

Dispensa del Corso di Complementi di Ricerca Operativa

Marco Trubian*,

Milano, Marzo 2008

*Dipartimento di Scienze dell'Informazione, Università degli Studi di Milano, Via Comelico 39, 20135
Milano, e-mail: trubian@dsi.unimi.it

Indice

1	Notazione	1
2	Introduzione	2
2.1	Esempi	3
2.2	Gli algoritmi di ottimizzazione	4
2.3	Nozioni di convessità	5
2.4	Minimi locali e globali	7
3	Ottimizzazione non vincolata	8
4	Condizioni analitiche di ottimalità	9
4.1	Condizioni analitiche di ottimalità per problemi convessi	14
4.2	Condizioni <i>empiriche</i> di ottimalità	16
4.3	Programmazione Quadratica	16
4.3.1	Q non è semidefinita positiva	16
4.3.2	Q è definita positiva	16
4.3.3	Q è semidefinita positiva	17
4.3.4	Esempi	17
5	Convergenza e rapidità di convergenza	18
6	Metodi basati su ricerca lineare	20
6.1	Determinazione del passo α_k	20
6.1.1	Backtracking e metodo di Armijo	24
6.1.2	Metodo di ricerca esatto	25
6.1.3	Metodo di interpolazione	26
6.1.4	Inizializzazione del passo α_0	26
6.1.5	Tecniche che non calcolano le derivate	27
6.2	Scelta della direzione \mathbf{d}_k	27
6.2.1	Metodo del gradiente	29
6.2.2	Analisi del metodo del gradiente	30
6.2.3	Metodo di Newton	34
6.2.4	Confronto fra i metodi del gradiente e di Newton	38
6.2.5	Metodi quasi Newton	38
6.2.6	Formula di aggiornamento di rango uno	40
6.2.7	Formula di aggiornamento di rango due (DFP)	41
6.2.8	Formula di aggiornamento di rango due inversa (BFGS)	42
6.2.9	Famiglia di Broyden	43
6.3	Metodi alle direzioni coniugate	43
6.3.1	Il metodo di gradiente coniugato	45
7	Metodi di Trust-Region	46
7.1	Il punto di Cauchy	49
7.2	Il metodo <i>Dogleg</i>	50

8	Problemi ai Minimi Quadrati	51
8.1	Il metodo di Gauss-Newton	53
8.2	Il metodo di Levenberg-Marquardt	54
9	Ottimizzazione vincolata	55
9.1	Condizioni analitiche: vincoli di uguaglianza	56
9.1.1	Funzione obiettivo quadratica e vincoli di uguaglianza lineari	59
9.1.2	Da vincoli di disuguaglianza a vincoli di uguaglianza	60
9.2	Il caso generale: le condizioni KKT	61
9.3	Condizioni di ottimalità del secondo ordine	65
9.4	Punti di sella e dualità	67
9.5	Programmazione quadratica con vincoli di disuguaglianza lineari	70
9.6	Metodi con funzione di penalità	70
9.7	Metodi di barriera	72
9.8	Metodo del gradiente proiettivo	72
9.9	Metodo dei lagrangiani aumentati	75
9.10	SQP (Sequential Quadratic Programming)	76
10	Appendice	78

Elenco delle figure

1	Esempio in due dimensioni	3
2	Curva con minimo globale non isolato	7
3	Direzione di discesa	11
4	Esempi di funzioni quadratiche	17
5	Convergenza a valori errati	22
6	Condizione di Armijo	22
7	Condizioni di Wolfe	23
8	Scelta del passo ottima	25
9	Andamento zigzagante del metodo del gradiente	32
10	Il metodo del gradiente applicato alla funzione di Rosenbrock.	37
11	Il metodo di Newton applicato alla funzione di Rosenbrock.	39
12	Esempio di modello quadratico per la funzione di Rosenbrock.	47
13	Esempio di funzione obiettivo convessa solo nell'insieme ammissibile.	55
14	Esempio di funzione obiettivo convessa e infiniti minimi locali.	56
15	Condizione di ottimalità per problemi con vincoli di uguaglianza.	57
16	Esempio di gradienti di vincoli di uguaglianza linearmente dipendenti.	58
17	Condizioni di ottimalità con vincoli di disuguaglianza.	62
18	Condizioni di non ottimalità con vincoli di disuguaglianza.	63
19	Condizioni di ottimalità con vincoli di disuguaglianza.	63
20	Esempio di gradienti di vincoli attivi linearmente dipendenti.	64
21	Direzioni critiche.	66
22	Esempio sulle condizioni del secondo ordine.	67
23	Punto di sella della funzione Lagrangiana	69
24	Passo di correzione nel gradiente proiettivo	75
25	Andamento dei grafici di $\Phi(\mathbf{y})$ e $\Psi(\mathbf{y})$	80

1 Notazione

\mathbb{R}^n = insieme dei vettori reali n -dimensionali

\mathbb{R}_+^n = insieme dei vettori reali n -dimensionali non negativi

\mathbb{Z}^n = insieme dei vettori interi n -dimensionali

\mathbb{Z}_+^n = insieme dei vettori interi n -dimensionali interi non negativi

$[a, b] = \{x \in \mathbb{R} : a \leq x \leq b\}$ = intervallo chiuso

$(a, b) = \{x \in \mathbb{R} : a < x < b\}$ = intervallo aperto

$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$ = vettore colonna ad n componenti ($\mathbf{x} \in \mathbb{R}^n$)

$\mathbf{x}^T = [x_1, x_2, \dots, x_n]$ = vettore riga ad n componenti ($\mathbf{x} \in \mathbb{R}^n$)

$e_i^T = (0, 0, \dots, 1, 0, \dots, 0)$ i -esimo versore: vettore con un 1 in posizione i -esima e 0 altrove

$\|\mathbf{x}\| = \sqrt{\sum_{j=1}^n x_j^2}$ = norma euclidea del vettore ($\mathbf{x} \in \mathbb{R}^n$)

$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$ = matrice $m \times n$

$\arg \min\{f(\mathbf{x}) : \mathbf{x} \in X\} = \mathbf{x}^* \in \mathbb{R}^n$ tale che $f(\mathbf{x}^*) = \min\{f(\mathbf{x}) : \mathbf{x} \in X\}$

$I(x, \varepsilon) = \{y : \|y - x\| < \varepsilon\}$ = intorno aperto di raggio $\varepsilon > 0$ di x

$\frac{\partial f}{\partial x_i}$ derivata parziale prima rispetto a x_i

∇ operatore di derivata parziale

$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \dots \quad \frac{\partial f}{\partial x_n} \right]^T$ vettore gradiente

∇^2 operatore di derivata seconda (elementi $\frac{\partial^2 f}{\partial x_i \partial x_j}$)

$\nabla^2 f(\mathbf{x}) = H(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$ matrice Hessiana

$\frac{\partial \mathbf{q}}{\partial \mathbf{x}} = \mathbf{J}$ = matrice Jacobiana $n \times k = [\nabla q_1, \nabla q_2, \dots, \nabla q_k]$ dove $q_j : \mathbb{R}^n \rightarrow \mathbb{R}$, $j = 1, \dots, k$

2 Introduzione

L'obiettivo di questa dispensa è quello di fornire una introduzione agli algoritmi per l'ottimizzazione numerica per problemi nel continuo.

Molte applicazioni possono essere formulate come problemi di ottimizzazione nel continuo; per esempio

- trovare la traiettoria ottima di un velivolo o del braccio di un robot;
- identificare le proprietà sismiche di un pezzo di crosta terrestre adattando un modello della regione studiata ad un insieme di valori campionati da una rete di stazioni di rilevamento;
- comporre un portafoglio di investimenti che massimizzi il ritorno atteso nel rispetto di un accettabile livello di rischio;
- calcolare la forma ottimale di un componente di una automobile, di un aereomobile o dello scafo di una imbarcazione;
- determinare la conformazione spaziale di una proteina a partire dalla sequenza di aminoacidi che la compongono.

Formulazione matematica Formalmente, l'ottimizzazione consiste nella massimizzazione o minimizzazione di una funzione soggetta a vincoli sulle sue variabili. La notazione utilizzata è la seguente:

- \mathbf{x} è il vettore delle *variabili* o delle *incognite*;
- $f(\mathbf{x})$ è la *funzione obiettivo* che vogliamo massimizzare o minimizzare;
- $g(\mathbf{x})$ sono *vincoli di disuguaglianza*, funzioni scalari di \mathbf{x} che definiscono delle *disequazioni* che il vettore delle variabili deve soddisfare;
- $h(\mathbf{x})$ sono *vincoli di uguaglianza*, funzioni scalari di \mathbf{x} che definiscono delle *equazioni* che il vettore delle variabili deve soddisfare;
- X è la *regione ammissibile*, cioè il luogo dei punti dove il vettore delle variabili soddisfa tutti i vincoli, di disuguaglianza o di uguaglianza.

