SCUOLA DI SCIENZE Corso di Laurea Triennale in Matematica

# Il problema della Nearest Correlation Matrix: il metodo delle proiezioni alternate e l'accelerazione di Anderson

Tesi di Laurea in Analisi Numerica

Relatore: Chiar.ma Prof.ssa VALERIA SIMONCINI Presentata da: FEDERICO CANDELETTI

Sessione Unica Anno Accademico 2021-2022

"Essendomi accontentato, ho trovato l'anima gemella..."

# Indice

Al	ostra	$\mathbf{ct}$	III
In	trod	uzione	$\mathbf{V}$
	Nota	azioni e definizioni	VII
1	La	matrice di correlazione più vicina	1
	1.1	La matrice di correlazione	2
	1.2	Il problema NCM	6
		1.2.1 Caratterizzazione della soluzione	8
<b>2</b>	Il n	netodo delle proiezioni alternate	15
	2.1	Proiezioni sugli insiemi convessi	15
	2.2	Implementazione del metodo	17
		2.2.1 Il caso di matrici non semidefinite positive	21
	2.3	Risultati numerici	23
3	Acc	elerazione di Anderson	25
	3.1	Accelerazione di Anderson per metodi di punto fisso	26
	3.2	Accelerazione del metodo delle proiezioni alternate	29
	3.3	Varianti del metodo	31
		3.3.1 Strategia di vincolo sugli elementi della matrice di correlazione	32
		3.3.2 Ricerca di una matrice di correlazione definita positiva	34
	3.4	Esperimenti numerici e confronto con il metodo non accelerato $\ . \ .$	36
Co	onclu	Isioni	47
	App	endice	49
Bi	bliog	grafia	51

## Abstract

Un problema frequente nell'analisi dei dati del mondo reale è quello di lavorare con dati "non coerenti", poiché raccolti a tempi asincroni o a causa di una successiva modifica "manuale" per ragioni che esulano da quelle matematiche. In particolare l'indagine dell'elaborato nasce da motivazioni di tipo finanziario, dove strumenti semplici come le matrici di correlazione, che sono utilizzate per capire le relazioni tra vari titoli o strategie, non rispettano delle caratteristiche cruciali a causa dell'incoerenza dei dati. A partire da queste matrici "invalide" si cerca la matrice di correlazione più vicina in norma, in modo da mantenere più informazioni originali possibili. Caratterizzando la soluzione del problema tramite analisi convessa, si utilizza il metodo delle proiezioni alternate, largamente utilizzato per la sua flessibilità anche se penalizzato dalla velocità di convergenza lineare. Viene quindi proposto l'utilizzo dell'accelerazione di Anderson, una tecnica per accelerare la convergenza dei metodi di punto fisso che, applicata al metodo di proiezione alternata, porta significativi miglioramenti in termini di tempo computazionale e numero di iterazioni. Si mostra inoltre come, nel caso di varianti del problema, l'applicazione dell'accelerazione di Anderson abbia un effetto maggiore rispetto al caso del problema "classico".

# Introduzione

L'elaborato si propone di indagare una strategia per trovare la matrice di correlazione più vicina ad una matrice dei dati di partenza (in inglese il problema è detto *Nearest Correlation Matrix*, abbreviato con "NCM"). Tale strategia si ispira alla trattazione del problema da parte di Higham [11, 12, 13].

Il primo capitolo introduce alcuni concetti di natura statistica per la definizione della matrice di correlazione e delle sue caratteristiche, e in seguito usa tali caratteristiche per formulare il problema NCM in maniera rigorosa.

Il secondo capitolo propone il metodo più largamente utilizzato per la risoluzione di questo problema: il metodo delle proiezioni alternate, formulato da Higham in [12]. Vengono gettate le basi per la teoria matematica necessaria al metodo e viene costruito l'algoritmo con la modifica di Dykstra in quanto si utilizzano insiemi convessi chiusi e non sottospazi.

Dopo la presentazione di qualche esempio e del caso particolare di matrici a rango basso, sono riportati gli esperimenti numerici, con l'accuratezza dell'algoritmo nei diversi casi.

Nel terzo capitolo viene proposta una tecnica per velocizzare l'algoritmo delle proiezioni alternate proveniente dai metodi di risoluzione dei metodi di punto fisso: l'accelerazione di Anderson. Per applicare tale accelerazione all'algoritmo precedentemente introdotto è necessario riformulare quest'ultimo come metodo di punto fisso, per poi adattare l'accelerazione al caso del problema NCM.

Insieme all'accelerazione di Anderson, vengono esplicitate due varianti del problema originario che, aggiungendo un vincolo al problema e rendendolo dunque più complesso, danno maggior senso alla discussione dell'accelerazione dell'algoritmo. Dopo una breve trattazione dell'assenza di garanzie di convergenza per la versione accelerata del metodo, vengono presentati i risultati numerici per mettere a confronto il metodo classico con quello accelerato appena costruito. Come si può vedere in [3] e [11], questo approccio non è l'unico possibile per la risoluzione del problema presentato. Esistono diverse variazioni del problema a seconda delle necessità, dalle matrici strutturate in modi precisi, al passaggio per il problema duale come viene fatto in [18]. Tale problema è molto legato a causa delle sue applicazioni al problema del completamento di matrice. Nonostante ciò, il metodo delle proiezioni alternate (accelerato o no) rimane uno dei più ampiamente utilizzati per la sua semplicità e la sua flessibilità a seconda della richiesta specifica del problema.

#### Notazioni e definizioni

Risulta necessario introdurre nozioni e concetti fondamentali che saranno utilizzati nell'elaborato.

**Definizione** (Prodotto interno). Sia X spazio vettoriale su  $\mathbb{C}$  (o su  $\mathbb{R}$ ). Si dice prodotto interno una funzione  $\langle \cdot, \cdot \rangle : X \times X \to \mathbb{C}$  tale che

1)  $\forall x, y \in X \quad \langle x, y \rangle = \langle y, x \rangle, \quad (su \mathbb{R} \ e \ la \ propriet a \ simmetrica)$ 

2)  $\forall x \in X \quad \langle x, x \rangle \ge 0 \ e \ \langle x, x \rangle = 0 \Leftrightarrow x = 0,$ 

3)  $\forall y \in X \text{ fissata, la funzione } x \mapsto \langle x, y \rangle \text{ è lineare.}$ 

Se X spazio vettoriale su  $\mathbb{R}$  dotato di prodotto interno e inoltre X è uno spazio di Banach (completo rispetto alla norma indotta dal prodotto interno) si dice che X è uno spazio di Hilbert.

Sia X uno spazio di Hilbert.

**Definizione** (Insieme chiuso). In termini topologici, un sottoinsieme C di uno spazio topologico X che contiene la sua frontiera è un insieme chiuso. Se X è uno spazio metrico, un insieme  $C \subset X$  è chiuso se e solo se per ogni

Se X e uno spazio metrico, un insieme  $C \subset X$  e chiuso se e solo se per ogni successione in C che converge in X, il limite della successione sta in C.

**Definizione** (Parte interna). La parte interna di un insieme S (o i punti interni di S) è l'insieme dei punti  $x \in S$  tali che  $\exists \epsilon > 0$  per cui  $B_{\epsilon}(x) \subset S$ . Viene indicata con  $\overset{\circ}{S}$ , int(S).

**Definizione** (Insieme convesso). Dato un sottoinsieme  $\mathcal{K} \subset X$  con  $x_1, x_2 \in \mathcal{K}$ , si dice che  $\mathcal{K}$  è convesso se

$$y = (1-t)x_1 + tx_2 \in \mathcal{K}, \quad \forall t \in [0,1].$$

Si consideri lo spazio vettoriale  $\mathbb{R}^{n \times n}$  della matrici quadrate di dimensione  $n \times n$ . Sia A una matrice nello spazio  $\mathbb{R}^{n \times n}$ , gli elementi della matrice in posizione ij si denotano con la lettera minuscola corrispondente a quella della matrice, ovvero  $a_{ij}$ .

Si indica con  $A^T$  la matrice trasposta di A. Se  $A = A^T$  la matrice è simmetrica; la matrice di correlazione e quella di covarianza sono matrici simmetriche.

$$A = Q\Lambda Q^T = Q\Lambda Q^{-1},$$

dove Q è una matrice ortogonale  $(Q^T Q = I)$  e  $\Lambda$  è una matrice diagonale i cui elementi sono gli autovalori di A.

**Definizione** (Matrice definita positiva). Una matrice A è definita positiva se  $\forall x \in \mathbb{R}^n$ ,  $x \neq 0$  vale  $x^T A x > 0$  (semidefinita positiva se  $x^T A x \geq 0$ ). Per indicare A definita positiva si scrive A > 0 ( $A \geq 0$  per A semidefinita positiva).

Se  $-A \ge 0$  allora A è semidefinita negativa e si indica con  $A \le 0$ . Una matrice è definita positiva se lo è la forma quadratica associata  $x^T A x$ .

**Definizione** (Traccia di matrice quadrata). *Si definisce la traccia di una matrice quadrata come la somma degli elementi diagonali. Si indica con* 

$$tr(A) = \sum_{i=1}^{n} a_{ii}.$$

Per una matrice simmetrica vale  $tr(A) = \sum_{i=1}^{n} a_{ii} = \sum_{i=1}^{n} \lambda_i$ .

Osservazione. La traccia di matrice gode della proprietà di invarianza per permutazione ciclica. Siano  $A, B, C, D \in \mathbb{R}^{n \times n}$ :

$$\operatorname{tr}(ABCD) = \operatorname{tr}(BCDA) = \operatorname{tr}(CDAB) = \operatorname{tr}(DABC).$$

Lo spazio vettoriale  $\mathbb{R}^{n \times n}$  è uno spazio di Hilbert. Essendo uno spazio a valori in  $\mathbb{R}$ , il prodotto interno di cui è dotato è una forma bilineare simmetrica definita positiva, cioè un prodotto scalare.

Il prodotto scalare tra matrici quadrate è dato da

$$\langle A, B \rangle = \operatorname{tr}(A^T B) = \operatorname{vec}(A)^T \operatorname{vec}(B).$$

**Notazione** (Operatore vec/unvec). Data  $A = (a_{ij})_{i,j} \in \mathbb{R}^{n \times n}$  si denota con vec(A) il vettore lungo  $n^2$  che possiede tutte le colonne della matrice una sopra l'altra:

$$A \stackrel{vec}{\longmapsto} [a_{11}, \ldots, a_{n1}, a_{12}, \ldots, \ldots, a_{nn}]^T$$
.

 $\acute{\mathrm{E}}$  possibile definire anche l'operazione inversa.

Dato un vettore  $\boldsymbol{a}$  lungo  $n^2$  si denota con unvec $(\boldsymbol{a})$  la preimmagine di  $\boldsymbol{a}$  attraverso l'operatore vec:

$$\boldsymbol{a} = [a_{11}, \dots, a_{n1}, a_{12}, \dots, \dots, a_{nn}]^T \xrightarrow{unvec} \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}$$

Fissando una matrice simmetrica definita positiva  $W \in \mathbb{R}^{n \times n}$ , è possibile definire un altro prodotto interno:

$$\langle A, B \rangle_W = \langle W^{1/2} A^T W^{1/2}, W^{1/2} B W^{1/2} \rangle, \qquad A, B \in \mathbb{R}^{n \times n}$$

(Anche la matrice  $W^{1/2}$  è simmetrica)

Questo prodotto interno induce su  $\mathbb{R}^{n\times n}$  la norma-W nel modo usuale:

$$\|A\|_W = \sqrt{\langle A, A \rangle_W}.$$

Se  $A \in \mathbb{R}^{n \times n}$  è simmetrica si può riscrivere  $\langle A, A \rangle_W$  come

$$\langle W^{1/2} A^T W^{1/2}, W^{1/2} A W^{1/2} \rangle =$$
  
= tr(W^{1/2} A W^{1/2} W^{1/2} A W^{1/2}) =  
= tr((W^{1/2} A W^{1/2})^T W^{1/2} A W^{1/2}).

Considerando che la norma Frobenius per matrici è

$$||A||_F^2 = \sum_{i,j} a_{ij}^2 = \operatorname{tr}(A^T A),$$

la norma-W risulta uguale alla norma Frobenius di una matrice congruente ad A:

$$\|A\|_W = \|W^{1/2}AW^{1/2}\|_F.$$

Il risultato della norma-W dipende dalla scelta della matrice W simmetrica definita positiva, ovvero cambia a seconda della trasformazione  $W^{1/2}AW^{1/2}$ . Per questo la norma-W viene vista come una versione pesata della norma Frobenius.

A partire dalla norma-W si definisce la distanza tra due matrici  $A, B \in \mathbb{R}^{n \times n}$  come  $||A - B||_W$ .

**Notazione** (Prodotto Hadamard). Date due matrici  $A \in B$  con le stesse dimensioni (ad es.  $A, B \in \mathbb{R}^{n \times n}$ ), di elementi  $a_{ij} \in b_{ij}$  rispettivamente, si denota il prodotto Hadamard (*element-wise product*) di  $A \in B$  con la matrice  $A \circ B$  i cui elementi sono del tipo  $a_{ij}b_{ij}$ .

**Notazione** (Proiezione su un insieme). Dato uno spazio X e un suo elemento  $x \in X$  si denota la proiezione di x su un sottoinsieme  $C \subset X$  con  $P_C(x)$ .

**Notazione** (Operatore diag). Dato un vettore  $\boldsymbol{v} = [v_1, \ldots, v_n]^T \in \mathbb{R}^n$  si denota con diag $(v_i)$  (oppure diag $(v_1, \ldots, v_n)$ ) la matrice diagonale che ha per elementi le componenti del vettore v,

$$[v_1, \dots, v_n]^T \xrightarrow{\text{diag}} \begin{bmatrix} v_1 & 0 & \cdots & 0 \\ 0 & v_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & v_n \end{bmatrix}$$

**Notazione** (Operatore Diag). Data una matrice  $A \in \mathbb{R}^{n \times n}$  si denota con Diag(A) il vettore che ha come componenti gli elementi della diagonale di A.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \xrightarrow{\text{Diag}} [a_{11}, a_{22}, \dots, a_{nn}]^T.$$

**Notazione** (Operatore DIAG). Data una matrice  $A \in \mathbb{R}^{n \times n}$  si denota con DIAG(A) la matrice diagonale con la stessa diagonale della matrice A.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \xrightarrow{\text{DIAG}} \begin{bmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & a_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{bmatrix}$$

## Capitolo 1

# La matrice di correlazione più vicina

La matrice di correlazione è una matrice simmetrica, semidefinita positiva e con diagonale unitaria:

$$X \in \mathbb{R}^{n \times n}$$
:  $X = X^T$ ,  $X \ge 0$ ,  $(X)_{ii} = 1$ ,  $\forall i = 1 \dots n$ .

