T)^\perp} \tilde{T}^{-1} P_R \tilde{T} x' + P_{N(T)^\perp} \tilde{T}^{-1} P_R \tilde{T} x''$$

$$= P_{N(T)^\perp} \tilde{T}^{-1} P_R y' + P_{N(T)^\perp} \tilde{T}^{-1} P_R y''$$

$$= P_{N(T)^\perp} \tilde{T}^{-1} P_R y' = P_{N(T)^\perp} \tilde{T}^{-1} y' = P_{N(T)^\perp} x'$$

$$= P_{N(T)^\perp} (x' + x'') = P_{N(T)^\perp} x$$

(iii) Utilizzando (ii) si ottiene
 $T^{\dagger}TT^{\dagger} = P_{N(T^{\dagger})}T^{\dagger} = P_{N(T^{\dagger})}P_{N(T^{\dagger})}\tilde{T}^{-1} = P_{N(T^{\dagger})}\tilde{T}^{-1} = T^{\dagger}$.

(iv) Utilizzando (i) si ottiene
 $TT^{\dagger}T = P_R\tilde{T} = T$

(v) - (vi) Per dimostrare le due relazioni, si osservi che, da (i), (ii) segue che $T^{\dagger}T$ e TT^{\dagger} sono operatori di proiezione. Mostriamo che ogni operatore di proiezione è autoaggiunto; la (v), (vi) ne saranno una conseguenza diretta.

Sia P un generico operatore di proiezione su un generico sottospazio chiuso R ,
 con $x = x' + x''$, $y = y' + y''$, $x', y' \in R$, $x'', y'' \in R^{\perp}$.
 $(P_Rx, y) = (x', y) = (x', y) = (x, y) = (x, P_Ry)$
 quindi $(P_R)^* = P_R$

(vii) Sapendo che $T^{\dagger} = P_{N(T^{\dagger})}\tilde{T}^{-1}$, si può notare che T^{\dagger} ammette una rappresentazione del tipo (1.2) con $(\tilde{T}^{\dagger}) = \tilde{T}^{-1}$ invertibile con inversa limitata.

Si osservi inoltre che $T^{\dagger}x = 0$ per tutti e soli gli $x \in R(T^{\dagger})^{\perp}$, infatti \tilde{T}^{-1} è invertibile e fa corrispondere $R(T^{\dagger})^{\perp}$ con $N(T)$. Quindi $N(T^{\dagger}) = R(T)^{\perp}$.

Applichiamo allora la definizione 1.3 all'operatore T^{\dagger} . Si ha
 $(T^{\dagger})^{\dagger} = P_{N(T^{\dagger})}(T^{\dagger})^{-1} = P_{(R(T^{\dagger})^{\perp})}(\tilde{T}^{-1})^{-1} = P_{R(T)}\tilde{T} = T$

Si osservi che se T è autoaggiunto, considerando $R(T) = N(T^*)^{\perp} = N(T)^{\perp}$ e le proprietà (i) - (ii), si può concludere che $TT^{\dagger} = T^{\dagger}T$.
 Altre proprietà si possono verificare facilmente facendo uso di quelle enunciate nella proposizione precedente.

A conclusione di questo paragrafo enunciamo altre definizioni di inversa generalizzata, limitandoci al caso in cui $R(T)$ è chiuso.

Queste definizioni, come già accennato, poggiano sulle proprietà appena verificate. È interessante notare che queste proprietà sono di notevole importanza poiché caratterizzano esattamente l'inverso generalizzato, ossia esso è l'unico operatore che le soddisfa.

Definizione 1.5 Moore

Se $T \in B(X, Y)$, $R(T)$ chiuso, allora T^{\dagger} è l'unico operatore in $B(Y, X)$ che soddisfi:

- (i) $TT^{\dagger} = P_{R(T)}$
- (ii) $T^{\dagger}T = P_{R(T^{\dagger})}$

Definizione 1.6 Penrose

Se $T \in B(X, Y)$, $R(T)$ chiuso, allora T^{\dagger} è l'unico operatore in $B(Y, X)$ che soddisfi:

- (i) $(TT^{\dagger})^* = (TT^{\dagger})$
- (ii) $(T^{\dagger}T)^* = (T^{\dagger}T)$
- (iii) $T^{\dagger}TT^{\dagger} = T^{\dagger}$
- (iv) $TT^{\dagger}T = T$

Definizione 1.7 Deoser and Whalen

Se $T \in B(X, Y)$, $R(T)$ chiuso, allora T^{\dagger} è l'unico operatore in $B(Y, X)$ che soddisfi:

- (i) $T^{\dagger}Tx = x \quad \forall x \in N(T)^{\perp}$
- (ii) $T^{\dagger}y = 0 \quad \forall y \in R(T)^{\perp}$

Nel paragrafo precedente è stato dimostrato che la definizione variazionale è equivalente a quella operatoriale. Si ha inoltre il seguente risultato:

Proposizione 1.4 Le definizioni

1.2 (variazionale), 1.3 (operatoriale), 1.5 (di Moore), 1.6 (di Penrose), 1.7 (di Deoser-Whalen) sono equivalenti.

Si osservi che l'ultima proposizione è valida solo se $R(T)$ è chiuso. Una dimostrazione di quest'ultima proposizione si può trovare in [13]. Per maggiori dettagli sulla definizione di Moore si consulti [31]; per quella di Penrose si consulti [37]; per quella di Deoser-Whalen si consulti [10].

1.5 Inversa generalizzata di operatori compatti e matrici

Nella maggior parte delle applicazioni, l'operatore T risulta compatto.

In questo caso, utilizzando il sistema singolare $\{\mu_n; u_n, v_n\}$ di T (vedi appendice A), è possibile determinare una semplice espressione dell'operatore inverso T^{\dagger} .

Teorema 1.10 Sia $T \in B(X, Y)$ un'operatore compatto con sistema singolare $\{\mu_n; u_n, v_n\}$. Se $y \in R(T) \oplus R(T)^{\perp}$, allora

$$T^{\dagger}y = \sum_{n=1}^{\infty} \frac{1}{\mu_n} (y, v_n) Y u_n \tag{1.7}$$

Dimostrazione. Sia P_R l'operatore di proiezione ortogonale su $R(T)$.

Allora, poichè $y \in R(T) \oplus R(T)^{\perp}$, $P_R y \in R(T)$.

Dal teorema di Picard (vedi appendice A) applicato a $P_R y$, si ottiene

$$\sum_{n=1}^{\infty} \frac{1}{\mu_n^2} |(P_R y, v_n)|^2 < +\infty$$

Si osservi che, poichè $v_n \in R(T)$, allora $(P_R y, v_n) = (y, P_R v_n) = (y, v_n)$. Utilizzando nuovamente il teorema di Picard, il quale stabilisce condizioni non solo sufficienti ma anche necessarie alla convergenza, si conclude che

$$\sum_{n=1}^{\infty} \frac{1}{\mu_n} (y, v_n) u_n \quad \text{è convergente in } X$$

Sia allora $v \in X$, $v = \sum_{n=1}^{\infty} \frac{1}{\mu_n} (y, v_n) u_n$.

Si osservi che $\{u_n\} \subset N(T)^{\perp}$; quindi, poichè $N(T)^{\perp}$ è chiuso, $v \in N(T)^{\perp}$.

Si ha

$$\begin{aligned} Tv &= \sum_{n=0}^{\infty} \frac{1}{\mu_n} (P_R y, v_n) T v_n \\ &= \sum_{n=0}^{\infty} \frac{1}{\mu_n} (y, v_n) \mu_n v_n \\ &= \sum_{n=0}^{\infty} (P_R y, v_n) v_n \\ &= P_R y \end{aligned}$$

Quindi, per il teorema 1.1 (i), il limite è una soluzione ai minimi quadrati; ricordando che $v \in N(T)^\perp$, necessariamente $v = T^\dagger y$.

Nel caso in cui gli spazi di Hilbert siano di dimensione finita, l'operatore T può essere rappresentato per mezzo di una matrice. Più precisamente, un generico operatore $T' \in \mathcal{L}(X, Y)$, con $\dim(X) = r$, $\dim(Y) = s$, viene rappresentato da una matrice $T \in M_{s,r}(\mathbb{R})$, dove $M_{s,r}(\mathbb{R})$ è lo spazio delle matrici $s \times r$ a coefficienti reali. Questo si verifica ogni volta che un problema inverso viene discretizzato e quindi rappresentato su uno spazio di dimensione finita.

Teorema 1.11 Sia $T \in M_{s,r}(\mathbb{R})$, $n = \text{rank}(T)$ allora esistono $M \in M_{r,n}(\mathbb{R})$ e $N \in M_{r,n}(\mathbb{R})$ con $\text{rank}(M) = \text{rank}(N) = n$, tali che

$$T = MN^t$$

Per una precisa dimostrazione si consulti [40].
 Diamo qui una spiegazione euristica di come si determinano M e N .
 Una volta determinato il rango di T , ad esempio per mezzo dell'uso dei minori principali, si isolino n colonne linearmente indipendenti. Queste n colonne formano la matrice M .
 Per ogni colonna c linearmente dipendente da quelle considerate in M , si determini la combinazione lineare di c in funzione delle colonne di M . I coefficienti trovati dalle singole combinazioni lineari formano la matrice N .
 Vediamo un semplice esempio.

$$\begin{pmatrix} 1 & 3 & 5 \\ 2 & 1 & 5 \\ 1 & 0 & 2 \\ 1 & 1 & 3 \end{pmatrix} \text{ con } T \in M_{4,3}(\mathbb{R}), \text{ rank}(T) = 2.$$

Si ha $M \in M_{4,2}(\mathbb{R})$, $N_{3,2}(\mathbb{R})$.

Indicando con c_1 la colonna iesima:

$$M = (c_1, c_2) = \begin{pmatrix} 1 & 3 \\ 2 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix}$$

Per determinare i coefficienti della combinazione lineare che esprime c_3 in funzione di c_1 e c_2 , si risolve il sistema

$$\alpha c_1 + \beta c_2 = c_3$$

Si ottiene $\alpha = 2$, $\beta = 1$. Quindi

$$N = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ \alpha & \beta \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 2 & 1 \end{pmatrix} \text{ con } T = MN^t$$

La decomposizione descritta permette di calcolare in maniera diretta T^\dagger .

Definizione 1.8 Sia $T \in M_{r,r}(\mathbb{R})$, $n = \text{rank}(T)$, $M \in M_{r,n}(\mathbb{R})$, $N \in M_{r,n}(\mathbb{R})$ con $n = \text{rank}(M) = \text{rank}(N)$, $T = MN^t$.
 Allora si definisce pseudoinversa di T la matrice

$$T^\dagger = N(N^t N)^{-1} (M^t M)^{-1} M^t$$

Si osservi che

(a) $T^\dagger \in M_{r,r}(\mathbb{R})$

(b) l'operatore corrispondente alla matrice T^\dagger è un operatore con range chiuso.

Teorema 1.12 Se $T \in M_{r,r}(\mathbb{R})$ è la rappresentazione dell'operatore $T' \in \mathcal{B}(X, Y)$, allora la matrice pseudoinversa $T^\dagger \in M_{r,r}(\mathbb{R})$ è la rappresentazione matriciale dell'operatore $T'^\dagger \in \mathcal{B}(Y, X)$.

Inoltre $\text{rank}(T') = \text{rank}(T^\dagger)$

Dimostrazione. La dimostrazione è immediata considerando la definizione 1.6 di Penrose e le usuali proprietà del calcolo matriciale. Occorre verificare le relazioni (i) - (iv) della definizione 1.6. Mostriamo la verifica di (iii); le altre si verificano in modo analogo.

(iii) $T^\dagger T T^\dagger = T^\dagger$
 $T^\dagger T T^\dagger = N(N^t N)^{-1} (M^t M)^{-1} M^t M N^t N (N^t N)^{-1} (M^t M)^{-1} M^t$
 $= N(N^t N)^{-1} (M^t M)^{-1} M^t M^{-1} M^{-1} M^t = T^\dagger$

Per dimostrare che $\text{rank}(T) = \text{rank}(T^\dagger)$, si osservi che, per il teorema 1.8 $R(T^\dagger) = R(T)$. La matrice che rappresenta T^\dagger è la matrice trasposta.

Quindi $\text{rank}(T^\dagger) = \text{rank}(T^t) = \text{rank}(T)$

Il metodo descritto è utile per la risoluzione di problemi semplici di piccola dimensione. Vediamo adesso un altro algoritmo per calcolare T^\dagger ; esso è l'applicazione numerica del procedimento che ha condotto alla definizione operatoriale 1.3 e risulta particolarmente efficace nel caso in cui la matrice T è di dimensione molto grande.

In moltissimi problemi numerici, come ad esempio la elaborazione della immagini, la matrice, oltre ad avere una precisa struttura, è simmetrica. Consideriamo quindi una generica matrice hermitiana.

Per non appesantire troppo le notazioni, d'ora in poi chiameremo con lo stesso nome l'operatore lineare e la matrice che lo rappresenta.

Sappiamo dall'algebra lineare che se $H \in M_{n,n}(\mathbb{R})$ hermitiana, allora $\exists C, F \in M_{n,n}(\mathbb{R})$ con C matrice unitaria, le cui colonne formano una base ortogonale di autovettori di H e $F = \text{diag}(\lambda_1, \dots, \lambda_2)$, dove λ_i rappresenta l'autovale corrispondente all'autovettore dell'iesima colonna, t.c.

$$H = C^* F C$$

Si può supporre che, posto $k = \text{rank}(H)$, con $k \leq n$, allora $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_k| > 0$ e $\lambda_{k+1} = \lambda_{k+2} = \dots = \lambda_n = 0$.

Un ragionamento analogo a quello fatto all'inizio del paragrafo per la determinazione delle matrici M e N , ci permette di osservare che, posto

$$G = \text{diag}(g_1, \dots, g_n) \text{ con } g_i = \begin{cases} 1 & i = 1, \dots, k \\ 0 & \text{altrimenti} \end{cases}$$

si ha che

$$P_R = C^* G C$$

dove R è lo spazio generato dalle colonne di H .

Poiché H è hermitiana $R(H) = R(H^*) = N(H)^\perp$.

Consideriamo ora la definizione 1.3 di inversa generalizzata.

Posto $\tilde{H} = C^*VC$ dove $V = \text{diag}(v_1, \dots, v_n)$ con $v_i = \begin{cases} \lambda_i & i = 1, \dots, k \\ 1 & \text{altrimenti} \end{cases}$ verificiamo che

$$H = P_R \tilde{H}$$

Si ha che \tilde{H} è invertibile e

$$P_R \tilde{H} = (C^*GC)(C^*VC) = C^*FC = H, \text{ infatti } GV = F.$$

Inoltre $\tilde{H}^{-1} = C^*WC$ con $W = \text{diag}(w_1, \dots, w_n)$, $w_i = \begin{cases} \frac{1}{\lambda_i} & i = 1, \dots, k \\ 1 & \text{altrimenti} \end{cases}$

Sappiamo infine che

$$H^1 = P_{N^+} \tilde{H}^{-1}$$

quindi

$$H^1 = (C^*GC)(C^*WC) = C^*GWC$$

Posto $Z = GW$, si ottiene $Z = \text{diag}(z_1, \dots, z_n)$ $z_i = \begin{cases} \frac{1}{\lambda_i} & i = 1, \dots, k \\ 0 & \text{altrimenti} \end{cases}$

La matrice H^1 può essere scritta come

$$H^1 = C^*ZC$$

Quest'ultimo procedimento è utile se implementato. Accenniamo a tal proposito che gli autovettori e gli autovalori di H possono essere calcolati per mezzo di metodi di approssimazione come, ad esempio, il metodo delle potenze e il metodo basato sulla decomposizione QR di H , dove Q è una matrice unitaria e R triangolare superiore (si consulti [8]).

Capitolo 2

Metodi di regolarizzazione per problemi mal posti

2.1 Introduzione

Ogni volta che si affronta un problema concreto, i dati a disposizione sono sempre affetti da errore; è chiaro infatti che a causa della tolleranza della apparecchiatura che effettua le rilevazioni, i dati acquisiti da fenomeni naturali sono noti entro un margine fissato. Nel caso in cui una piccola variazione del dato comporti una grande variazione della soluzione, non si hanno garanzie sulla affidabilità del risultato.

Scopo di questo paragrafo è analizzare metodologie che permettano di affrontare questo problema che, come vedremo, rientra nella categoria dei problemi mal posti.

In quest'ottica verrà sviluppata la teoria generale dei metodi di regolarizzazione per la risoluzione di equazioni lineari in spazi di Hilbert. Verrà illustrato il noto algoritmo regolarizzante di A.N. Tikhonov e l'algoritmo iterativo di Landweber-Fridman. Questi, per la loro semplicità, consentono una rapida analisi delle proprietà caratteristiche dei metodi di regolarizzazione.

2.2 Problemi mal posti e problemi inversi

Consideriamo una generica equazione

$$Kx = y \tag{2.1}$$

dove K è un'applicazione tra gli spazi di topologici X e Y , e $y \in Y$ è il dato assegnato. Agli inizi di questo secolo J. Hadamard classificò la categoria dei problemi mal posti.

Definizione 2.1 Il problema (2.1) è detto ben posto se sono soddisfatte le condizioni seguenti

(a) $\forall y \in Y$ la soluzione $x \in X$ esiste

(b) la soluzione è unica

(c) la soluzione dipende con continuità da y , cioè l'applicazione $y \rightarrow x$ è continua.

Ogni problema non ben posto è detto mal posto.

Risulta evidente che l'essere ben posto è una proprietà legata non solo all'operatore K , ma dipende dagli spazi X, Y e dalle loro topologie.

La condizione (a) equivale alla surgettività dell'operatore K , ed è quindi soddisfatta se $Y = K(X)$, mentre la condizione (b) equivale all'injectività. Affinchè entrambe le condizioni (a), (b) siano soddisfatte è necessario e sufficiente che l'applicazione sia invertibile.

Come accennato nell'introduzione, la condizione (c) è quella di maggior interesse e sta alla base dello studio dei problemi non ben posti. Si osservi che questa condizione dipende dalle topologie introdotte in X e in Y , le quali possono rendere l'applicazione $y \rightarrow x$ continua o meno.

Nasce quindi l'idea di rendere ben posto qualsiasi problema modificando la terna (K, X, Y) . Ad esempio nel caso in cui K sia injectiva e continua in un generico insieme compatto

$\bar{X} \subset X$, si può considerare il problema rispetto alla terna $(K, \bar{X}, K(\bar{X}))$. In questo caso l'applicazione risulta ovviamente surgettiva e l'inversa K^{-1} , ristretta all'insieme chiuso $K(\bar{X})$, risulta continua (cfr. [47]).

Il problema così ridefinito soddisfa le condizioni (a) e (c) ed è quindi ben posto.

Si osservi che nel caso in cui l'applicazione K è lineare e gli spazi X, Y sono di Hilbert, è sempre possibile attivare una strategia del genere, qui enunciata nel caso in cui K è injectivo e continuo su un sottospazio compatto \bar{X} (si consulti [2] pag. 19).