Con questa notazione il problema di ottimizzazione può venir così formulato:

$$\begin{array}{ll} \min & f(\mathbf{x}) \\ \text{t.c.} & \mathbf{x} \in X \end{array}$$

dove $X = \{\mathbf{x} \in \mathbb{R}^n : g_i(\mathbf{x}) \leq 0 \quad i = 1, \dots, k; \quad h_j(\mathbf{x}) = 0 \quad j = 1, \dots, h\}$

Nel caso in cui $X \subset \mathbb{R}^n$, cioè X è un sottoinsieme proprio di \mathbb{R}^n , si parla di *ottimizzazione vincolata*, mentre quando vale $X \equiv \mathbb{R}^n$ si parla di *ottimizzazione non vincolata*.

Quando almeno una delle funzioni $f(\mathbf{x}), g(\mathbf{x}), h(\mathbf{x})$ è non lineare si parla di *Programmazione Non Lineare* (PNL).

2.1 Esempi

Consideriamo il problema di determinare il punto $P = (x_1, x_2)$ più vicino al punto di coordinate $(2, 1)$, con la richiesta che P deve appartenere all'intersezione della superficie interna alla parabola di equazione $x_2 = x_1^2$ con l'area del triangolo di estremi $(0, 0)$, $(2, 0)$ e $(0, 2)$. Il problema può essere modellizzato nel seguente modo

$$\begin{aligned} \min \quad & (x_1 - 2)^2 + (x_2 - 1)^2 \\ \text{t.c.} \quad & x_1^2 - x_2 \leq 0 \\ & x_1 + x_2 \leq 2 \\ & x_1, x_2 \geq 0 \end{aligned}$$

Esso può essere visto come un problema di PNL leggendo $f(\mathbf{x}) = (x_1 - 2)^2 + (x_2 - 1)^2$ ed i vincoli di disuguaglianza come: $g_1(\mathbf{x}) = x_1^2 - x_2$, $g_2(\mathbf{x}) = x_1 + x_2 - 2$, $g_3(\mathbf{x}) = -x_1$ e $g_4(\mathbf{x}) = -x_2$. La Figura 2.1 rappresenta in chiaro la regione ammissibile del problema, e,

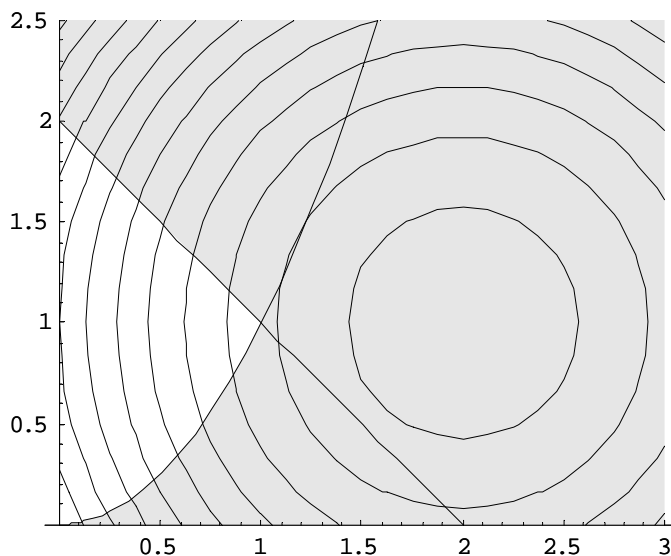


Figura 1: Esempio in due dimensioni

come cerchi concentrici, le curve di livello della funzione obiettivo.

Anche i problemi di Programmazione Lineare Intera (PLI) possono essere visti come problemi di PNL. Cominciamo con la Programmazione Lineare Binaria (PLB):

$$\begin{aligned} 1) \quad \min \quad & \mathbf{c}^T \mathbf{x} \\ \text{t.c.} \quad & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \in \{0, 1\}^n \end{aligned}$$

Essa è riconducibile alla PNL mediante, ad esempio, la seguente trasformazione

$$\begin{aligned} 2) \quad \min \quad & \mathbf{c}^T \mathbf{x} + M\mathbf{x}^T(\mathbf{1} - \mathbf{x}) \\ \text{t.c.} \quad & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \in [0, 1]^n \end{aligned}$$

I problemi 1) e 2) sono equivalenti. Nella funzione obiettivo del problema 2) abbiamo aggiunto un valore di penalità, $M \gg 0$, grande a sufficienza. Ciascuna componente $Mx_i(1-x_i)$ assume il suo valore massimo in corrispondenza del valore $x_i = 0.5$, mentre tende a zero al tendere della variabile x_i a 0 oppure a 1.

Il problema 2) è un caso particolare di programmazione quadratica, dove le variabili hanno esponente di grado 2 o 1 ed i vincoli sono lineari. La programmazione lineare binaria è un caso particolare di lineare intera ma è anche un caso particolare di programmazione quadratica. In alternativa si può scrivere

$$\begin{aligned} 3) \quad & \min \quad \mathbf{c}^T \mathbf{x} \\ & t.c. \quad A\mathbf{x} = \mathbf{b} \\ & \quad \mathbf{x}^T(\mathbf{1} - \mathbf{x}) = \mathbf{0} \\ & \quad \mathbf{x} \in [0, 1]^n \end{aligned}$$

I problemi 1) e 3) sono equivalenti. Poiché le variabili sono non negative l'unico modo per soddisfare il vincolo non lineare $\mathbf{x}^T(\mathbf{1} - \mathbf{x}) = \mathbf{0}$, è richiedere che ogni componente $x_i(1-x_i)$ sia uguale a zero. Questo si verifica solo quando ciascuna variabile x_i assume il valore 0 oppure 1. Notiamo, infine, che anche la Programmazione Lineare Intera (PLI):

$$\begin{aligned} 4) \quad & \min \quad \mathbf{c}^T \mathbf{x} \\ & t.c. \quad A\mathbf{x} = \mathbf{b} \\ & \quad \mathbf{x} \in Z^n \end{aligned}$$

è riconducibile alla PNL mediante, ad esempio, la seguente trasformazione

$$\begin{aligned} 5) \quad & \min \quad \mathbf{c}^T \mathbf{x} \\ & t.c. \quad A\mathbf{x} = \mathbf{b} \\ & \quad \sin(\pi x_j) = 0, \quad j = 1, \dots, n \end{aligned}$$

Queste semplici trasformazioni mettono in evidenza come i problemi di PNL costituiscano una generalizzazione dei problemi incontrati nell'ambito della programmazione lineare e lineare intera e che pertanto siano in generale più difficili da risolvere (frequentemente, molto più difficili da risolvere) di quanto non lo siano i pur difficili problemi di PLI.

2.2 Gli algoritmi di ottimizzazione

In generale, non è possibile identificare per via analitica le soluzioni ottime dei problemi di PNL. È possibile, come vedremo nelle sezioni 4 e 4.1, fornire, in qualche caso, le condizioni analitiche che devono essere soddisfatte dalle soluzioni ottime, ed in alcuni, fortunati casi, determinare una soluzione ottima proprio imponendo tali condizioni. Di conseguenza, si faranno di volta in volta delle assunzioni per specializzare le funzioni $f(\mathbf{x})$, $g_i(\mathbf{x})$, e $h_i(\mathbf{x})$, (continuità, differenziabilità, convessità, ecc.) in modo da dividere i problemi in classi, applicando ai problemi di ciascuna classe i metodi più opportuni. Nella maggior parte dei casi quindi, la

via analitica non permette di fornire una soluzione, anche se le proprietà analitiche possono ancora tornare utili per dire se una data soluzione, in qualche modo generata, è o meno una soluzione ottima. In generale quindi, è necessario ricorrere a tecniche di tipo algoritmico.

Gli algoritmi di ottimizzazione sono iterativi. Iniziano generando un vettore di variabili \mathbf{x}_0 , non necessariamente appartenente alla regione ammissibile X , e generano iterativamente una sequenza di vettori $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ che sono tutte (sperabilmente) approssimazioni migliori di una soluzione ottima.

Abbiamo già incontrato questa tecnica quando abbiamo risolto i problemi di Programmazione Lineare. L'*algoritmo del Simplex* è un tipico algoritmo di tipo iterativo che, partendo da una soluzione ammissibile, si sposta in soluzioni vicine alla soluzione corrente, arrestandosi quando nessuna di tali soluzioni migliora il valore della funzione obiettivo corrente e sono soddisfatte particolari condizioni analitiche, che garantiscono di trovarsi nella soluzione ottima. Uno schema analogo viene adottato dagli algoritmi che vedremo nel seguito, con una importante differenza però: nella maggior parte dei casi tali algoritmi si arrestano identificando solo un *minimo locale* del problema, una soluzione che è la migliore nella porzione di regione ammissibile che la circonda, ma che potrebbe essere anche molto diversa (e peggiore) dalla soluzione ottima del problema (cfr. la Sezione 2.4).