È uno strumento centrale dell'analisi multivariata discreta, principalmente trova applicazioni di tipo statistico nello studio di fenomeni del mondo reale, da cui provengono grandi quantità di dati sperimentali, ma viene utilizzata anche in ambiti più matematici, come in Algebra Lineare Numerica, nell'analisi dell'errore della fattorizzazione di Cholesky o del metodo di Jacobi, nel precondizionamento di metodi iterativi per la risoluzione di sistemi lineari.

Spesso però, quando i dati provengono dal mondo reale, si ha a che fare con matrici di *quasi-correlazione* (o anche matrici di correlazione *invalide*), vale a dire matrici che non rispettano tutte le condizioni descritte sopra, quali matrici che non sono semidefinite positive. Ciò accade perchè la matrice viene costruita a partire da dati incompleti o non coerenti tra loro, da osservazioni asincrone poichè magari non disponibili tutte allo stesso modo in un momento preciso.

In ambito finanziario, è comune modificare "manualmente" porzioni di matrice o singoli elementi a seguito di un giudizio esterno o quando si vuole che i dati rientrino all'interno di un modello prestabilito, come accade nell'aggregazione del rischio o nella valutazione di strategie e di capitale in situazioni "estreme" (in inglese, lo *stress testing*): ancora, questo genera una matrice non definita (indefinita).

Sorge quindi la necessità di cercare una matrice di correlazione più vicina possi-

bile a quella disponibile di *quasi-correlazione* per poter svolgere un'analisi corretta ed evitare risultati non interpretabili, come varianze negative. Contemporaneamente, cercando la matrice di correlazione che dista meno da quella dei dati *invalida* si auspica che la quantità di informazioni mantenuta permetta di lavorare ancora sul problema originale.

Il problema della matrice di correlazione più vicina si è rivelato di largo utilizzo in diverse discipline: dal mondo della finanza ad applicazioni più pratiche come la ricostruzione dei livelli del mare, la modellizzazione di dati sanitari o la modellizzazione dei danni causati da tempeste alle infrastrutture.

#### 1.1 La matrice di correlazione

Per comprendere meglio l'utilizzo della matrice di correlazione è necessario introdurla tramite concetti di statistica descrittiva.

Siano  $X_1 \dots X_n$  vettori *m*-dimensionali ognuno contenente *m* campioni di una variabile (aleatoria). Tali variabili possono essere raggruppate in una matrice dei dati

$$\mathbf{X}_{(m \times n)} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{bmatrix}$$

dove l'elemento  $x_{ik}$  corrisponde alla *j*-esima osservazione della *k*-esima variabile.

Per ogni variabile k si possono definire la **media campionaria**  $\bar{x}_k$  e la **varianza campionaria**  $s_{kk}$  (una misura della dispersione dei dati):

$$\bar{x}_k = \frac{1}{m} \sum_{j=1}^m x_{jk}, \quad s_{kk} = \frac{1}{m} \sum_{j=1}^m (x_{jk} - \bar{x}_k)^2, \quad \forall k = 1 \dots n$$

Osservazione. In statistica, il valore della varianza campionaria descritto sopra viene detto stimatore distorto (biased estimator) della varianza di popolazione; per ragioni teoriche deve essere quindi corretto ponendo al denominatore m-1 al posto di m, in modo da avere uno stimatore imparziale della varianza di popolazione. Questa correzione prende il norme di **correzione di Bessel**, ma in casi in cui il numero di osservazioni sia molto alto non ha un effetto rilevante e dunque non è argomento di discussione dell'elaborato.

La deviazione standard è la radice quadrata della varianza campionaria:  $\sqrt{s_{kk}}$ . Date due variabili, si può definire una misura dell'associazione lineare tra le due tramite la **covarianza campionaria** 

$$s_{ik} = \frac{1}{m} \sum_{j=1}^{m} (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k), \quad \forall i, k = 1 \dots n.$$

Per una misura di associazione lineare (o dipendenza lineare) che però non dipenda dalle unità di misura e dall'ordine di grandezza dei dati ci si affida al **coefficiente di correlazione campionaria** 

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}}\sqrt{s_{kk}}}, \quad \forall i, k = 1 \dots n.$$

Il coefficiente di correlazione tra due variabili può essere visto come la covarianza tra le corrispondenti variabili standardizzate, ovvero in cui le osservazioni  $x_{jk}$  sono sostituite con i dati standardizzati  $(x_{jk} - \bar{x}_k)/\sqrt{s_{kk}}$ . Il valore di  $r_{ik}$  non cambia se la correzione di Bessel viene applicata o meno.

Anche se i valori di covarianza e correlazione sembrano esprimere lo stesso concetto, tra i due è presente una differenza cruciale: la covarianza tra due variabili deve essere rapportata alle unità di misura dei dati in questione e risente dell'ordine di grandezza di questi, la correlazione come detto sopra è indipendente da tutto questo, è una quantità standardizzata e **adimensionale** ed è spesso preferibile alla covarianza in quanto limitata tra -1 e 1; se il valore di  $r_{ik}$  è vicino a 1 (o -1) c'è dipendenza lineare tra le variabili i e k, se è vicino allo zero indica assenza di dipendenza lineare tra le due. É importante notare che il coefficiente di correlazione tra due variabili esprime il grado di dipendenza lineare tra le due, ma non permette di dedurre legami di altro tipo: se tra due variabili ci fosse una dipendenza di tipo quadratico, non si potrebbe dedurre alcun legame tra esse (a partire da queste quantità statistiche) e il coefficiente  $r_{ik}$  tra le due potrebbe essere basso.

Passando al caso multivariato, le quantità  $s_{ik}$  e  $r_{ik}$  diventano rispettivamente

matrice di covarianza campionaria S

$$\mathbf{S}_{(n \times n)} = \begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1n} \\ s_{21} & s_{22} & \cdots & s_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ s_{n1} & s_{n2} & \cdots & s_{nn} \end{bmatrix}$$

e matrice di correlazione campionaria R

$$\mathbf{R}_{(n \times n)} = \begin{bmatrix} 1 & r_{12} & \cdots & r_{1n} \\ r_{21} & 1 & \cdots & r_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ r_{n1} & r_{n2} & \cdots & 1 \end{bmatrix}$$

Dal momento che  $s_{ik} = s_{ki}$ ,  $r_{ik} = r_{ki}$  le matrici **S** e **R** sono simmetriche.

A partire dalla matrice di covarianza, indicando con  $\mathbf{D}^{1/2}$  la matrice che ha sulla diagonale le deviazioni standard  $\sqrt{s_{kk}}$ , si può ottenere la matrice di correlazione:  $\mathbf{R} = \mathbf{D}^{-1/2} \mathbf{S} \mathbf{D}^{-1/2}$ .

La matrice di covarianza  $\mathbf{S}$  viene utilizzata nella regressione lineare, in strategie di riduzione delle variabili come l'Analisi delle Componenti Principali o nelle tecniche di normalizzazione e standardizzazione dei campioni di dati, per ricondursi allo studio di modelli probabilistici con distribuzione multinormale; le proprietà spettrali di  $\mathbf{S}$  hanno grande importanza nello studio e nella visualizzazione dei dati in analisi multivariata.

Essendo una matrice simmetrica semidefinita positiva, la forma quadratica associata è anch'essa semidefinita positiva. Da un punto di vista statistico questa forma quadratica è interpretabile come distanza di un punto  $x \in \mathbb{R}^n$  dall'origine dello spazio  $\mathbb{R}^n$  o da un punto  $\mu \in \mathbb{R}^n$ : la forma quadratica è espressa tramite la matrice inversa di **S** come

$$x^T S^{-1} x \ge 0.$$

Nel caso in cui la matrice di covarianza sia l'identità, questa distanza statistica (distanza di Mahalanobis) coincide con la distanza euclidea.

Fissato un  $c^2$ , i punti  $x \in \mathbb{R}^n$  che soddisfano l'equazione  $(x-\mu)^T S^{-1}(x-\mu) = c^2$ giacciono su un *iperellissoide* nello spazio  $\mathbb{R}^n$ , di centro  $\mu$ , i cui assi sono dati dagli autovettori della matrice e la loro lunghezza è proporzionale agli autovalori della matrice. La dominanza di questi ultimi è fondamentale in diverse discipline, dalla semplice riduzione della dimensione dei dati a tecniche che possono avere applicazione nel riconoscimento facciale. In ciò che concerne lo studio di modelli probabilistici mediante questi strumenti, questi iperellissoidi sono detti *contorni di densità di probabilità costante* [14].



Figura 1.1: Ellisse dei dati

Discorsi analoghi valgono per la matrice di correlazione, essendo questa la matrice di covarianza dei dati standardizzati.

Osservazione. L'espressione della forma quadratica  $(x-\mu)^T S^{-1}(x-\mu)$  è lecita solo nel momento in cui la matrice S è definita positiva e non semidefinita positiva, per cui anche  $S^{-1}$  lo è.

I concetti di matrice di covarianza  $\mathbf{S}$  e correlazione  $\mathbf{R}$  sono indubbiamente legati; per coerenza con gli studi citati nell'elaborato si farà riferimento alla matrice di correlazione: matrice simmetrica, semidefinita positiva, con diagonale unitaria.

- In quanto  ${f R}$  è una matrice simmetrica, i suoi autovalori sono tutti reali.

- La matrice **R** è semidefinita positiva dal momento che, come detto sopra, la matrice di correlazione coincide con la matrice di covarianza dei dati standardizzati. Di conseguenza tutti i suoi autovalori soddisfano  $\lambda_i \geq 0 \forall i$ .

- La diagonale unitaria si spiega con il fatto che una variabile è linearmente correlata con se stessa. Come conseguenza si ha che tr $(\mathbf{R}) = n$  e quindi che  $\sum_{k=1}^{n} \lambda_k = n$ . Si può dunque concludere che  $0 \le \lambda_i \le n \quad \forall i = 1 \dots n$ .

Vale inoltre il seguente risultato per matrici semidefinite positive.

**Proposizione.** Se A è una matrice simmetrica semidefinita positiva vale  $|a_{ij}| \leq \sqrt{a_{ii}a_{jj}}$ .

Dimostrazione. Dal momento che la matrice A è semidefinita positiva vale che i suoi minori principali sono tutti non negativi.

In particolare, considerando i minori principali di ordine 2, vale  $a_{ii}a_{jj} - a_{ij}a_{ji} \ge 0$ . Siccome la matrice è simmetrica si ha  $a_{ii}a_{jj} \ge a_{ij}^2$ , da cui la tesi.

# 1.2 Il problema NCM: approssimazione della matrice di correlazione più vicina

Come già accennato all'inizio del capitolo, non sempre è possibile reperire dei dati sincronizzati e consistenti, o talvolta è necessario eliminare delle osservazioni sperimentali perchè non corrette; nel momento in cui si costruisce la matrice delle correlazioni tra le variabili si ottiene quindi una matrice di correlazione *approssimata* o *invalida* (matrice di *quasi-correlazione*), che non gode delle proprietà di matrice semidefinita positiva e/o di diagonale unitaria. Per fare un esempio, in ricerca finanziaria, le matrici di correlazione campionaria sono usate con scopi predittivi: è chiaro che se i dati utilizzati sono asincroni l'utilizzo di una matrice ottenuta tramite quei dati può portare ad analisi errate. Per giustificare una corretta analisi si richiede dunque il calcolo della matrice di correlazione più vicina in norma Frobenius a quella *approssimata* ottenuta dai dati inconsistenti, che mantenga il più possibile le informazioni originali.

Problema. Si cerca, per un matrice simmetrica arbitraria  $A \in \mathbb{R}^{n \times n}$ , di approssimare la matrice di correlazione X corrispondente alla quantità

 $\gamma(A, W) = \min \left\{ \|A - X\|_W : X \text{ è una matrice di correlazione} \right\}, \tag{1.1}$ 

che indica la distanza minima di A da una qualsiasi matrice di correlazione, fissata la matrice W relativa alla norma- $W \|\cdot\|_W$ .

La trasformazione  $A \to W^{1/2}AW^{1/2}$  della norma  $\|\cdot\|_W$  è una congruenza quindi preserva la segnatura della matrice A. Ciò tornerà utile nel caso di matrici A a rango basso con molti autovalori nulli, frequenti nell'analisi di tipo finanziario che motiva la ricerca della NCM.

Definiamo gli insiemi, rispettivamente, delle matrici simmetriche semidefinite positive e delle matrici simmetriche con diagonale unitaria:

$$\mathcal{S} = \left\{ Y = Y^T \in \mathbb{R}^{n \times n} : Y \ge 0 \right\},$$
$$\mathcal{U} = \left\{ Y = Y^T \in \mathbb{R}^{n \times n} : y_{ii} = 1, \ i = 1 \dots n \right\}.$$

La soluzione del problema (1.1) si trova nell'intersezione dei due insiemi  $S \in \mathcal{U}$ ; siccome entrambi gli insiemi sono chiusi e convessi, lo è anche la loro intersezione. É dato il seguente teorema da [16]:

**Teorema 1.1** (*Distanza minima da un insieme convesso chiuso*). Sia H uno spazio di Hilbert e sia  $\mathcal{K} \neq \emptyset$  un sottoinsieme chiuso e convesso di H. Dato  $x \in H$  esiste ed è unico l'elemento  $k_0 \in \mathcal{K}$  t.c.  $||x - k_0||_H \leq ||x - k||_H \quad \forall k \in \mathcal{K}$ , dove  $||\cdot||_H$  è la norma indotta dal prodotto interno  $\langle \cdot, \cdot \rangle_H$  su H. Inoltre, una condizione necessaria e sufficiente affinché  $k_0$  sia l'unico elemento con questa caratteristica è che  $\langle x - k_0, k - k_0 \rangle_H \leq 0 \quad \forall k \in \mathcal{K}$ .

Dal momento che lo spazio delle matrici  $\mathbb{R}^{n \times n}$  è uno spazio di Hilbert e l'intersezione  $S \cap U$  è chiusa e convessa, segue dal teorema che la soluzione X del problema (1.1) esiste ed è unica.

Sono degni di nota due casi particolari, per cui viene fissata W diagonale: - se la matrice di partenza A è diagonale allora la soluzione è data da X = I, - se A è semidefinita positiva con gli elementi  $a_{ii} \leq 1 \quad \forall i$  allora X è ottenuta semplicemente sostituendo agli elementi diagonali il valore 1.