Una riformulazione di questo tipo è però priva di interesse pratico poichè modifica l'insieme dei dati Y . Si supponga che i dati in ingresso abbiano una componente che non appartiene a $K(\bar{X})$. In questo caso il problema ristretto alla terna $(K, \bar{X}, K(\bar{X}))$ non ammette soluzioni. Un esempio tipico è quello in cui l'operatore K è un operatore di convoluzione. È noto che se il nucleo ha certe proprietà di regolarità, l'insieme $K(X)$ risulta essere formato da funzioni "liscie" (addirittura C^∞). Spesso i dati, affetti da errore, non appartengono alla stessa classe di funzioni C^∞ ; non è possibile quindi affrontare il problema rispetto alla terna $(K, X, K(X))$.

Consideriamo nuovamente le condizioni (a) e (b).

L'invertibilità di K è una richiesta molto restrittiva che in svariati problemi non viene rispettata. Alla luce dei risultati enunciati nel capitolo precedente, è possibile considerare la definizione 2.1 nel senso della soluzione generalizzata e dare così un significato più ampio e concreto al concetto di problema ben posto.

Al fine di operare con la pseudoinversa, d'ora in poi considereremo X, Y spazi di Hilbert e K operatore lineare e continuo; per utilizzare le stesse notazioni del capitolo 1, indicheremo l'operatore con il simbolo T .

La soluzione generalizzata x^1 se esiste è unica; l'inversa generalizzata T^1 garantisce quindi l'unicità, ma non l'esistenza della soluzione. Sappiamo infatti che, nel caso in cui $R(T)$ non è chiuso, si ha $D(T^1) = R(T) \oplus R(T)^\perp \neq Y$ (cfr. proposizione 1.1). Viceversa, nel caso in cui $R(T)$ è chiuso, oltre a fornire una unica soluzione $\forall y \in Y$, l'operatore lineare T^1 è anche limitato; l'applicazione $y \rightarrow x^1 = T^1 y$ è continua e soddisfa il punto (c) della definizione di problema ben posto. Abbiamo così dimostrato la seguente proposizione.

Proposizione 2.1 *Se T è un operatore lineare limitato tra spazi di Hilbert, il problema (2.1), nel senso della soluzione generalizzata, è ben posto se e solo se $R(T)$ è chiuso.*

Questo risultato è di notevole importanza per la classificazione dei problemi ben posti. Si consideri l'operatore di Fredholm di prima specie con range di dimensione non finita, nella forma

$$\int_a^b k(r, s) u(s) ds = g(r) \quad c \leq r \leq d \quad (2.2)$$

dove il nucleo $k : L^2([c, d] \times [a, b]) \rightarrow \mathbb{R}$ e $g : L^2([a, b]) \rightarrow \mathbb{R}$ sono funzioni assegnate, e $u : L^2([c, d]) \rightarrow \mathbb{R}$ è la soluzione da determinare.

Si osservi che tale operatore, essendo compatto, ha range non chiuso (infatti ogni operatore compatto possiede range chiuso se e solo se questo ha dimensione finita; si consulti [2], pag. 17 osservazione 3.1) e dà luogo a un problema mal posto.

Problemi legati all'esistenza della soluzione si hanno ad esempio nel caso in cui il nucleo assume la forma $k(r, s) = m(r)n(s)$ con $g(r) \neq C m(r)$ $C \in \mathbb{R}$

Inoltre, nel caso in cui la soluzione esiste, questa non è necessariamente unica; nell'esempio precedente, se u è una soluzione, qualunque altra funzione

$$u' = u + \bar{u} \quad \text{con} \quad \int_a^b n(s) \bar{u}(s) ds = 0$$

risolve l'equazione (2.2).

Più interessante è lo studio della dipendenza continua dai dati.

Il lemma di Riemann-Lebesgue (cfr. [6] pag. 159) asserisce che

$$\int_a^b k(r, s) \sin ns ds \rightarrow 0 \quad \text{per} \quad n \rightarrow +\infty.$$

Consideriamo quindi il dato perturbato

$$\bar{g}_n(r) = g(r) + A \int_a^b k(r, s) \sin ns ds \quad \text{con} \quad A \in \mathbb{R} \text{ costante arbitraria.}$$

Ovviamente al dato \bar{g}_n corrisponde la soluzione $\bar{u}_n(s) = u(s) + A \sin ns$.

Quindi, sebbene $\bar{g}_n \rightarrow g$, la differenza tra le soluzioni corrispondenti può essere molto grande.

L'operatore di Fredholm di prima specie (2.2) rientra nella categoria degli operatori il cui range è non chiuso e, come ci aspettavamo, dà luogo a un problema mal posto.

Lo studio dei problemi mal posti ha acquistato notevole importanza con la necessità di affrontare i cosiddetti problemi inversi.

Il concetto di problema inverso non ammette una caratterizzazione precisa. Supponiamo di aver fissato per "problema diretto" la determinazione di una certa soluzione avendo a disposizione un modello matematico e dati in ingresso. Uno scambio di ruolo tra ciò che intendiamo per dato in ingresso e soluzione in uscita fa nascere un nuovo problema che chiameremo "problema inverso". Questo esempio chiarirà i termini della questione.

Supponiamo di avere a disposizione un trasduttore che, posto a contatto con una sorgente, restituisca un opportuno valore numerico α . Considerando quest'ultimo come problema diretto, il problema inverso si pone nel modo seguente: "se α è il valore restituito dal trasduttore, in quale configurazione si trova la sorgente?".

Generalmente in tutti i casi in cui si ha una sorgente che dà luogo a un campo (particelle che emettono onde elettromagnetiche, una corda che vibrando produce onde sonore, ecc.), per problema inverso si intende la determinazione di informazioni sulle sorgenti avente a disposizione informazioni sul campo generato.

Generalmente i problemi inversi sono mal posti; vediamo perchè.

Consideriamo un modello matematico riassunto da un operatore lineare $T : X \rightarrow Y$ con X, Y spazi di Hilbert.

Per problema diretto intendiamo la determinazione, a partire dal dato $x \in X$, del risultato $y = T(x)$.

Supponiamo ulteriormente che l'operatore T sia compatto, con sistema singolare $\{\mu_n, u_n, v_n\}$ (cfr. appendice A), ossia

$$Tx = \sum_{n=1}^{\infty} \mu_n(x, u_n) v_n \quad \forall x \in X$$

Il problema inverso sarà la determinazione del valore x la cui immagine secondo T sia y . Utilizziamo l'espressione dell'inversa generalizzata calcolata nel capitolo precedente:

$$T^1 y = \sum_{n=1}^{\infty} \frac{1}{\mu_n} (y, v_n) \gamma u_n$$

Si può osservare che, mentre nel problema diretto i coefficienti $\mu_n \rightarrow 0$ per $n \rightarrow +\infty$, nel problema inverso i corrispondenti coefficienti μ_n^{-1} divergono.

Inoltre, rispetto al sistema ortonormale v_n , l'errore che perturba il dato y ha componenti principalmente localizzate in corrispondenza dei valori grandi di n (questo è causato dal fatto che l'errore è solitamente una funzione altamente oscillante e le funzioni singolari più oscillanti sono quelle corrispondenti a valori singolari piccoli; si pensi all'analogia con le armoniche di Fourier; per maggiori dettagli si consulti [2]). Si ottiene quindi che in un problema inverso le componenti dovute all'errore vengono moltiplicate per i coefficienti μ_n^{-1} divergenti, mentre in un problema diretto le eventuali componenti di errore vengono moltiplicate per valori sempre più piccoli.

È utile osservare che si incontrano problemi mal posti in tutte le branche della fisica e dell'ingegneria. Una serie di esempi di problemi inversi può essere trovata in [2].

Una generica metodologia che permetta di affrontare un problema non ben posto viene detta "di regolarizzazione".

Come suggerisce il nome, il suo scopo è quello di approssimare la soluzione cercando di fare in modo che essa soddisfi a certe proprietà di regolarità che hanno in genere le soluzioni corrispondenti a dati non perturbati. Una prima analisi di tali metodi viene fatta nel paragrafo seguente.

2.3 Metodi di regolarizzazione

Si consideri nuovamente l'equazione

$$Tx = y \quad (2.3)$$

con X, Y spazi di Hilbert $T \in \mathcal{B}(X, Y)$, $x \in X$, $y \in Y$.

Supponiamo inoltre $R(T)$ non chiuso; abbiamo già visto che l'equazione (2.3) dà così luogo ad un problema mal posto. Cerchiamo di risolvere l'equazione nel senso della soluzione generalizzata. Il teorema 1.9 - (ii) asserisce che T^1 è un operatore non limitato. Questa caratteristica comporta la mancata dipendenza continua dai dati.

L'idea che sta alla base dei metodi di regolarizzazione è di definire una famiglia di operatori R_α limitati che approssimino con continuità T^1 .

Definizione 2.2 Una famiglia di operatori $\{R_\alpha\}$ dipendenti dal parametro reale $\alpha > 0$, dove $R_\alpha : Y \rightarrow X$ $\forall \alpha > 0$ è detto algoritmo regolarizzante per il problema (2.3) nel senso della soluzione generalizzata se sono soddisfatte le condizioni seguenti:

- (i) R_α è lineare e continuo $\forall \alpha > 0$
- (ii) $\forall y \in R(T) \oplus R(T)^\perp$ $\lim_{\alpha \rightarrow 0} R_\alpha y = x^1$ dove $x^1 = T^1 y$

La definizione precedente considera esclusivamente operatori R_α lineari. Questa condizione semplifica parecchio lo studio successivo, ma non riveste carattere di necessità (esistono infatti metodi di regolarizzazione basati su famiglie di operatori non limitati). A tal proposito si osservi che nel capitolo successivo verrà enunciata una definizione di caratteripiù generale, che non richiede l'ipotesi di linearità (cfr. definizione 3.6).

Come si può vedere dalla definizione precedente, ogni algoritmo regolarizzante converge alla soluzione generalizzata qualora questa esista, ossia nel caso in cui $y \in D(T^1)$. È importante però sottolineare che il dominio di ogni operatore R_α è tutto l'insieme Y . Nel caso in cui il dato sia affetto da errore può succedere che y non appartenga a $D(T^1)$. In questo caso l'inversa generalizzata non ha alcuna utilità, mentre l'uso di un algoritmo regolarizzante permette di ottenere una soluzione approssimata del dato non perturbato.

L'idea centrale nella realizzazione pratica di un metodo di regolarizzazione, cioè nella definizione di come deve essere fatta la famiglia $\{R_\alpha\}$, è quella di utilizzare informazioni aggiuntive che permettano di gestire il problema mal posto in una maniera più "regolare". Queste informazioni vanno dedotte dalla formulazione del problema. Ad esempio si può fare in modo che le soluzioni approssimate soddisfino a certe proprietà tipiche delle soluzioni corrispondenti a dati "esatti" (si supponga di avere a disposizione l'ampiezza massima E delle soluzioni corrispondenti a dati non perturbati, $\|x^1\| \leq E$).

Nei paragrafi successivi verranno studiate alcune particolari famiglie regolarizzanti $\{R_\alpha\}$; qui procediamo ancora in maniera del tutto generale.

Utilizzando la rappresentazione spettrale dell'operatore T^*T (cfr. Appendice B), ogni algoritmo regolarizzante ammette una rappresentazione semplice.

Sappiamo che x^1 è una soluzione nel senso dei minimi quadrati, quindi $T^*Tx^1 = T^*y$.

Nel caso in cui T^*T è invertibile si ha chiaramente $x^1 = (T^*T)^{-1}T^*y$.

Sappiamo che un algoritmo regolarizzante R_α "approssima" la soluzione x^1 . Scriviamo allora

$$x^1_\alpha = \tilde{R}_\alpha(T^*T)^{-1}T^*y \quad (2.4)$$

Qui \tilde{R}_α è una funzione reale definita su $\sigma(T^*T) \subseteq [0, \|T\|^2]$, continua e $\tilde{R}_\alpha(T^*T)$ va letta nel contesto della risoluzione spettrale.

La famiglia di funzioni $\{\tilde{R}_\alpha\}$ ha il ruolo di approssimare, per $\alpha \rightarrow 0$, la funzione $f(t) = t^{-1}$. Il teorema seguente aiuta a chiarire il significato dell'operatore $\tilde{R}_\alpha(T^*T)$ e stabilisce condizioni sufficienti affinché $x^1_\alpha \rightarrow x^1$ per $\alpha \rightarrow 0$.

Il teorema è enunciato nel caso in cui l'operatore è compatto; per una formulazione generale si consulti [13] - teorema (3.2.2).

Teorema 2.1 Si consideri il problema (2.3), con T compatto. Sia $\{\tilde{R}_\alpha\}_{\alpha > 0}$ una famiglia di funzioni continue di variabile reale, definite su $[0, \|T\|^2]$. Si supponga inoltre che

$$(i) \tilde{R}_\alpha(t) \rightarrow t^{-1} \text{ per } \alpha \rightarrow 0, \forall t \in [0, \|T\|^2]$$

(ii) $(t\tilde{R}_\alpha(t))$ sia uniformemente limitato

Allora, se $y \in D(T^1)$, si ha

$$x^1_\alpha = \tilde{R}_\alpha(T^*T)^{-1}T^*y \rightarrow x^1 = T^1y \text{ per } \alpha \rightarrow 0 \quad (2.5)$$

Dimostrazione. Si osservi che, se P è un polinomio reale, si ha $P(T^*T)T^* = T^*P(T^*T)$; la verifica può essere fatta per sostituzione diretta. Una generalizzazione di questo risultato, dovuta a Weierstrass, assicura che questo è valido per ogni funzione reale e continua definita su $\sigma(T^*T) = \sigma(TT^*)$. Si ottiene

$$x^1_\alpha = \tilde{R}_\alpha(T^*T)^{-1}T^*y = T^*\tilde{R}_\alpha(TT^*)y \in R(T^*) \quad (2.6)$$

Possiamo allora sviluppare x^1_α rispetto al sistema singolare $\{\mu_n, u_n, v_n\}$ dell'operatore T .

$$x^1_\alpha = \sum_{n=1}^{\infty} \tilde{R}_\alpha(\mu_n^{-2})(T^*y, u_n)u_n = \sum_{n=1}^{\infty} \tilde{R}_\alpha(\mu_n^{-2})(y, T u_n)u_n$$

$$= \sum_{n=1}^{\infty} \tilde{R}_\alpha(\mu_n^2) \mu_n (y, v_n) u_n = \sum_{n=1}^{\infty} \mu_n^2 \tilde{R}_\alpha(\mu_n^2) \mu_n^{-1} (y, v_n) u_n$$

Le ipotesi (i), (ii) garantiscono che $\mu_n^2 \tilde{R}_\alpha(\mu_n^2) \rightarrow 1$, quindi, ricordando la rappresentazione di inversa generalizzata per operatori compatti vista nel capitolo precedente, si ottiene

$$x^1_\alpha \rightarrow \sum_{n=1}^{\infty} \mu_n^{-1} (y, v_n) u_n = x^1$$

Il teorema 2.1 garantisce la convergenza all'inversa generalizzata nel caso in cui $y \in D(T^1)$. Inoltre, per la continuità della funzione \tilde{R}_α , l'operatore $\tilde{R}_\alpha(T^*T)^*$ è lineare e continuo (cfr. appendice B). Si ha così che questo metodo, basato sulla famiglia di funzioni continue $\{\tilde{R}_\alpha\}_{\alpha>0}$ rientra nella categoria degli algoritmi regolarizzanti lineari. Si osservi che gli algoritmi che considereremo nei prossimi paragrafi ammettono una rappresentazione di questo tipo, pur di scegliere un'opportuna famiglia $\{\tilde{R}_\alpha\}$.
 Alla luce di quest'ultima considerazione, spesso indicheremo con la stessa notazione R_α , sia l'algoritmo regolarizzante della definizione 2.2, che l'operatore (2.4); risulterà dal contesto a cosa ci si riferisce.

Il significato della condizione (i) è stato accennato. La condizione (ii) è una condizione di regolarità sulla famiglia di funzioni. Vedremo nel prossimo paragrafo che la famiglia di funzioni che definisce la regolarizzazione di Tikhonov verifica queste due condizioni.

Abbiamo già detto che nella maggior parte dei casi di interesse pratico y non appartiene a $D(T^1)$. In questo caso la soluzione ottenuta mediante l'algoritmo regolarizzante diverge.

Teorema 2.2 Nelle ipotesi del teorema precedente, se y non appartiene a $D(T^1)$, allora, scelta qualsiasi successione $\{\alpha_n\}_{n \in \mathbb{N}}$, con $\alpha_n \rightarrow 0$ per $n \rightarrow +\infty$, si ha che $R_{\alpha_n}(T^*T)^* y$ non converge (puntualmente).

Dimostrazione. Sia P l'operatore di proiezione su $\overline{R(T^1)} = N(T^*)^\perp$. Supponiamo per assurdo che si abbia la convergenza, cioè $\exists z \in X$ t.c.

$$R_{\alpha_n}(T^*T)^* y = R_{\alpha_n}(T^*T)^* z \rightarrow z \quad \text{per } n \rightarrow +\infty.$$

Allora

$$T R_{\alpha_n}(T^*T)^* y = T T^* R_{\alpha_n}(T^*T)^* y = T T^* R_{\alpha_n}(T T^*) P y.$$

Ma per le ipotesi (i) e (ii) del teorema precedente si ha che $T T^* R_{\alpha_n}(T T^*) P y \rightarrow P y$, per $n \rightarrow +\infty$.

Quindi $P y = T z$, e necessariamente $y \in D(T^1)$, che è in contraddizione con l'ipotesi.

Corollario 2.1 Se y non appartiene a $D(T^1)$, allora

$$\lim_{\alpha \rightarrow 0} \|R_\alpha(T^*T)^* y\| = +\infty$$

Il corollario segue dal fatto che ogni successione non divergente (quindi limitata) in uno spazio di Hilbert ammette sottosuccessione convergente, mentre il teorema precedente assicura la divergenza ad ogni successione.

Gli ultimi due teoremi enunciati mostrano che la condizione $y \in D(T^1)$ è necessaria e sufficiente affinché si abbia la convergenza di un algoritmo di tipo (2.4). Si intuisce così

il comportamento e l'uso di un algoritmo regolarizzante nel caso di un dato perturbato $y = y' + y''$ con $y' \in D(T^1)$, $y'' \in Y \setminus D(T^1)$.

Utilizzare l'operatore R_α , con α piccolo, produce una approssimazione di T^1 troppo "fine", la quale, anche se riproduce bene la soluzione corrispondente al dato y' , comporta un grosso errore a causa della divergenza della componente dovuta a y'' .

Occorre quindi scegliere un opportuno parametro $\alpha > 0$, tale che l'operatore corrispondente approssimi in maniera soddisfacente T^1 , ma nello stesso tempo non risenta dell'errore sul dato. Questo discorso, ora accennato, verrà chiarito nel resto del capitolo e riassume il significato della regolarizzazione dei problemi mal posti.

2.4 Regolarizzazione di Tikhonov

Nella risoluzione di problemi mal posti, le strategie comunemente usate possono essere raggruppate in due categorie.

La prima consiste nel considerare unicamente soluzioni appartenenti a una particolare classe (ad esempio una sfera di raggio E). Tra queste soluzioni si cercherà quella che meglio approssima i dati, ossia quella che rende minima la distanza $\|Tx - y\|$.

La seconda metodologia è in un certo senso duale della precedente. Si tratta di cercare tra le soluzioni che approssimano il dato entro un margine di errore prefissato, la più "regolare" (ad esempio quella di norma minima, oppure, nel caso di spazi di funzioni, quella meno oscillante). Una particolare famiglia di operatori regolarizzanti proposta da Tikhonov (cfr. [42]) permette di gestire in maniera precisa entrambi i metodi.