Le strategie adottate per passare da un vettore \mathbf{x}_k al successivo \mathbf{x}_{k+1} fanno la differenza fra un algoritmo ed un altro. La maggior parte degli algoritmi usano i valori della funzione f , delle funzioni g_i ed h_i e, quando conveniente, delle derivate prime e seconde di tali funzioni. Alcuni algoritmi usano l'informazione raccolta durante il processo iterativo di ricerca, mentre altri usano solo informazione locale al vettore \mathbf{x} corrente. In ogni caso i buoni algoritmi dovrebbero possedere le seguenti caratteristiche:

- **Robustezza.** Dovrebbero comportarsi bene su una varietà di problemi della classe per cui sono stati progettati a partire da qualsiasi ragionevole vettore iniziale \mathbf{x}_0 .
- **Efficienza.** Non dovrebbero richiedere un eccessivo tempo di calcolo o spazio di memoria.
- **Accuratezza.** Dovrebbero identificare una soluzione con precisione, senza risultare eccessivamente sensibili ad errori nei dati o agli inevitabili errori di arrotondamento che si verificano quando l'algoritmo viene implementato.

Queste caratteristiche sono spesso conflittuali, come vedremo, e occorrerà scendere a compromessi.

2.3 Nozioni di convessità

In ottimizzazione la nozione di convessità è fondamentale. Molti problemi possiedono tale proprietà che in genere li rende più facili da risolvere sia in teoria che in pratica. Il termine convesso può essere applicato sia a insiemi che a funzioni.

Definizione 1 *Insieme convesso.* Un insieme $X \subset \mathbb{R}^n$ è convesso se comunque presi due punti $\mathbf{x}, \mathbf{y} \in X$, allora $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in X$, per ogni $\lambda \in [0, 1]$.

Cioè ogni punto nel segmento che congiunge due qualsiasi punti in X è contenuto in X .

Definizione 2 *Funzione convessa.* Una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$ è convessa se il suo dominio è un insieme convesso $X \subseteq \mathbb{R}^n$ e comunque presi due punti $\mathbf{x}, \mathbf{y} \in X$ vale la relazione:

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \text{ per ogni } \lambda \in [0, 1].$$

In altre parole, il valore della funzione sui punti del segmento che unisce \mathbf{x} a \mathbf{y} è minore o uguale alla combinazione convessa dei valori che assume la funzione nei punti \mathbf{x} e \mathbf{y} .

Le funzioni convesse e gli insiemi convessi rivestono un ruolo importante nella PNL ed in particolare lo rivestono i problemi dove si minimizza una funzione obiettivo convessa su una regione ammissibile convessa.

Definizione 3 *Un problema di ottimizzazione con funzione obiettivo e regione ammissibile entrambe convesse viene detto problema convesso.*

I problemi convessi sono importanti perché, in genere, più semplici da risolvere. Quindi è importante capire come si può stabilire se una funzione o un insieme sono convessi.

Relativamente agli insiemi convessi valgono le seguenti proprietà:

Proprietà 1 *Siano X e Y due insiemi convessi ed $\alpha > 0$ uno scalare allora*

- a) *l'insieme αX è convesso;*
- b) *l'insieme $X + Y$ è convesso;*
- c) *l'insieme $X \cap Y$ è convesso.*

Relativamente alle funzioni convesse valgono le seguenti proprietà:

Proprietà 2 *Siano $f()$ e $g()$ due funzioni convesse ed $\alpha > 0$ uno scalare allora*

- a) *la funzione $\alpha f()$ è convessa;*
- b) *la combinazione lineare $f() + g()$ è convessa;*
- c) *la funzione $\max\{f(), g()\}$ è convessa;*
- d) *il luogo dei punti \mathbf{x} per i quali vale $f(\mathbf{x}) \leq \alpha$ è convesso.*

Tali proprietà ci aiutano ad identificare la convessità di una funzione complessa dall'analisi delle sue componenti.

Una immediata e importante conseguenza della Proprietà 1.c e della 2.d è la seguente

Proprietà 3 *Un problema di programmazione non lineare:*

$$\begin{aligned} \min & && f(\mathbf{x}) \\ \text{t.c.} & && g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, k \\ & && h_i(\mathbf{x}) = 0, \quad i = 1, \dots, h \end{aligned}$$

è convesso se $f(\mathbf{x})$ è convessa, le $g_i(\mathbf{x})$ sono convesse e le $h_i(\mathbf{x})$ sono lineari.

Osserviamo infine, che la funzione $\max\{f(), g()\}$ è convessa se $f()$ e $g()$ sono convesse, ma la funzione $\max\{f(), g()\}$ non è in genere differenziabile anche se $f()$ e $g()$ lo sono. La richiesta che la funzione obiettivo $f()$ sia differenziabile è fondamentale per la maggior parte dei metodi che vedremo e la trattazione di funzioni non differenziabili risulta molto più complessa del caso di funzioni differenziabili.

2.4 Minimi locali e globali

Data una funzione $f : X \rightarrow \mathbb{R}$, con $X \subseteq \mathbb{R}^n$, in generale, ciò che desideriamo è trovare un *ottimo globale* di $f(\mathbf{x})$, un punto in cui la funzione assume il suo minimo valore. Formalmente,

Definizione 4 Un punto $\mathbf{x}^* \in X$ è un punto di minimo globale di $f(\mathbf{x})$

se:

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \forall \mathbf{x} \in X$$

Poiché i nostri algoritmi iterativi non visitano (chiaramente) tutti i punti di X , né (sperabilmente) nemmeno molti, non si è sicuri che la funzione $f(\mathbf{x})$ non prenda una ripida discesa in una delle regioni che non sono state visitate dall'algoritmo. La maggior parte degli algoritmi sono quindi in grado di trovare solo *minimi locali*, cioè punti che assumono il valore minimo fra tutti quelli del loro *intorno*. Formalmente,

Definizione 5 Un punto $\mathbf{x}^* \in X$ è un punto di minimo locale di $f(\mathbf{x})$ se esiste un intorno aperto $I(\mathbf{x}^*, \varepsilon)$ di \mathbf{x}^* avente raggio $\varepsilon > 0$ tale che:

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \forall \mathbf{x} \in X \cap I(\mathbf{x}^*, \varepsilon)$$

Un punto è detto di minimo *in senso stretto* se al posto di \leq abbiamo $<$ nelle precedenti definizioni.

Come detto la maggior parte delle proprietà analitiche e delle tecniche algoritmiche hanno come finalità quella di determinare punti di minimo locale della funzione $f(\mathbf{x})$.

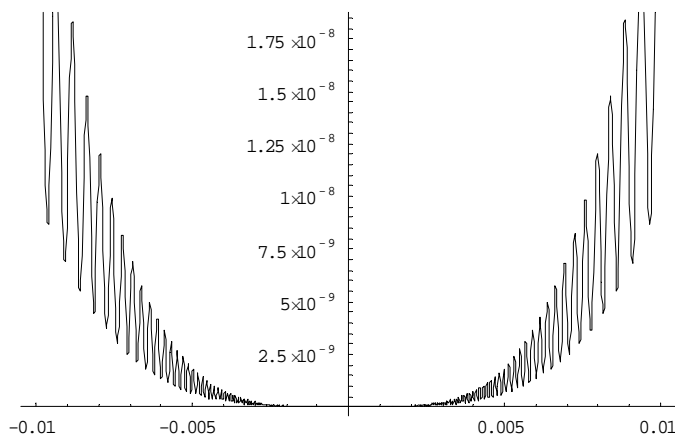


Figura 2: Curva con minimo globale non isolato

La figura illustra l'andamento della funzione

$$f(x) = x^4 \cos(1/x) + 2x^4$$

due volte differenziabile con continuità e che ha un minimo globale in $x^* = 0$ e minimi locali in senso stretto in molti punti prossimi ad esso. Questo esempio, benché costruito ad arte, non è però forzato in quanto, nei problemi di ottimizzazione associati con la determinazione di conformazioni molecolari, la funzione potenziale, che deve essere minimizzata, può avere

un numero di minimi locali che è esponenziale rispetto al numero di variabili. Consideriamo il seguente esempio di ottimizzazione non vincolata definito in \mathbb{R}^3 .

$$\min f(\mathbf{x}) = \sum_{i=1}^N \sum_{j=1}^{i-1} \left(\frac{1}{\|\mathbf{x}_i - \mathbf{x}_j\|^{12}} - \frac{2}{\|\mathbf{x}_i - \mathbf{x}_j\|^6} \right)$$

con $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathbb{R}^3$. La funzione $f(\mathbf{x})$ rappresenta l'energia potenziale di un raggruppamento di atomi identici, centrati nei punti $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, quando esista una forza di attrazione/repulsione fra ogni loro coppia. Si è potuto verificare sperimentalmente che questo modello rappresenta con particolare accuratezza gli stati di energia di alcuni gas ideali, quali il *kripto* (dal greco *kryptos*, nascosto, per la difficoltà di ottenere l'elemento in presenza di aria) e l'*argo* (dal greco *argos*, inattivo, poiché l'elemento non è reattivo) e di alcuni metalli. Si ritiene che la configurazione a minima energia, cioè la soluzione ottima del problema sia quella più stabile in natura.