Un'interessante particolarità del problema (1.1) è che, mentre la definizione positiva di una matrice è indipendente dalla scelta della base per lo spazio, la diagonale unitaria è una caratteristica che dipende dalla base dello spazio: questo indubbiamente modifica la posizione dell'intersezione nella quale si trova la soluzione all'interno dello spazio  $\mathbb{R}^{n \times n}$ , e dunque rende difficile l'espressione di una soluzione in forma esplicita.

Nel seguente enunciato sono fornite delle stime dall'alto e dal basso per la distanza  $\gamma(A, W)$ .

**Proposizione.** Sia  $A \in \mathbb{R}^{n \times n}$  matrice simmetrica, sia W matrice simmetrica definita positiva fissata. Allora max  $\{\alpha_1, \alpha_2\} \leq \gamma(A, W) \leq \max\{\beta_1, \beta_2, \beta_3\}$ , dove

$$\alpha_1^2 = \sum_{i=1}^n w_i^2 (a_{ii} - 1)^2 + \sum_{\substack{i \neq j \\ |a_{ij}| > 1}} w_i w_j (1 - |a_{ij}|)^2,$$
  
$$\alpha_2^2 = \min \{ \|A - X\|_W : X \in S \},$$
  
$$\beta_1 = \|A - I\|_W,$$
  
$$\beta_2 = \min \{ \|A - zz^T\|_W : z_i = \pm 1, i = 1 \dots n \},$$
  
$$\beta_3 = \min_{0 \le \rho \le 1} \|A - (\rho^{|i-j|})\|_W.$$

Infine, fissando la matrice W, la distanza soddisfa la seguente proprietà

$$|\gamma(A, W) - \gamma(B, W)| \le ||A - B||_W.$$

Questa stima suggerisce che se la distanza  $||A - B||_W$  è sufficientemente piccola, allora la matrice di correlazione più vicina ad A è una buona approssimazione della matrice di correlazione più vicina a B.

#### 1.2.1 Caratterizzazione della soluzione

Per sviluppare una strategia che conduca alla soluzione del problema (1.1) si cerca di caratterizzare la soluzione lavorando con uno specifico prodotto interno e usando le proprietà degli insiemi convessi.

Fissata una matrice W simmetrica definita positiva, sia  $\|\cdot\|_W$  la norma-W indotta su  $\mathbb{R}^{n \times n}$  dal prodotto interno

$$\langle A, B \rangle_W = \operatorname{tr}(W^{1/2}AW^{1/2}W^{1/2}BW^{1/2}) = \operatorname{tr}(AWBW).$$

Siccome la matrice di partenza A viene scelta simmetrica, si può considerare di lavorare solo con matrici simmetriche. Di conseguenza vale  $\langle A, B \rangle_W = \langle B, A \rangle_W$ .

Dato un insieme convesso  $\mathcal{K} \subset \mathbb{R}^{n \times n}$ , si definisce il **cono normale** all'insieme convesso nel punto  $B \in \mathcal{K}$  come:

$$\partial \mathcal{K}(B) = \left\{ Y = Y^T \in \mathbb{R}^{n \times n} : \langle Y, B \rangle_W = \sup_{Z \in \mathcal{K}} \langle Y, Z \rangle_W \right\}$$
(1.2)  
=  $\left\{ Y = Y^T \in \mathbb{R}^{n \times n} : \langle Z - B, Y \rangle_W \le 0, \ \forall Z \in \mathcal{K} \right\}$ 

La seconda parte del Teorema 1.1 afferma che la soluzione X del problema (1.1) soddisfa la seguente condizione:

$$\langle A - X, Z - X \rangle_W \leq 0 \quad \forall Z \in \mathcal{S} \cap \mathcal{U}.$$

Questa condizione può essere riscritta come  $A - X \in \partial(\mathcal{S} \cap \mathcal{U})(X)$ , in cui  $\partial(\mathcal{S} \cap \mathcal{U})(X)$  è il cono normale a  $\mathcal{S} \cap \mathcal{U}$  in X.

Vale il seguente risultato da [19] per coni normali a insiemi convessi:

**Proposizione.** Siano  $\mathcal{K}_1 \in \mathcal{K}_2$  due insiemi convessi tali che  $\overset{\circ}{\mathcal{K}_1} \cap \overset{\circ}{\mathcal{K}_2} \neq \emptyset$ . Dato  $P \in \mathcal{K}_1 \cap \mathcal{K}_2$  si ha  $\partial(\mathcal{K}_1 \cap \mathcal{K}_2)(P) = \partial\mathcal{K}_1(P) + \partial\mathcal{K}_2(P)$ .

Una qualsiasi matrice di correlazione appartiene ai punti interni sia di S che di  $\mathcal{U}$ , per cui si può concludere che la soluzione X soddisfa

$$A - X \in \partial \mathcal{S}(X) + \partial \mathcal{U}(X). \tag{1.3}$$

É dunque importante determinare delle espressioni esplicite per  $\partial S$  e  $\partial U$ .

#### Lemma 1.2. Data $A \in \mathcal{U}$ ,

$$\partial \mathcal{U}(A) = \left\{ W^{-1} \operatorname{diag}(\theta_i) W^{-1} : \boldsymbol{\theta} = [\theta_1, \dots, \theta_n]^T \text{ vettore arbitrario} \right\}.$$

Dimostrazione. Si ha

$$\partial \mathcal{U}(A) = \left\{ Y = Y^T \in \mathbb{R}^{n \times n} : \langle Y, A \rangle_W = \sup_{Z \in \mathcal{U}} \langle Y, Z \rangle_W \right\}$$

Le matrici in questione sono tutte simmetriche e quindi il vincolo può essere riscritto come  $\langle A, Y \rangle_W = \sup_{Z \in \mathcal{U}} \langle Z, Y \rangle_W$ , ovvero,

$$\operatorname{tr}(A\widetilde{Y}) = \sup_{Z \in \mathcal{U}} \operatorname{tr}(Z\widetilde{Y}) \Leftrightarrow \sum_{i} (A\widetilde{Y})_{ii} = \sup_{Z \in \mathcal{U}} \sum_{i} (Z\widetilde{Y})_{ii} \Leftrightarrow$$
$$\Leftrightarrow \sum_{i,j} a_{ij}\widetilde{y}_{ji} = \sup_{Z \in \mathcal{U}} \sum_{i,j} z_{ij}\widetilde{y}_{ji} \Leftrightarrow \sum_{i,j} \widetilde{y}_{ij}a_{ij} = \sup_{Z \in \mathcal{U}} \sum_{i,j} \widetilde{y}_{ij}z_{ij}$$

con  $\widetilde{Y} = WYW$  simmetrica.

Se  $\tilde{Y}$  non fosse diagonale, una scelta di  $z_{ij}$  abbastanza grande violerebbe la condizione del sup. Di conseguenza  $\tilde{Y}$  deve essere diagonale, e qualsiasi matrice della forma  $Y = W^{-1} \operatorname{diag}(\theta_i) W^{-1}$  soddisfa la condizione richiesta.

Lemma 1.3. Data  $A \in \mathcal{S}$ ,

$$\partial \mathcal{S}(A) = \left\{ Y = Y^T \in \mathbb{R}^{n \times n} : \langle Y, A \rangle_W = 0, \ Y \le 0 \right\}.$$

Dimostrazione. Si ha

$$\partial \mathcal{S}(A) = \left\{ Y = Y^T \in \mathbb{R}^{n \times n} : \langle Y, A \rangle_W = \sup_{Z \in \mathcal{S}} \langle Y, Z \rangle_W \right\}$$

Sia  $Z \in \mathcal{S}$  con decomposizione spettrale  $Z = Q\Lambda Q^T$ , dove Q è ortogonale e  $\Lambda = \operatorname{diag}(\lambda_i) \geq 0$ . Sia  $C = Q^T W Y W Q$ .

$$\sup_{Z \in \mathcal{S}} \langle Y, Z \rangle_{W} = \sup_{\Lambda \ge 0, \ Q^{T}Q = I} \langle Y, Q\Lambda Q^{T} \rangle_{W}$$
$$= \sup_{\Lambda \ge 0, \ Q^{T}Q = I} \operatorname{tr}(YWQ\Lambda Q^{T}W)$$
$$= \sup_{\Lambda \ge 0, \ Q^{T}Q = I} \operatorname{tr}(C\Lambda)$$
$$= \sup_{\lambda_{i} \ge 0, \ Q^{T}Q = I} \sum_{i} \lambda_{i}c_{ii}$$
$$= \begin{cases} 0 \quad \text{se } Y \le 0 \\ \infty \quad \text{altrimenti.} \end{cases}$$

Siccome deve valere l'uguaglianza  $\langle Y, A \rangle_W = \sup_{Z \in S} \langle Y, Z \rangle_W$  e il prodotto interno  $\langle \cdot, \cdot \rangle_W$  è a valori in  $\mathbb{R}$ , il sup calcolato sopra deve essere finito. Dunque l'uguaglianza vale per  $\langle Y, A \rangle_W = 0$  e  $Y \leq 0$ .

Si può dedurre da questo il seguente corollario.

**Corollario 1.4.** Per  $A \in S$ , sia p la dimensione di Ker(A),

$$\partial \mathcal{S}(A) = \{Y : WYW = -VDV^T, \text{dove } V \in \mathbb{R}^{n \times p} \text{ ha colonne ortonormali} \\ \text{che generano } \text{Ker}(A) \in D = \text{diag}(d_i) \ge 0 \}.$$

Dimostrazione. Sia  $A = Q\Lambda Q^T$  la decomposizione spettrale della matrice  $A \in S$ dove  $\Lambda = \operatorname{diag}(\lambda_i) \operatorname{con} \lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_{n-p} > 0 = \lambda_{n-p+1} = \ldots = \lambda_n$ . Sia  $\Lambda_1 = \operatorname{diag}(\lambda_1, \ldots, \lambda_{n-p}), \ Q = [Q_1, Q_2] \operatorname{con} Q_1 \in \mathbb{R}^{n \times (n-p)}$  e  $Q_2$  con colonne ortonormali che generano Ker(A). Sia  $G = Q_1^T W Y W Q_1 \in \mathbb{R}^{(n-p) \times (n-p)}$ . Se  $Y \in \partial \mathcal{S}(A)$ , per il Lemma 1.3,

$$0 = \langle Y, A \rangle_W = \operatorname{tr}(AWYW) = \operatorname{tr}(Q_1 \Lambda_1 Q_1^T WYW) = \operatorname{tr}(\Lambda_1 G) = \sum_{i=1}^{n-p} \lambda_i g_{ii},$$

dove  $\Lambda_1 > 0$  e  $Y \leq 0$  (e quindi  $Q_1^T W Y W Q_1 \leq 0$ ). Di conseguenza vale che  $\text{Diag}(Q_1^T W Y W Q_1) = [0, \ldots, 0]^T = \mathbf{0}$  (Poichè gli elementi  $g_{ii}$  sono tutti negativi o tutti nulli).

Siccome dal Lemma 1.3 si ha  $Y \leq 0$ , si può costruire

$$F = \begin{bmatrix} G & H \\ H^T & M \end{bmatrix} := Q^T (WYW)Q = \begin{bmatrix} Q_1^T (WYW)Q_1 & Q_1^T (WYW)Q_2 \\ Q_2^T (WYW)Q_1 & Q_2^T (WYW)Q_2 \end{bmatrix} \le 0.$$

Ne consegue che i blocchi  $G \in \mathbb{R}^{(n-p)\times(n-p)}, M \in \mathbb{R}^{p\times p}$  siano semidefiniti negativi. Se Diag $(G) = \mathbf{0}$ , come conseguenza la matrice G è indefinita. Si ha quindi che G è la matrice nulla.

La matrice F può essere fattorizzata in questo modo:

$$\begin{bmatrix} 0 & H \\ H^T & M \end{bmatrix} := L \begin{bmatrix} B \\ & M \end{bmatrix} L^T,$$

dove L è una matrice triangolare inferiore e il blocco  $B \in \mathbb{R}^{(n-p) \times (n-p)}$ , che dipende da H, ha autovalori di segno opposto rispetto a quelli di M. Se quindi H non fosse una matrice nulla, la matrice F risulterebbe allo stesso tempo indefinita e semidefinita negativa. Per questo, H = 0 e  $M \leq 0$ .

Allora  $0 \ge WYW = Q_2MQ_2^T = -VDV^T$ , usando la decomposizione spettrale di *M* per generare  $V \in \mathbb{R}^{n \times p}$  con colonne ortonormali e  $D \in \mathbb{R}^{p \times p}$ ,  $D \ge 0$ , diagonale.

É possibile ora caratterizzare esplicitamente la soluzione del problema NCM.

**Teorema 1.5.** La matrice di correlazione X risolve il problema (1.1) se e solo se

$$X = A + W^{-1} \left( V D V^T - \operatorname{diag}(\theta_i) \right) W^{-1}, \qquad (1.4)$$

dove  $V \in \mathbb{R}^{n \times p}$  ha colonne ortonormali che generano Ker(X),  $D = \text{diag}(d_i) \ge 0$ , e i  $\theta_i$  sono arbitrari.

Dimostrazione. Il risultato segue dalla condizione (1.3) applicando il Lemma 1.2 e il Corollario 1.4.  $\hfill \Box$ 

Osservazione. Il Teorema 1.5 vale anche per V con colonne non ortonormali.

Dimostrazione. Si prova allo stesso modo del Teorema 1.5 ma usando una versione del Corollario 1.4 in cui la decomposizione spettrale di A sia  $A = Q\Lambda Q^T$ , dove  $\Lambda_1 = diag(\lambda_1, \ldots, \lambda_{n-p}), \ Q = [Q_1, Q_2] \text{ con } Q_1 \in \mathbb{R}^{n \times (n-p)}$  e  $Q_2$  con colonne non ortonormali ma che comunque generano Ker(A).

In questo modo si ha  $0 \ge WYW = Q_2MQ_2^T = -Q_2XDX^TQ_2^T = -VDV^T$ , dove  $V = Q_2X$  è una matrice che non ha colonne ortonormali. Si utilizza poi questa decomposizione per provare la (1.4).

Si può utilizzare il Teorema 1.5 in maniera immediata nei seguenti casi particolari.

- Se W viene scelta diagonale la matrice X generalmente sarà singolare. Nel caso in cui X sia invece non singolare, allora risulta V = 0 e  $X = A - W^{-1} \operatorname{diag}(\theta_i) W^{-1}$ , vale a dire che X è ottenuta aggiustando gli elementi diagonali di A ad 1.

- Se W viene scelta diagonale e inoltre gli elementi diagonali di A sono almeno 1, si può dire di più attraverso il seguente risultato.