Tikhonov cerca una soluzione che soddisfi entrambi i vincoli

$$(1) \|Tx - y\| \leq \epsilon$$

$$(2) \|x\| \leq E$$

minimizzando il funzionale, dipendente dal parametro α ,

$$\Phi_\alpha[x, y] = \|Tx - y\|_Y^2 + \alpha \|x\|_X^2 \quad (2.7)$$

Si osservi innanzitutto che non è sempre possibile ottenere soluzioni che soddisfino i due vincoli; si consideri ad esempio $X = Y$, $T = I$, $\|y\| = 3E$, $\epsilon < E$.

Lo studio del funzionale $\Phi_\alpha[x, y]$ permette di stabilire l'esistenza di una soluzione e , nel caso affermativo, di determinarla.

Teorema 2.3 Per ogni $y \in Y$, $\alpha > 0$ il funzionale $\Phi_\alpha[x, y]$ ammette uno e un solo punto di minimo $x_\alpha \in N(T)^\perp$, soluzione dell'equazione di Eulero

$$(T^*T + \alpha I)x_\alpha = T^*y \quad (2.8)$$

Dimostrazione. Supponiamo che x_α sia un punto di minimo di $\Phi_\alpha[x, y]$.

Allora $\forall w \in X$, ponendo $f(t) = \Phi_\alpha[x_\alpha + tw, y]$, si ha $f'(0) = 0$.

Si osservi che per la continuità e la linearità dell'operatore T , $f(t)$ risulta continua e derivabile. Inoltre

$$\begin{aligned} f(t) &= \Phi_\alpha[x_\alpha + tw, y] = \|Tx_\alpha - y + tTw\|_Y^2 + \alpha \|x_\alpha + tw\|_X^2 \\ &= \|Tx_\alpha - y\|_Y^2 + t^2 \|Tw\|_Y^2 + 2t(Tx_\alpha - y, Tw)_Y + \alpha \|x_\alpha\|_X^2 + t^2 \alpha \|w\|_X^2 + 2t\alpha(x_\alpha, w)_X \\ &= f(0) + 2t(T^*Tx_\alpha - T^*y + \alpha x_\alpha, w)_X + t^2 (\|Tw\|_Y^2 + \alpha \|w\|_X^2). \end{aligned}$$

Si ha $f'(0) = 0 \Leftrightarrow (T^*Tx_\alpha - T^*y + \alpha x_\alpha, w)_X = 0$.

L'unico elemento ortogonale ad ogni $w \in X$ è l'elemento nullo, quindi se x_α è punto di

minimo, allora
 $T^*T x_\alpha - T^*y + \alpha x_\alpha = 0$, equivalentemente,
 $(T^*T + \alpha I)x_\alpha = T^*y$.

Dimostriamo adesso che tale punto di minimo esiste sempre ed è unico. Ciò segue dal fatto che l'operatore $T^*T + \alpha I$ è definito positivo (i suoi autovalori sono contenuti nell'insieme $[\alpha, \|T\|^2 + \alpha]$) e quindi ha inverso continuo. Ne segue l'esistenza e l'unicità di $x_\alpha = (T^*T + \alpha I)^{-1}T^*y$.

Concludiamo dimostrando che tale soluzione x_α appartiene a $N(T)^\perp$. Come già visto nella dimostrazione del teorema 2.1, per ogni funzione reale e continua $h(t)$,
 $h(T^*T)T^*y = T^*h(TT^*y)$.

Ponendo $h(s) = (s + \alpha)^{-1}$, $s \in \mathbb{R}$, $s > 0$, si ottiene
 $(T^*T + \alpha I)^{-1}T^*y = T^*(TT^* + \alpha I)^{-1}T^*y$.
 Abbiamo così mostrato che $x_\alpha \in R(T^*) \subseteq N(T)^\perp$.

Nel prossimo paragrafo vedremo come viene trattato il problema dei vincoli esposto all'inizio, utilizzando il funzionale $\Phi_\alpha[x, y]$ (e quindi l'equazione di Eulero (2.8)). Dimostriamo adesso che l'equazione di Eulero, al variare di $\alpha > 0$, definisce un algoritmo regolarizzante e studiamone la convergenza.

Utilizzando le notazioni introdotte nel paragrafo precedente, consideriamo la funzione reale

$$\tilde{R}_\alpha(t) = (t + \alpha)^{-1} \quad (2.9)$$

Si osservi che

- (i) fissato $\alpha > 0$
 \tilde{R}_α è continua su $[0, +\infty)$
- (ii) $\tilde{R}_\alpha(t) \rightarrow t^{-1}$ ($\alpha \rightarrow 0$)
- (iii) $|\tilde{R}_\alpha(t)| = |\frac{t}{t+\alpha}|$ è uniformemente limitata, infatti

$$|\frac{t}{t+\alpha}| = \frac{1}{1+\frac{\alpha}{t}} \leq 1, \quad t > 0.$$

Allora la famiglia di operatori $\{\tilde{R}_\alpha\}_{\alpha>0}$, con

$$R_\alpha = \tilde{R}_\alpha(T^*T)T^* = (T^*T + \alpha I)^{-1}T^* \quad (2.10)$$

definisce un algoritmo regolarizzante lineare detto algoritmo regolarizzante di Tikhonov del primo ordine.

I risultati del paragrafo precedente sulla convergenza alla soluzione generalizzata garantiscono che

- (i) $x^1_\alpha = R_\alpha y \rightarrow x^1 = T^1y$, se $y \in R(T) \oplus R(T)^\perp$
- (ii) $\{x^1_\alpha\}$ diverge se y non appartiene a $R(T) \oplus R(T)^\perp$.

È possibile anche ottenere una maggiorazione sull'ordine di convergenza dell'algoritmo (2.10). A questo scopo si osservi che nel teorema 2.1, condizione sufficiente alla convergenza è che $y \in R(T) \oplus R(T)^\perp$, ossia $Py \in R(T)$, dove P è l'operatore di proiezione di Y su $\tilde{R}(T)$. È ragionevole pensare che più il dato soddisfi a certe condizioni di regolarità, più la successione $\{x^1_\alpha\}$ converga rapidamente alla soluzione.

Vediamo infatti che la condizione $Py \in R(T(T^*T)^\gamma)$ $0 < \gamma \leq 1$, permette di stimare il raggio di convergenza. Si osservi che gli insiemi $R(T(T^*T)^\gamma)$, al crescere di γ , formano una famiglia di insiemi decrescenti, ognuno costituito da elementi sempre più regolari.

Teorema 2.4 Se $\exists w \in X$ tale che $Py = T(T^*T)^\gamma w$, $0 < \gamma \leq 1$, allora

$$\|T^1y - x^1_\alpha\| \leq \alpha^\gamma \|w\|$$

dove $\{x^1_\alpha\}$ è la successione determinata con l'algoritmo di Tikhonov (2.10).

Prima di procedere, enunciamo due risultati (lemma 2.1 e 2.2) che saranno utilizzati nel corso della dimostrazione del teorema.

Lemma 2.1 $R((T^*T)^\gamma) \subseteq N(T)^\perp$

Dimostrazione. Per semplicità, verifichiamo la relazione solo nel caso in cui l'operatore T è compatto.

Utilizziamo il sistema singolare $\{\mu_n; u_n, v_n\}$ per l'operatore T ; sia $x \in X$.

Allora
 $(T^*T)^\gamma x = \sum_{n=1}^{\infty} \mu_n^{2\gamma} (x, u_n) X u_n$.

Inoltre $\{u_n\}$ è un sistema ortonormale completo per $\tilde{R}(T^*)$.

Quind $(T^*T)^\gamma x \in R(T^*) = N(T)^\perp$.

Lemma 2.2 Si consideri la famiglia di funzioni reali $\{\tilde{R}_\alpha\}$ (2.9),

$$\tilde{R}_\alpha(t) = (t + \alpha)^{-1}, \quad \alpha > 0.$$

Allora, fissato $\gamma \in (0, 1]$,

$$t^\gamma (1 - t\tilde{R}_\alpha(t)) < \alpha^\gamma$$

per ogni $t \in [0, +\infty)$.

Dimostrazione. Si osservi che, essendo $t \geq 0$ e $\alpha > 0$, segue

$$\begin{aligned} t^\gamma |1 - t\tilde{R}_\alpha(t)| &= t^\gamma \left| 1 - \frac{t}{t+\alpha} \right| \\ &= \frac{t^\gamma \alpha}{t+\alpha} \\ &= \frac{1}{t^{1-\gamma} \alpha^{-1} + t^{-\gamma}} \end{aligned}$$

Quindi:

- (i) Se $0 < t \leq \alpha$
 - (ii) Se $t > \alpha$
- $$\frac{1}{t^{1-\gamma} \alpha^{-1} + t^{-\gamma}} \leq \frac{1}{t^{-\gamma}} = t^\gamma \leq \alpha^\gamma$$
- $$\frac{1}{t^{1-\gamma} \alpha^{-1} + t^{-\gamma}} \leq \frac{1}{t^{1-\gamma} \alpha^{-1}} = t^{\gamma-1} \alpha \leq \alpha^{\gamma-1} \alpha = \alpha^\gamma$$

Possiamo ora dimostrare il teorema 2.4.

Dimostrazione. (teorema 2.4)

$$\begin{aligned} x^1_\alpha &= \tilde{R}_\alpha(T^*T)T^*y = \tilde{R}_\alpha(T^*T)T^*(Py + y - Py) \\ &= \tilde{R}_\alpha(T^*T)T^*Py = \tilde{R}_\alpha(T^*T)(T^*T)^\gamma w \end{aligned}$$

Si osservi che $TT^1y = Py = T(T^*T)^{\gamma}w$, quindi $T(T^1y - (T^*T)^{\gamma}w) = 0$ cioè $T^1y - (T^*T)^{\gamma}w \in N(T)$. Dal lemma 2.1, sappiamo che $(T^*T)^{\gamma}w \in N(T)^{\perp}$. Ricordando che anche $T^1y \in N(T)^{\perp}$, si ha $T^1y = (T^*T)^{\gamma}w$. Quindi

$$\begin{aligned} \|T^1y - x^1_{\alpha}\| &= \|(T^*T)^{\gamma}w - \tilde{R}_{\alpha}(T^*T)(T^*T)^{\gamma+1}w\| \\ &= \|(T^*T)^{\gamma}(I - (T^*T)\tilde{R}_{\alpha}(T^*T))w\| \\ &\leq \|(T^*T)^{\gamma}(I - (T^*T)\tilde{R}_{\alpha}(T^*T))\| \|w\| \\ &\leq \alpha^{\gamma}\|w\|. \end{aligned}$$

Per l'ultima maggioranza si è utilizzato il lemma 2.2 (si consideri come al solito la risoluzione spettrale dell'operatore T^*T , i cui autovalori sono contenuti in $[0, \|T\|^2] \subset [0, +\infty)$).

Il teorema ora dimostrato asserisce che l'algoritmo regolarizzante di Tikhonov, nel caso migliore contemplato, converge almeno con ordine α . Questo si ottiene se $\gamma = 1$, cioè se $T^1y \in R(T^*T)$.

Si può dimostrare che in realtà $O(\alpha)$ è il miglior raggio di convergenza ottenibile dall'algoritmo di Tikhonov e si ottiene necessariamente nel caso in cui $T^1y \in R(T^*T)$.

La convergenza dell'algoritmo di Tikhonov non è la proprietà più importante per la risoluzione di problemi inversi. Più che sapere se e come un algoritmo converge nel caso di un dato y , è importante conoscere le qualità dell'algoritmo nel caso in cui si abbiano a disposizione dati affetti da errore e quindi non conosciuti esattamente. Vedremo nel prossimo paragrafo che l'algoritmo di Tikhonov garantisce soluzioni affidabili (stabili e sufficientemente approssimate) anche in questo caso.

2.5 Soluzioni vincolate , approssimate regolarizzate e metodo della discrepanza

Si considerino le due funzioni introdotte nel paragrafo precedente

$$\begin{aligned} (1) \quad \alpha \rightarrow \epsilon(\alpha, y) &= \|Tx_{\alpha} - y\|_Y \\ (2) \quad \alpha \rightarrow E(\alpha, y) &= \|x_{\alpha}\|_X \end{aligned}$$

La prima distanza, detta discrepanza, permette di stimare l'accuratezza della soluzione x_{α} , ossia la precisione con cui approssima il dato; qualora l'inversa generalizzata esista, si ottengono valori piccoli di $\epsilon(\alpha, y)$ in corrispondenza di soluzioni prossime a T^1y . La seconda in un certo senso misura la "regolarità" della soluzione x_{α} ; se ad esempio $X = C^1(I)$ e $\|x\|_X^2 = \int (x'(t))^2 dt$, soluzioni poco oscillanti corrispondono a valori piccoli di $E(\alpha, y)$.

Nel caso dell'algoritmo di Tikhonov del primo ordine si possono avere utili indicazioni sull'andamento di queste due funzioni.

Teorema 2.5 Per ogni $y \in Y$, considerando l'algoritmo regolarizzante di Tikhonov 2.9, (i) la funzione $(\epsilon(\alpha, y))^2$ è monotona non decrescente a valori in $(\|y^{\perp}\|^2, \|y\|^2)$, dove $y^{\perp} = (I - P)y$ con P operatore di proiezione ortogonale di Y su $R(T)$,

(ii) la funzione $(E(\alpha, y))^2$ è monotona non crescente, a valori in $(0, \|T^1y\|^2)$ se $y \in R(T) \oplus R(T)^{\perp} = D(T^1)$, o a valori in $(0, +\infty)$ se $y \notin R(T) \oplus R(T)^{\perp}$,

Dimostrazione. Dalla definizione si ha

$$x_{\alpha} = R_{\alpha}y = (T^*T + \alpha I)^{-1}T^*y.$$

Si osservi che $R_{\alpha}y = R_{\alpha}P_{N(T^*T)^{\perp}}y = R_{\alpha}P_{R(T)^{\perp}}y$. Inoltre, come si è visto nel teorema 2.3, l'algoritmo di Tikhonov si può scrivere nella forma equivalente

$$x_{\alpha} = T^*(TT^* + \alpha I)^{-1}y$$

D'ora in poi spesso indicheremo con lo stesso simbolo sia il prodotto scalare in X che quello in Y , lo stesso dicasi per le norme. Possiamo allora scrivere:

$$\begin{aligned} E^2(\alpha, y) &= \|x_{\alpha}\|_X^2 = \|(T^*T + \alpha I)^{-1}T^*y\|_Y^2 \\ &= \|(T^*T + \alpha I)^{-1}T^*Py\|_Y^2 = \|T^*(TT^* + \alpha I)^{-1}Py\|_Y^2 \\ &= (TT^*(TT^* + \alpha I)^{-2}Py, Py) \end{aligned}$$

$$\begin{aligned} \epsilon^2(\alpha, y) &= \|Tx_{\alpha} - y\|_Y^2 = \|TR_{\alpha}y - y\|_Y^2 \\ &= \|TR_{\alpha}Py - Py - (I - P)y\|_Y^2 \\ &= \|(TR_{\alpha} - I)Py\|_Y^2 + \|(I - P)y\|_Y^2 \\ &= \|(TT^* - (TT^* + \alpha I)(TT^* + \alpha I)^{-1}Py)\|_Y^2 + \|y^{\perp}\|_Y^2 \\ &= \alpha^2\|(TT^* + \alpha I)^{-1}Py\|_Y^2 + \|y^{\perp}\|_Y^2 \end{aligned}$$

Indicando con $\{E_{\lambda}\}_{\lambda > 0}$ la famiglia spettrale associata all'operatore TT^* (cfr. appendice B) e ricordando che $\sigma(T^*T) \subset [0, \|T\|^2]$ (cfr. [15], teorema 3.3.1), si ottengono le relazioni seguenti

$$\begin{aligned} E^2(\alpha, y) &= \int_0^{\|T\|^2} \frac{\lambda}{(\lambda + \alpha)^2} d\|E_{\lambda}Py\|_Y^2 \\ \epsilon^2(\alpha, y) &= \alpha^2 \int_0^{\|T\|^2} \frac{1}{(\lambda + \alpha)^2} d\|E_{\lambda}Py\|_Y^2 + \|y^{\perp}\|_Y^2 \end{aligned}$$

Consideriamo le funzioni

$$\xi(\alpha, \lambda) = \frac{\lambda}{(\lambda + \alpha)^2} \leq \lambda^{-1}$$

$$\psi(\alpha, \lambda) = \frac{\alpha^2}{(\lambda + \alpha)^2} \leq 1$$

definite su $\{(0, +\infty) \times (0, \|T\|^2)\}$.

Fissato $\bar{\lambda} \in (0, \|T\|^2)$, $\xi(\alpha, \bar{\lambda})$ è monotona decrescente dal valore $\bar{\lambda}^{-1}$ per $\alpha = 0$, al valore 0 per $\alpha \rightarrow +\infty$, $\psi(\alpha, \bar{\lambda})$ è monotona crescente dal valore 0 per $\alpha = 0$, al valore 1 per $\alpha \rightarrow +\infty$.

Inoltre

$$(a) \int_0^{\|T\|^2} \lambda^{-1} d\|E_\lambda P y\|_Y^2 = \|T^{-1} y\|^2$$

$$(b) \int_0^{\|T\|^2} 1 d\|E_\lambda P y\|_Y^2 = \|P y\|^2$$

dove nel caso in cui $y \in R(T) \oplus R(T)^\perp$, $T^{-1} y$ rappresenta l'inversa generalizzata, mentre nel caso in cui y non appartiene a $R(T) \oplus R(T)^\perp$, con abuso di linguaggio indichiamo $\|T^{-1} y\| = +\infty$ (si confronti il corollario 2.1).

Utilizzando il teorema di convergenza dominata per le famiglie di funzioni reali $\Xi_\alpha(\lambda) = \xi(\alpha, \lambda)$ e $\Psi_\alpha(\lambda) = \psi(\alpha, \lambda)$, la monotonia rispetto a α delle funzioni $\xi(\alpha, \lambda)$ e $\psi(\alpha, \lambda)$ e i limiti (a), (b) si ottiene la tesi.

Rappresentiamo graficamente i risultati del teorema.

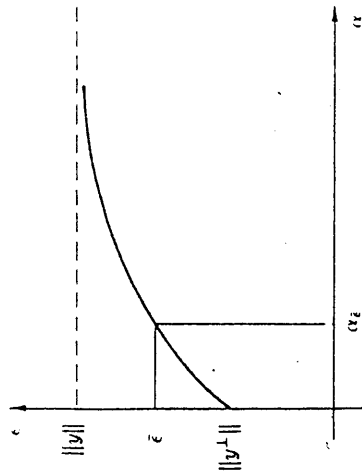
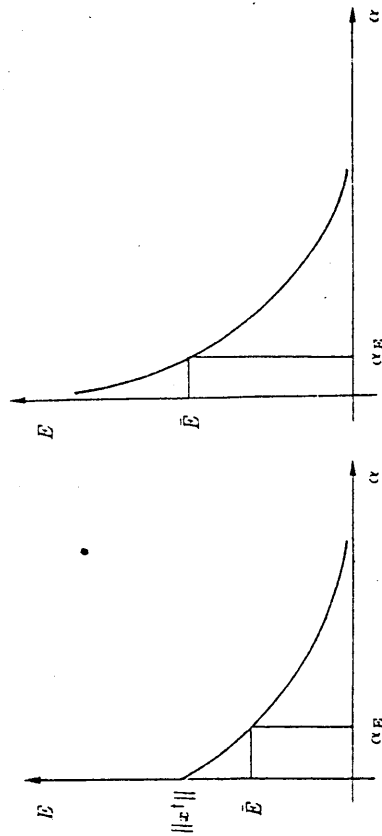


Figura A



$y \in R(T) \oplus R(T)^\perp$

$y \notin R(T) \oplus R(T)^\perp$

La monotonia delle funzioni ϵ e E rende particolarmente semplice il trattamento dei due problemi introdotti all'inizio del paragrafo precedente:

- (1) Tra le soluzioni x che approssimano il dato a livello ϵ (ossia $\|Tx - y\| \leq \epsilon$), determinare la più regolare ($\|x\|$ minimo)
- (2) Tra le soluzioni x soddisfacenti a certe condizioni di regolarità (ossia $\|x\| \leq E$), determinare quella che meglio approssima il dato ($\|Tx - y\|$ minimo).