La soluzione di questo, apparentemente semplice problema, è in realtà estremamente complessa, tenuto conto che si stima esistano $O(e^N)$ configurazioni che corrispondono a minimi locali, ma non globali, di f .

In questi casi la tecnica più frequentemente adottata è quella di campionare sistematicamente, o casualmente, lo spazio delle soluzioni ammissibili, di generare cioè un insieme di soluzioni ammissibili distinte e possibilmente ben ripartite nello spazio delle soluzioni X , e di portare, mediante un algoritmo di ricerca iterativo, ciascun punto campionato al suo minimo locale. Il miglior minimo locale viene scelto come approssimazione della soluzione ottima del problema. La disciplina che si occupa, nell'ambito della PNL, di determinare i punti di minimo globale prende il nome di *Ottimizzazione Globale* (OG). Le tecniche dell'OG si basano comunque sulla disponibilità di algoritmi per l'individuazione di punti di minimo locale ed è perciò su queste tecniche che noi concentriamo la nostra attenzione.

Osserviamo infine che se la funzione $f(\mathbf{x})$ e l'insieme che definisce la regione ammissibile X , godono di alcune proprietà molto generali è possibile affermare l'esistenza di almeno un punto di minimo e giustificarne quindi la ricerca.

Proprietà 4 *Se una funzione $f(\mathbf{x})$ è continua su un insieme $X \subseteq \mathbb{R}^n$ chiuso e limitato, allora $f(\mathbf{x})$ ammette un punto di minimo globale in X .*

Proprietà 5 *Se per una funzione $f(\mathbf{x})$ definita su tutto \mathbb{R}^n vale*

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

allora $f(\mathbf{x})$ ammette minimo globale.

3 Ottimizzazione non vincolata

Nei problemi di ottimizzazione non vincolata si minimizza la funzione obiettivo senza porre restrizioni al valore che le variabili possono assumere.

Di solito non vi è una prospettiva globale dell'andamento di tale funzione. Tutto quello che sappiamo è il valore di $f(\mathbf{x})$ e (forse) di qualche sua derivata in un insieme di punti $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots$. Fortunatamente, questo è spesso tutto quello che serve ai nostri algoritmi per identificare una soluzione affidabile senza impiegare troppo tempo o troppa memoria.

Esempio Supponiamo di cercare una curva che si adatti ad un insieme di dati sperimentali, corrispondenti alle misure y_1, y_2, \dots, y_m di un segnale campionato agli istanti t_1, t_2, \dots, t_m . Dai dati in nostro possesso e da conoscenze relative all'applicazione, deduciamo che il segnale ha un andamento esponenziale ed oscillatorio, e scegliamo di modellarlo con la funzione

$$\phi(t, \mathbf{x}) = x_1 + x_2 e^{-(x_3 - t)^2 / x_4} + x_5 \cos(x_6 t)$$

I numeri reali x_1, x_2, \dots, x_6 sono i parametri del modello. Desideriamo sceglierli in modo da far aderire il più possibile i valori $\phi(t, \mathbf{x})$, calcolati dal modello, ai dati osservati y_1, y_2, \dots, y_m . Per trasformare il nostro obiettivo in un problema di ottimizzazione, definiamo i residui

$$r_j = y_j - \phi(t_j, \mathbf{x}), \quad j = 1, \dots, m,$$

che misurano la discrepanza fra i dati predetti dal modello e quelli misurati. La nostra stima del vettore \mathbf{x} verrà ottenuta risolvendo il problema

$$\min_{\mathbf{x} \in \mathbb{R}^6} f(\mathbf{x}) = r_1^2(\mathbf{x}) + r_2^2(\mathbf{x}) + \dots + r_m^2(\mathbf{x}).$$

Questo è un *problema ai minimi quadrati non lineare*, un caso speciale di ottimizzazione non vincolata (cfr. Sezione 8).

Si osservi che il calcolo della sola funzione obiettivo, per un dato valore delle variabili \mathbf{x} , può essere computazionalmente oneroso, anche se le variabili sono solo 6, dipendendo dal numero m di osservazioni, che può essere anche molto elevato.

Supponiamo ora che il valore $f(\mathbf{x}^*)$ della soluzione ottima \mathbf{x}^* sia diverso da zero, poiché, come spesso accade, il modello non ha riprodotto esattamente il valore di tutti i punti campionati.

Come possiamo verificare che \mathbf{x}^* è effettivamente un minimo per la funzione $f(\mathbf{x})$?

Per rispondere dobbiamo definire quali condizioni analitiche sono soddisfatte dai punti di minimo e spiegare come fare a riconoscerle.

4 Condizioni analitiche di ottimalità

La prima questione alla quale occorre dare risposta è la seguente:

Come è possibile riconoscere un minimo locale senza esplorare tutti gli infiniti punti che appartengono al suo intorno?

È possibile ottenere questo risultato quando la funzione non è *spigolosa*, ad esempio quando è due volte differenziabile con continuità o più formalmente se è di classe C^2 .

Definizione 6 Una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$ è detta di classe $C^k(X)$ se e solo se essa è differenziabile con continuità almeno k volte in tutti i punti $\mathbf{x} \in X \subseteq \mathbb{R}^n$.

Possiamo dire che un punto \mathbf{x}^* è un minimo locale, magari in senso stretto, esaminando il vettore gradiente e la matrice hessiana, calcolati in \mathbf{x}^* , $\nabla f(\mathbf{x}^*)$ e $H(\mathbf{x}^*)$. Lo strumento che ci permette di effettuare questa analisi è lo sviluppo in serie di Taylor di $f(\mathbf{x})$, arrestato al più al secondo ordine.

Proprietà 6 Sia $f : \mathbb{R}^n \rightarrow \mathbb{R}$ una funzione differenziabile con continuità. Dati un punto $\mathbf{x} \in \mathbb{R}^n$ ed un vettore $\mathbf{h} \in \mathbb{R}^n$, vale la seguente relazione

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \beta(\mathbf{x}, \mathbf{h})$$

dove $\beta(\mathbf{x}, \mathbf{h})$ è un infinitesimo che tende a zero più velocemente della norma di \mathbf{h} , o, in altri termini, vale

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x} + t\mathbf{h})^T \mathbf{h}$$

per qualche $t \in (0, 1)$. Inoltre, se $f(\mathbf{x})$ è una funzione di classe $C^2(\mathbb{R}^n)$ si ha

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T H(\mathbf{x}) \mathbf{h} + \beta_2(\mathbf{x}, \mathbf{h})$$

dove $\beta_2(\mathbf{x}, \mathbf{h})$ è un infinitesimo che tende a zero più velocemente della norma al quadrato di \mathbf{h} , o, in altri termini, valgono

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T H(\mathbf{x} + t\mathbf{h}) \mathbf{h}$$

e

$$\nabla f(\mathbf{x} + \mathbf{h}) = \nabla f(\mathbf{x}) + \int_0^1 H(\mathbf{x} + t\mathbf{h}) \mathbf{h} dt$$

per qualche $t \in (0, 1)$.

Mediante lo sviluppo in serie di Taylor si possono ricavare *condizioni necessarie* di ottimalità: si assume che \mathbf{x}^* sia un minimo locale e si studiano proprietà di $\nabla f(\mathbf{x}^*)$ e di $H(\mathbf{x}^*)$. Tali proprietà risultano più semplici da ricavare se prima introduciamo i concetti di *direzione di discesa* e di *derivata direzionale*.

Definizione 7 Data una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$, un vettore $\mathbf{d} \in \mathbb{R}^n$ si dice *direzione di discesa* per f in \mathbf{x} se esiste $\bar{\lambda} > 0$ tale che:

$$f(\mathbf{x} + \lambda \mathbf{d}) < f(\mathbf{x})$$

per ogni $0 < \lambda < \bar{\lambda}$.

In pratica, se siamo in un punto \mathbf{x} e ci spostiamo da \mathbf{x} lungo \mathbf{d} , la funzione, almeno per un tratto, decresce. Quindi, le direzioni di discesa sono quelle da seguire per migliorare il valore della funzione obiettivo. Si intuisce quindi che se in un punto \mathbf{x} esistono direzioni di discesa, quel punto non può essere un minimo locale. Per sapere se un punto è un minimo locale è necessario scoprire che non esistono direzioni di discesa in quel punto.