**Teorema 1.6.** Sia  $A = A^T$  con elementi diagonali  $a_{ii} \ge 1$  e sia W diagonale. Allora nel Teorema 1.5 vale  $\theta_i \ge 0$  per ogni i. Inoltre, se A ha t autovalori non positivi allora la matrice di correlazione ad essa più vicina ha almeno t autovalori nulli.

Dimostrazione. La (1.4) può essere riscritta come

$$X = A + W^{-1} (V D V^T) W^{-1} - W^{-1} \operatorname{diag}(\theta_i) W^{-1}.$$

Siccome la matrice A ha elementi diagonali  $a_{ii} \ge 1$  e  $W^{-1}(VDV^T)W^{-1} \ge 0$ , gli elementi diagonali di  $A + W^{-1}(VDV^T)W^{-1}$  sono almeno 1. Affinché X abbia diagonale unitaria serve che  $\theta_i \ge 0 \forall i$ .

Parallelamente si nota che il termine  $W^{-1}(VDV^T)W^{-1}$  porta i t autovalori non positivi di  $A - W^{-1}$ diag $(\theta_i)W^{-1}$  a diventare non negativi (poichè  $X \in S$ ). La perturbazione  $W^{-1}(VDV^T)W^{-1}$  ha rango al massimo p, quindi si ha che  $p \ge t$ . Per ipotesi p è la dimensione di Ker(X), dunque X ha almeno t autovalori nulli.

Il Teorema 1.5 ed il Teorema 1.6 non permettono di calcolare effettivamente una soluzione del problema (1.1), ma possono essere usati per verificare un possibile candidato.

Per esempio, se W = I, data la matrice indefinita

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix},$$

un potenziale candidato per la matrice di correlazione più vicina è  $X = ee^T$ , in cui  $e = [1, 1, 1]^T$ , con  $||A - X||_F = \sqrt{2}$ . Ker(X) è generato dalle colonne di

$$V = \begin{bmatrix} 1 & -1 \\ -1 & 1 \\ 0 & 1 \end{bmatrix},$$

per cui

$$X = A + \begin{bmatrix} d_1 + d_2 & -d_1 & -d_2 \\ -d_1 & d_1 & 0 \\ -d_2 & 0 & d_2 \end{bmatrix} + \operatorname{diag}(\theta_i).$$

Per il Teorema 1.5 deve essere  $d_i \ge 0$  ma l'equazione implica  $d_2 = -1$ , dunque X non può essere una soluzione.

La soluzione corretta è infatti

$$X = \begin{bmatrix} 1 & 0.7607 & 0.1573 \\ 0.7607 & 1 & 0.7607 \\ 0.1573 & 0.7607 & 1 \end{bmatrix},$$

con  $\left\|A-X\right\|_{F}=0.5278,$ trovata con l'Algoritmo 1 descritto nel prossimo capitolo.

### Capitolo 2

# Il metodo delle proiezioni alternate

Tra i metodi iterativi disponibili per risolvere il problema (1.1) con garanzia di convergenza all'unica soluzione, quello più facile da implementare è quello delle proiezioni alternate sugli insiemi convessi introdotti nel Capitolo 1, proposto da Higham in [12].

Un altro algoritmo più veloce e basato su un metodo quasi-Newton è stato sviluppato da Qi e Sun in [18] e successivamente migliorato da Bosdorf e Higham in [3], ma il metodo di proiezione alternata rimane più largamente utilizzato in diverse discipline, per popolarità e semplicità.

Oltre ad essere di semplice implementazione, è facilmente modificabile nel caso si vogliano aggiungere ulteriori vincoli alla matrice da calcolare, in particolare se è necessario fissare elementi o cercare una matrice strettamente definita positiva.

La strategia del metodo consiste nel proiettare alternativamente su dei sottoinsiemi l'iterata corrente. A tal fine, l'algoritmo richiede ad ogni iterazione una decomposizione spettrale e per questo la velocità di convergenza è lineare. Queste caratteristiche rendono il metodo piuttosto lento.

#### 2.1 Proiezioni sugli insiemi convessi

Sviluppata la teoria necessaria nel Capitolo 1, è di interesse capire come proiettare le matrici simmetriche sulle matrici di correlazione. Si considerano le proiezioni sui due insiemi convessi  $\mathcal{U}$  ed  $\mathcal{S}$  individualmente.

La proiezione su  $\mathcal{U}$  è piuttosto semplice.

**Teorema 2.1.** Data  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$ , in norma-W,

$$P_{\mathcal{U}}(A) = A - W^{-1} \operatorname{diag}(\theta_i) W^{-1},$$

dove  $\boldsymbol{\theta} = [\theta_1, \dots, \theta_n]^T$  è la soluzione del sistema lineare

$$(W^{-1} \circ W^{-1})\boldsymbol{\theta} = \operatorname{Diag}(A - I).$$

*Dimostrazione.* La proiezione  $H = P_{\mathcal{U}}(A)$  è caratterizzata da  $A - H \in \partial \mathcal{U}(H)$ , che per il Lemma 1.2 si scrive

$$A - H = W^{-1} \operatorname{diag}(\theta_i) W^{-1}.$$

Uguagliando gli elementi diagonali e ponendo  $W^{-1} = (w_{ij})_{i,j}$  (se W è diagonale anche  $W^{-1}$  lo è), si hanno le equazioni

$$\sum_{j=1}^{n} w_{ij}^2 \theta_j = a_{ii} - 1, \ \forall i = 1, \dots, n.$$

Ne risulta il sistema lineare  $(W^{-1} \circ W^{-1})\boldsymbol{\theta} = \text{Diag}(A - I)$ , dove  $\boldsymbol{\theta} = [\theta_1, \dots, \theta_n]^T$ . Essendo W definita positiva, lo è anche la matrice  $(W^{-1} \circ W^{-1})$  e il sistema ha soluzione unica.

Osservazione. Nel caso particolare in cui W sia diagonale, si ha

$$A - H = W^{-1} \operatorname{diag}(\theta_i) W^{-1} = D$$

con D diagonale. Di conseguenza  $H = P_{\mathcal{U}}(A) = A - D$  e si può dire semplicemente:

$$h_{ij} = \begin{cases} a_{ij}, & i \neq j, \\ 1, & i = j. \end{cases}$$

Per quanto riguarda la proiezione su S, il discorso non è così immediato. Sia  $A \in \mathbb{R}^{n \times n}$  simmetrica con decomposizione spettrale  $A = QDQ^T$ , dove D =  $\operatorname{diag}(\lambda_i) \in Q$  ortogonale. Si definiscono:

$$A_{+} = Q \operatorname{diag}(\max(\lambda_{i}, 0))Q^{T}, \quad A_{-} = Q(\min(\lambda_{i}, 0))Q^{T}$$

Valgono  $A = A_{+} + A_{-} e A_{+} A_{-} = A_{-}A_{+} = 0.$ 

Teorema 2.2. Data 
$$A \in \mathbb{R}^{n \times n}$$
,  $A = A^T$ , in norma- $W$ ,  
 $P_{\mathcal{S}}(A) = W^{-1/2} ((W^{1/2}AW^{1/2})_+)W^{-1/2}$ .  
In più, vale  $DIAG(P_{\mathcal{S}}(A)) \ge DIAG(A)$ .

Dimostrazione. Si dimostra la correttezza dell'espressione fornita per  $H = P_{\mathcal{S}}(A)$ provando che soddisfa la richiesta  $A - H \in \partial \mathcal{S}(H)$ , cioè, dal Lemma 1.3, che siano soddisfatte:

$$A - H \le 0, \ \operatorname{tr}((A - H)WHW) = 0.$$

Si vede prima

$$A - H = W^{-1/2} (W^{1/2} A W^{1/2} - (W^{1/2} A W^{1/2})_{+}) W^{-1/2} = W^{-1/2} (W^{1/2} A W^{1/2})_{-} W^{-1/2} \le 0$$

e poi

$$(A - H)WHW = W^{-1/2}(W^{1/2}AW^{1/2})_{-}W^{-1/2} \cdot W^{1/2}((W^{1/2}AW^{1/2})_{+})W^{1/2}$$
$$= W^{-1/2}(W^{1/2}AW^{1/2})_{-}(W^{1/2}AW^{1/2})_{+}W^{1/2} = 0.$$

Infine, detta  $M = W^{1/2}AW^{1/2}$  si ha  $M_+ - M = -M_- \ge 0$ , ovvero

$$(W^{1/2}AW^{1/2})_+ - W^{1/2}AW^{1/2} \ge 0$$

dunque, moltiplicando a destra e sinistra per  $W^{-1/2}$  si ottiene una congruenza che preserva la disuguaglianza, e prendendo le parti diagonali si ha la tesi.

#### 2.2 Implementazione del metodo

Per trovare la matrice di correlazione più vicina nell'intersezione di  $S \in U$  si proietta iterativamente ripetendo ad ogni ciclo l'operazione

$$A \leftarrow P_{\mathcal{U}}(P_{\mathcal{S}}(A)).$$

L'idea di proiettare iterativamente su sottospazi è stata analizzata in spazi di Hilbert da von Neumann, il quale ha provato la convergenza del metodo al punto dell'intersezione più vicino al dato iniziale. In questo caso  $S \in U$  non sono sottospazi ma insiemi chiusi e convessi, di conseguenza il risultato di von Neumann non si applica e l'iterazione può convergere a punti non ottimali. Viene quindi adottata un'iterazione modificata che incorpora per ogni proiezione **la correzione di Dykstra**, interpretabile come un vettore normale all'insieme convesso sul quale si proietta.

La procedura descritta da Dykstra in [8] è relativa al caso generale in cui la soluzione  $g^*$  si trova nell'intersezione di r insiemi convessi chiusi  $\mathcal{K}_1, \ldots, \mathcal{K}_r$ . Ad ogni ciclo t vengono eseguite le r proiezioni individualmente su ogni  $\mathcal{K}_i$ , con  $i = 1, \ldots, r$ ; se dunque  $g_{t,i-1}$  è l'iterata corrente e  $P_i$  è la proiezione sull'insieme  $\mathcal{K}_i$ , si scrive  $g_{t,i} = P_i(g_{t,i-1})$ .

La modifica di Dykstra consiste nell'applicare la proiezione sull'insieme  $\mathcal{K}_i$  del ciclo t non all'iterata corrente  $g_{t,i-1}$  ma al termine  $(g_{t,i-1} - I_{t-1,i})$ , ovvero all'iterata corrente diminuita dell'incremento  $I_{t-1,i}$  corrispondente alla proiezione su  $\mathcal{K}_i$  nel ciclo precedente t-1.

Dunque, se g è il dato iniziale si ha

$$P_{1}(g) = g_{1,1} = g + I_{1,1}$$

$$P_{2}(g_{1,1}) = g_{1,2} = g_{1,1} + I_{1,2}$$

$$\vdots$$

$$P_{r}(g_{1,r-1}) = g_{1,r} = g_{1,r-1} + I_{1,r}$$

Poi per il secondo ciclo,

$$P_1(g_{1,r} - I_{1,1}) = g_{2,1} = g_{1,r} - I_{1,1} + I_{2,1}$$
$$P_2(g_{2,1} - I_{1,2}) = g_{2,2} = g_{2,1} - I_{1,2} + I_{2,2}$$
$$\vdots$$

Continuando in questo modo si genera una successione infinita di  $(g_{t,i})$  che converge alla soluzione  $g^*$ , come provato in [8] e [4].

Si riporta in Figura 2.1 un'idea dell'algoritmo appena descritto. In figura gli insiemi  $\mathcal{K}_i^*$  corrispondono ai coni normali  $\partial \mathcal{K}_i$ .



Figura 2.1: Disegno schematico dell'algoritmo [8]

Come affermato da Boyle e Dykstra in [4], se l'insieme sul quale si proietta è un sottospazio affine, essendo la proiezione un'operazione lineare, non è necessario applicare la correzione corrispondente: siccome  $\mathcal{U}$  è un traslato di un sottospazio

$$\mathcal{U} = \{ Y = N + I \in \mathbb{R}^{n \times n} : N = N^T, n_{ii} = 0, \forall i = 1, ..., n \},\$$

la correzione della proiezione su  $\mathcal{U}$  viene omessa.

É presentato l'algoritmo del metodo di proiezione alternata, con la correzione di Dykstra.

Algoritmo 1: (Metodo delle proiezioni alternate) Data una matrice simmetrica  $A \in \mathbb{R}^{n \times n}$  si calcola la matrice di correlazione più vicina ad A in norma-W.

**Dati:**  $\Delta S_0 = 0$ ,  $Y_0 = A$ **Risultato:** X soluzione di (1.1)1 for  $k = 1, \ldots$ , convergenza do  $R_k = Y_{k-1} - \Delta S_{k-1}$  % $\Delta S_{k-1}$  è la correzione di Dykstra  $\mathbf{2}$  $X_k = P_{\mathcal{S}}(R_k)$ 3  $\Delta S_k = X_k - R_k$ 4  $Y_k = P_{\mathcal{U}}(X_k)$  $\mathbf{5}$ if test di convergenza then 6  $X = Y_k$ 7 quit 8 end 9 10 end

Il seguente teorema da [6] (Teorema 9.24) e da [4] (Teorema 2) assicura che sia  $X_k$  che  $Y_k$  convergono alla matrice di correlazione cercata per  $k \to \infty$ .

**Teorema 2.3.** Siano  $\mathcal{K}_1, \mathcal{K}_2, \ldots, \mathcal{K}_r$  sottoinsiemi convessi chiusi di uno spazio di Hilbert H, tali che  $\mathcal{K} := \bigcap_{i=1}^r \mathcal{K}_i \neq \emptyset$ . Dato  $x \in H$ , sia  $\{x_n\}_{n \in \mathbb{N}}$  la successione generata dall'Algoritmo 1, allora  $\lim_{n\to\infty} ||x_n - P_{\mathcal{K}}(x)||_H = 0$ , con  $||\cdot||_H$  la norma dello spazio di Hilbert H.

La velocità di convergenza è lineare quando gli insiemi di proiezione sono sottospazi, con la costante proporzionale all'angolo tra i due sottospazi [7]. In questo caso, non avendo sottospazi ma insiemi chiusi convessi, la convergenza lineare è una prospettiva ottimistica.

È inoltre evidente che il calcolo della decomposizione spettrale della matrice ad ogni iterazione (cioè la proiezione su S) rallenta significativamente l'algoritmo: risulterà infatti dagli esperimenti numerici che la maggior parte del tempo necessario all'esecuzione dell'Algoritmo 1 sarà impiegata per la decomposizione spettrale dell'iterata.