Nel caso (1), supponendo $\|y^\perp\| \leq \epsilon \leq \|y\|$, i valori del parametro di regolarizzazione α per cui $\|Tx_\alpha - y\| \leq \epsilon$, appartengono all'insieme $I_1 = (0, \alpha_\epsilon)$ dove α_ϵ è il valore tale che $\|Tx_{\alpha_\epsilon} - y\| = \epsilon$, (vedi fig. A).
Nell'insieme I_1 , la soluzione più regolare si ottiene per $\alpha = \alpha_\epsilon$; (vedi fig. B).

Nel caso (2), supponendo $E < \|x^\perp\|$, i valori del parametro di regolarizzazione α per cui $\|x_\alpha\| \leq E$ appartengono all'insieme $I_2 = (\alpha_E, \infty)$, dove α_E è il valore tale che $\|x_{\alpha_E}\| = E$ (vedi fig. B).
Nell'insieme I_2 , la soluzione che meglio approssima il dato si ottiene per $\alpha = \alpha_E$, (vedi fig. A).

Considerando la osservazioni fin qui svolte, la risoluzione del problema (1) può essere affrontata nel seguente modo.

Per ogni valore del parametro di regolarizzazione $\alpha > 0$, utilizzando l'operatore $R_\alpha y$ (2.10), si calcoli il minimo z_α del funzionale $\Phi_\alpha[x, y]$. Si determini quindi il valore $\bar{\alpha}$ tale che

$$\|Tx_{\bar{\alpha}} - y\| = \epsilon$$

Questo valore esiste ed è unico se $\|y^\perp\| < \epsilon < \|y\|$.

La soluzione $z_{\bar{\alpha}}$ ottenuta è quella di norma minima tra tutte quelle ammesse.

Per quanto riguarda il problema (2) si procede nel seguente modo. Per ogni valore del parametro $\alpha > 0$, si calcoli come nel caso precedente la soluzione regolarizzata $R_\alpha y$. Si determini poi il valore $\bar{\alpha}$ tale che $\|z_{\bar{\alpha}}\| = E$. Questo valore esiste ed è unico se $E < \|x^\perp\|$. La soluzione $z_{\bar{\alpha}}$ ottenuta è quella che meglio approssima il dato tra quelle ammesse.

La soluzione del problema (1) è detta soluzione approssimata regolarizzata ed è la soluzione del problema variazionale

$$\|x_\alpha\|_X = \inf\{\|x\|_X : \|Tx - y\|_Y \leq \epsilon\}$$

La soluzione del problema (2) è detta soluzione vincolata a livello E ed è la soluzione del problema variazionale

$$\|Tx_\alpha - y\|_Y = \inf\{\|Tx - y\|_Y : \|x\|_X \leq E\}$$

Nel caso in cui vi siano indicazioni sia sull'errore che sulla norma della soluzione, cioè nel caso in cui esistano e siano ammesse soluzioni tali che $\|Tx - y\| \leq \epsilon$ e $\|x\| \leq E$, si può scegliere direttamente $\bar{\alpha} = (\frac{\epsilon}{E})^2$. Questa scelta, che va sotto il nome di metodo di Miller, ha valide basi teoriche. La soluzione che si ottiene è un buon compromesso tra approssimazione del dato e regolarità. Ovviamente tale soluzione ha una norma maggiore

di quella del problema (1) e una discrepanza maggiore di quella del problema (2). Per una ampia trattazione si può consultare [2] o [29].

È già stato più volte ricordato che nei problemi concreti il dato a disposizione è contaminato da una componente di errore detto noise dovuto a cause sistematiche (tra le quali la precisione con cui l'elaboratore maneggia i dati) e accidentali. In pratica non abbiamo mai a disposizione il dato $y \in Y$ ma una sua approssimazione e, nei casi migliori, una stima dell'errore massimo con cui viene fornito il dato. Nella parte restante di questo paragrafo vedremo come i risultati ottenuti per i problemi (1) e (2), i quali sono stati formulati nella ipotesi di conoscere "esattamente" il dato $y \in Y$, siano la base su cui lavorare per affrontare il problema concreto.

Indichiamo con y^δ una approssimazione del dato "esatto" $y \in R(T)$, e con y_{err}^δ la componente in $CR(T)$ di errore sul dato.

$$y^\delta = y + y_{err}^\delta$$

$$\|y - y^\delta\| = \|y_{err}^\delta\| \leq \delta$$

Si tenga presente che considereremo come dati non affetti da errore solo quelli in $R(T)$. Questa scelta viene fatta al fine di rendere più semplice i calcoli (non si dovrà fare continuamente ricorso all'operatore di proiezione su $R(T)$, altrimenti necessario), ma ha una motivazione più forte. Sappiamo che l'operatore T rappresenta una generica macchina che, ricevendo in ingresso un dato x , fornisce in uscita il corrispondente risultato y . Possiamo allora considerare un dato y^δ non appartenente a $R(T)$ come affetto da errore poiché non è tra quelli che la macchina, con un funzionamento corretto, è in grado di restituire.

Vediamo come agisce un generico algoritmo regolarizzante lineare in presenza di un dato perturbato y^δ .
Indicando con x^\dagger la soluzione corrispondente al dato $y \in R(T)$, si ottiene

$$x_\alpha = R_\alpha y^\delta = R_\alpha y + R_\alpha y_{err}^\delta$$

$$= R_\alpha T^\dagger x^\dagger + R_\alpha y_{err}^\delta$$

Valutiamo allora l'andamento della soluzione x_α in funzione del parametro α , confrontandola con la soluzione corrispondente al dato esatto x^\dagger .

$$x_\alpha - x^\dagger = R_\alpha y^\delta - x^\dagger = R_\alpha T x^\dagger + R_\alpha y_{err}^\delta - x^\dagger$$

$$= (R_\alpha T x^\dagger - x^\dagger) + R_\alpha y_{err}^\delta$$

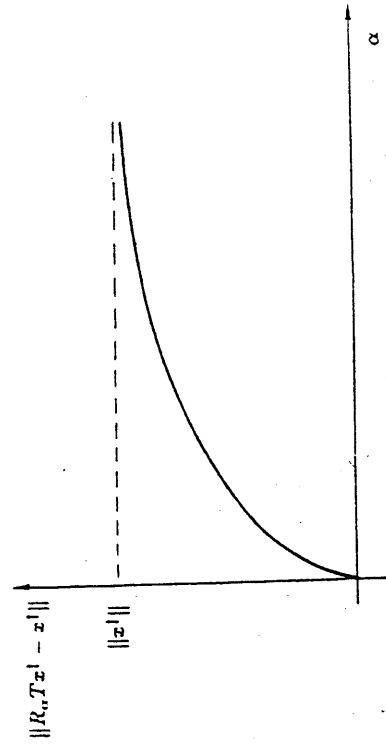
L'errore che si commette è la somma dei due termini $(R_\alpha T x^\dagger - x^\dagger)$ e $R_\alpha y_{err}^\delta$. Il primo rappresenta l'errore dovuto all'algoritmo regolarizzante R_α il quale approssima l'operatore T^\dagger (si osservi che questo errore è indipendente dall'errore sul dato). Chiamiamo questa componente errore di approssimazione.
Il secondo termine è dovuto alla perturbazione y_{err}^δ e rappresenta come l'errore sul dato si propaga sulla soluzione.

Studiamo, al variare di α , il comportamento dei due addendi. Analogamente a quanto è stato fatto nella dimostrazione precedente, indicando con $\{F_\lambda\}_{\lambda>0}$ la famiglia spettrale associata all'operatore $T^\dagger T$, si ottiene

$$\|R_\alpha T x^\dagger - x^\dagger\|^2 = \|(T^\dagger T + \alpha I)^{-1} T^\dagger T x^\dagger - x^\dagger\|^2 = \int_0^{\|T\|^2} \left(\frac{\lambda}{\lambda + \alpha} - 1 \right)^2 d\|F_\lambda x^\dagger\|^2$$

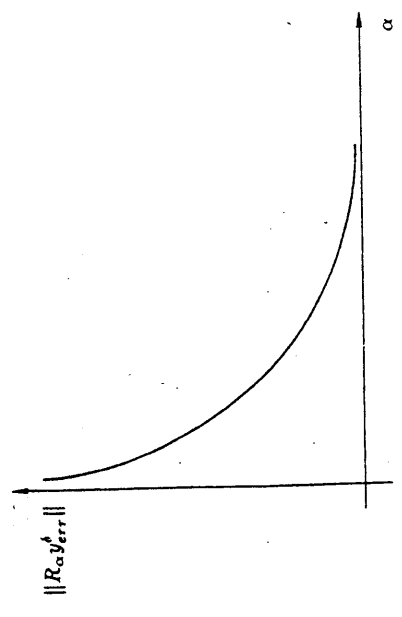
$$= \int_0^{\|T\|^2} \frac{\alpha^2}{(\lambda + \alpha)^2} d\|F_\lambda x^\dagger\|^2$$

Considerando nuovamente l'andamento della funzione reale $\frac{\alpha^2}{(\lambda + \alpha)^2}$ vista nel teorema precedente, si ha che l'errore di approssimazione è una funzione crescente, dal valore 0 per $\alpha = 0$ al valore $\|x^\dagger\|$ per $\alpha \rightarrow +\infty$



Errore di approssimazione

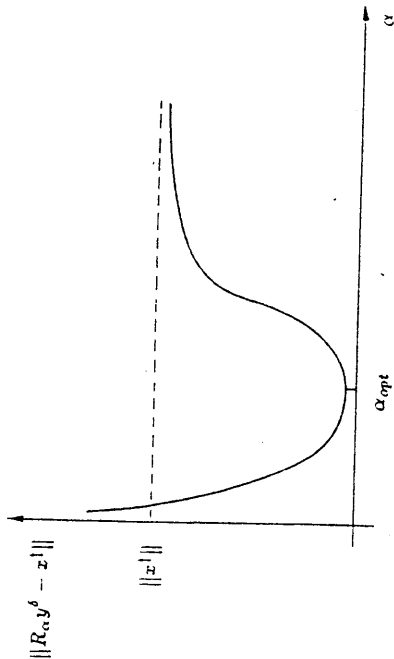
Per quanto riguarda l'errore generato dall'errore sul dato $R_\alpha y_{err}^\delta$, questo può essere visto come se fosse la soluzione ottenuta dal dato "non perturbato" y_{err}^δ . Si ottiene quindi un andamento di $\|R_\alpha y_{err}^\delta\|^2$ analogo a quello di $E = \|x_\alpha\|^2$ visto nel teorema 2.5, cioè decrescente dal valore $\|T^\dagger y_{err}^\delta\|^2$ (nella maggior parte dei casi $+\infty$ poiché solitamente y_{err}^δ non appartiene a $R(T) \ominus R(T)^\perp$) per $\alpha \rightarrow 0$, al valore 0 per $\alpha \rightarrow +\infty$.



Possiamo allora ritornare alla valutazione dell'errore globale

$$\|x_\alpha - x^d\|^2 \leq \|R_\alpha T x^d - x^d\|^2 + \|R_\alpha y^d - y^d\|^2$$

Si ottiene un andamento che può essere rappresentato nel seguente modo:



Esiste un parametro di regolarizzazione ottimale in corrispondenza del quale l'operatore R_α fornisce la migliore approssimazione della soluzione x^d . Il funzionamento di un algoritmo regolarizzante si può riassumere nel seguente modo:

» valori piccoli del parametro α riducono l'errore di approssimazione, mentre l'errore prodotto dalla perturbazione sul dato viene controllato da valori grandi di α ; il parametro ottimale, detto α_{opt} , rappresenta il miglior compromesso tra approssimazione e propagazione degli errori »

Dal punto di vista qualitativo, la situazione esposta nella figura precedente si verifica in svariati problemi numerici. Ad esempio nelle approssimazioni per discretizzazioni di problemi ben posti (equazioni differenziali, integrazione numerica, ecc.), esiste un valore ottimale del passo di discretizzazione; discretizzazioni più fini, sebbene possano far supporre una soluzione più precisa del problema, in realtà comportano un errore maggiore. Le cause di questo fenomeno sono gli errori di calcolo dovuti alla precisione limitata del calcolatore (numero di macchina). Ogni volta che la discretizzazione diventa più fine, la mole di calcoli aumenta come anche l'errore ad essi associato. Se si sceglie un passo di discretizzazione minore di quello ottimale, il miglioramento ottenuto dalla maggiore precisione con cui si maneggia il problema non compensa l'aumento dell'errore dovuto ai calcoli e la soluzione risulta peggiore.

Vi è però una sostanziale differenza tra il valore ottimale del parametro di regolarizzazione e il valore ottimale del passo di discretizzazione di un problema ben posto, infatti l'esistenza di α_{opt} non è dovuta a cause numeriche, ma unicamente alla natura del problema (che è mal posto) e al dato y^d (che è perturbato). Si osservi che una risoluzione simbolica di un problema ben posto per mezzo di discretizzazione comporta invece risultati migliori per passi sempre più piccoli.

Si intuisce che la risoluzione del problema è ricondotta alla determinazione del parametro di regolarizzazione ottimale.

Analizziamo un importante criterio per la scelta del parametro α proposto da Morosov, che va sotto il nome di **metodo della discrepanza** (cfr. [32]).

Il criterio di Morosov nasce dall'osservazione intuitiva che, conoscendo il dato a meno di δ , non è ragionevole cercare una soluzione che approssimi il dato con una precisione maggiore di δ .

Il metodo della discrepanza consiste appunto nella scelta del parametro $\bar{\alpha}$ tale che

$$\|Tx_{\bar{\alpha}} - y^d\| = \delta \tag{2.11}$$

Utilizzando i risultati del problema (1) considerato precedentemente, si possono trarre indicazioni sulla soluzione $x_{\bar{\alpha}}$ che rafforzano la validità del criterio.

Riconsideriamo il problema nella forma seguente:

» avendo a disposizione il dato approssimato y^d , si determini una soluzione approssimata di norma minima ».

In virtù di quanto detto circa la motivazione intuitiva, appare naturale cercare la soluzione di norma minima nell'insieme

$$C_\delta = \{x \in X : \|Tx - y^d\| \leq \delta\}$$

Il problema ora formulato è equivalente al problema (1) (si osservi che la tolleranza qui indicata con δ , nel problema (1) è indicata con ϵ). Dall'analisi svolta sul problema (1), si conclude che la soluzione di norma minima si ottiene in corrispondenza del valore del parametro di regolarizzazione $\bar{\alpha}$ tale che

$$\|Tx_{\bar{\alpha}} - y^d\| = \delta$$

La soluzione così ottenuta è quella di norma minima, ossia più regolare, tra quelle che consideriamo affidabili, cioè tra quelle che approssimano il dato con la stessa precisione con cui esso è fornito.

La proprietà di minimo della soluzione $x_{\bar{\alpha}}$, può essere anche verificata direttamente. Si ha il seguente risultato.

Lemma 2.3 Si consideri l'insieme chiuso e convesso

$$C_\delta = \{x \in X : \|Tx - y^d\| \leq \delta\}$$

e sia \bar{x} l'unico elemento di norma minima in C_δ .

Allora

$$\|T\bar{x} - y^d\| = \delta$$

cioè \bar{x} corrisponde alla soluzione $x_{\bar{\alpha}}$ ottenuta col metodo della discrepanza.

Dimostrazione. Supponiamo per assurdo che l'elemento di norma minima \bar{x} non appartenga al bordo di C_δ , cioè

$$\|T\bar{x} - y^d\| < \delta$$

Allora esiste $t \in (0, 1)$ tale che, per t sufficientemente prossimo a 1, si ha ancora

$$\|T(t\bar{x}) - y^d\| < \delta, \text{ quindi } t\bar{x} \in C_\delta.$$

Ma $\|t\bar{x}\| = t\|\bar{x}\| < \|\bar{x}\|$, in contraddizione con l'ipotesi.

Quindi necessariamente $\|T\bar{x} - y^d\| = \delta$.

La scelta del parametro $\bar{\alpha}$ basata sul criterio di Morozov nella pratica si realizza determinando la soluzione per una successione di valori positivi $\{\alpha_n\}_{n=1}^N$. Tra questi valori si sceglie quello tale che $\|Tx_{\alpha_n} - y^\delta\|$ sia più vicino possibile a δ .

La scelta è detta "a posteriori" poiché necessita della determinazione delle varie soluzioni x_{α_n} , ossia dipende dal particolare dato y^δ e non solo dal livello di errore δ . Si osservi che il metodo della discrepanza risulta particolarmente efficace nel caso di algoritmi regolarizzanti iterativi, poiché la soluzione regolarizzata ad ogni passo è immediatamente disponibile (questo discorso verrà chiarito nel prossimo capitolo).

Nel caso dell'algoritmo di Tikhonov è possibile avere una maggioranza "a priori", cioè indipendente dal particolare y^δ , del parametro $\bar{\alpha}$. Questa stima è dovuta a Vinokurov (si veda [48]).

Teorema 2.6 Si consideri l'algoritmo regolarizzante di Tikhonov (2.10) e sia $\bar{\alpha} > 0$ il parametro di regolarizzazione corrispondente al metodo della discrepanza. Si supponga inoltre $\delta < \|y^\delta\|$. Allora

$$\bar{\alpha} \leq \frac{\delta \|T\|^2}{\|y^\delta\| - \delta} \quad (2.12)$$

Dimostrazione.

$$\begin{aligned} \|y^\delta\| - \delta &= \|y^\delta\| - \|Tx_{\bar{\alpha}} - y^\delta\| \leq \|Tx_{\bar{\alpha}}\| \\ &= \|T(T^*T + \bar{\alpha}I)^{-1}T^*y^\delta\| \\ &= \|T\left(\frac{1}{\bar{\alpha}} - \frac{1}{T^*T + \bar{\alpha}I}\right)T^*y^\delta\| \\ &= \frac{1}{\bar{\alpha}} \|T[1 - T^*T(T^*T + \bar{\alpha}I)^{-1}]T^*y^\delta\| \\ &= \frac{1}{\bar{\alpha}} \|TT^*y^\delta - TT^*Tx_{\bar{\alpha}}\| \\ &= \frac{1}{\bar{\alpha}} \|TT^*(y^\delta - Tx_{\bar{\alpha}})\| \\ &\leq \frac{1}{\bar{\alpha}} \|T\|^2 \|y^\delta - Tx_{\bar{\alpha}}\| \\ &\leq \frac{1}{\bar{\alpha}} \|T\|^2 \delta \end{aligned}$$

Da cui

$$\bar{\alpha} \leq \frac{\delta \|T\|^2}{\|y^\delta\| - \delta}$$

Corollario 2.2 Si consideri la funzione $\alpha = \alpha(\delta)$ dove $\alpha(\delta)$ è il valore del parametro di regolarizzazione corrispondente al metodo della discrepanza, cioè

$$\alpha(\delta) = \{\alpha > 0 : \|Tx_{\alpha} - y^\delta\| = \delta\}$$

Allora

$$\alpha(\delta) \rightarrow 0^+ \quad , \quad \text{per } \delta \rightarrow 0^+$$

di ordine ≥ 1 .