Ma come si possono individuare le direzioni di discesa?

Definizione 8 Sia data una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$, un vettore $\mathbf{d} \in \mathbb{R}^n$ e un punto \mathbf{x} dove f è definita. Se esiste il limite:

$$\lim_{\lambda \rightarrow 0^+} \frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda}$$

allora tale limite prende il nome di *derivata direzionale* della funzione f nel punto \mathbf{x} lungo la direzione \mathbf{d} .

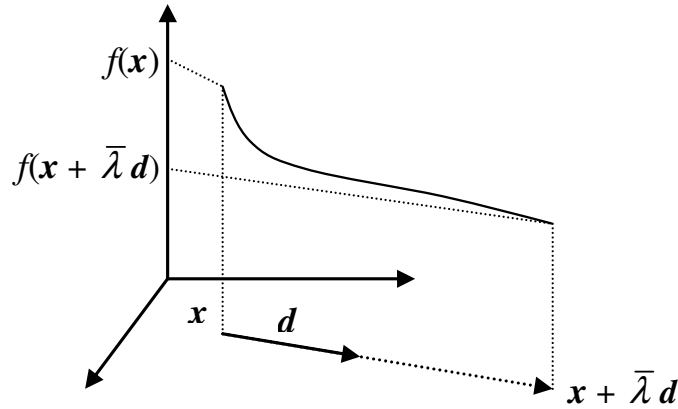


Figura 3: Direzione di discesa

Il significato geometrico della derivata direzionale è il seguente. Dato un punto $\mathbf{x} \in X$ e una direzione $\mathbf{d} \in \mathbb{R}^n$ si considerano i punti $\mathbf{x} + \lambda \mathbf{d}$, con $\lambda > 0$, cioè i punti che si incontrano partendo da \mathbf{x} e muovendosi lungo la direzione \mathbf{d} . Il valore della derivata direzionale rappresenta il tasso di variazione del valore della funzione obiettivo, calcolato in \mathbf{x} , lungo la direzione \mathbf{d} . In altre parole, come varia il valore di $f(\mathbf{x})$ se ci si allontana da \mathbf{x} lungo la direzione \mathbf{d} . Se $\mathbf{d} = \mathbf{e}_i$, cioè la direzione coincide con quella dell' i -esimo versore degli assi coordinati, allora la derivata direzionale coincide con la derivata parziale della funzione obiettivo nella componente i -esima: $\frac{\partial f}{\partial x_i}$. La derivata direzionale è utile in quanto può essere calcolata anche in punti nei quali una funzione non è differenziabile. Nel caso in cui invece sia possibile derivare la funzione $f(\mathbf{x})$, si ricava una interessante relazione fra derivata direzionale e vettore gradiente $\nabla f(\mathbf{x})$.

Proprietà 7 *Data una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$, derivabile con continuità in $\mathbf{x} \in \mathbb{R}^n$ la derivata direzionale di f nel punto \mathbf{x} lungo la direzione \mathbf{d} è data da $\nabla f(\mathbf{x})^T \mathbf{d}$.*

Dimostrazione Applicando lo sviluppo in serie di Taylor arrestato al primo ordine ad f in \mathbf{x} ricaviamo: $f(\mathbf{x} + \lambda \mathbf{d}) = f(\mathbf{x}) + \lambda \nabla f(\mathbf{x})^T \mathbf{d} + \beta(\mathbf{x}, \lambda \mathbf{d})$. Dividendo per λ e riarrangiando i termini si ottiene: $\frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda} = \nabla f(\mathbf{x})^T \mathbf{d} + \frac{\beta(\mathbf{x}, \lambda \mathbf{d})}{\lambda}$. Il limite per $\lambda \rightarrow 0$ del termine a sinistra del segno di uguale porta al risultato. \square

Il legame fra derivata direzionale e gradiente è proprio quello che ci serve per legare le informazioni che ci dà il vettore gradiente calcolato in un punto \mathbf{x} con la presenza di direzioni di discesa in \mathbf{x} .

Proprietà 8 *Sia data una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$, derivabile con continuità in \mathbb{R}^n . Dati un punto \mathbf{x} ed una direzione $\mathbf{d} \in \mathbb{R}^n$, se la derivata direzionale di f in \mathbf{x} lungo \mathbf{d} è negativa, allora la direzione \mathbf{d} in \mathbf{x} è di discesa.*

Dimostrazione Sappiamo che $\lim_{\lambda \rightarrow 0^+} \frac{f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x})}{\lambda} = \nabla f(\mathbf{x})^T \mathbf{d}$ e che per ipotesi $\nabla f(\mathbf{x})^T \mathbf{d} < 0$. Quindi, per λ sufficientemente piccolo, vale $f(\mathbf{x} + \lambda \mathbf{d}) - f(\mathbf{x}) < 0$, che è il significato di direzione di discesa. \square

Quindi, qualunque direzione formi un angolo ottuso ($> 90^\circ$) con il gradiente della funzione obiettivo calcolato in \mathbf{x} è una direzione di discesa, e tutti i punti sufficientemente vicini a \mathbf{x}

lungo tale direzione sono caratterizzati dall'avere un valore della funzione f inferiore a quello che ha f in \mathbf{x} . Inoltre, l'antigradiente $-\nabla f(\mathbf{x})$ è esso stesso una direzione di discesa (mentre $\nabla f(\mathbf{x})$ è una direzione di salita).

Questo significa che se in un punto \mathbf{x} il vettore $\nabla f(\mathbf{x})$ non è nullo, allora esistono direzioni di discesa in \mathbf{x} e tale punto non può essere un minimo locale.

Teorema 1 *Data una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$, derivabile con continuità in $\mathbf{x}^* \in \mathbb{R}^n$, condizione necessaria affinché il punto \mathbf{x}^* sia un minimo locale per f è che il gradiente della funzione calcolato in \mathbf{x}^* sia nullo: $\nabla f(\mathbf{x}^*) = \mathbf{0}$.*

Dimostrazione Se fosse $\nabla f(\mathbf{x}^*) \neq \mathbf{0}$ allora $-\nabla f(\mathbf{x}^*)$ sarebbe una direzione di discesa ed esisterebbe un punto vicino a \mathbf{x}^* , $\mathbf{x}^* - \lambda \nabla f(\mathbf{x}^*)$, in cui il valore della funzione obiettivo sarebbe $f(\mathbf{x}^* - \lambda \nabla f(\mathbf{x}^*)) < f(\mathbf{x}^*)$, contraddicendo il fatto che \mathbf{x}^* sia un minimo locale. \square

Il Teorema 1 fornisce delle condizioni molto generali, dette condizioni necessarie di ottimalità del 1° ordine. Un punto che soddisfa tali condizioni si dice *punto stazionario*, e PS indica l'insieme dei punti stazionari per f , ossia: $PS = \{\mathbf{x} : \mathbf{x} \in \mathbb{R}^n, \nabla f(\mathbf{x}) = \mathbf{0}\}$

Naturalmente, i punti stazionari possono essere sia punti di minimo locale, sia punti di ...massimo locale! E, purtroppo, anche punti che non sono né punti di massimo, né punti di minimo. Per dirimere questa nuova questione (punti di minimo, punti di massimo, altri punti) abbiamo bisogno di usare le informazioni che ci può dare la matrice hessiana.

Teorema 2 *Data una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$ di classe $C^2(\mathbf{x}^*)$, condizioni necessarie affinché il punto \mathbf{x}^* sia un minimo locale per f sono che il gradiente della funzione calcolato in \mathbf{x}^* , $\nabla f(\mathbf{x}^*)$, sia nullo e che valga la relazione: $\mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} \geq 0$ per ogni $\mathbf{d} \in \mathbb{R}^n$.*

Dimostrazione Per lo sviluppo in serie di Taylor arrestato al secondo termine:

$$f(\mathbf{x}^* + \lambda \mathbf{d}) = f(\mathbf{x}^*) + \lambda \nabla f(\mathbf{x}^*)^T \mathbf{d} + \frac{1}{2} \lambda^2 \mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} + \beta_2(\mathbf{x}^*, \lambda \mathbf{d})$$

e poiché $\nabla f(\mathbf{x}^*) = \mathbf{0}$ si può ricavare

$$\frac{f(\mathbf{x}^* + \lambda \mathbf{d}) - f(\mathbf{x}^*)}{\lambda^2} = \frac{1}{2} \mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} + \frac{\beta_2(\mathbf{x}^*, \lambda \mathbf{d})}{\lambda^2}$$

Poiché per ipotesi \mathbf{x}^* è minimo locale allora, per λ sufficientemente piccolo, il termine a sinistra è non negativo. Risulta quindi $\frac{1}{2} \mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} + \frac{\beta_2(\mathbf{x}^*, \lambda \mathbf{d})}{\lambda^2} \geq 0$. Passando al limite per $\lambda \rightarrow 0$, e osservando che \mathbf{d} è una direzione qualsiasi, segue la tesi. \square

Le condizioni espresse nel Teorema 2 sono dette condizioni necessarie di ottimalità del 2° ordine. Le condizioni appena espresse non sono però sufficienti. Se il punto $\mathbf{x}^* \in PS$ possiamo calcolare la matrice hessiana $H(\mathbf{x}^*)$ e se la matrice soddisfa la relazione $\mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} \geq 0$ per ogni $\mathbf{d} \in \mathbb{R}^n$ possiamo solo escludere che il punto sia un massimo locale. Tuttavia, se la funzione obiettivo è due volte differenziabile con continuità è possibile enunciare anche (finalmente) delle condizioni sufficienti di ottimalità del 2° ordine. Cioè condizioni che, se soddisfatte da un punto, garantiscono che quel punto sia di minimo.