Per quanto riguarda il test di convergenza dell'Algoritmo 1, il criterio d'arresto

utilizzato è

$$\max\left\{\frac{\|X_k - X_{k-1}\|_F}{\|X_k\|_F}, \frac{\|Y_k - Y_{k-1}\|_F}{\|Y_k\|_F}, \frac{\|Y_k - X_k\|_F}{\|Y_k\|_F}\right\} \le \text{tol}$$
(2.1)

dove tol è un valore di tolleranza fissato.

Le tre quantità nel test sono spesso dello stesso ordine di grandezza, quindi una qualsiasi tra esse può essere usata nel criterio d'arresto. Negli esperimenti del prossimo capitolo verrà usata solo la terza quantità in quanto analoga ad un criterio d'arresto per l'algoritmo di Dykstra provato da Birgin e Raydan in [2].

#### 2.2.1 Il caso di matrici non semidefinite positive

**Teorema 2.4.** È data  $A = A^T$  con elementi diagonali  $a_{ii} \ge 1$  e W diagonale. Sia  $Y_k = P_{\mathcal{U}}(X_k) = X_k + D_k$  dove  $D_k$  è diagonale (per il Teorema 2.1). Allora nell'Algoritmo 1

$$R_k = A + \Delta_k,$$

dove  $\Delta_k = \sum_{i=1}^{k-1} D_i$  è una matrice semidefinita negativa. Dimostrazione. Sapendo che  $R_1 = A$  si ha

 $r_{\rm III} = r_{\rm III} = r_{\rm III}$ 

$$R_{k+1} = Y_k - \Delta S_k = P_{\mathcal{U}}(X_k) - \Delta S_k = X_k + D_k - (X_k - R_k)$$
$$= R_k + D_k = \dots = R_1 + D_1 + \dots + D_k,$$

da cui la tesi. Inoltre

$$I = \text{DIAG}(Y_k) = \text{DIAG}(X_k + D_k)$$
  
=  $\text{DIAG}(P_S(R_k) + D_k)$   
 $\geq \text{DIAG}(R_k + D_k) \quad (\text{per il Teorema 2.2})$   
=  $\text{DIAG}(A + \Delta_{k+1})$   
 $\geq I + \Delta_{k+1}$   
 $\Rightarrow \Delta_{k+1} \leq 0.$ 

Il teorema precedente mostra che  $R_k$  è uguale ad A diminuita di una matrice diagonale semidefinita positiva; questo ha un'importante conseguenza nel caso in cui A sia "altamente" definita non positiva o anche a rango basso. **Corollario 2.5.** Sia  $A = A^T$  con elementi diagonali  $a_{ii} \ge 1$ , t autovalori non positivi e sia W diagonale. Allora nell'Algoritmo 1  $R_k$  ha almeno t autovalori non positivi e  $X_k$  ha almeno t autovalori nulli, per ogni k.

Dimostrazione. Segue dal Teorema 2.4 che  $R_k = A - (-\Delta_k)$ . Se  $A = Q\Lambda Q^T$  con autovalori  $\lambda_1 \ge \ldots \ge \lambda_{n-t} > 0 \ge \lambda_{n-t+1} \ge \ldots \ge \lambda_n$  e  $(-\Delta_k)$  semidefinita positiva segue che  $R_k$  ha  $p \ge t$  autovalori non positivi. Essendo  $X_k$  la proiezione di  $R_k$  sulle matrici semidefinite positive, avrà  $p \ge t$  autovalori nulli.

Osservazione. Mandando  $k \to \infty$ nel Corollario 2.5 si ritrova la seconda parte del Teorema 1.6.

Il Corollario 2.5 permette di distinguere due casi:

- Se  $t \gg 1$  (se è grande rispetto alla dimensione dei dati) allora è possibile ottenere  $P_{\mathcal{S}}(R_k)$  senza risolvere l'intero problema di autovalori e autovettori della matrice  $W^{1/2}R_kW^{1/2}$  (vedere Teorema 2.2). È sufficiente calcolare gli  $n - t \ll n$ autovalori più grandi  $\lambda_j$  e i corrispondenti autovettori ortonormali  $q_j$  e inserire nel Teorema 2.2 l'espressione

$$(W^{1/2}R_kW^{1/2})_+ = \sum_{\lambda_j>0} \lambda_j q_j q_j^T.$$

Questo può essere fatto in maniera efficiente riducendo ortogonalmente  $R_k$  a forma tridiagonale e applicando il metodo di bisezione seguito dal metodo delle potenze inverse per gli autovettori [20];

- Se  $t \ll n$ , cioè A ha pochi autovalori non positivi, è probabile che lo stesso valga anche per  $R_k$ . È dunque possibile operare in termini di risparmio computazionale calcolando autovalori e autovettori non positivi:8

$$(W^{1/2}R_kW^{1/2})_+ = W^{1/2}R_kW^{1/2} - \left(\sum_{\lambda_j \le 0} \lambda_j q_j q_j^T\right),$$

dove il numero di autovalori positivi di  $W^{1/2}R_kW^{1/2}$  è stimato ad ogni iterazione.

#### 2.3 Risultati numerici

Gli esperimenti numerici sono eseguiti su MATLAB R2022a, con unità di roundoff  $u = 2^{-53} \approx 1.1 \times 10^{-16}$ , su PC con processore Intel Core i5-4430S @2.70GHz e 16GB di RAM. Per gli esperimenti riportati si ha W = I, di conseguenza la norma-W coincide con la norma Frobenius.

• Esperimento 2.1. È data la matrice definita positiva

$$A = \begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix},$$

con tol= $10^{-8}$ . L'Algoritmo 1 converge in 19 iterazioni e la soluzione trovata è

$$X = \begin{bmatrix} 1 & -0.8084 & 0.1916 & 0.1068 \\ -0.8084 & 1 & -0.6562 & 0.1916 \\ 0.1916 & -0.6562 & 1 & -0.8084 \\ 0.1068 & 0.1916 & -0.8084 & 1 \end{bmatrix}.$$

La matrice X ha rango 3 e vale  $||A - X||_W = 2.1337$ .

• Esperimento 2.2. A partire dalla raccolta di matrici di correlazione (randcorr) disponibile all'interno della funzione gallery di MATLAB, sono fissate 10 matrici di dimensione 500 × 500.

In un primo caso le matrici di correlazione vengono perturbate tramite una matrice simmetrica casuale E con elementi dell'ordine di  $10^{-4}$ , per simulare casi in cui la matrice dei dati viene interamente modificata. Con tol=  $10^{-6}$ , l'Algoritmo 1 converge mediamente (tra le 10 matrici) in 3 iterazioni, con una distanza media  $||A - X||_W = 2.2 \cdot 10^{-3}$ . Inoltre la distanza media tra la soluzione X e la matrice di correlazione originale è  $||R - X||_W = 3.5 \cdot 10^{-2}$ .

Nel secondo caso, ogni elemento delle matrici di correlazione che in modulo è minore o pari di una soglia fissata a 0.1 viene sostituito con uno zero. In media l'algoritmo converge in una sola iterazione, con  $||A - X||_W = 1.7 \cdot 10^{-15}$ . La distanza media tra la matrice di correlazione originale R e la soluzione trovata è però

 $||R - X||_W = 12.7$ , vale a dire che la soluzione X è molto lontana dalla matrice che contiene i dati corretti.

L'ipotesi più probabile è che la perturbazione effettuata nel secondo caso vada a "spostare" significativamente la matrice e che quindi l'algoritmo converga ad una nuova matrice di correlazione, molto vicina alla matrice perturbata A ma drasticamente diversa da quella originaria R. Osservando nello specifico la matrice A si nota infatti che quest'ultima è una matrice sparsa, e suggerisce l'idea che la matrice originale R contenga molti elementi che soddisfano la condizione  $|r_{ij}| \leq 0.1$ . Se infine si abbassa la soglia a  $10^{-2}$ , le iterazioni necessarie all'Algoritmo 1 per convergere diventano 3, si ha  $||A - X||_W = 1.0 \cdot 10^{-3}$  e  $||R - X||_W = 1.57$ , per cui la soluzione X diventa una buona approssimazione della matrice di correlazione originale R.

• Esperimento 2.3. Sia A una matrice del pacchetto di matrici di correlazione "invalide" fornito da Higham [10], di dimensione  $1399 \times 1399$ . A è una matrice di dati azionari resa disponibile da una compagnia di gestione di fondi finanziari, con diagonale unitaria, elementi in modulo minori o uguali a 1, ma non semidefinita positiva. Inoltre, siccome ha 1245 autovalori non positivi, per il Teorema 1.6 la soluzione X del problema 1.1 deve avere almeno 1245 autovalori nulli, e dunque rango 154 al massimo. Una matrice di questo tipo è frequente in applicazioni di tipo finanziario, soprattutto quando viene costruita a partire da poche osservazioni.

Viene applicato l'Algoritmo 1 con tol= $10^{-4}$ . Si ha convergenza dopo 67 iterazioni, con  $||A - X||_W = 20.9621$  e una soluzione X con 154 autovalori positivi (si intende maggiori di  $10^{-2}$ , gli altri autovalori non nulli risultano dell'ordine di  $10^{-14}$ , dunque trascurabili). Si può sfruttare il largo numero di autovalori non positivi della matrice A per evitare di fare una completa decomposizione spettrale ad ogni ciclo, calcolando solo i 154 autovalori più grandi in modulo (è utile aumentare leggermente il numero 154 per avere un margine di sicurezza e controllare di aver ottenuto tutti gli autovalori non trascurabili).

In questo modo infatti l'Algoritmo 1 converge alla soluzione X in 29 iterazioni, mantenendo una distanza in norma-W dalla matrice A accettabile,  $||A - X||_W =$ 24.1683. In più, con questa strategia il tempo computazionale per la convergenza passa da 39.79 a 11.52 secondi, mostrando come la maggior parte del tempo di esecuzione dell'algoritmo sia impiegato per la decomposizione spettrale.

### Capitolo 3

# Accelerazione di Anderson

Il metodo delle proiezioni alternate descritto nel Capitolo 2 è forse l'approccio più largamente utilizzato per la risoluzione del problema NCM. Purtroppo, come già detto, l'Algoritmo 1 ha velocità di convergenza solo lineare. Può quindi impiegare un grande numero di iterazioni per convergere entro una certa tolleranza fissata.

In questo Capitolo si cerca di velocizzare l'algoritmo delle proiezioni alternate adottando l'accelerazione di Anderson, una tecnica per aumentare la velocità di convergenza delle iterazioni di punto fisso, che porterà ad una significativa riduzione del numero di iterazioni e del tempo computazionale necessario alla convergenza dell'algoritmo.

In breve, l'accelerazione di Anderson per un metodo di punto fisso, non utilizza solo l'iterata corrente per determinare quella successiva ma aggiunge anche l'informazione delle ultime  $m_k$  iterazioni. É nota in chimica quantistica come DIIS (*Direct Inversion in the Iterative Subspace*) dove ha prodotto ottimi risultati nel calcolo delle strutture elettroniche degli atomi; in Analisi Numerica, per sistemi lineari coincide con il metodo generalizzato dei minimi residui (GMRES) se  $m_k = k$ (k numero dell'iterazione corrente), mentre per problemi non lineari si può dimostrare che converge localmente con velocità lineare. Sebbene non ci siano garanzie generali di convergenza, l'accelerazione di Anderson ha in generale un buon potenziale per migliorare la convergenza del metodo delle proiezioni alternate per il problema NCM, come si evince dai risultati sperimentali.

#### 3.1 Accelerazione di Anderson per metodi di punto fisso

Problema (Problema di punto fisso). Data la funzione  $g: \mathbb{R}^n \to \mathbb{R}^n$ , determinare x tale che si abbia

$$g(x) = x. \tag{3.1}$$

I problemi di punto fisso sono spesso formulati per la risoluzione di equazioni non lineari del tipo f(x) = g(x) - x = 0.

Un metodo per trovare la soluzione del problema (3.1) è l'iterazione di punto fisso (o metodo delle iterate successive): dato  $x_0 \in \mathbb{R}^n$  punto iniziale, l'iterazione è data da

$$x_{k+1} = g(x_k), \quad k \ge 0.$$
 (3.2)

La convergenza dell'iterazione (3.1) è garantita sotto certe ipotesi di regolarità della funzione g (vedi Teorema 4.2.1 in [15]).

Un modo per aumentare la velocità di convergenza dell'iterazione di punto fisso consiste nell'usare l'accelerazione di Anderson, che calcola  $x_{k+1}$  utilizzando anche le  $m_k$  iterate precedenti. Il parametro  $m_k$  è fissato a m una volta fatte le prime miterazioni.

Si introduce l'accelerazione di Anderson per metodi di punto fisso.

Algoritmo 2: (Accelerazione di Anderson originale) Dato  $x_0$  punto iniziale l'algoritmo produce la successione  $\{x_k\}$  convergente al punto fisso della funzione  $g: \mathbb{R}^n \to \mathbb{R}^n$ 

**Dati:**  $x_0 \in \mathbb{R}^n$ ,  $m \in \mathbb{Z}$  t.c.  $m \ge 1$ ,  $g : \mathbb{R}^n \to \mathbb{R}^n$  **Risultato:**  $x^*$  punto fisso di g1  $x_1 = g(x_0)$ 2 for  $k = 1, 2, \dots$ , convergenza do

- $\mathbf{s} \quad m_k = \min(m, k)$
- 4 Determina  $\theta^{(k)} = (\theta_1^{(k)}, \dots, \theta_{m_k}^{(k)})^T \in \mathbb{R}^{m_k}$  che minimizza  $||u_k v_k||_2^2$ , dove

$$\iota_k = x_k + \sum_{j=1}^{m_k} \theta_j(x_{k-j} - x_k) \in v_k = g(x_k) + \sum_{j=1}^{m_k} \theta_j \big( g(x_{k-j}) - g(x_k) \big)$$

5  $x_{k+1} = v_k$  usando i parametri in  $\theta^{(k)}$ .

6 end

Nella prima versione dell'algoritmo l'ultimo passo dell'iterazione era  $x_{k+1} = u_k + \beta_k (v_k - u_k)$ , dove  $u_k$  e  $v_k$  sono calcolati con  $\theta^{(k)}$  e  $\beta_k > 0$  è determinato empiricamente; spesso, come in questo caso, il valore più usato è  $\beta_k = 1$ .