La scelta di Morozov non è l'unica che offra buoni risultati. Un metodo alternativo è stato proposto da Arcangeli e consiste nella scelta del parametro α tale che

$$\sqrt{\alpha} \|Tx_{\alpha} - y^\delta\| = \delta$$

Anche il criterio di Arcangeli, come quello di Morozov, necessita della conoscenza del livello di errore δ . Per un'analisi delle proprietà di questa scelta si consulti [1].

2.6 Finestre spettrali e regolarizzazione di ordine superiore

L'algoritmo regolarizzante di Tikhonov del primo ordine (2.10) è definito come

$$R_{\alpha} = (T^*T + \alpha I)^{-1}T^*$$

Supponiamo per il momento che l'operatore T , oltre ad essere lineare e limitato, sia autoaggiunto. Indichiamo con $\{E_{\lambda}\}_{\lambda \geq 0}$ la famiglia spettrale associata all'operatore T . L'operatore R_{α} ammette la seguente rappresentazione :

$$\begin{aligned} R_{\alpha} &= \int_0^{\|T\|} \frac{\lambda}{\lambda^2 + \alpha} dE_{\lambda} \\ &= \int_0^{\|T\|} \frac{\lambda^2}{\lambda^2 + \alpha} \frac{1}{\lambda} dE_{\lambda} \\ &= \int_0^{\|T\|} W_{\alpha}(\lambda) \frac{1}{\lambda} dE_{\lambda} \end{aligned}$$

dove

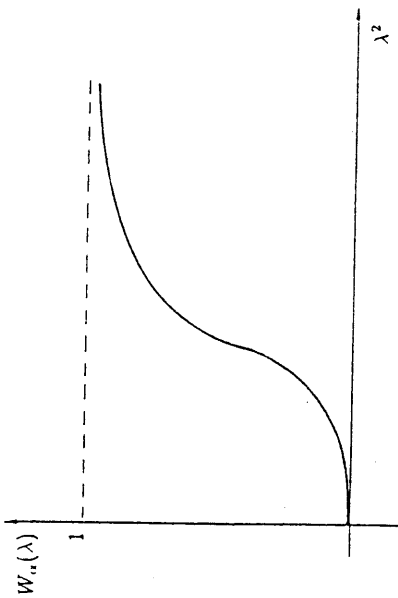
$$W_{\alpha}(\lambda) = \frac{\lambda^2}{\lambda^2 + \alpha}$$

Ricordando che qualora l'operatore autoaggiunto T sia invertibile, l'inversa è data da

$$T^{-1} = \int_0^{\|T\|} \frac{1}{\lambda} dE_{\lambda}$$

è facile notare che l'operatore di Tikhonov è una versione filtrata dell'operatore inverso T^{-1} . La famiglia di funzioni reali $\{W_{\alpha}\}_{\alpha > 0}$ definisce una famiglia di finestre spettrali con le seguenti proprietà

- $W_{\alpha}(\lambda) \simeq 1$ per $\lambda^2 \gg \alpha$
- $W_{\alpha}(\lambda) \simeq \frac{\lambda^2}{\alpha}$ per $\alpha \gg \lambda^2$
- $W_{\alpha}(\lambda) \rightarrow 0$ per $\alpha \rightarrow +\infty$



Essa riduce gli effetti delle componenti corrispondenti agli autovalori più piccoli. Queste componenti sono quelle che risentono maggiormente dell'errore sui dati (si confronti nel paragrafo 2.2 la proprietà analoga per i valori singolari μ_n corrispondenti a n grandi). Proprio in questo senso l'operatore R_α offre soluzioni stabilizzate, meno suscettibili cioè, da perturbazioni sul dato.

Si osservi che la richiesta di considerare T autoaggiunto è sorta al fine di poter utilizzare la rappresentazione spettrale, ma non riveste carattere di necessità. Ritorniamo al caso di generico operatore lineare limitato.

Si consideri l'operatore di Fredholm di prima specie a nucleo convolutivo

$$T : L^2(\mathbb{R}^n) \rightarrow L^2(\mathbb{R}^n), \quad k \in L^1(\mathbb{R}^n)$$

$$(Tx)(t) = \int_{\mathbb{R}^n} k(t-s)x(s)ds$$

e l'equazione

$$Tx = y \quad (2.13)$$

con $y \in L^2(\mathbb{R}^n)$.

In questo caso, utilizzando la teoria delle trasformate di Fourier, è possibile caratterizzare le finestre spettrali in una maniera più semplice. Sia $x \in L^2(\mathbb{R}^n)$. Possiamo scrivere

$$\begin{aligned} (Tx)(\xi) &= (k * x)(\xi) \hat{x}(\xi) \\ &= \int_{\mathbb{R}^n} k(t) e^{-2\pi i t \xi} dt \int_{\mathbb{R}^n} x(t) e^{-2\pi i t \xi} dt \\ (Tx)(t) &= \int_{\mathbb{R}^n} (\hat{Tx})(\xi) e^{+2\pi i t \xi} d\xi = \int_{\mathbb{R}^n} \hat{k}(\xi) \hat{x}(\xi) e^{+2\pi i t \xi} d\xi \end{aligned}$$

dove $\mathcal{A} \subset \mathbb{R}^n$ è la banda, che spesso considereremo limitata, della funzione $k(t)$. Inoltre, ricordando che $\bar{\bar{k}} = k$,

$$\begin{aligned} (T^*x)(t) &= \int_{\mathbb{R}^n} \overline{\hat{k}(\xi)} \hat{x}(\xi) e^{+2\pi i t \xi} d\xi \\ (T^*T)(t) &= (T^*(Tx))(t) \end{aligned}$$

$$\begin{aligned} &= \int_{\mathcal{A}} \overline{\hat{k}(\xi)} \hat{k}(\xi) \hat{x}(\xi) e^{+2\pi i t \xi} d\xi \\ &= \int_{\mathcal{A}} |\hat{k}(\xi)|^2 \hat{x}(\xi) e^{+2\pi i t \xi} d\xi \end{aligned}$$

Utilizzando la linearità dell'integrale, l'equazione di Eulero (2.8) assume la forma seguente

$$\int_{\mathcal{A}} (|\hat{k}(\xi)|^2 + \alpha) \hat{x}_\alpha(\xi) e^{+2\pi i t \xi} d\xi + \alpha \int_{CA} \hat{x}_\alpha(\xi) e^{+2\pi i t \xi} d\xi = \int_{\mathcal{A}} \overline{\hat{k}(\xi)} \hat{y}(\xi) e^{+2\pi i t \xi} d\xi \quad (2.14)$$

dove CA è l'insieme di frequenze fuori banda. Si ottiene

$$\hat{x}_\alpha(\xi) = \begin{cases} \frac{\overline{\hat{k}(\xi)} \hat{y}(\xi)}{|\hat{k}(\xi)|^2 + \alpha} & \xi \in \mathcal{A} \\ 0 & \xi \in CA \end{cases}$$

da cui

$$\begin{aligned} x_\alpha(t) &= \int_{\mathcal{A}} \frac{\overline{\hat{k}(\xi)} \hat{y}(\xi)}{|\hat{k}(\xi)|^2 + \alpha} e^{+2\pi i t \xi} d\xi \\ &= \int_{\mathcal{A}} \frac{|\hat{k}(\xi)|^2}{|\hat{k}(\xi)|^2 + \alpha} \frac{\hat{y}(\xi)}{\hat{k}(\xi)} e^{+2\pi i t \xi} d\xi \\ &= \int_{\mathcal{A}} \overline{W_\alpha(\xi)} \frac{\hat{y}(\xi)}{\hat{k}(\xi)} e^{+2\pi i t \xi} d\xi \end{aligned}$$

dove si è posto

$$\overline{W_\alpha(\xi)} = \frac{|\hat{k}(\xi)|^2}{|\hat{k}(\xi)|^2 + \alpha} \quad (2.15)$$

Si osservi che, qualora esista la soluzione generalizzata, questa assume la forma

$$x^!(t) = \int_{\mathcal{A}} \frac{\hat{y}(\xi)}{\hat{k}(\xi)} e^{+2\pi i t \xi} d\xi \quad (2.16)$$

Alla finestra spettrale $W_\alpha(\lambda)$ ottenuta utilizzando la risoluzione spettrale dell'operatore R_α , corrisponde il filtro in frequenza $\overline{W_\alpha(\xi)}$.

$\overline{W_\alpha(\xi)}$ ha le proprietà di un filtro "passa basso" ed elimina le armoniche di frequenza elevata, contaminate maggiormente da errore. La famiglia $\{\overline{W_\alpha(\xi)}\}_{\alpha > 0}$ caratterizza completamente l'algoritmo regolarizzante di Tikhonov per l'operatore di Fredholm.

Si possono costruire algoritmi regolarizzanti definendo opportunamente famiglie di filtri. Più precisamente:

Teorema 2.7 Si consideri il problema (2.13). Se la famiglia di funzioni $\{W_\alpha(\xi)\}_{\alpha > 0}$ soddisfa alle seguenti proprietà

- (a) $\forall \alpha > 0, W_\alpha$ è continua q.o.
- (b) $\forall \alpha > 0, 0 \leq W_\alpha(\xi) \leq 1$
- (c) $\lim_{\alpha \rightarrow 0} W_\alpha(\xi) = 1 \quad \forall \xi \in \mathcal{A}$

$$(d) \quad \forall \alpha > 0, \quad \frac{W_\alpha(\xi)}{k(\xi)} \in L^\infty(\mathbb{R}^n)$$

allora la famiglia di operatori definita da

$$(R_\alpha)(y) = \int_{\mathbb{A}} W_\alpha(\xi) \frac{\hat{y}(\xi)}{k(\xi)} e^{+2\pi i \xi t} d\xi$$

è un algoritmo regolarizzante per il problema (2.13).

Dimostrazione. Verifichiamo che la famiglia $\{W_\alpha(\xi)\}$ soddisfa le condizioni (i) - (ii) della definizione 2.2 di algoritmo regolarizzante.

Per (i), si osservi che R_α è lineare per le proprietà di linearità della trasformata di Fourier; inoltre è continuo:

infatti, utilizzando l'uguaglianza di Parseval,

$$\begin{aligned} \|R_\alpha y\|_Y^2 &= \int_{\mathbb{A}} \frac{W_\alpha(\xi)}{k(\xi)} |\hat{y}(\xi)|^2 d\xi \\ &\leq C(\alpha) \int_{\mathbb{A}} |\hat{y}(\xi)|^2 d\xi \\ &= C(\alpha) \|y\|_X^2 \end{aligned}$$

dove $|\frac{W_\alpha(\xi)}{k(\xi)}| \leq C(\alpha) < +\infty$ per (d).

Per quanto riguarda (ii)

$$\begin{aligned} \|R_\alpha y - x\|^2 &= \|R_\alpha T x - x\|^2 \\ &= \int_{\mathbb{A}} \frac{W_\alpha(\xi)}{k(\xi)} k(\xi) \hat{x}(\xi) - \hat{x}(\xi) \|^2 d\xi \\ &= \int_{\mathbb{A}} |W_\alpha(\xi) - 1|^2 |\hat{x}(\xi)|^2 d\xi \end{aligned}$$

Per (b) si ha

$$(W_\alpha(\xi) - 1)\hat{x}(\xi) \leq \hat{x}(\xi) \in L^2(\mathbb{R}^n)$$

allora, utilizzando il teorema di convergenza dominata e (c) si ottiene

$$\|R_\alpha y - x\|^2 \rightarrow 0, \quad \alpha \rightarrow 0, \quad \forall y \in L^2(\mathbb{R}^n).$$

La proprietà (d) è verificata ogni volta che $W_\alpha(\xi)$ tende a zero almeno come $k(\xi)$ o $\text{Supp}(W_\alpha(\xi))$ è contenuto in un intervallo limitato di \mathbb{R}^n , cioè ogni volta che la funzione antitrasformata di $W_\alpha(\xi)$ è a banda limitata. Ad esempio, nel caso in cui $k(\xi)$ non si annulli mai, la finestra rettangolare di semiampiezza $\frac{1}{\alpha}$ è il filtro più semplice che soddisfi (d) (si osservi però che tale filtro non offre buoni risultati a causa della forte discontinuità nei punti $\{-\frac{1}{\alpha}, +\frac{1}{\alpha}\}$).

In spazi di funzioni i metodi di regolarizzazione si caratterizzano quindi come particolari metodi di filtraggio in frequenza; in relazione alla particolare classe di dati in ingresso, la teoria del filtraggio suggerisce la scelta dell'algoritmo più opportuno.

Nel caso in cui l'operatore T sia compatto, ricorrendo al sistema singolare $\{\mu_n; u_n, v_n\}$, la rappresentazione spettrale di $R_\alpha y$ assume una forma più semplice. Con un procedimento analogo a quello sviluppato nel cap. 1.5 per ottenere l'inversa generalizzata

$$T^{-1}y = \sum_{n=1}^{\infty} \frac{1}{\mu_n} (y, v_n) \gamma u_n,$$

applicato a $R_\alpha = (T^*T + \alpha I)^{-1}T^*$, si ottiene

$$\begin{aligned} R_\alpha y &= \sum_{n=1}^{\infty} \frac{\mu_n}{\mu_n^2 + \alpha} (y, v_n) \gamma u_n \\ &= \sum_{n=1}^{\infty} \left(\frac{\mu_n^2}{\mu_n^2 + \alpha} \right) \frac{1}{\mu_n} (y, v_n) \gamma u_n \\ &= \sum_{n=1}^{\infty} W_{\alpha, n} \frac{1}{\mu_n} (y, v_n) \gamma u_n \end{aligned}$$

Analogamente all'equazione di Fredholm, si determina la famiglia di successioni filtro $\{W_{\alpha, n}\}_{\alpha > 0}$. Ogni successione ha un andamento qualitativo illustrato in figura (b).

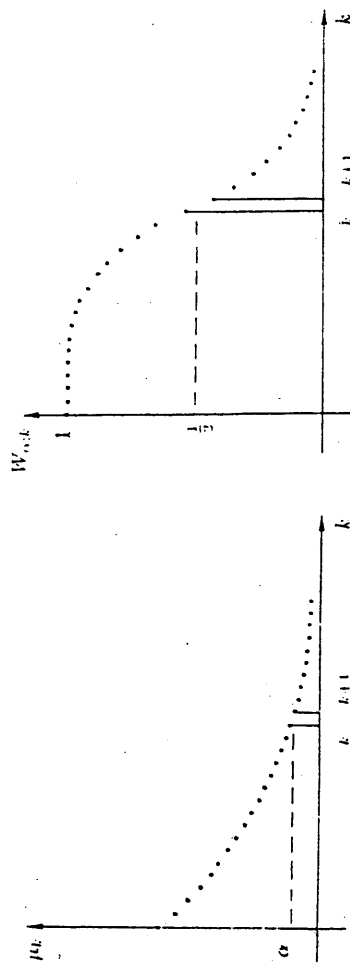


Figura (a)

Figura (b)

Per ogni fissato $\alpha > 0$, il filtro $W_{\alpha, n}$ ha la funzione di attenuare le componenti corrispondenti agli n grandi, le quali, come è già stato detto, risentono maggiormente delle perturbazioni sul dato (infatti queste componenti sono moltiplicate per $\mu_n^{-1} \rightarrow +\infty$). Ovviamente anche in questo caso è possibile definire successioni di filtri alternative a quella ottenuta con l'algoritmo di Tikhonov.

Le proprietà enunciate nel teorema 2.7 si caratterizzano nel modo seguente

- (i) $\forall \alpha > 0 \quad 0 \leq W_{\alpha, n} \leq 1$
- (ii) $\lim_{\alpha \rightarrow 0} W_{\alpha, n} = 1 \quad \forall n \in \mathbb{N}$
- (iii) $\lim_{n \rightarrow +\infty} W_{\alpha, n} = 0 \quad \forall \alpha > 0$

Ad esempio, il filtro

$$\tilde{W}_{\alpha, n} = \begin{cases} 1 & n \leq \frac{1}{\alpha} \\ 0 & n > \frac{1}{\alpha} \end{cases},$$

corrisponde nel caso continuo al filtraggio con finestra rettangolare e conduce alla nota decomposizione ai valori singolari troncata (Truncated S.V.D.) che verrà analizzata numericamente nel paragrafo 2.8.

2.7 Metodi iterativi di regolarizzazione

In alternativa alla regolarizzazione di Tikhonov, l'equazione (2.3) può essere risolta con metodi iterativi.

Per metodo iterativo intendiamo un algoritmo che determini approssimazioni successive della soluzione, ognuna delle quali calcolata per mezzo delle precedenti.

Semplificando le notazioni, considerando un arbitrario punto iniziale $x_0 \in X$, un algoritmo iterativo A , che possiamo immaginare come una mappa $X \rightarrow X$, genera una sequenza $\{x_k\}_{k \in \mathbb{N}} \subset X$ definita dalla relazione seguente

$$x_{k+1} = A(x_k)$$

Si osservi che x_{k+1} dipende in realtà da x_0, x_1, \dots, x_k .

La proprietà principale che deve avere un metodo iterativo è la convergenza alla soluzione esatta x^* , cioè $x_k \rightarrow x^*$ per $k \rightarrow +\infty$. In alcuni casi la soluzione esatta viene raggiunta necessariamente in un numero finito di passi; il metodo iterativo diventa così diretto, intendendo per diretto ogni metodo che permetta di ottenere la soluzione in un numero finito di operazioni.

Un'altra caratteristica che un metodo iterativo deve avere è la stazionarietà, cioè, supponendo che l'iterata n -esima risolva l'equazione (ossia $x_n = x^*$), si deve verificare che $x_m = x_n \quad \forall m > n$, o equivalentemente $A(x^*) = x^*$.

Importantissima è la velocità di convergenza. Un metodo iterativo che, sebbene converga alla soluzione, richieda una mole enorme di calcoli non è in pratica utilizzabile.

I metodi iterativi hanno acquistato notevole interesse con l'introduzione dei calcolatori. L'implementazione di tali metodi risulta particolarmente semplice e la velocità di calcolo in continua crescita offerta dalla tecnologia li rende sempre più efficienti. Si osservi inoltre che, dal punto di vista numerico, sono generalmente più stabili dei metodi diretti.