Teorema 3 *Data una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$ di classe $C^2(\mathbf{x}^*)$, condizioni sufficienti affinché il punto \mathbf{x}^* sia un minimo locale in senso stretto per f sono che il gradiente della funzione calcolato in \mathbf{x}^* , $\nabla f(\mathbf{x}^*)$, sia nullo e che valga la relazione: $\mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} > 0$ per ogni $\mathbf{d} \in \mathbb{R}^n$.*

Dimostrazione Per lo sviluppo in serie di Taylor arrestato al secondo termine:

$$f(\mathbf{x}^* + \lambda \mathbf{d}) = f(\mathbf{x}^*) + \lambda \nabla f(\mathbf{x}^*)^T \mathbf{d} + \frac{1}{2} \lambda^2 \mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} + \beta_2(\mathbf{x}^*, \lambda \mathbf{d})$$

e poiché $\nabla f(\mathbf{x}^*) = \mathbf{0}$ si può ricavare

$$\frac{f(\mathbf{x}^* + \lambda \mathbf{d}) - f(\mathbf{x}^*)}{\lambda^2} = \frac{1}{2} \mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} + \frac{\beta_2(\mathbf{x}^*, \lambda \mathbf{d})}{\lambda^2}$$

Poiché per λ sufficientemente piccolo, il termine a destra è dominato dalla componente $\frac{1}{2} \mathbf{d}^T H(\mathbf{x}^*) \mathbf{d}$, risulta $f(\mathbf{x}^* + \lambda \mathbf{d}) - f(\mathbf{x}^*) > 0$. Osservando che \mathbf{d} è una direzione qualsiasi, segue la tesi. \square

Si noti che le condizioni sufficienti del secondo ordine appena enunciate garantiscono qualcosa in più delle condizioni necessarie discusse precedentemente, e cioè che il punto in questione è un minimo locale *in senso stretto*. Inoltre, le condizioni appena enunciate sono sufficienti ma non necessarie, cioè un punto può essere un minimo locale in senso stretto e non soddisfare le condizioni sufficienti del secondo ordine. Ad esempio, il punto $x^* = 0$ è un minimo in senso stretto della funzione $f(x) = x^4$, ma in tale punto l'hessiano si annulla e non è quindi definito positivo.

Apparentemente però abbiamo solo spostato il problema. Infatti, non sembra arduo valutare se il vettore gradiente ha o meno tutte le componenti nulle, ma sembra invece computazionalmente complicato stabilire se una matrice quadrata H soddisfi o meno la relazione $\mathbf{d}^T H \mathbf{d} > 0$ per ogni $\mathbf{d} \in \mathbb{R}^n$.

Definizione 9 *Una matrice H quadrata di ordine n , si dice (semi)-definita positiva su un insieme $X \subseteq \mathbb{R}^n$ se per ogni $\mathbf{d} \in X$, $\mathbf{d} \neq \mathbf{0}$, vale*

$$\begin{aligned} \mathbf{d}^T H \mathbf{d} &> 0 \text{ definita positiva} \\ \mathbf{d}^T H \mathbf{d} &\geq 0 \text{ semi-definita positiva} \end{aligned}$$

Una matrice H è (semi)-definita negativa se $-H$ è (semi)-definita positiva. Se H non è né (semi)-definita positiva, né (semi)-definita negativa, allora la matrice è indefinita.

Ora, se una matrice è *simmetrica*, come è il caso di ogni matrice hessiana, è possibile stabilire se essa è (semi)-definita positiva osservando il segno dei determinanti delle n sottomatrici quadrate che si ottengono considerando le matrici formate dalle sue prime i righe ed i colonne.

Proprietà 9 *Una matrice H simmetrica è (semi)-definita positiva se e solo se i determinanti di tutti i minori principali della matrice sono $(\geq) > 0$*

Un'altra caratterizzazione della (semi)-definita positività di una matrice H è data dal segno dei suoi autovalori.

Proprietà 10 Una matrice H simmetrica è (semi)-definita positiva se e solo se ha tutti gli autovalori (\geq) > 0 .

Inoltre gli autovalori di H simmetrica possono fornire un limite superiore o inferiore al valore di $\mathbf{d}^T H \mathbf{d}$ per ogni valore di \mathbf{d} : $\lambda_{\min} \|\mathbf{d}\|^2 \leq \mathbf{d}^T H \mathbf{d} \leq \lambda_{\max} \|\mathbf{d}\|^2$, dove $\lambda_{\min} :=$ autovalore minimo e $\lambda_{\max} :=$ autovalore massimo di H .

In pratica quindi, in un punto di stazionarietà \mathbf{x}^* , si deve poter analizzare efficientemente il segno degli autovalori della matrice hessiana $H(\mathbf{x}^*)$, o di matrici che la approssimano, spesso sfruttando una loro rappresentazione come prodotti di matrici triangolari.

Osserviamo infine che, se in un punto \mathbf{x}^* si annulla il gradiente, ma l'Hessiana è indefinita allora \mathbf{x}^* non è né un punto di minimo né un punto di massimo.

4.1 Condizioni analitiche di ottimalità per problemi convessi

Come abbiamo già detto i problemi convessi rivestono un ruolo particolare in ottimizzazione. La ragione principale è che nel caso di problemi convessi saper risolvere il problema di ottimizzazione locale (determinare cioè un qualsiasi minimo locale) implica saper risolvere il problema di ottimizzazione globale. Infatti, nel caso dei problemi convessi, i minimi locali coincidono con i minimi globali.

Proprietà 11 In un problema di ottimizzazione convesso (cioè con funzione obiettivo e regione ammissibile entrambe convesse) ogni minimo locale è anche minimo globale.

Quindi, nel caso di problemi convessi l'identificazione di un minimo locale è sufficiente per la loro soluzione. Di conseguenza è preliminare ad ogni approccio risolutivo stabilire se il problema che si intende risolvere è convesso o meno. Di seguito richiamiamo utili proprietà che permettono di stabilire, dall'analisi del vettore gradiente o della matrice hessiana calcolati in un punto, quando una funzione è convessa nell'intorno di quel punto. Di seguito poi, presenteremo proprietà che risulteranno utili nell'identificare i minimi locali (e quindi globali) di funzioni convesse.

Proprietà 12 Una funzione $f : X \rightarrow \mathbb{R}$, derivabile con continuità nell'insieme convesso X , è convessa su X se e solo se, comunque presi due punti \mathbf{x} e $\mathbf{y} \in X$, vale la relazione:

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x})$$

Dimostrazione Se f è convessa allora vale la relazione $\lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x}) \geq f(\lambda \mathbf{y} + (1 - \lambda)\mathbf{x})$ che possiamo riscrivere come

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \frac{f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\lambda}$$

e passando la limite per $\lambda \rightarrow 0$ si ha la tesi.

Viceversa, se vale $f(\mathbf{y}) - f(\mathbf{x}) \geq \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x})$, allora per $\mathbf{z} = \lambda \mathbf{y} + (1 - \lambda)\mathbf{x} \in X$, con $\lambda \in [0, 1]$ valgono $f(\mathbf{y}) - f(\mathbf{z}) \geq \nabla f(\mathbf{z})^T (\mathbf{y} - \mathbf{z})$ e $f(\mathbf{x}) - f(\mathbf{z}) \geq \nabla f(\mathbf{z})^T (\mathbf{x} - \mathbf{z})$. Moltiplicando la prima espressione per λ , la seconda per $(1 - \lambda)$ e sommando si ottiene

$$\lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x}) - f(\mathbf{z}) \geq \nabla f(\mathbf{z})^T (\lambda(\mathbf{y} - \mathbf{z}) + (1 - \lambda)(\mathbf{x} - \mathbf{z})) = 0,$$

poiché $\lambda(\mathbf{y} - \mathbf{z}) + (1 - \lambda)(\mathbf{x} - \mathbf{z}) = \mathbf{0}$, e quindi

$$\lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x}) \geq f(\mathbf{z}) = f(\lambda\mathbf{y} + (1 - \lambda)\mathbf{x})$$

□

Proprietà 13 *Sia X un insieme convesso e sia $f : X \rightarrow \mathbb{R}$ di classe $C^2(X)$, allora f è convessa se e solo se la matrice hessiana $H(\mathbf{x})$ è semidefinita positiva in \mathbf{x} , per ogni punto $\mathbf{x} \in X$.*

In pratica quindi, anche per stabilire la convessità della funzione obiettivo si deve poter analizzare efficientemente il segno degli autovalori della matrice hessiana H .