Si possono riscrivere  $u_k \in v_k$  in questo modo:

$$u_{k} = \left(1 - \sum_{j=1}^{m_{k}} \theta_{j}^{(k)}\right) x_{k} + \sum_{j=1}^{m_{k}} \theta_{j}^{(k)} x_{k-j} = \sum_{j=0}^{m_{k}} w_{j} x_{k-j},$$
$$v_{k} = \left(1 - \sum_{j=1}^{m_{k}} \theta_{j}^{(k)}\right) g(x_{k}) + \sum_{j=1}^{m_{k}} \theta_{j}^{(k)} g(x_{k-j}) = \sum_{j=0}^{m_{k}} w_{j} g(x_{k-j}),$$
$$\cos w_{j} = \begin{cases} 1 - \sum_{j=1}^{m_{k}} \theta_{j} & j = 0\\ \theta_{j} & j \ge 1. \end{cases}$$

Inoltre

$$\sum_{j=0}^{m_k} w_j = w_0 + \sum_{j=1}^{m_k} w_j = \left(1 - \sum_{j=1}^{m_k} \theta_j\right) + \sum_{j=1}^{m_k} \theta_j = 1.$$

In conclusione l'Algoritmo 2 cerca di minimizzare la quantità  $||u_k - v_k||_2^2$  con il vincolo  $\sum_{j=0}^{m_k} w_j = 1$ .

Se g è lineare allora la funzione da minimizzare diventa

$$||u_k - g(u_k)||_2^2$$
,

in quanto  $v_k = \sum_j w_j g(x_{k-j}) = g(\sum_j w_j x_{k-j}) = g(u_k)$ . Di conseguenza  $v_k$  è un vettore del sottospazio affine generato dalle ultime  $m_k + 1$  iterazioni, che minimizza il residuo dell'equazione di punto fisso.

É utile considerare una forma equivalente dell'Algoritmo 2 che memorizza in due matrici rispettivamente le differenze tra iterate successive e le differenze tra le corrispondenti immagini delle iterate attraverso la funzione f, con f(x) = g(x) - x. **Algoritmo 3:** (Accelerazione di Anderson) Dato  $x_0$  punto iniziale l'algoritmo produce la successione  $\{x_k\}$  convergente ad uno zero della funzione  $f : \mathbb{R}^n \to \mathbb{R}^n$ 

Notazioni:  $f_i = f(x_i), \quad \Delta x_i = x_{i+1} - x_i, \quad \mathcal{X}_k = [\Delta x_{k-m_k} \dots \Delta x_{k-1}],$  $\Delta f_i = f_{i+1} - f_i, \ \mathcal{F}_k = \left[\Delta f_{k-m_k} \dots \Delta f_{k-1}\right]$ **Dati:**  $x_0 \in \mathbb{R}^n$ ,  $m \in \mathbb{Z}$  t.c.  $m \ge 1$ ,  $f : \mathbb{R}^n \to \mathbb{R}^n$ **Risultato:**  $x^*$  zero di f1  $f_0 = f(x_0)$ **2**  $x_1 = x_0 + f_0$ **3**  $f_1 = f(x_1)$  $4 \Delta x_0 = x_1 - x_0$ **5**  $\Delta f_0 = f_1 - f_0$ 6 for  $k = 1, 2, \ldots$ , convergenza do  $m_k = \min(m, k)$ 7  $\mathcal{X}_k = [\Delta x_{k-m_k} \dots \Delta x_{k-1}]$ (aggiornamento e riduzione di  $\mathcal{X}_k$ ) 8  $\mathcal{F}_k = [\Delta f_{k-m_k} \dots \Delta f_{k-1}]$ (aggiornamento e riduzione di  $\mathcal{F}_k$ ) 9 Determina  $\mathbf{10}$  $\gamma^{(k)} = (\gamma_{k-m_k}^{(k)}, \dots, \gamma_{k-1}^{(k)})^T \in \mathbb{R}^{m_k} \text{ che risolve } \min_{\gamma \in \mathbb{R}^{m_k}} \|f_k - \mathcal{F}_k \gamma\|_2^2$  $\bar{x}_k = x_k - \sum_{i=k-m_k}^{k-1} \gamma_i^{(k)} \Delta x_i = x_k - \mathcal{X}_k \gamma^{(k)}$ 11  $\bar{f}_k = f_k - \sum_{i=k-m_k}^{k-1} \gamma_i^{(k)} \Delta f_i = f_k - \mathcal{F}_k \gamma^{(k)}$ 12 $x_{k+1} = \bar{x}_k + \bar{f}_k$ 13  $f_{k+1} = f(x_{k+1})$  $\mathbf{14}$  $\Delta x_k = x_{k+1} - x_k$ 15 $\Delta f_k = f_{k+1} - f_k$  $\mathbf{16}$ 17 end

La riga 10 dell'Algoritmo 3 contiene la parte più consistente dei calcoli.

Si suppone che  $\mathcal{F}_k$  abbia rango massimo e il problema ai minimi quadrati sia risolto tramite fattorizzazione QR.

La fattorizzazione di  $\mathcal{F}_k$  si ottiene a partire da quella di  $\mathcal{F}_{k-1}$ ; siccome  $\mathcal{F}_k$  è ottenuta da  $\mathcal{F}_{k-1}$  rimuovendo la prima colonna (se  $k \ge m$ ) e aggiungendo una colonna alla fine, si tratta di un problema di aggiornamento di fattorizzazione QR. Fatto l'aggiornamento, si risolve il problema  $\min_{\gamma \in \mathbb{R}^{m_k}} \|f_k - \mathcal{F}_k\gamma\|_2^2$ . Sia  $k \ge m$ . La matrice  $\mathcal{F}_k$  ha dimensione  $n \times m$  e il costo per l'aggiornamento di  $Q \in R$  dopo l'eliminazione della prima colonna e l'aggiunta dell'ultima è di  $m^2/2 + 10mn + 3n$  flops. Il costo della risoluzione dei minimi quadrati tramite sostituzione all'indietro del sistema triangolare con R è di  $2mn + m^2$  flops. Di conseguenza il metodo accelerato richiede  $3m^2/2 + 12mn + 3n$  operazioni aggiuntive per ciclo, rispetto all'iterazione non accelerata.

#### 3.2 Accelerazione del metodo delle proiezioni alternate

Introdotta l'idea dell'accelerazione di Anderson, questa viene applicata al metodo delle proiezioni alternate presentato nel precedente capitolo, come in [13]. Il test di convergenza utilizzato per l'Algoritmo 1 si riduce solo al terzo termine

$$\frac{\|Y_k - X_k\|_F}{\|Y_k\|_F},$$

in quanto corrisponde ad un robusto criterio d'arresto per l'algoritmo di Dykstra, come indicato in [2].

Per velocizzare l'Algoritmo 1 è però necessario che questo venga riformulato come metodo di punto fisso. Bisogna quindi definire la funzione g del problema (3.1). Le variabili che ricorrono sono  $Y_k \in \Delta S_k$ , mentre  $X_k$  è usata per il test di convergenza.

Sia quindi l'Algoritmo 1 in forma di punto fisso:

**Algoritmo 4:** (Forma di punto fisso dell'Algoritmo 1) Data una matrice simmetrica  $A \in \mathbb{R}^{n \times n}$  si calcola la matrice di correlazione più vicina ad A in norma-W. È richiesta una tolleranza tol per il test di convergenza.

Dati:  $\Delta S_0 = 0$ ,  $Y_0 = A$ , tol Risultato: X soluzione del problema (1.1) 1 for k = 1, ..., convergenza do 2  $| [X_k, Y_k, \Delta S_k] = g(Y_{k-1}, \Delta S_{k-1})$ 3  $| \mathbf{if} || Y_k - X_k ||_F \le \text{tol} || Y_k ||_F$  then 4  $| X = Y_k$ 5 | quit6 | end7 end La funzione g contiene precisamente i calcoli dell'Algoritmo 1:

 $R_k = Y_{k-1} - \Delta S_{k-1}$  $X_k = P_S(R_k)$  $\Delta S_k = X_k - R_k$  $Y_k = P_U(X_k)$ 

Si applica quindi l'accelerazione di Anderson al metodo delle proiezioni alternate come punto fisso, nella forma equivalente dell'Algoritmo 3, in cui al posto delle matrici si utilizza la loro immagine attraverso l'operatore vec. Segue l'algoritmo completo.

Algoritmo 5: (Accelerazione di Anderson del metodo delle proiezioni alternate in forma di punto fisso) Data una matrice simmetrica  $A \in \mathbb{R}^{n \times n}$  l'algoritmo calcola la matrice di correlazione più vicina ad A tramite le proiezioni alternate con l'accelerazione di Anderson. Richiede una tolleranza tol per il test di convergenza.

Notazioni:  $z_k = \operatorname{vec}(Z_k), \quad Z_k = [Y_k, \Delta S_k] \in \mathbb{R}^{n \times 2n}$ Dati:  $\Delta S_0 = 0, \quad Y_0 = A, \quad z_0 \in \mathbb{R}^{2n^2}, \quad tol,$  $f : \mathbb{R}^{2n^2} \to \mathbb{R}^{2n^2}, \quad g : \mathbb{R}^{2n^2} \to \mathbb{R}^{n \times 3n}$ 

**Risultato:** X soluzione del problema (1.1)

1 Eseguire l'Algoritmo 2 sulla funzione

$$f(z) = \operatorname{vec}(\widetilde{g}(Z) - Z),$$

con  $\tilde{g}(Z_k) = Z_{k+1}$  e la funzione  $g(Z_k) = [X_k, \tilde{g}(Z_k)]$  definita nell'Algoritmo 4. Utilizzare il test di convergenza  $||Y_k - X_k||_2 \le \text{tol } ||Y_k||_2$ . Denotare il risultato con  $x_*$ .

2  $X = \operatorname{unvec}(x_*)$ 

Il test di convergenza utilizzato per l'Algoritmo 5 è equivalente a quello dell'Algoritmo 4.

Convergenza del metodo accelerato: a differenza degli Algoritmi 1 e 4 che convergono per i risultati dei capitoli precedenti, l'Algoritmo 5 potrebbe non arrivare a convergenza in quanto non ci sono risultati sulla convergenza dell'accelerazione di Anderson. La convergenza di un metodo di punto fisso con accelerazione di Anderson sotto determinate ipotesi è ancora un problema aperto.

Il costo computazionale per iterazione dell'Algoritmo 1 è dominato dalla decomposizione spettrale, che è  $10n^3$  flops per la decomposizione completa e  $14n^3/3$ flops usando la procedura di tridiagonalizzazione seguita dalla bisezione e dalle potenze inverse per matrici con pochi autovalori positivi. Un ciclo dell'accelerazione di Anderson applicata al metodo di proiezione alternata come punto fisso ha un costo computazionale di  $3m^2/2 + 24mn^2 + 6n^2$  flops, in quanto usa vettori lunghi  $2n^2$ . Poichè sperimentalmente si verifica che è sufficiente prendere  $m \leq 5$ , il costo aggiuntivo dell'accelerazione è un  $O(n^2)$ , trascurabile per n grande.

L'accelerazione comporta anche un aumento di  $2mn^2$  elementi da memorizzare.

#### 3.3 Varianti del metodo

Vengono presentate ora due modifiche dell'algoritmo delle proiezioni alternate per il problema NCM. La prima in cui specifici elementi della matrice A sono fissati, la seconda in cui viene richiesto che l'autovalore minimo superi una tolleranza positiva  $\delta$ .

La possibilità di imporre vincoli di questo genere, di variare l'algoritmo a seconda di diverse necessità e di applicare una strategia come quella dell'accelerazione di Anderson rende il metodo delle proiezioni alternate l'approccio più utilizzato per la risoluzione del problema NCM, a discapito di metodi con convergenza più veloce ma poco flessibili e poco pratiche nel caso di varianti del problema (ad esempio, il metodo quasi-Newton utilizzato da Qi e Sun in [18]).

Gli esperimenti numerici della Sezione 3.4 mostrano che l'aggiunta dei vincoli sopra citati aumenta di molto il numero di iterazioni necessarie alla convergenza, rispetto alla risoluzione del problema "*classico*". Inoltre è importante notare che le due modifiche possono essere integrate simultaneamente nell'Algoritmo 1, ovviamente incrementando significativamente il numero di iterazioni e il tempo di esecuzione.

Dal momento che le varianti appena descritte sono comunemente usate, aumenta l'interesse nell'applicazione frequente dell'accelerazione di Anderson al metodo delle proiezioni alternate, grazie alla sua versatilità nella risoluzione del problema NCM.

#### 3.3.1 Strategia di vincolo sugli elementi della matrice di correlazione

L'esigenza di fissare determinati elementi della matrice nasce da analisi di tipo statistico nelle quali mancano alcune osservazioni.

Si assuma che, in una matrice dei dati  $m \times n$ , siano presenti tutte le osservazioni delle prime  $n_1$  colonne (sempre possibile a meno di permutare l'ordine delle variabili).

Un metodo per calcolare i coefficienti di correlazione tra le variabili è quello dell'eliminazione a coppie, dove il coefficiente tra due variabili viene calcolato usando solo le componenti presenti in entrambi i vettori. In questo modo si ottiene una matrice simmetrica della forma

$$C = \begin{bmatrix} A & Y \\ Y^T & B \end{bmatrix} \in \mathbb{R}^{n \times n},$$

dove  $A \in \mathbb{R}^{n_1 \times n_1}$  è effettivamente una matrice di correlazione in quanto costruita a partire da dati senza osservazioni mancanti. Non c'è però alcuna garanzia che la matrice C sia semidefinita positiva. Si cerca quindi di calcolare la matrice di correlazione più vicina a C, ma dal momento che i coefficienti presenti in A sono considerati esatti, si richiede che il blocco A rimanga invariato.

In finanza, per valutare la resistenza di precise strategie a situazioni estreme, o per vedere gli effetti di una crisi economica, spesso si modifica solo una parte degli strumenti finanziari analizzati; siccome la modifica rende la matrice indefinita, si cerca la matrice di correlazione più vicina a quella modificata, richiedendo che il blocco della matrice di correlazione che corrisponde al gruppo di strumenti finanziari non modificati rimanga immutato.

Si denota con E la matrice degli elementi vincolati (matrice di *pattern*), i cui elementi non nulli sono quelli che devono rimanere invariati nella ricerca della NCM. Questa variante del problema (1.1) cerca dunque la matrice più vicina ad A (in norma Frobenius) che appartenga all'intersezione tra l'insieme S e l'insieme

$$\mathcal{E} = \left\{ Y = Y^T \in \mathbb{R}^{n \times n} : y_{ii} = 1, \ i = 1, \dots, n \in y_{ij} = e_{ij} \ \forall (i,j) \in \mathcal{N} \right\},\$$

dove  $\mathcal{N}$  denota l'insieme degli indici relativi agli elementi non diagonali fissati. Chiaramente se  $(i, j) \in \mathcal{N}$  allora  $(j, i) \in \mathcal{N}$ .