Per una breve introduzione sugli algoritmi iterativi e sulla varie definizioni riguardanti la velocità di convergenza (raggio di convergenza, convergenza lineare, superlineare, ecc.) si consulti [25]. Per un'analisi e un confronto di diversi metodi può essere utile consultare [36]

I metodi iterativi possiedono proprietà regolarizzanti; essi risultano così un valido strumento per il trattamento di problemi inversi. Vedremo come il numero di iterazioni gioca il ruolo del parametro di regolarizzazione α della definizione 2.2 di algoritmo regolarizzante. Più precisamente possiamo immaginare di definire la famiglia di operatori regolarizzanti $\{R_\alpha\}_{\alpha > 0}$ nel seguente modo

$$R_\alpha y = S_n(\alpha, y)$$

dove S_n è l'operatore ottenuto considerando le prime n iterazioni del metodo e l'applicazione $n(\alpha, y)$ rappresenta la relazione che determina il numero di iterazioni necessarie in funzione del dato y e del "livello di accuratezza" α . Si osservi che l'applicazione n dipende dal valore $y \in Y$; ciò significa che in genere non è possibile considerare una semplice funzione $n(\alpha)$, ossia il numero di iterazioni necessarie per ottenere un operatore R_α prossimo all'operatore inverso, dipende fortemente dal dato del problema.

Un metodo iterativo che per la sua semplicità permette una verifica delle proprietà di regolarizzazione è il metodo di Landweber-Fridman. L'analisi di tale algoritmo aiuta a capire in che modo un metodo iterativo regolarizza un problema mal posto.

Il metodo viene costruito iterando l'equazione ai minimi quadrati (iii) del teorema 1.1. In generale la ricerca delle radici di una generica equazione $f(x) = 0$ può essere ricondotta allo studio dei punti fissi dell'equazione $g(x) = x - f(x)$, cioè alla ricerca dei punti \bar{x} tali

L'algoritmo regolarizzante di Tikhonov è stato introdotto per mezzo del funzionale $\Phi_\alpha(x, y)$ (2.7), il quale stima la regolarità della soluzione per mezzo della norma $\|\cdot\|_X$. Consideriamo adesso un operatore lineare $L: X \rightarrow Y$ che permetta di controllare la regolarità della soluzione secondo criteri diversi. Ad esempio, nel caso in cui X sia uno spazio di funzioni, l'operatore L può rappresentare la derivata prima, ossia $Lx = x'$, in modo da controllare le oscillazioni della soluzione più che la sua ampiezza.

Per mezzo di tale operatore è possibile definire il seguente funzionale

$$\Phi_\alpha[x, y] = \|Tx - y\|_Y^2 + \alpha \|Lx\|_X^2$$

dove la stima della regolarità di x è demandata all'operatore L .

Un teorema del tutto analogo al 2.3 asserisce che, fissato α , il minimo del funzionale esiste e corrisponde alla soluzione dell'equazione di Eulero generalizzata

$$(T^*T + \alpha L^*L)x_\alpha = T^*y \quad (2.17)$$

Va precisato che l'unicità della soluzione è garantita solo nel caso in cui

$$N(T) \cap N(L) = \{0\}, \text{ relazione che è banalmente verificata nel caso dell'equazione (2.8)}$$

(dove $L = I$).

Nel caso in cui $N(T) \cap N(L) = \{0\}$, esiste un elemento $u \in X$, $u \neq 0$, tale che $Tu = Lu = 0$; allora se x_α è soluzione dell'equazione (2.17), ogni altro elemento $\tilde{x}_\alpha = x_\alpha + \beta u$, con $\beta \in \mathbb{R}$, risulta anch'esso soluzione.

Nel caso già considerato in cui $X = L^2(\mathbb{R}^n)$, oltre alla derivata prima vengono spesso utilizzate derivate successive o combinazioni lineari di esse. L'operatore L assume la forma

$$L = \sum_{i=0}^N a_i x^{(i)}$$

con $a_i \in \mathbb{R}$, $a_i \geq 0$ e il funzionale

$$\Phi_\alpha[x, y] = \|Tx - y\|_Y^2 + \sum_{i=0}^N a_i \|x^{(i)}\|^2$$

L'equazione di Eulero viene risolta nello spazio dove le derivate N -esime sono calcolabili; più precisamente

$$U \subset L^2(\mathbb{R}^n) \\ U = \{u: u^{(i)} \in C(\mathbb{R}^n) \cap L^2(\mathbb{R}^n) \quad i = 0, \dots, N-2, \quad u^{N-1} \text{ assolutamente continua}\}$$

I coefficienti a_i , che hanno la funzione di dare peso differente alle varie derivate, possono essere ulteriormente generalizzati considerando funzioni peso $a_i \in C^i(\mathbb{R}^n) \cap L^2(\mathbb{R}^n)$. L'algoritmo così costruito viene detto algoritmo regolarizzante di Tikhonov di ordine N .

Con semplici calcoli ([2]), si può dimostrare che all'algoritmo di Tikhonov di ordine N , corrisponde il filtro in frequenza seguente

$$W_\alpha(\xi) = \frac{|\hat{k}(\xi)|^2}{|\hat{k}(\xi)|^2 + \alpha \sum_{i=0}^N a_i |\xi|^{2i}}$$

che $\bar{x} = g(\bar{z})$.

La ricerca di tali punti può essere fatta per mezzo del seguente procedimento iterativo $x_{k+1} = g(x_k)$, che, se convergente, risulta essere banalmente consistente con l'equazione iniziale $f(x) = 0$ (si consulti [8]).

Applicando questo risultato all'equazione (iii) del teorema 1.1, si ottiene l'algoritmo di Landweber-Fridman:

$$x_{n+1} = x_n + \tau(T^n y - T^n T x_n)$$

dove τ è un valore reale, detto parametro di rilassamento, che permette di controllare la convergenza.

Con un semplice calcolo per induzione si verifica che il metodo assume la forma seguente

$$x_{n+1} = \tau [I + (I - \tau T^n T) + (I - \tau T^n T)^2 + \dots + (I - \tau T^n T)^n] T^n y + (I - \tau T^n T)^{n+1} x_0 \quad (2.18)$$

È immediato così constatare che l'algoritmo di Landweber-Fridman risulta lineare.

Proposizione 2.2 Il metodo iterativo lineare di Landweber-Fridman è un algoritmo lineare di regolarizzazione

Dimostrazione. Consideriamo per semplicità $x_0 = 0$. Si osservi innanzitutto che, in accordo con il teorema 2.1, opportunamente generalizzato al caso iterativo, la famiglia di funzioni reali continue $\{R_n\}_{n \in \mathbb{N}}$, definite da

$$\bar{R}_n(t) = \tau \sum_{j=1}^n (1 - \tau t)^j,$$

genera il metodo (il parametro $\alpha \rightarrow 0$ ora va interpretato come $n \rightarrow +\infty$). Una verifica delle condizioni (i) - (ii) del teorema 2.1 consente di affermare che il metodo è effettivamente un algoritmo di regolarizzazione.

Per quanto riguarda la condizione (i), si ha

$$\begin{aligned} \lambda \tau \sum_{n=1}^{\infty} (1 - \tau t)^n &= \tau \frac{1}{1 - (1 - \tau t)} = \frac{1}{\tau} \quad \forall t \in (0, \|T\|^2) \text{ se } \tau \in (0, \frac{2}{\|T\|^2}) \\ \text{Inoltre} \\ |\tau \sum_{n=1}^{\infty} (1 - \tau t)^n| &= \left| \frac{\tau t (1 - (1 - \tau t)^n)}{1 - (1 - \tau t)} \right| \\ &= \left| \frac{\tau t (1 - (1 - \tau t)^n)}{\tau} \right| \leq 1 \quad \forall t \in (0, \|T\|^2) \text{ se } \tau \in (0, \|T\|^{-2}) \end{aligned}$$

che verifica la condizione (ii) e che conclude la dimostrazione. ■

Come importante esempio, consideriamo l'equazione integrale a nucleo convolitivo (2.2). È interessante verificare che l'itema n-esima ha l'effetto di un filtro "passa basso" (cfr. paragrafo precedente). La relazione (2.18) diventa

$$\begin{aligned} x_{n+1}(\xi) &= (1 - \tau |\hat{k}(\xi)|^2)^{n+1} \hat{x}_0(\xi) \\ &\quad + \tau \left[1 + (1 - \tau |\hat{k}(\xi)|^2) + (1 - \tau |\hat{k}(\xi)|^2)^2 + \dots + (1 - \tau |\hat{k}(\xi)|^2)^n \right] \hat{k}(\xi) \hat{y}(\xi) \end{aligned}$$

Osservando che

$$\begin{aligned} 1 + (1 - \tau |\hat{k}(\xi)|^2) + (1 - \tau |\hat{k}(\xi)|^2)^2 + \dots + (1 - \tau |\hat{k}(\xi)|^2)^n &= \\ = \frac{1 - (1 - \tau |\hat{k}(\xi)|^2)^{n+1}}{1 - (1 - \tau |\hat{k}(\xi)|^2)} &= \frac{1 - (1 - \tau |\hat{k}(\xi)|^2)^{n+1}}{\tau |\hat{k}(\xi)|^2} \end{aligned}$$

si ha:

$$x_{n+1}(\xi) = (1 - \tau |\hat{k}(\xi)|^2)^{n+1} \hat{x}_0(\xi) + \left[1 - (1 - \tau |\hat{k}(\xi)|^2)^{n+1} \right] \frac{\hat{y}(\xi)}{\hat{k}(\xi)}$$

Possiamo così scrivere l'espressione dell'iterata n-esima:

$$\begin{aligned} x_{n+1}(t) &= \int_{\mathcal{A}} \hat{x}_0(\xi) e^{+2\pi i \xi t} d\xi + \\ &\quad \int_{\mathcal{A}} (1 - \tau |\hat{k}(\xi)|^2)^{n+1} \hat{x}_0(\xi) e^{+2\pi i \xi t} d\xi + \\ &\quad \int_{\mathcal{A}} \left[1 - (1 - \tau |\hat{k}(\xi)|^2)^{n+1} \right] \frac{\hat{y}(\xi)}{\hat{k}(\xi)} e^{+2\pi i \xi t} d\xi \end{aligned}$$

A parte i primi due integrali, che indichiamo con $U_n(x_0(\tau))$, nulli se $x_0(\tau) \equiv 0$, l'iterata (n+1)-esima può essere scritta nel modo seguente

$$(\hat{S}_{n+1} y)(t) = U_n(x_0(\tau))(t) + \int_{\mathcal{A}} W_{n+1}(\xi) \frac{\hat{y}(\xi)}{\hat{k}(\xi)} e^{+2\pi i \xi t} d\xi$$

dove

$$W_{n+1}(\xi) = 1 - (1 - \tau |\hat{k}(\xi)|^2)^{n+1}$$

è un filtro in frequenza dipendente da n, cioè dal numero di iterate effettuate.

Il filtro soddisfa alle condizioni del teorema 2.7, infatti:

- (a) $\forall n \in \mathbb{N} \quad W_n(\xi)$ è continua q.o.
 (b) $\forall n \in \mathbb{N} \quad 0 \leq W_n(\xi) \leq 1 \quad \forall \xi \in \mathcal{A}$,
 che è verificata se e solo se
 $0 \leq (1 - \tau |\hat{k}(\xi)|^2) \leq 1 \Leftrightarrow 0 \leq \tau \leq \frac{1}{|\hat{k}(\xi)|^2} \quad \forall \xi \in \mathcal{A}$
 (c) $\lim_{n \rightarrow +\infty} W_n(\xi) = 1 \quad \forall \xi \in \mathcal{A}$, condizione che è verificata se e solo se
 $|1 - \tau |\hat{k}(\xi)|^2| < 1$
 che dà luogo, ripercorrendo passaggi analoghi al punto (ii), a
 $0 < \tau < \frac{2}{|\hat{k}(\xi)|^2} \quad \forall \xi \in \mathcal{A}$

- (d) $\frac{W_n(\xi)}{\hat{k}(\xi)} \in L^\infty(\mathbb{R}^n)$

Il numeratore, per (b), è limitato; verificiamo che il rapporto è limitato nei punti in cui $\hat{k}(\xi) \rightarrow 0$.

Sviluppando secondo Taylor, possiamo scrivere che $W_n(\xi) \simeq \tau n |\hat{k}(\xi)|^2$ per $\xi \in \mathcal{A}$ dove $|\hat{k}(\xi)|$ è piccolo.
 Quindi $\frac{W_n(\xi)}{\hat{k}(\xi)} \simeq \tau n \hat{k}(\xi) \in L^\infty(\mathbb{R}^n)$.

La famiglia di funzioni $\{W_n(\xi)\}_{n \in \mathbb{N}}$ definite così un algoritmo regolarizzante dove il parametro di regolarizzazione è il numero di iterazioni.
 Le condizioni (b)-(c) imposte sul parametro di rilassamento τ forniscono $0 < \tau \leq |\hat{k}(\xi)|^{-2}$.
 La convergenza richiesta in (c) è assicurata anche per $\tau \in \left(|\hat{k}(\xi)|^{-2}, \frac{2}{|\hat{k}(\xi)|^2} \right)$; per questi valori la condizione (b) non è soddisfatta poiché $W_n(\xi) \in (0, 2)$. Si osservi che anche per questi $\tau \in \left(|\hat{k}(\xi)|^{-2}, \frac{2}{|\hat{k}(\xi)|^2} \right)$ si ottiene un filtro che, sebbene raddoppi l'ampiezza di alcune armoniche, converge alla funzione $f(t) \equiv 1$.
 Consideriamo quindi

$$\tau \in \left(0, \frac{2}{\|K\|_\infty^2} \right)$$

dove $\|K\|_\infty = \sup_{\xi \in \mathbb{R}_+} |\hat{k}(\xi)|$. Si ottiene

$$\lim_{n \rightarrow +\infty} x_n(t) = x^1(t) + x_0^{(0)}(t)$$

dove $x_0^{(0)}(t) = \int_{c_A} x(\xi) e^{+2\pi i \xi t} d\xi$ è la componente di x_0 parallela a $N(T)$.

Il limite è la soluzione ai minimi quadrati che ha la distanza minima da x_0 .

I risultati del paragrafo 2.5 sulla regolarizzazione in presenza di dati perturbati, ottenuti in corrispondenza dell'algoritmo di Tikhonov, possono essere facilmente generalizzati al caso iterativo.

Indichiamo nuovamente con y^δ una approssimazione a livello δ del dato esatto $y \in R(T)$ e con y_{err}^δ la componente di errore sul dato.

Allora:

$$x_n - x^1 = S_n y^\delta - x^1 = (S_n T x^1 - x^1) + S_n y_{err}^\delta$$

Analogamente a quanto già visto, la differenza tra la soluzione esatta e l'iterata n-esima è la somma dell'errore di approssimazione, dovuto esclusivamente al numero di iterazioni effettuate, e dell'errore dovuto alla perturbazione y_{err}^δ . Utilizzando l'uguaglianza di Parseval, nel caso in cui $x_0 = 0$, si ottiene

$$\begin{aligned} \|S_n T x^1 - x^1\|^2 &= \|\tau [I + (I - \tau T^* T) + (I - \tau T^* T)^2 + \dots + (I - \tau T^* T)^{n-1}] T^* T x^1 - x^1\|^2 \\ &= \int_{\lambda} |1 - (1 - \tau |\hat{k}(\xi)|^2)^n - 1|^2 \frac{|\hat{k}(\xi)|^2 |x^1(\xi)|^2}{|\hat{k}(\xi)|^2} d\xi \\ &= \int_{\lambda} |1 - \tau |\hat{k}(\xi)|^2|^{2n} |x^1(\xi)|^2 d\xi \end{aligned}$$

L'errore di approssimazione è una funzione decrescente di n dal valore $\|x^1\|^2$ per $n = 0$, al valore 0 per $n \rightarrow +\infty$ (cfr. figura A seguente).

L'errore dovuto alla perturbazione sul dato assume la forma seguente

$$\|S_n y_{err}^\delta\|^2 = \int_{\lambda} |1 - (1 - \tau |\hat{k}(\xi)|^2)^n|^2 \frac{|y_{err}^\delta(\xi)|^2}{|\hat{k}(\xi)|^2} d\xi$$

ed è una funzione di n crescente dal valore 0 per $n = 0$ ad un valore grande per $n \rightarrow +\infty$ (spesso $+\infty$ poiché y_{err}^δ solitamente non appartiene a $R(T) \oplus R(T)^\perp$, cfr. figura B seguente).

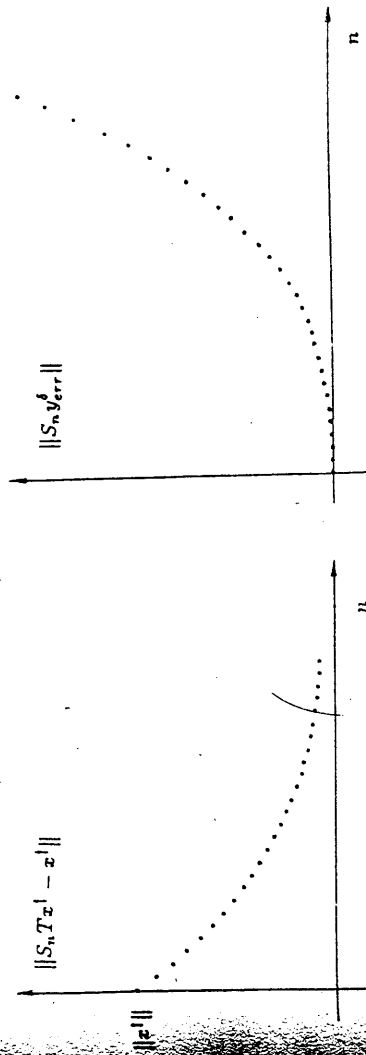


Figura A

L'errore globale $\|x_n - x^1\|^2$ ha le caratteristiche rappresentate in figura seguente:

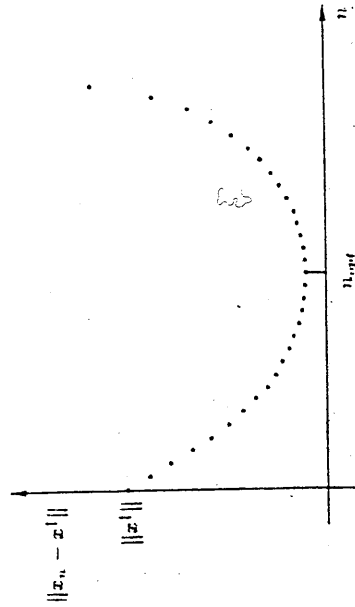


Figura B

Esiste un valore ottimale n_{opt} in corrispondenza del quale si ottengono i risultati migliori, e che rappresenta il miglior compromesso tra approssimazione dell'operatore inverso e propagazione degli errori.

Generalizziamo infine il criterio della discrepanza al caso iterativo.
 Si supponga di decomporre il dato perturbato y^δ nel seguente modo

$$y^\delta = y' + y^1$$

con $y' \in R(T)$ e $y^1 \in R(T)^\perp$.

Poiché consideriamo come dati esatti solo quelli appartenenti a $R(T)$, allora si ha

$$\|y - y^\delta\| \leq \|y - y'\| + \|y^1\| \leq \delta$$

che implica

$$\|y^1\| \leq \delta$$

Valutiamo la discrepanza in presenza del dato perturbato.

$$\begin{aligned} \epsilon^2(n, y^\delta) &= \|Tx_n - y^\delta\|^2 \\ &= \int_A |k(\xi)[1 - (1 - \tau k(\xi)^2)^n] \frac{y^\delta(\xi)}{k(\xi)} - y^\delta(\xi)|^2 d\xi + \int_{cA} |y^\delta(\xi)|^2 d\xi \\ &= \int_A |1 - \tau k(\xi)|^{2n} |y^\delta(\xi)|^2 d\xi + \int_{cA} |y^\delta(\xi)|^2 d\xi \end{aligned}$$

La discrepanza $\epsilon(n, y^\delta)$ è una funzione decrescente di n dal valore $\|y^\delta\|$ per $n = 0$, al valore $\|y^\delta\|$ per $n \rightarrow +\infty$.