Inoltre, quando analizzeremo i problemi di ottimizzazione vincolata, sfruttando le proprietà 1, 2 e 3, saremo interessati allo studio delle matrici hessiane delle funzioni $g_i(\mathbf{x})$ allo scopo di stabilire se la regione ammissibile sia convessa o meno.

Proprietà 14 *Data una funzione $f : \mathbb{R}^n \rightarrow \mathbb{R}$ di classe $C^2(\mathbb{R}^n)$, sia \mathbf{x} un punto in cui la matrice hessiana $H(\mathbf{x})$ è definita positiva. Allora, almeno in un intorno di \mathbf{x} , f è strettamente convessa.*

Dimostrazione Dati due punti \mathbf{x} e \mathbf{y} , approssimiamo f con Taylor, nell'intorno di \mathbf{x} :

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T H(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \beta_2(\mathbf{x}, \mathbf{y})$$

Poiché $H(\mathbf{x})$ è definita positiva vale $(\mathbf{y} - \mathbf{x})^T H(\mathbf{x})(\mathbf{y} - \mathbf{x}) > 0$ e, almeno in un intorno di \mathbf{x} , questa espressione domina su $\beta_2(\mathbf{x}, \mathbf{y})$. Quindi, se al secondo membro sottraiamo questa quantità positiva, otteniamo $f(\mathbf{y}) - f(\mathbf{x}) > \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})$. □

Come detto le funzioni convesse rivestono un ruolo importante nella PNL poiché sono le uniche per le quali il problema di ottimizzazione globale è, almeno in linea di principio risolubile. In particolare, le condizioni analitiche si semplificano sensibilmente.

Teorema 4 *Data una funzione convessa $f : \mathbb{R}^n \rightarrow \mathbb{R}$ derivabile con continuità, condizione necessaria e sufficiente affinché il punto \mathbf{x}^* sia un minimo locale per f è che $\nabla f(\mathbf{x}^*) = \mathbf{0}$.*

Dimostrazione Per la convessità di f per ogni coppia di punti \mathbf{x} e $\mathbf{y} \in \mathbb{R}^n$ vale $f(\mathbf{y}) - f(\mathbf{x}) \geq \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})$ e poiché nel punto di minimo locale \mathbf{x}^* il gradiente si annulla si ricava $f(\mathbf{y}) \geq f(\mathbf{x})$. □

Nella funzione convessa ogni punto di minimo locale è anche di minimo globale (cfr. proprietà 11) quindi se abbiamo due punti distinti a gradiente nullo, allora abbiamo infiniti punti di stazionarietà (a gradiente nullo).

Teorema 5 *Data una funzione convessa $f : \mathbb{R}^n \rightarrow \mathbb{R}$, l'insieme dei punti di minimo della funzione è convesso.*

Dimostrazione Immediata conseguenza del punto .d della proprietà 2. □

Purtroppo la maggior parte dei problemi di ottimizzazione che si incontrano in pratica non ammettono modellizzazioni convesse. Come conseguenza l'individuazione di un punto di minimo globale può risultare un'operazione molto costosa, in termini computazionali, se non addirittura impossibile. In queste situazioni il meglio che si può fare è individuare un insieme di minimi locali e da essi dedurre stime, più o meno buone, del valore della soluzione ottima.

4.2 Condizioni empiriche di ottimalità

Dal punto di vista pratico un algoritmo di tipo iterativo viene fatto terminare se soddisfa un qualche criterio di convergenza di tipo empirico. Si tratta in pratica di trovare delle relazioni di tipo quantitativo che ci permettano di riconoscere il punto corrente, nella successione $\{\mathbf{x}_k\}$ di punti visitati dall'algoritmo, come un punto stazionario. In genere la terminazione viene imposta qualora all'iterazione k uno, o una combinazione dei seguenti criteri viene soddisfatta

$$\begin{aligned}\|\mathbf{x}_k - \mathbf{x}_{k-1}\| &< \varepsilon_1 \\ \|\nabla f(\mathbf{x}_k)\| &< \varepsilon_2 \\ \|f(\mathbf{x}_k) - f(\mathbf{x}_{k-1})\| &< \varepsilon_3\end{aligned}$$

dove $\varepsilon_1, \varepsilon_2$ e ε_3 sono prescritti valori positivi di tolleranza.

4.3 Programmazione Quadratica

Analizziamo ora una classe di problemi che riveste una importanza particolare nell'ambito delle tecniche di ottimizzazione. Si tratta dei problemi di tipo *quadratico*:

$$\begin{aligned}\min \quad & f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x} \\ \text{t.c.} \quad & \mathbf{x} \in \mathbb{R}^n\end{aligned}$$

dove Q è una matrice quadrata simmetrica di ordine n .

Nel caso di funzioni obiettivo quadratiche¹ si ricava immediatamente $\nabla f(\mathbf{x}) = Q\mathbf{x} - \mathbf{b}$ e $H(\mathbf{x}) = Q$. In particolare la matrice hessiana è costante e il vettore gradiente è lineare.

Nell'ambito della programmazione quadratica la condizione necessaria di ottimalità del primo ordine diventa: $\nabla f(\mathbf{x}) = \mathbf{0}$, che si riduce al sistema lineare in n equazioni ed n incognite $Q\mathbf{x} = \mathbf{b}$. In pratica, quando la matrice Q è invertibile, possiamo, *imponendo* il soddisfacimento della condizione necessaria di ottimalità del primo ordine, ricavare l'unico punto stazionario di $f(\mathbf{x})$:

$$\mathbf{x}^* = Q^{-1}\mathbf{b}$$

Ora, per distinguere se tale punto sia di minimo o di massimo è sufficiente analizzare la matrice hessiana $H = Q$. Come conseguenza della Proprietà 14, se Q è definita positiva allora la funzione quadratica f è convessa ovunque e, in conseguenza del teorema 4, il punto \mathbf{x}^* è un minimo globale.

Analizziamo ora i casi che si possono presentare.

4.3.1 Q non è semidefinita positiva

In questo caso non è vero che $\forall \mathbf{x}$ vale $\mathbf{x}^T Q \mathbf{x} \geq 0$ e quindi non possiamo limitare la nostra attenzione al solo vettore gradiente. La funzione obiettivo non ha punti di minimo.

4.3.2 Q è definita positiva

In questo caso $\mathbf{x}^* = Q^{-1}\mathbf{b}$ è l'unico ottimo globale del problema.

¹Si osservi che tali problemi non hanno nulla a che fare con i *problemi ai minimi quadrati* introdotti nella Sezione 8.

4.3.3 Q è semidefinita positiva

In questo caso se Q è invertibile ricadiamo nel caso precedente, altrimenti in uno dei due sottocasi seguenti

- non esistono soluzioni, cioè non esiste minimo
- ci sono infinite soluzioni.

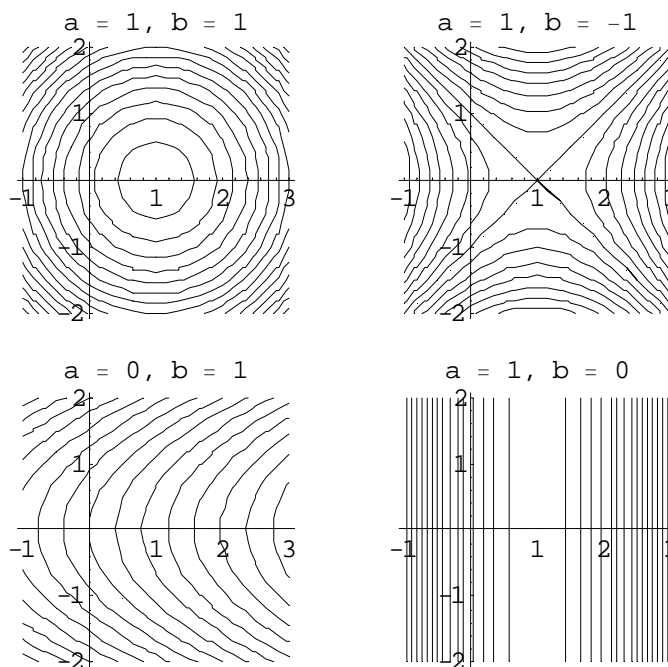


Figura 4: Esempi di funzioni quadratiche

4.3.4 Esempi

Consideriamo la funzione in due variabili $f(x, y) = \frac{1}{2} (ax^2 + by^2) - x$ che possiamo riscrivere come

$$f(\mathbf{x}) = \frac{1}{2} (x, y) \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} - (x, y) \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

Si osservi che in questo esempio, a e b coincidono con gli autovalori di Q , poiché essa è diagonale.