Per visualizzare meglio il concetto, si riportano in Figura 3.1 degli esempi di matrice di *pattern*, relativi a due delle matrici dell'Esperimento 3.3 della Sezione 3.4. Il numero degli elementi fissati di ogni matrice (indicati in blu) è evidenziato sotto ogni grafico



Figura 3.1: Pattern delle matrici di ordine 70 e 94 dell'Esperimento 3.3

Similmente al ragionamento fatto nel Capitolo 1, dal momento che l'intersezione  $S \cap \mathcal{E}$  è non vuota il problema ha soluzione unica, ma questo accade solo se l'insieme  $\mathcal{N}$  è scelto in maniera opportuna. Infatti non è scontato che, dati gli indici degli elementi da fissare, esista sempre una matrice di correlazione con gli elementi fissati richiesti.

Esempio. Sia

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

e  $\mathcal{N} = \{(2,3), (3,2), (2,4), (4,2), (3,4), (4,3)\}$ . Non è possibile sostituire A con una matrice di correlazione che mantenga fissati gli elementi richiesti in  $\mathcal{N}$ , poichè questi corrispondono ad un blocco  $3 \times 3$  di A che è indefinito.

Il metodo delle proiezioni alternate incorpora facilmente il vincolo degli elementi fissati sostituendo la proiezione sull'insieme  $\mathcal{U}$  con quella sull'insieme  $\mathcal{E}$  al passo 5 dell'Algoritmo 1.

Data una matrice simmetrica A, con E matrice di pattern, la proiezione  $H = P_{\mathcal{E}}(A)$ è data da

$$h_{ij} = \begin{cases} 1 & i = j, \\ e_{ij} & (i,j) \in \mathcal{N}, \\ a_{ij} & \text{altrimenti.} \end{cases}$$

Dal momento che si considera  $\mathcal{N}$  privo degli indici relativi agli elementi diagonali,  $P_{\mathcal{E}}$  rimane ben definita anche se A non ha diagonale unitaria.

Il fatto che la soluzione per questa variante del problema NCM possa non esistere si deve riflettere anche sul test di convergenza.

Nell'esempio precedente si vede facilmente che, applicando l'Algoritmo 1, le matrici  $X_k$  e  $Y_k$  sono costanti per  $k \ge 1$ , per cui i primi due termini del test in (2.1) sono costantemente nulli, mentre il terzo termine è non nullo per ogni k. L'utilizzo del terzo termine è quindi adatto alla variante del problema anche nel caso in cui non esista una soluzione.

#### 3.3.2 Ricerca di una matrice di correlazione definita positiva

Il Teorema 1.6 dice che se la matrice A di partenza ha t autovalori nulli, anche la matrice di correlazione X soluzione del problema (1.1) avrà almeno t autovalori nulli, ovvero sarà una matrice singolare. La singolarità è un problema in analisi multivariata, ad esempio quando è necessaria la matrice inversa di una matrice di correlazione, per la visualizzazione della dispersione dei dati.

É dunque di grande utilità calcolare la matrice di correlazione X più vicina ad una matrice data A con il vincolo  $\lambda_{min}(X) \geq \delta$ , dove  $\lambda_{min}(X)$  è l'autovalore minimo della matrice X e  $\delta$  è una tolleranza positiva fissata. Siccome per una matrice di correlazione vale tr $(X) = n = \sum_i \lambda_i(X)$  segue che  $\delta \leq 1$ .

Applicando questa modifica al problema la soluzione che si ottiene è una matrice

di correlazione strettamente definita positiva.

Dato  $0 \le \delta \le 1$ , si definisce l'insieme

$$\mathcal{S}^{\delta} = \left\{ Y = Y^T \in \mathbb{R}^{n \times n} : \lambda_{min}(Y) \ge \delta \right\}.$$

Se  $\delta = 0$ , si ha  $\mathcal{S}^0 = \mathcal{S}$ .

Questa variante del problema (1.1) cerca la matrice più vicina ad A (in norma Frobenius) che stia nell'intersezione  $S^{\delta} \cap \mathcal{U}$ .

L'intersezione  $S^{\delta} \cap U$  è non vuota, chiusa e convessa in quanto  $I \in S^{\delta} \quad \forall \delta \in [0, 1]$ . Di conseguenza questa variante del problema NCM ha soluzione unica.

Anche in questo caso, il metodo delle proiezioni alternate incorpora in maniera semplice la modifica appena descritta, sostituendo la proiezione su S con quella su  $S^{\delta}$  al passo 3 dell'Algoritmo 1.

Una formula esplicita per la proiezione di una matrice simmetrica  $A \in \mathbb{R}^{n \times n}$ sull'insieme  $\mathcal{S}^{\delta}$  è data dal seguente risultato di Cheng e Higham [5].

**Teorema 3.1.** Sia  $A \in \mathbb{R}^{n \times n}$  matrice simmetrica con decomposizione spettrale  $A = QDQ^T$ , dove  $D = \text{diag}(\lambda_1, \ldots, \lambda_n)$ . Sia  $\delta \ge 0$ . Allora la matrice più vicina ad A che abbia autovalore minimo almeno  $\delta$  è data data

$$P_{\mathcal{S}^{\delta}}(A) = Q \operatorname{diag}(\tau_i) Q^T, \quad \tau_i = \begin{cases} \lambda_i & \text{se } \lambda_i \ge \delta, \\ \delta & \text{se } \lambda_i < \delta. \end{cases} \quad i = 1, \dots, n$$

Infine, l'implementazione del vincolo con il  $\delta$  permette di affrontare un problema relativo alla matrice di correlazione più vicina. Nelle proiezioni alternate, a seconda dell'ordine delle proiezioni individuali, si può ottenere una matrice con diagonale unitaria ma indefinita o, contrariamente, una matrice semidefinita positiva che non ha diagonale unitaria.

Probabilmente la soluzione migliore a questo problema è imporre un  $\delta$  sufficientemente grande in modo che la definizione della matrice non sia influenzata da perturbazioni dell'ordine di grandezza della tolleranza tol (ad esempio, tol  $\approx 10^{-16}$ e  $\delta \approx 10^{-8}$ ).

### 3.4 Esperimenti numerici e confronto con il metodo non accelerato

Si riportano ora gli esperimenti utili a verificare l'efficacia dell'accelerazione di Anderson nella riduzione del numero di iterazioni e del tempo totale di esecuzione dell'algoritmo delle proiezioni alternate per il problema NCM.

Gli esperimenti sono eseguiti su MATLAB R2022a, su PC con processore Intel Core i5-4430S @2.70GHz e 16GB di RAM.

Vengono utilizzati i seguenti algoritmi.

- **nearcorr\_new**: il metodo delle proiezioni alternate descritto dall'Algoritmo 1, modificato per incorporare le varianti della Sezione 3.3.

- nearcorr\_AA: il metodo delle proiezioni alternate accelerato descritto dall'Algoritmo 5. Si usa l'implementazione di Anderson presentata in [22, 23], che sfrutta l'aggiornamento della fattorizzazione QR, come descritto nella Sezione 3.1.

La tolleranza per la convergenza è fissata a  $n \cdot u$ , dove n è l'ordine della matrice e  $u \approx 1.1 \times 10^{-16}$  è l'unità di round-off. La matrice W è sempre fissata come W = I per comodità, per cui la norma-W coincide con la norma Frobenius. La convergenza è assicurata solo per l'Algoritmo 1 (senza le modifiche delle varianti), ma per la versione con l'accelerazione di Anderson potrebbe non essere garantita per mancanza di risultati teorici di convergenza. Nonostante ciò, in caso di convergenza si può verificare che in ogni esperimento i due algoritmi convergono circa alla stessa soluzione X.

Si<br/>a $\delta=0,$ senza vincolo sugli elementi. Si cerca di risolvere il problema NCM standard, senza varianti.

• Esperimento 3.1. Si confronta il numero di iterazioni necessarie alla convergenza degli Algoritmi 1 e 5 al variare del parametro *m*, utilizzando 4 matrici invalide di piccole dimensioni, con dati provenienti da portafogli di investimento e strategie di gestione del rischio (vedi Appendice). Sono riportati i risultati in Tabella 3.1, dove it è il numero di iterazioni per nearcorr\_new, mentre itAA è relativo a nearcorr\_AA.

n	it	itAA						
		m = 1	m = 2	m = 3	m = 4	m = 5	m = 6	
4	39	15	10	9	9	9	9	
5	27	17	14	12	11	10	10	
6	804	319	225	100	56	43	30	
7	33	15	10	10	10	9	9	

Tabella 3.1: Esperimento 3.1

Si osserva che l'accelerazione di Anderson porta ad una significativa riduzione del numero di iterazioni, già per m = 1, fino ad un fattore di riduzione di 26.80 nel caso m = 6 per la matrice  $6 \times 6$ .

Il numero di iterazioni tende a stabilizzarsi per m che cresce, il che è compatibile con il comportamento dell'accelerazione di Anderson riportato in letteratura. Per questo motivo, a parte qualche caso, per gli esperimenti successivi verrà fissato il valore di m.

Si riporta in Figura 3.2 la storia dell'errore relativo

$$\frac{\|X_k - Y_k\|_F}{\|Y_k\|_F}$$

per le matrici di ordine n = 5 e n = 6, in scala logaritmica. L'espressione "NCMaa<sub>j</sub>" indica l'algoritmo **nearcorr\_AA** con m = j. Oltre al decremento nel numero di iterazioni si può notare la differenza nell'andamento dell'errore tra i casi accelerati e il caso non accelerato, a riprova del fatto che non è stato ancora dimostrato alcun risultato che garantisca la convergenza dell'Algoritmo 5.

• Esperimento 3.2. Si confrontano ora il numero di iterazioni ed il tempo di esecuzione (espresso in secondi) per i due algoritmi nearcorr\_new e nearcorr\_AA. Sia m = 2. Si utilizzano 3 matrici di grandi dimensioni. Le prime due sono fornite da una compagnia di fondi di investimento, la prima ha dimensione  $1399 \times 1399$  ed è a rango basso (è la stessa dell'Esperimento 2.3) mentre la seconda ha dimensione  $3120 \times 3120$  ed ha rango pieno. La terza contiene dati provenienti da 3250 banche di 27 stati europei ed ha dimensione  $3250 \times 3250$ .

In Tabella 3.2 vengono riportate le iterazioni necessarie alla convergenza insieme al tempo computazionale t, con l'aggiunta di t\_ap e t\_AA relativi a nearcorr\_AA,



Figura 3.2: Storia dell'errore per le matrici di ordine 5 e 6, Esperimento 3.1

che sono il tempo totale per le chiamate della funzione g dell'Algoritmo 4 insieme al test di convergenza e il tempo necessario alla risoluzione del problema ai minimi quadrati, rispettivamente. Il tempo mancante rappresenta il tempo necessario per le operazioni generali che MATLAB svolge, ad esempio l'operazione vec.

L'accelerazione indicata alla fine corrisponde al rapporto tra il tempo impiegato da nearcorr\_new e quello impiegato da nearcorr\_AA.

n	nearcorr_new		nearcori		acceler		
	it	t	itAA	t	t_ap	t_AA	- acceler.
1399	476	265.1	124	90.1	70.9	16.9	2.9
3120	559	2734.5	175	1020.6	874.4	130.5	2.7
3250	7	36.9	6	35.9	32.6	2.6	1.03

Si osserva un forte decremento del numero di iterazioni, con quelle di **nearcorr\_AA** circa uguali ad un terzo di quelle di **nearcorr\_new**.

Per quanto riguarda il tempo computazionale si assiste ad un'accelerazione significativa dell'Algoritmo 5 (nearcorr\_AA), circa 3 volte più veloce dell'Algoritmo 1 (nearcorr\_new).

Così come per il metodo non accelerato, il tempo computazionale dell'Algoritmo 5 è principalmente costituito dalla parte di proiezione alternata e quindi di decomposizione spettrale, piuttosto che dal problema ai minimi quadrati.

É importante sottolineare che il numero di iterazioni e il tempo di esecuzione non sempre dipendono dalla dimensione della matrice di partenza, come è possibile vedere dalla terza matrice della Tabella 3.2.

Si considera ora il problema NCM nella variante con elementi fissati (ma con  $\delta=0).$ 

• Esperimento 3.3. Si utilizzano tre matrici. La prima è la matrice  $7 \times 7$  utilizzata nell'Esperimento 3.1, per cui è richiesto di calcolare la matrice di correlazione più vicina mantenendo inalterato il primo blocco principale  $3 \times 3$ . La seconda è una

matrice indefinita di  $7 \times 7$  blocchi, ognuno di ordine 10, contenente dati azionari; la NCM deve preservare il blocco in posizione (1,1), la diagonale principale e tutte le diagonali dei blocchi nella prima riga e prima colonna. La terza matrice, proveniente da valutazioni di stoccaggio di diossido di carbonio nella regione Rocky Mountains degli Stati Uniti, ha dimensione  $94 \times 94$  ed ha una struttura a blocchi; i 12 blocchi diagonali da cui è costituita (tutti di dimensioni diverse) sono definiti positivi e devono rimanere fissati (i "pattern" sono in Figura 3.1).

In Tabella 3.3 sono riportati il numero di iterazioni per nearcorr\_new senza elementi fissati (it), il numero di iterazioni per nearcorr\_new con elementi fissati (it\_fe) e il numero di iterazioni per nearcorr\_AA con elementi fissati (itAA\_fe), al variare di m = 1, ..., 5.

n	it	it fe	itAA_fe				
	10	10_10	m = 1	m = 2	m = 3	m = 4	m = 5
7	33	34	14	11	10	9	9
70	100	355	179	102	85	71	58
94	18	40	15	14	12	12	12

Tabella 3.3: Esperimento 3.3

Si può vedere chiaramente che l'aggiunta del vincolo degli elementi fissati aumenta (non di poco in alcuni casi) il numero di iterazioni necessarie alla convergenza, relativamente al problema standard.

In ogni caso, l'accelerazione di Anderson abbassa il numero di iterazioni con un fattore di riduzione simile ai casi precedenti anche in presenza del vincolo degli elementi fissati.

• Esperimento 3.4. Come altro esperimento per la variante NCM con elementi fissati si prendono matrici di correlazione invalide di dimensione rispettivamente 200, 400, 600, 800 e si confronta il tempo computazionale necessario alla convergenza di nearcorr\_new e nearcorr\_AA, al variare di m. Per ogni matrice si fissa il blocco principale di ordine n/2, che per costruzione è una matrice di correlazione. La tolleranza per il test di convergenza viene aumentata a  $n \cdot \sqrt{u}$ , dal momento che la quantità precedente  $n \cdot u$  non porta l'algoritmo a convergenza entro le 10000 iterazioni.