Con analogia a quanto visto per l'algoritmo di Tikhonov nel paragrafo 2.5, è possibile definire una regola d'arresto delle iterazioni che chiamiamo ancora criterio della discrepanza. Essa si caratterizza nel seguente modo:

in presenza di dati perturbati a livello δ , accettiamo come soluzione quella ottenuta dalla prima iterata tale che

$$\|Tx_n - y\| \leq \delta$$

Poiché, come abbiamo visto, $\|x_n\|$ è una funzione crescente di n , ne segue che la soluzione determinata è quella di norma minima tra quelle che consideriamo affidabili, cioè tra quelle la cui discrepanza è minore o uguale alla precisione δ con cui si ha a disposizione il dato.

Si osservi che il criterio, nel caso dell'algoritmo di Landweber-Fridman, è ben definito poiché la discrepanza è monotonicamente decrescente, a $\|y^\delta\| \leq \delta$.

Spesso il criterio formulato conduce ad un numero di iterazioni troppo elevato. Nelle applicazioni pratiche si utilizza allora una costante positiva $C > 1$, la quale permette di fissare l'accuratezza della soluzione.

Il criterio della discrepanza risulta una efficace regola di arresto anche per algoritmi iterativi di regolarizzazione non lineari. Nel prossimo capitolo verrà ampiamente trattato in relazione al metodo del gradiente coniugato.

2.8 Trattamento numerico dei metodi di regolarizzazione

Nei paragrafi precedenti sono stati introdotti metodi analitici per problemi mal posti. La risoluzione numerica di tali problemi, che solitamente sono formulati in spazi continui di dimensione infinita, necessita di una discretizzazione, ossia di una approssimazione in spazi di dimensione finita. Il problema (2.3) si caratterizza quindi come sistema lineare

$$\bar{T}\bar{x} = \bar{y} \tag{2.19}$$

dove la matrice \bar{T} è l'approssimazione numerica dell'operatore T , la cui dimensione dipende dal tipo di discretizzazione adottata.

Ovviamente il problema (2.19) è ben posto, poiché formulato in spazi di dimensione finita (cfr. proposizione 2.1). Da un punto di vista numerico però, l'essere ben posto non è sufficiente a garantire stabilità e affidabilità della soluzione. La proprietà che invece caratterizza ogni problema del tipo (2.19) è il condizionamento della matrice \bar{T} :

un'equazione lineare (2.3) mal posta dà luogo a una matrice mal condizionata.

Si è tentati di ovviare al cattivo condizionamento della matrice, semplicemente utilizzando discretizzazioni più precise, ossia di dimensione maggiore. Questa strategia, nel caso

di problemi mal posti, non è adatta, poiché discretizzazioni più fini danno luogo a matrici inverse generalizzate sempre peggio condizionate.

Nel caso in cui T è compatto, si può fornire una semplice spiegazione di questo fenomeno. Sia $\{\mu_n; u_n, v_n\}$ il sistema singolare dell'operatore T (vedi appendice A).

Supponiamo di aver discretizzato T per mezzo di una matrice $\bar{T} \in M_p(\mathbb{R})$ invertibile; inoltre si consideri la decomposizione in valori singolari della matrice \bar{T} , ossia

$$\bar{T}\bar{x} = \sum_{i=1}^p \bar{\mu}_i(\bar{x}, \bar{u}_i)\bar{v}_i$$

con $\bar{\mu}_1 \geq \bar{\mu}_2 \geq \dots \geq \bar{\mu}_p > 0$.

Si verifica che i p valori singolari di \bar{T} risultano essere una approssimazione dei primi p valori singolari dell'operatore T . Una spiegazione euristica di questo fatto è la seguente: ricordando che i primi p valori singolari sono quelli corrispondenti alle componenti più significative e meno oscillanti, si può concludere che queste sono le più facili da approssimare, ossia sono quelle rispetto alle quali la discretizzazione adottata risulta più sensibile.

Si supponga ora di realizzare una discretizzazione più fine, che dia luogo a una matrice $\bar{T}' \in M_{p'}(\mathbb{R})$, $p' > p$. In questo caso i p' valori singolari di \bar{T}' sono una approssimazione dei primi p' valori singolari di T .

Consideriamo allora le rispettive inverse generalizzate. I valori singolari di \bar{T}' sono

$$\bar{\mu}'_1 \geq \bar{\mu}'_2 \geq \dots \geq \bar{\mu}'_{p-1}$$

mentre quelli di \bar{T} sono

$$\bar{\mu}^{-1}_1 \geq \bar{\mu}^{-1}_2 \geq \dots \geq \bar{\mu}^{-1}_{p-1}$$

Si ottiene

$$\text{cond}(\bar{T}') = \frac{\frac{1}{\bar{\mu}'_p}}{\frac{1}{\bar{\mu}'_1}} = \frac{\bar{\mu}_1}{\bar{\mu}_p}$$

e

$$\text{cond}(\bar{T}) = \frac{\frac{1}{\bar{\mu}_p}}{\frac{1}{\bar{\mu}_1}} = \frac{\bar{\mu}'_1}{\bar{\mu}'_p}$$

Considerando nuovamente le proprietà delle componenti associate ai primi valori singolari, si osservi che le due approssimazioni di μ_1 non differiscono di molto (la approssimazione ottenuta per mezzo della prima discretizzazione è già buona), mentre, per quanto riguarda i valori singolari più piccoli, in genere $\bar{\mu}'_p \gg \bar{\mu}_p$. Si è così ottenuto che

$$\text{cond}(\bar{T}') \gg \text{cond}(\bar{T})$$

Si osservi inoltre che $\text{cond}(\bar{T}) = \text{cond}(\bar{T}')$; l'instabilità in presenza di perturbazione sui dati si manifesta però maggiormente nel problema inverso per le ragioni già evidenziate nel paragrafo 2.2.

Prima di procedere con l'analisi numerica dei metodi di regolarizzazione, vediamo come discretizzare l'equazione di Fredholm di prima specie (2.2)

$$\int_a^b k(r, s)x(r)dr = y(s) \quad c \leq s \leq d$$

con $x \in L^2([a, b])$, $y \in L^2([c, d])$, $k \in L^2([a, b] \times [c, d])$.

L'idea più semplice, che va sotto il nome di metodo di collocazione, è richiedere che l'equazione sia verificata su un numero finito di punti $s_1, s_2, \dots, s_M \in [c, d]$, detti punti di collocazione. Si ottengono così M equazioni integrali

$$\int_a^b k(\tau, s_i)x(\tau)dr = y(s_i) \quad i = 1, \dots, M \tag{2.20}$$

Esaminiamo due diversi metodi per la risoluzione del sistema ottenuto.

(1) Le equazioni possono essere risolte numericamente per mezzo di qualsiasi formula di quadratura (cfr. [8]) con nodi $\tau_1, \tau_2, \dots, \tau_N \in [a, b]$ e pesi $\beta_1, \beta_2, \dots, \beta_N$, ossia

$$\int_a^b k(\tau, s_i)x(\tau)dr \approx \sum_{j=1}^N \beta_j k(\tau_j, s_i)x(\tau_j).$$

Si ottiene il seguente sistema lineare

$$\sum_{j=1}^N \beta_j k(\tau_j, s_i)x(\tau_j) = y(s_i) \quad i = 1, \dots, M$$

nell'incognita $(x(\tau_1), \dots, x(\tau_N))$, o equivalentemente,

$$\tilde{T}\underline{x} = \underline{y} \quad \text{dove } \tilde{T} \in M_{M,N}(\mathbb{R}) \tag{2.21}$$

$$\begin{aligned} \tilde{T}_{i,j} &= \beta_j k(\tau_j, s_i) \\ x_j &= x(\tau_j) \quad j = 1, \dots, N \\ y_i &= y(s_i) \quad i = 1, \dots, M. \end{aligned}$$

Il sistema (2.21), risolto nel senso della soluzione generalizzata con metodi di regolarizzazione numerici (che vedremo più avanti), dà la soluzione cercata. Si osservi che nel caso in cui $[a, b] = [c, d]$, si possono scegliere i punti di collocazione uguali ai nodi di integrazione; in questo caso si ottiene una matrice quadrata.

(2) Più efficiente è il metodo basato sulle funzioni di rappresentazione $\{K_i\}_{i=1}^M$, definite come

$$K_i \in L^2([a, b]) \quad K_i(\tau) = k(\tau, s_i)$$

Indichiamo con $\langle \bullet, \bullet \rangle$, $L^2([a, b]) \times L^2([a, b]) \rightarrow \mathbb{R}$ il prodotto scalare in $L^2([a, b])$ definito nel modo seguente

$$\langle a, b \rangle = \int_a^b a(\tau)b(\tau)dr$$

e con $L : L^2([a, b]) \rightarrow \mathbb{R}^M$ l'operatore lineare definito da

$$Lf = \begin{pmatrix} \langle K_1, f \rangle \\ \langle K_2, f \rangle \\ \dots \\ \langle K_M, f \rangle \end{pmatrix}$$

L'insieme di equazioni (2.20) può così essere espresso nella forma seguente

$$Lx = \underline{y} \tag{2.22}$$

dove \underline{y} è il vettore come definito in (1).

Dimostriamo che la soluzione generalizzata di (2.22), che si osservi appartiene a $L^2([a, b])$, è combinazione lineare delle funzioni K_i (da cui il nome "funzioni di rappresentazione").

Sappiamo che $x^1 \in N(L)^\perp = R(L^*) = R(L^*)$, infatti L^* è un operatore discreto. Valutiamo allora l'espressione dell'operatore L^* . Occorre determinare l'operatore L^* tale che

$$(Lf, \phi)_{\mathbb{R}^M} = \langle f, L^*\phi \rangle \quad \forall f \in L^2([a, b]) \quad \forall \phi \in \mathbb{R}^M$$

$$\text{Si ha } (Lf, \phi)_{\mathbb{R}^M} = \sum_{i=1}^M \phi_i \langle K_i, f \rangle = \langle \sum_{i=1}^M \phi_i K_i, f \rangle,$$

quindi

$$L^*\phi = \sum_{i=1}^M \phi_i K_i \quad \forall \phi \in \mathbb{R}^M$$

Ne segue che $x^1 \in L^2[a, b]$ ammette la rappresentazione

$$x^1 = \sum_{i=1}^M \alpha_i K_i = L^*\bar{\alpha}$$

per $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_M)'$ opportuno.

Il problema è così ricondotto alla determinazione del vettore dei coefficienti $\bar{\alpha}$.

Verifichiamo che questo è la soluzione del sistema lineare

$$G\bar{\alpha} = \underline{y} \tag{2.23}$$

dove $G \in M_{M,M}(\mathbb{R})$, detta matrice di Gram, è simmetrica semidefinita positiva con $(G)_{i,j} = \langle K_i, K_j \rangle$.

Sappiamo che $L^*\phi = \sum_{i=1}^M \phi_i K_i \quad \forall \phi \in \mathbb{R}^M$, pertanto

$$LL^*\phi = L\left(\sum_{i=1}^M \phi_i K_i\right) = \begin{bmatrix} \langle \sum_{i=1}^M \phi_i K_i, K_1 \rangle \\ \langle \sum_{i=1}^M \phi_i K_i, K_2 \rangle \\ \dots \\ \langle \sum_{i=1}^M \phi_i K_i, K_M \rangle \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^M \phi_i \langle K_i, K_1 \rangle \\ \sum_{i=1}^M \phi_i \langle K_i, K_2 \rangle \\ \dots \\ \sum_{i=1}^M \phi_i \langle K_i, K_M \rangle \end{bmatrix} = G\underline{\phi}$$

$$\Leftrightarrow G = LL^*$$

La matrice di Gram risulta autoaggiunta semidefinita positiva. Inoltre

$$Lx^1 = \underline{y} \Leftrightarrow LL^*\bar{\alpha} = \underline{y} \Leftrightarrow G\bar{\alpha} = \underline{y}$$

quindi effettivamente $\bar{\alpha}$ è soluzione del sistema (2.23).

La matrice G può essere decomposta nel prodotto

$$G = UDU^t$$

con $U \in M_{M,M}(\mathbb{R})$ matrice ortogonale le cui colonne sono gli autovettori u_1, u_2, \dots, u_M associati rispettivamente agli autovalori $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M \geq 0$ e $D \in M_{M,M}(\mathbb{R})$ matrice diagonale contenente gli autovalori, $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_M)$. Si ottiene (cfr. cap. 1.7)

$$\text{con } D^l = \begin{cases} \lambda_i & \text{se } \lambda_i > 0 \\ 0 & \text{se } \lambda_i = 0 \end{cases}$$

Esplicitando le componenti, la soluzione assume la forma seguente

$$\begin{aligned} x^l &= \sum_{i=1}^M \alpha_i K_i \\ &= \sum_{i=1}^M \left(\sum_{j=1, \lambda_j \neq 0}^M \frac{(u_j, y)}{\lambda_j} u_{i,j} \right) K_i \\ &= \sum_{j=1, \lambda_j \neq 0}^M \left(\frac{(u_j, y)}{\lambda_j} \right) \eta_j \end{aligned}$$

con

$$\eta_j = \sum_{i=1}^M u_{i,j} K_i$$

Si osservi che anche in questo caso per la determinazione della matrice di Gram, occorre utilizzare una formula di quadratura.

La soluzione ottenuta in funzione delle K_i può essere valutata nei punti $\{x_l\}_{l=1}^N$; si determina la soluzione numerica seguente

$$x_l = \sum_{j=1, \lambda_j \neq 0}^M \left(\frac{(u_j, y)}{\lambda_j} \right) \sum_{i=1}^M u_{i,j} k(\tau_i, s_i) \quad l = 1, \dots, N$$

Procediamo come accennato con l'illustrazione di alcuni metodi numerici per il trattamento di problemi mal posti, ossia metodi per la regolarizzazione del sistema lineare derivante dalla discretizzazione. È opportuno precisare che le considerazioni e i risultati della restante parte del paragrafo possono essere applicati semplicemente al problema algebrico di risoluzione di un sistema mal condizionato, senza necessariamente fare riferimento al problema continuo da cui esso deriva.

Indichiamo d'ora in poi la matrice \bar{T} semplicemente con T . Supponiamo T quadrata (eventualmente la si può rendere quadrata aggiungendo gli opportuni zeri). Riformuliamo il problema nel modo seguente:

" si risolva il sistema lineare $Tx = y$, dove la matrice $T \in M_n(\mathbb{R})$ è mal condizionata "

Una soluzione numerica diretta di tale problema darebbe luogo a soluzioni inaffidabili, spesso molto oscillanti rispetto alla soluzione analitica. Si ricordi inoltre che abbiamo a disposizione dati comunque affetti da errore (nel migliore dei casi va comunque contemplato il numero di macchina dell'elaboratore).

Lo scopo è ottenere una soluzione

$$z = \sum_{i=1}^k c_i z_i$$

dove il valore $k \leq n$ e il sistema di vettori $\{z_i\}_{i=1}^k$ va scelto in relazione al problema considerato e al metodo utilizzato.

Si osservi che, mentre il sistema ha dimensione n (generalmente molto grande), la soluzione che cerchiamo è combinazione lineare di $k \leq n$ vettori; in questo modo risolviamo il sistema in uno spazio di dimensione minore che garantisce però maggiore stabilità. Nei nostri esempi il sistema $\{z_i\}_{i=1}^k$ sarà sempre formato da vettori ortonormali.

L'efficienza del metodo dipende dalla scelta della base $\{z_i\}_{i=1}^k$; la base canonica $z_i = e_i = (0, \dots, 0, 1, 0, \dots, 0)$, con 1 all' i -esimo posto, comunemente utilizzata, spesso non offre i risultati migliori. Inoltre, come già osservato, una base può risultare adatta per la risoluzione di un sistema $Kz = y$, mentre può non esserlo per un altro sistema $K'z = y$.

L'opportunità di utilizzare una particolare base può essere "misurata" per mezzo del numero di condizionamento di T rispetto ai vettori $\{z_i\}_{i=1}^k$ che indicheremo con \mathcal{K} , e che è definito nel seguente modo

$$\mathcal{K} = \frac{\mu_1(T)}{\mu_k(TX_k)} \quad (2.24)$$

dove X_k è la matrice $n \times k$ formata dai vettori colonna $\{z_i\}_{i=1}^k$ e $\mu_k(T)$ è l' i -esimo dei valori singolari, ordinati in senso decrescente, della matrice T . Più il numero di condizionamento è piccolo, minore è la sensibilità della soluzione a perturbazioni sui dati, quindi migliore è la stabilità.

Fissata la base $\{z_i\}_{i=1}^k$, il numero di condizionamento in effetti dipende dal più piccolo valore singolare della matrice TX_k . \mathcal{K} può essere reso piccolo considerando un numero minore di vettori, cioè considerando la base $\{z_i\}_{i=1}^{k'}$ $k' < k$. Anche qui si presenta il compromesso tra approssimazione del dato e stabilità: se una base di cardinalità minore permette maggiore regolarità, occorre però controllare che la distanza $\|Tz - y\|$ sia entro la tolleranza richiesta.

Elenchiamo ora alcuni semplici metodi numerici; in seguito verrà brevemente affrontato il problema della stabilità in presenza di dati perturbati.

1 Fattorizzazione QR troncata (TQR)

La fattorizzazione QR, ossia decomposizione della matrice T nel prodotto di una matrice unitaria Q per una triangolare superiore R , è un procedimento comunemente usato per la risoluzione di sistemi lineari (cfr. [8]). Esso ha buone proprietà di stabilità e di costo. Nel caso in cui T sia mal condizionata, il metodo si caratterizza nel modo seguente.

Indicata nuovamente con X la matrice la cui i -esima colonna corrisponde al vettore i -esimo della base $\{z_i\}_{i=1}^n$, si determini la decomposizione QR di TX_m , dove X_m è la matrice $(n \times m)$ formata dalle prime m colonne di X , con $m \leq n$. In formule

$$TX_m = QR$$

con $Q \in M_{n,m}$ e $R \in M_{m,m}$.

Il parametro m corrisponde al numero massimo di vettori rispetto ai quali vogliamo sia espressa la soluzione $z = \sum_{i=1}^m c_i z_i$, e va scelto sulla base di considerazioni tra costo e affidabilità.

Si risolvano poi, con "sostituzione all'indietro", le prime $k \leq m$ equazioni (da cui il nome di fattorizzazione "troncata") del sistema

$$R\bar{z} = Q^t y$$

Otteniamo così la soluzione approssimata

$$\bar{z}^{(k)} = \sum_{i=1}^k c_i \bar{z}_i \quad (2.25)$$

Il valore k gioca il ruolo di parametro di regolarizzazione. Come si intuisce, esso deve essere il più piccolo tra quelli che permettono una buona approssimazione del dato, in modo da ottenere un numero di condizionamento \mathcal{K} minore possibile.

Il residuo $\|T\bar{z}^{(k)} - y\|_2$, che misura l'accuratezza della soluzione, può essere facilmente stimato calcolando la norma l^2 della ultime $(n - k)$ componenti di $Q^t y$. Un calcolo iterativo di tale norma per diversi valori di k , dà indicazioni sulla scelta del valore ottimale.

2 Regolarizzazione di Tikhonov (del primo ordine)

Analogamente a quanto visto nei paragrafi precedenti, la soluzione del sistema $T\bar{x} = y$ viene determinata risolvendo il problema di minimo

$$\min_{\bar{x}} (\|T\bar{x} - y\|_2^2 + \alpha \|\bar{x}\|_2^2)$$

dove $\|\bullet\|_2$ indica la norma l^2 , e il problema di minimo è equivalente all'equazione di Eulero

$$(T^t T + \alpha I)\bar{x} = T^t y$$

Si osservi che la soluzione dell'equazione di Eulero è la soluzione ai minimi quadrati del sistema sovradeterminato $(2n \times n)$

$$\begin{pmatrix} T \\ \alpha \end{pmatrix} \bar{x} = \begin{pmatrix} y \\ 0 \end{pmatrix}$$

Per la soluzione dell'equazione di Eulero può essere utilizzata la decomposizione ai valori singolari (S.V.D.) della matrice T (vedi [30]). Si determina così la seguente decomposizione

$$T = VDU^t$$

dove D è matrice diagonale contenente i valori singolari di T , ossia $d_{k,k} = \mu_k \geq 0$, e U, V sono ortogonali le cui colonne sono i corrispondenti vettori singolari \bar{u}_i, \bar{v}_i (si osservi che $\bar{u}_i = \mu_i^{-1} T \bar{v}_i$ e $\bar{v}_i = \mu_i^{-1} T^t \bar{u}_i$, come da definizione di sistema singolare - vedi appendice A).

Occorre quindi esprimere il vettore dei dati y come combinazione lineare dei vettori singolari \bar{u}_i

$$\bar{y} = \sum_{i=1}^n \beta_i \bar{u}_i$$

il vettore dei coefficienti $\bar{b} = (\beta_1, \beta_2, \dots, \beta_n)^t$ può essere determinato immediatamente facendo uso della matrice ortogonale V

$$\bar{b} = V^t \bar{y}$$

Si ottiene, posto $\bar{z}_\alpha = \sum_{i=1}^n c_i \bar{z}_i$,

$$(T^t T + \alpha I)\bar{z}_\alpha = T^t y \Leftrightarrow \sum_{i=1}^n c_i (T^t T + \alpha I)\bar{z}_i = \sum_{i=1}^n \beta_i T^t \bar{u}_i$$

$$\Leftrightarrow \sum_{i=1}^n c_i (\mu_i^2 + \alpha)\bar{z}_i = \sum_{i=1}^n \beta_i \mu_i \bar{z}_i$$

$$\Leftrightarrow c_i = \frac{\beta_i \mu_i}{\mu_i^2 + \alpha}$$

Si ottiene la soluzione

$$\bar{z}_\alpha = \sum_{i=1}^n \left(\frac{\beta_i}{\mu_i + \frac{\alpha}{\mu_i}} \right) \bar{z}_i \quad (2.26)$$

o equivalentemente, posto $\bar{z}_\alpha = (c_1, c_2, \dots, c_n)^t$ $c_i = \frac{\beta_i}{\mu_i + \frac{\alpha}{\mu_i}}$,

$$\bar{z}_\alpha = U \bar{c}_\alpha$$

Secondo la terminologia adottata precedentemente, il sistema $\{\bar{z}_i\}$ in questo caso corrisponde al sistema di vettori singolari $\{\bar{u}_i\}$, ossia $X = U$. È importante sottolineare che i valori β_i, μ_i e i vettori \bar{u}_i e \bar{v}_i $i = 1 \dots n$, sono indipendenti dal parametro di regolarizzazione α e quindi vengono calcolati una sola volta.

3 Decomposizione ai valori singolari troncata (T.S.V.D.)

Questo metodo è la versione numerica dell'omonimo sviluppato nel paragrafo 2.6. Si decomponga la matrice T come già visto al punto 2:

$$T = VDU^t$$

Ripetendo nuovamente il vettore \bar{y} in funzione del sistema $\{\bar{u}_i\}_{i=1}^n$, $\bar{y} = \sum_{i=1}^n \beta_i \bar{u}_i$,

e calcolando i coefficienti $\bar{b} = (\beta_1, \beta_2, \dots, \beta_n)^t$, con $\bar{b} = V^t \bar{y}$,

la soluzione del sistema $T\bar{x} = \bar{y}$ assume la forma $DU^t \bar{x} = V^t \bar{y} \Leftrightarrow U^t \bar{x} = D^{-1} V^t \bar{y}$, ossia

$$\bar{x} = \sum_{i=1}^n \left(\frac{\beta_i}{\mu_i} \right) \bar{u}_i$$

La regolarizzazione della soluzione si ottiene considerando uno sviluppo di \bar{x} troncato al k -esimo termine

$$\bar{x}^{(k)} = \sum_{i=1}^k \left(\frac{\beta_i}{\mu_i} \right) \bar{u}_i \quad (2.27)$$

Come per il metodo TQR , k assume il ruolo di parametro di regolarizzazione; se uno sviluppo con pochi termini offre buoni risultati, allora il metodo risulta particolarmente utile. Questo dipende, oltre che dalla matrice T , anche dal dato \bar{y} :

se i coefficienti β_i tendono a zero più velocemente dei μ_i , allora la differenza $\|\bar{x} - \bar{x}^{(k)}\|_2^2 = \sum_{i=k+1}^n \left(\frac{\beta_i}{\mu_i} \right)^2$ risulta piccola e il metodo efficace.

Si osservi che, nel caso in cui questa situazione non si verifichi, la T.S.V.D. non deve essere utilizzata, ma occorre cercare una soluzione rispetto a una base diversa da $\{\underline{u}_i\}$. Come per il metodo TQR, la ricerca del valore δ ottimale deve essere fatta iterativamente, ossia attraverso prove successive (utilizzando strategie di tipo dicotomico, per esempio).

4 Regolarizzazione di Tikhonov di ordine superiore

Come evidenziato nel paragrafo 2.6, si tratta di risolvere il seguente problema di minimo

$$\min_{\underline{x}} (\|T\underline{x} - \underline{y}\|_2^2 + \alpha \|L\underline{x}\|_2^2)$$

dove la matrice L rappresenta l'operatore che stima la regolarità della soluzione. Nel caso in cui la matrice T sia la discretizzazione dell'operatore di Fredholm di prima specie, L generalmente è la discretizzazione di un opportuno operatore differenziale, solitamente alle differenze finite. Ad esempio, si consideri il metodo di collocazione con reticolo formato da nodi equidistanti $\tau_j = a + j\delta \quad j = 0, \dots, N-1, \quad \delta = \frac{b-a}{N-1}$.

L'approssimazione dell'operatore $\tilde{L} = x'$ con differenze finite assume la forma

$$x''(\tau_j) = \frac{1}{\delta} [x(\tau_{j+1}) - x(\tau_j)] \quad j = 0, \dots, N-2$$

e l'approssimazione dell'operatore $\tilde{L} = x''$

$$\begin{aligned} x''(\tau_j) &= \frac{1}{\delta} [x'(\tau_{j+1}) - x'(\tau_j)] \\ &= \frac{1}{\delta} \left\{ \frac{1}{\delta} [x(\tau_{j+1}) - x(\tau_j)] - \frac{1}{\delta} [x(\tau_j) - x(\tau_{j-1})] \right\} \\ &= \frac{1}{\delta^2} \{ x(\tau_{j+1}) - 2x(\tau_j) + x(\tau_{j-1}) \} \quad j = 1, \dots, N-2 \end{aligned}$$

Questi, a meno di fattori di scala, assumono rispettivamente la forma

$$\begin{aligned} \tilde{L}x = x' \rightarrow L &= \begin{pmatrix} 1 & -1 & & & \\ & 1 & -1 & & \\ & & 1 & -1 & \\ & & & \dots & \dots \\ & & & & \dots & \dots \\ & & & & & 1 & -1 \end{pmatrix}_{(n-1) \times n} \\ \tilde{L}x = x'' \rightarrow L &= \begin{pmatrix} 1 & -2 & 1 & & & \\ & 1 & -2 & 1 & & \\ & & 1 & -2 & 1 & \\ & & & \dots & \dots & \dots \\ & & & & \dots & \dots \\ & & & & & 1 & -2 & 1 \end{pmatrix}_{(n-2) \times n} \end{aligned}$$

In generale è possibile utilizzare ordini di derivazione superiori, e combinazioni lineari delle rispettive matrici, ma queste non sempre comportano miglioramenti apprezzabili. Si osservi che, nel caso in cui L sia l'approssimazione della derivata p -esima, la matrice corrispondente è di dimensione $(n-p) \times n$ e rango $(n-p)$ (ad esempio la derivata 4ª ha una approssimazione alle differenze finite di dimensione $(n-4) \times n$ la cui prima riga è $(1, -4, 6, -4, 1, 0, \dots, 0)$).

Ritornando al problema di minimo enunciato all'inizio, analogamente a quanto fatto al punto 2, questo è equivalente all'equazione di Eulero modificata

$$(T^t T + \alpha L^t L)\underline{x} = T^t \underline{y}$$

la quale a sua volta è equivalente alla risoluzione ai minimi quadrati del sistema sovradeterminato $(2n \times n)$

$$\begin{pmatrix} T \\ \alpha L \end{pmatrix} \underline{x} = \begin{pmatrix} \underline{y} \\ 0 \end{pmatrix}$$

Indichiamo con K la matrice le cui colonne sono gli autovettori generalizzati di (T, L) , cioè i vettori \underline{k}_i soluzione del sistema

$$\lambda_i (T^t T)\underline{k}_i = (L^t L)\underline{k}_i \quad i = 1, \dots, n$$

(si osservi che $L^t L$ ha rango $(n-p)$, quindi necessariamente esistono $(n-p)$ autovettori generalizzati λ_i non nulli).

Per mezzo della matrice K è possibile determinare una particolare decomposizione per T e per L detta decomposizione ai valori singolari generalizzati (Generalised S.V.D.) (si consulti [45]). Si ha così

$$T = V D_a K^{-1} \quad L = U D_b K^{-1}$$

dove $V (n \times n)$ e $U ((n-p) \times (n-p))$ sono ortogonali,

$$D_a = \text{diag}(a_1, a_2, \dots, a_n),$$

$$D_b = \begin{pmatrix} b_1 & & & & \\ & b_2 & & & \\ & & \dots & & \\ & & & \dots & b_{n-p} \end{pmatrix}_{(n-p) \times n}$$

Inoltre, posto $b_i = 0 \quad i = n-p+1, \dots, n$, si può dimostrare che $\lambda_i = \frac{a_i}{b_i} \quad i = 1, \dots, n$. In questo modo le ultime p colonne di K formano una base per il nucleo di L .

Con l'ausilio della G.S.V.D., l'equazione di Eulero viene risolta nel seguente modo.

$$\begin{aligned} (T^t T + \alpha L^t L)\underline{x}_\alpha = T^t \underline{y} &\Leftrightarrow (K^{-1} D_a^2 K^{-1} + \alpha K^{-1} D_b^2 D_b K^{-1})\underline{x}_\alpha = K^{-1} D_a V^t \underline{y} \\ &\Leftrightarrow (D_a^2 + \alpha D_b^2) K^{-1} \underline{x}_\alpha = D_a V^t \underline{y} \end{aligned}$$

Posto $\underline{y} = \sum_{i=1}^n \beta_i \underline{v}_i$ con $\underline{b} = (\beta_1, \dots, \beta_n)$ e $\underline{b} = V^t \underline{y}$, si ottiene la soluzione

$$\underline{x}_\alpha = \sum_{i=1}^n \left(\frac{a_i \beta_i}{a_i^2 + \alpha b_i^2} \right) \underline{k}_i \quad (2.28)$$

Il problema è così ricondotto alla determinazione del parametro di regolarizzazione α più opportuno. Anche in questo caso, analogamente al punto 2, si osservi che

$i = 1, \dots, n$ sono indipendenti dal parametro α e vengono calcolati una sola volta. Si osservi inoltre che questo metodo rientra nella categoria vista all'inizio del paragrafo, ponendo $c_i = \frac{a_i d_i}{a_i^2 + \alpha b_i^2}$ e $z_i = k_i$. È importante sottolineare che il calcolo numerico della G.S.V.D. non è affatto semplice e, a tutt'oggi, non esistono pacchetti software appositamente dedicati; in pratica il metodo è difficilmente utilizzabile.

Sappiamo che l'obiettivo principale dei metodi esposti è rendere la soluzione poco sensibile a perturbazioni del dato (regolarizzazione). Si possono trarre importanti indicazioni analizzando il comportamento di tali metodi in presenza di dati affetti da errore. A tal proposito indichiamo con \underline{y} il dato non perturbato, con \underline{y}_j il dato a disposizione, e consideriamo $\|\underline{y} - \underline{y}_j\|_2 = \delta$.

1 T.Q.R.

Sia $\underline{z} = \sum_{i=1}^n d_i \underline{z}_i$ la soluzione corrispondente al dato esatto e $\underline{z}_j = \sum_{i=1}^k d_i \underline{z}_i$ la soluzione corrispondente al dato perturbato, troncata al k -esimo termine. Poiché la base $\{\underline{z}_i\}_{i=1}^n$ è scelta ortogonale, si ha

$$\|\underline{z} - \underline{z}_j\|_2^2 = \sum_{i=1}^k (d_i - \hat{d}_i)^2 + \sum_{i=k+1}^n d_i^2$$

Si osservi che, anche in assenza di perturbazioni, la soluzione troncata non è "buona" nel caso in cui i primi k vettori non riescano ad approssimare sufficientemente la soluzione (cioè, per $i > k$, i coefficienti d_i non siano piccoli). Si può ulteriormente migliorare la prima sommatoria.

Si consideri la matrice R della decomposizione $TX = QR$, partizionata nel seguente modo

$$R = \begin{pmatrix} R_1 & S \\ 0 & R_2 \end{pmatrix} \text{ con } R_1 \text{ di dimensione } (k \times k) \text{ e } R_2 \text{ di dimensione } ((n-k) \times (n-k)).$$

Posto $\underline{d} = (d_1, \dots, d_n)$, $\hat{\underline{d}} = (\hat{d}_1, \dots, \hat{d}_k)$, si indichi con \underline{d}_1 il vettore formato dalle prime k componenti di \underline{d} , con \underline{d}_2 il vettore formato dalle seguenti $(n-k)$ componenti, con Q_1 la matrice formata dalle prime k colonne di Q , e con Q_2 la matrice formata dalle ultime $(n-k)$ colonne.

Si ha così

$$Q \begin{pmatrix} R_1 & S \\ 0 & R_2 \end{pmatrix} \begin{pmatrix} \underline{d}_1 \\ \underline{d}_2 \end{pmatrix} = \underline{y} \Leftrightarrow \begin{pmatrix} R_1 & S \\ 0 & R_2 \end{pmatrix} \begin{pmatrix} \underline{d}_1 \\ \underline{d}_2 \end{pmatrix} = \begin{pmatrix} Q_1^T \underline{y} \\ Q_2^T \underline{y} \end{pmatrix}$$

Si ottiene

$$Q_1 R_1 \hat{\underline{d}} = \underline{y}_j \Leftrightarrow R_1 \hat{\underline{d}} = Q_1^T \underline{y}_j$$

$$R_1 (\underline{d}_1 - \hat{\underline{d}}) = Q_1^T \underline{y} - Q_1^T \underline{y}_j = S \underline{d}_2$$

$$\|\underline{d}_1 - \hat{\underline{d}}\|_2 \leq \|R_1^{-1}\|_2 (\|Q_1^T\|_2 \|\underline{y} - \underline{y}_j\|_2 + \|S\|_2 \|\underline{d}_2\|_2)$$

$$= \|R_1^{-1}\|_2 (\delta + \|S\|_2 \|\underline{d}_2\|_2)$$

Dall'ultima espressione si può notare che, sebbene k crescenti riducano $\|\underline{d}_2\|_2$, la differenza $\|\underline{d}_1 - \hat{\underline{d}}\|_2$ può crescere molto; infatti, a causa del peggior condizionamento di $(KX)_k$, il termine $\|R_1^{-1}\|_2$ può diventare molto grande. Questo risultato è coerente con quanto detto nella prima parte del paragrafo circa il significato di k .

2 Regolarizzazione di Tikhonov

Consideriamo l'algoritmo sviluppato con l'ausilio della S.V.D.; in questo caso la soluzione corrispondente al dato non perturbato (senza regolarizzazione) è data da

$$\underline{z} = \sum_{i=1}^n \left(\frac{\hat{\beta}_i}{\mu_i}\right) \underline{u}_i \text{ dove } \underline{y} = \sum_{i=1}^n \beta_i \underline{v}_i.$$

La soluzione regolarizzata in presenza del dato perturbato è invece data da

$$\underline{z}_j = \sum_{i=1}^n \left(\frac{\hat{\beta}_i}{\mu_i + \alpha}\right) \underline{u}_i \text{ dove } \underline{y}_j = \sum_{i=1}^n \hat{\beta}_i \underline{v}_i.$$

Possiamo così scrivere

$$\underline{z} - \underline{z}_j = \sum_{i=1}^n \left(\frac{(\beta_i - \hat{\beta}_i) \mu_i + \alpha \frac{\beta_i}{\mu_i}}{\mu_i^2 + \alpha} \right) \underline{u}_i$$

Anche in questo caso, per $\alpha \rightarrow 0$, l'errore cresce a causa delle ultime componenti \underline{u}_i che hanno i valori singolari μ_i sempre più piccoli (si osservi che l'unica informazione di cui disponiamo è $|\beta_i - \hat{\beta}_i| < \delta$), mentre per α grande l'errore cresce a causa del termine $\frac{\beta_i}{\mu_i}$ al numeratore. La determinazione del valore di α che rende minima la distanza $\|\underline{z} - \underline{z}_j\|_2$ non è semplice.

3 T.S.V.D.

Utilizzando le notazioni già introdotte, la soluzione corrispondente al dato perturbato troncata al k -esimo termine, è data da $\underline{z}_j = \sum_{i=1}^k \frac{\hat{\beta}_i}{\mu_i} \underline{u}_i$, e l'errore

$$\begin{aligned} \underline{z} - \underline{z}_j &= \sum_{i=1}^k \left(\frac{\beta_i - \hat{\beta}_i}{\mu_i}\right) \underline{u}_i + \sum_{i=k+1}^n \frac{\beta_i}{\mu_i} \underline{u}_i \\ \|\underline{z} - \underline{z}_j\|_2^2 &\leq \delta^2 \sum_{i=1}^k \left(\frac{1}{\mu_i^2}\right) + \sum_{i=k+1}^n \left(\frac{\beta_i}{\mu_i}\right)^2 \\ &= \delta^2 M(k) + N(k) \end{aligned}$$

Si constata nuovamente la presenza di un termine $M(k)$ crescente al crescere di k e di un termine $N(k)$ decrescente, che rappresenta il compromesso tra approssimazione del dato e regolarità della soluzione.

Accenniamo brevemente che, nel caso della regolarizzazione di Tikhonov, esistono algoritmi per la determinazione del parametro α che non necessitano della conoscenza del livello di errore δ sui dati, informazione che spesso non è conosciuta.

Il più famoso è il metodo della "generalized cross-validation". L'idea che sta alla base di tale metodo è di scegliere il parametro in modo che, escludendo un valore y_j dal vettore \underline{y} dei dati e considerando solo i rimanenti, l'operatore ottenuto riesca ad approssimare al meglio il dato y_j non utilizzato. Questo procedimento di tipo statistico è ricorrente in calcolo numerico. Si rimanda a [50] per maggiori dettagli.