I risultati in Tabella 3.4 mostrano un decremento importante del tempo compu-

n	time_fe	time_fe_AA						
		m = 1	m = 2	m = 3	m = 4	m = 5		
200	30.73	3.34	4.45	2.37	1.55	1.82		
400	68.11	22.00	20.65	11.43	11.14	7.57		
600	114.28	33.09	45.07	15.26	30.78	14.25		
800	507.39	88.22	121.74	88.56	66.96	36.60		

Tabella 3.4: Esperimento 3.4

Il successivo esperimento è relativo alla variante del problema NCM con $\delta>0,$ senza elementi fissati.

• Esperimento 3.5. Si esegue l'algoritmo sulle matrici dell'Esperimento 3.1, per  $\delta = 10^{-8}$  e  $\delta = 0.1$ , al variare di m. I risultati sono riportati in Tabella 3.5.

Caso $\delta = 10^{-8}$									
n	it	itAA	itAA						
	10	m = 1	m = 2	m = 3	m = 4	m = 5	m = 6		
4	39	15	10	9	9	9	10		
5	27	17	14	12	11	10	10		
6	804	309	191	101	61	40	30		
7	33	15	10	10	10	9	9		
Caso $\delta$	= 0.1								
n	i+	itAA	itAA						
	10	m = 1	m = 2	m = 3	m = 4	m = 5	m = 6		
4	65	31	19	16	13	14	13		
5	34	23	15	14	13	12	12		
6	894	304	168	136	59	133	47		
7	55	31	24	15	15	14	13		

Tabella 3.5: Esperimento 3.5 con matrici Esperimento 3.1

Per il primo caso il numero di iterazioni è sostanzialmente identico all'Esperimento 3.1, ma qui si ha la certezza che la soluzione X sia definita positiva.

Per il secondo caso si osserva un aumento del numero di iterazioni necessarie alla convergenza del metodo nel caso standard. Come nel caso degli elementi fissati, l'accelerazione di Anderson riduce le iterazioni di un fattore simile al caso standard, per cui si può affermare che il vincolo non influenza l'effetto dell'accelerazione.

Per osservare in maniera più evidente i benefici dell'accelerazione di Anderson si utilizzano le matrici di grandi dimensioni dell'Esperimento 3.2. Sia m = 2 e  $\delta = 0.1$ . I risultati sono riportati in Tabella 3.6.

	nearcorr_new		nearcorr_AA				acceler
	it	t	itAA	t	t_ap	t_AA	
1399	555	322.2	112	83.8	65.1	16.5	3.8
3120	1142	5542.2	345	2012.9	1726.6	253.9	2.8
3250	7	36.6	6	36.2	33.1	2.4	1.01

Tabella 3.6: Esperimento 3.5 con matrici Esperimento 3.2

Si nota che con il vincolo  $\delta > 0$  l'accelerazione di Anderson riduce il numero di iterazioni di un fattore più alto rispetto a quello del problema standard, il che mostra che l'accelerazione di Anderson può avere molto più effetto nel caso vincolato che nel caso standard. Stesso discorso per il tempo computazionale, che subisce un'accelerazione maggiore rispetto al caso standard.

Dopo aver sperimentato la maggior efficacia dell'accelerazione di Anderson nel caso vincolato rispetto al caso standard, si combinano entrambi i vincoli di elementi fissati e  $\delta > 0$ .

• Esperimento 3.6. Si utilizzano nuovamente le matrici dell'Esperimento 3.3, con gli stessi elementi fissati e  $\delta = 0.1$ .

É importante notare che in questo caso non si hanno garanzie di convergenza. Risulta infatti che per la matrice di ordine 70 l'algoritmo non converge entro le prime 100000 iterazioni, né per la tolleranza fissata in precedenza, né per una tolleranza più alta (come  $4 \times 10^{-3}$ ). Per questo non ne viene riportato alcun dato nelle Tabelle 3.7 e 3.8.

In Tabella 3.7 vengono messi a confronto il numero di iterazioni necessarie alla convergenza dell'algoritmo nearcorr\_new nel caso del problema standard senza vincoli (it\_old), poi nel caso vincolato (it) e infine le iterazioni dell'algoritmo nearcorr\_AA nel caso vincolato (itAA).

In Tabella 3.8 sono riportati i tempi di calcolo degli algoritmi nearcorr\_new e nearcorr\_AA nel caso vincolato, con m = 1, ..., 5.

n	it_old	it	itAA				
		10	m = 1	m = 2	m = 3	m = 4	m = 5
7	33	55	31	25	16	15	15
70	-	-	-	-	-	-	-
94	18	128	36	25	24	20	19

Tabella 3.7: Iterazioni Esperimento 3.6

n	t.	t_AA						
	0	m = 1	m = 2	m = 3	m = 4	m = 5		
7	1.55e-1	5.29e-2	2.36e-2	9.46e-3	3.08e-3	6.97e-3		
70	-	-	-	-	-	-		
94	2.01e-1	6.69e-2	5.09e-2	5.31e-2	5.39e-2	5.78e-2		

Tabella 3.8: Tempo computazionale Esperimento 3.6

Si nota ovviamente l'incremento nel numero di iterazioni rispetto al problema con il solo vincolo degli elementi fissati, ma in questo caso l'accelerazione di Anderson riduce questo numero di un fattore uguale a 3.7 per la matrice  $7 \times 7$  e uguale a 6.4 per la matrice di ordine 94, mentre nell'esperimento con un solo vincolo i fattori erano 3.6 e 3.3 rispettivamente. Per quanto riguarda il tempo computazionale, il miglioramento non è così evidente come nei casi precedenti vista la dimensione ridotta delle matrici in esame.

Si include in Figura 3.3 la storia dell'errore relativo anche nel caso vincolato. Le curve riportate mostrano un ottimo miglioramento rispetto al caso vincolato non accelerato nonostante l'assenza di risultati di convergenza per l'algoritmo accelerato. Di nuovo, l'espressione "NCMaa<sub>j</sub>" indica l'algoritmo nearcorr\_AA con m = j.

Si può dunque affermare che l'accelerazione di Anderson applicata al metodo delle proiezioni alternate per il problema NCM porta un notevole decremento nel numero di iterazioni, con un fattore di riduzione di almeno 2 nel caso non vincolato e almeno 3 nel caso vincolato (ci sono dei casi limite in cui il fattore di riduzione è di oltre 20). Anche il tempo di esecuzione dell'algoritmo risente della modifica



Figura 3.3: Storia dell'errore relativo per l'Esperimento 3.6

di Anderson, maggiormente nel caso vincolato.

Sembra infatti che l'accelerazione tenda a produrre miglioramenti di entità maggiore per i problemi che il metodo delle proiezioni alternate standard trova più ardui.

## Conclusioni

L'interesse riguardo al problema NCM in Analisi Numerica è motivato soprattutto dall'analisi dei dati del mondo reale, che in questo elaborato sono dati di tipo finanziario.

L'approccio delle proiezioni alternate permette di formulare un problema che abbia garanzie di convergenza, tramite le proprietà degli spazi di Hilbert e le caratteristiche degli insiemi convessi e chiusi. L'Algoritmo 1 permette di sfruttare le proprietà spettrali della matrice A di partenza e concede la possibilità di analizzare problemi diversi come le varianti descritte nella Sezione 3.3, a differenza di algoritmi meno flessibili come quello di tipo Newton formulato da Qi e Sun [18]. La principale debolezza dell'Algoritmo 1 sta nella velocità di convergenza, lineare e dipendente dall'angolo tra gli insiemi di proiezione.

Subentra quindi la strategia dell'accelerazione di Anderson, tipica dei problemi di punto fisso. Nonostante non vi siano risultati certi che garantiscano la convergenza dell'accelerazione di Anderson applicata al metodo delle proiezioni alternate, i risultati numerici mostrano che un valore di m fissato uguale a 2 o 3 sembra velocizzare il metodo nel calcolo della matrice di correlazione più vicina, sia nella forma standard sia nelle varianti con vincolo sugli elementi e sull'autovalore minore. L'applicazione alle varianti suscita molto interesse soprattutto rispetto a certi metodi più veloci delle proiezioni alternate ma che non sono in grado di incorporare le modifiche necessarie (ad esempio il metodo in [18, 3]). Da ciò viene la raccomandazione di tenere il parametro m "relativamente basso", in quanto l'aumento del valore m non si riflette sempre in un'accelerazione del metodo, anzi potrebbe complicare ulteriormente i calcoli.

Il successo dell'accelerazione di Anderson nel contesto del problema NCM suggerisce la possibilità di applicarla ad altri algoritmi di proiezione, in modo da rendere competitivi i relativi metodi di proiezione, spesso lenti a causa della loro natura iterativa.

### Appendice

Si elencano le 4 matrici invalide utilizzate nell'Esperimento 3.1.

1) Turkay, Epperlein e Christofides, p. 86 in [21]

	1	-0.55	-0.15	-0.10
4 —	-0.55	1	0.90	0.90
A –	-0.15	0.90	1	0.90
	-0.10	0.90	0.90	1

2) Bhansali e Wise, Sezione 2 in [1]

$$A = \begin{bmatrix} 1 & -0.50 & -0.30 & -0.25 & -0.70 \\ -0.50 & 1 & 0.90 & 0.30 & 0.70 \\ -0.30 & 0.90 & 1 & 0.25 & 0.20 \\ -0.25 & 0.30 & 0.25 & 1 & 0.75 \\ -0.70 & 0.70 & 0.20 & 0.75 & 1 \end{bmatrix}$$

3) Minabutdinov, Manaev e Bouev, p. 36 in [17]

$$A = \begin{bmatrix} 1 & 3.1595 & 0.2009 & 1.0020 & -0.5809 & 0.1953 \\ 3.1595 & 1 & 12.6826 & 3.3505 & 8.1755 & 16.9360 \\ 0.2009 & 12.6826 & 1 & -0.0471 & 0.6807 & 1.0424 \\ 1.0020 & 3.3505 & -0.0471 & 1 & -0.7884 & -0.1409 \\ -0.5809 & 8.1755 & 0.6807 & -0.7884 & 1 & 0.7201 \\ 0.1953 & 16.9360 & 1.0424 & -0.1409 & 0.7201 & 1 \end{bmatrix}$$

4) Finger, Tabella 4 in [9]

$$A = \begin{bmatrix} 1 & 0.18 & -0.13 & -0.26 & 0.19 & -0.25 & -0.12 \\ 0.18 & 1 & 0.22 & -0.14 & 0.31 & 0.16 & 0.09 \\ -0.13 & 0.22 & 1 & 0.06 & -0.08 & 0.04 & 0.04 \\ -0.26 & -0.14 & 0.06 & 1 & 0.85 & 0.85 & 0.85 \\ 0.19 & 0.31 & -0.08 & 0.85 & 1 & 0.85 & 0.85 \\ -0.25 & 0.16 & 0.04 & 0.85 & 0.85 & 1 & 0.85 \\ -0.12 & 0.09 & 0.04 & 0.85 & 0.85 & 0.85 & 1 \end{bmatrix}$$

## Bibliografia

- BHANSALI, V., AND WISE, M. B. Forecasting portfolio risk in normal and stressed markets. arXiv preprint nlin/0108022 (2001).
- [2] BIRGIN, E. G., AND RAYDAN, M. Robust stopping criteria for Dykstra's algorithm. SIAM Journal on Scientific Computing 26, 4 (2005), 1405–1414.
- [3] BORSDORF, R., AND HIGHAM, N. J. A preconditioned Newton algorithm for the nearest correlation matrix. *IMA Journal of Numerical Analysis 30*, 1 (2010), 94–107.
- [4] BOYLE, J. P., AND DYKSTRA, R. L. A method for finding projections onto the intersection of convex sets in Hilbert spaces. In Advances in order restricted statistical inference. Springer, 1986, pp. 28–47.
- [5] CHENG, S. H., AND HIGHAM, N. J. A modified Cholesky algorithm based on a symmetric indefinite factorization. SIAM Journal on Matrix Analysis and Applications 19, 4 (1998), 1097–1110.
- [6] DEUTSCH, F. Best approximation in inner product spaces, vol. 7. Springer, 2001.
- [7] DEUTSCH, F., AND HUNDAL, H. The rate of convergence for the method of alternating projections, ii. *Journal of Mathematical Analysis and Applications* 205, 2 (1997), 381–405.
- [8] DYKSTRA, R. L. An algorithm for restricted least squares regression. *Journal* of the American Statistical Association 78, 384 (1983), 837–842.
- [9] FINGER, C. C. A methodology to stress correlation. *RiskMetrics Monitor Fourth Quarter* (1997), 3–11.

- [10] HIGHAM, N. J. A Collection of Invalid Correlation Matrices. https://nhigham.com/2016/03/23/ a-collection-of-invalid-correlation-matrices/.
- [11] HIGHAM, N. J. The Nearest Correlation Matrix. https://nhigham.com/ 2013/02/13/the-nearest-correlation-matrix/#fn.7.
- [12] HIGHAM, N. J. Computing the Nearest Correlation Matrix a problem from finance. IMA Journal of Numerical Analysis 22, 3 (2002), 329–343.
- [13] HIGHAM, N. J., AND STRABIĆ, N. Anderson acceleration of the alternating projections method for computing the Nearest Correlation Matrix. *Numerical Algorithms* 72, 4 (2016), 1021–1042.
- [14] JOHNSON, R. A., WICHERN, D. W., ET AL. Applied multivariate statistical analysis, vol. 6. Pearson London, UK:, 2014.
- [15] KELLEY, C. T. Iterative methods for linear and nonlinear equations. SIAM, 1995.
- [16] LUENBERGER, D. G. Optimization by vector space methods. John Wiley & Sons, 1997.
- [17] MINABUTDINOV, A., MANAEV, I., BOUEV, M., ET AL. Finding the nearest valid covariance matrix: A fx market case. European University at St. Petersburg Department of Economics Working Paper Ec-07/13 (2014).
- [18] QI, H., AND SUN, D. A quadratically convergent Newton method for computing the Nearest Correlation Matrix. SIAM Journal on matrix analysis and applications 28, 2 (2006), 360–385.
- [19] ROCKAFELLAR, R. T. Convex analysis. In *Convex analysis*. Princeton University Press, 2015.
- [20] SIMONCINI, V. Metodo delle potenze. In *Dispense del corso di Calcolo Numerico*. Università di Bologna, 2021.
- [21] TURKAY, S., EPPERLEIN, E., AND CHRISTOFIDES, N. Correlation stress testing for value-at-risk. *Journal of Risk* 5, 4 (2003).
- [22] WALKER, H. F. Anderson acceleration: Algorithms and implementations. WPI Math. Sciences Dept. Report MS-6-15-50 (2011).

[23] WALKER, H. F., AND NI, P. Anderson acceleration for fixed-point iterations. SIAM Journal on Numerical Analysis 49, 4 (2011), 1715–1735.