Primo caso: $a > 0$, $b > 0$.

$$\nabla f(\mathbf{x}) = \begin{pmatrix} ax - 1 \\ by \end{pmatrix} = \mathbf{0}, \quad \mathbf{x}^* = \begin{pmatrix} 1/a \\ 0 \end{pmatrix} = \mathbf{0}$$

Q è invertibile e definita positiva. \mathbf{x}^* è ottimo globale.

Secondo caso: $a > 0$, $b < 0$

$$\nabla f(\mathbf{x}) = \begin{pmatrix} ax - 1 \\ by \end{pmatrix} = \mathbf{0}, \quad \mathbf{x}^* = \begin{pmatrix} 1/a \\ 0 \end{pmatrix} = \mathbf{0}$$

Q è invertibile ma non semidefinita positiva. \mathbf{x}^* è un punto di sella.

Terzo caso: $a = 0$, $b \neq 0$

$\nabla f(\mathbf{x}) = \begin{pmatrix} -1 \\ by \end{pmatrix}$, sempre diverso da $\mathbf{0}$, Q non è invertibile e non esiste soluzione.

Se $b > 0$ ($b < 0$) allora Q è semidefinita positiva (negativa).

Quarto caso: $a \neq 0$, $b = 0$

$\nabla f(\mathbf{x}) = \begin{pmatrix} ax - 1 \\ 0 \end{pmatrix} = \mathbf{0}$. Caso di ∞ soluzioni: i punti in cui il gradiente si annulla sono tutti i punti $(1/a, \gamma)$ con $\gamma \in \mathbb{R}$. Q non è invertibile.

Se $a > 0$ ($a < 0$) allora Q è semidefinita positiva (negativa) e gli infiniti punti stazionari sono tutti punti di minimo (massimo). La Figura 4 illustra i quattro casi.

5 Convergenza e rapidità di convergenza

In tutti i casi, e sono la maggioranza, in cui non si è in grado di calcolare analiticamente l'insieme dei punti stazionari PS , si deve ricorrere ad algoritmi di tipo iterativo. Come abbiamo detto all'inizio, tali algoritmi, a partire da un vettore iniziale di variabili \mathbf{x}_0 generano iterativamente una sequenza di vettori, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$, quali approssimazioni via a via migliori della soluzione ottima. Prima di descrivere, nella prossima sezione, tali algoritmi in maggior dettaglio, proviamo a rispondere qui alla seguente domanda:

In che modo possiamo classificare la bontà di un algoritmo iterativo nell'ambito della PNL?

Nell'analizzare le prestazioni di un algoritmo di discesa, una prima importante distinzione è stabilire se, e in che modo, esso *converga*.

In particolar modo, quando si parla di *convergenza* nell'ambito della PNL, si intende sempre *convergenza ad un minimo locale*.

È quindi nell'ambito della ricerca di un minimo locale, che faremo distinzione fra *convergenza locale* e *convergenza globale*.

Definizione 10 *Algoritmo globalmente convergente.* Un algoritmo è globalmente convergente se è convergente ad un punto di PS per qualunque $\mathbf{x}_0 \in \mathbb{R}^n$ (termina in un minimo locale da qualunque punto parta).

Definizione 11 *Algoritmo localmente convergente.* Un algoritmo è localmente convergente se è convergente solo per \mathbf{x}_0 tale che $\mathbf{x}_0 \in I(\mathbf{x}^*, \varepsilon)$, $\mathbf{x}^* \in PS$, (termina in un minimo locale solo partendo da un intorno del minimo locale).

Quindi un algoritmo *globalmente convergente* converge ad un minimo locale da qualunque punto parta, mentre un algoritmo *localmente convergente*, può convergere ad un minimo locale, o non convergere affatto, al variare del punto di partenza.

Ma, anche supponendo che le condizioni di convergenza di un algoritmo di discesa siano soddisfatte, e che esso, a partire dal punto \mathbf{x}_0 , converga ad un punto di minimo locale \mathbf{x}^* , è necessario chiedersi in che modo, cioè con che *velocità*, esso converge. Si tratta cioè di caratterizzare la *rapidità* con cui tale convergenza avviene.

Nell'ambito della PL, PLI e dei problemi di ottimizzazione su grafo, tutti gli algoritmi studiati sono in grado di calcolare la soluzione ottima in un numero finito, eventualmente esponenziale,

di iterazioni. In quell'ambito, è possibile fondare la definizione di *efficienza* di un algoritmo di ottimizzazione sulla base del *numero di iterazioni* necessarie, nel caso peggiore, per giungere alla soluzione ottima.

Nel caso della *programmazione non lineare*, tranne che in casi particolari (ad es. il metodo di Newton per funzioni quadratiche, cfr. la Sezione 6.2.3), gli algoritmi risolutivi producono, addirittura, una successione *infinita* di punti \mathbf{x}_k . Per misurare il grado di convergenza di un algoritmo è quindi necessario introdurre nuovi criteri.

I metodi più utilizzati per misurare la convergenza adottano come criterio di convergenza il rapporto tra gli scostamenti esistenti, ad un'iterazione ed alla successiva, tra la soluzione corrente \mathbf{x}_k ed il punto limite \mathbf{x}^* al quale l'algoritmo dovrebbe convergere, cioè:

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|}$$

Definizione 12 Sia $\{\mathbf{x}_k\}$ una successione in \mathbb{R}^n che converge a \mathbf{x}^* . Si dice che la convergenza è *Q-lineare* se esiste un valore costante $r \in (0, 1)$ tale che per ogni $k \geq \bar{k}$ vale

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq r$$

Questo significa che la distanza dalla soluzione \mathbf{x}^* decresce ad ogni iterazione di almeno un fattore costante limitato da 1. Per esempio, la successione $1 + (0.5)^k$ converge Q-linearmente ad 1, con tasso $r = 0.5$. Il prefisso Q significa quoziente, poiché questo tipo di convergenza è definita in termini di quozienti di errori successivi.

Se vale la relazione $r \geq 1$ si dice che l'algoritmo ha convergenza *sublineare* e, aggiungiamo, non è di alcun interesse! Molto più interessanti invece gli algoritmi che presentano una convergenza *superlineare*

Definizione 13 Sia $\{\mathbf{x}_k\}$ una successione in \mathbb{R}^n che converge a \mathbf{x}^* . Si dice che la convergenza è *Q-superlineare* se

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0$$

Si osservi che l'esistenza del limite implica che:

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq r_k \text{ con } r_k \rightarrow 0 \text{ quando } k \rightarrow \infty.$$

Ma si può far di meglio.

Definizione 14 Sia $\{\mathbf{x}_k\}$ una successione in \mathbb{R}^n che converge a \mathbf{x}^* . Si dice che la convergenza è *Q-quadratica* se esiste un valore costante $C > 0$ tale che per ogni $k \geq \bar{k}$ vale

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} \leq C$$

La rapidità di convergenza dipende da r e da C , i cui valori dipendono non solo dall'algoritmo ma anche dalle proprietà del particolare problema. Indipendentemente comunque da tali valori, una successione che converge in modo quadratico convergerà sempre più velocemente di una che converge in modo lineare.

Naturalmente, ogni successione che converge quadraticamente converge anche superlinearmente ed ogni successione che converge superlinearmente converge anche linearmente. È possibile definire anche tassi di convergenza superiori (cubici, o oltre) ma non sono di reale interesse pratico. In generale, diciamo che il Q -ordine di convergenza è $p > 1$ se esiste una costante positiva M tale che per ogni $k \geq \bar{k}$ vale

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^p} \leq M$$

Le seguenti successioni hanno la rapidità di convergenza riportata a fianco.

$$\begin{aligned} \{\mathbf{x}_k\} &= \frac{1}{k} \text{ (sublineare);} & \{\mathbf{x}_k\} &= \frac{1}{2^k} \text{ (lineare);} \\ \{\mathbf{x}_k\} &= \frac{1}{k!} \text{ (superlineare);} & \{\mathbf{x}_k\} &= \frac{1}{2^{2^k}} \text{ (quadratica);} \end{aligned}$$

6 Metodi basati su ricerca lineare

La struttura degli algoritmi iterativi di discesa basati su ricerca lineare per problemi di minimizzazione non vincolata è semplice.

```

Metodo di discesa;
{
  Scegli  $\mathbf{x}_0 \in \mathbb{R}^n$ ;  $k := 0$ ;
  While  $\nabla f(\mathbf{x}_k) \neq \mathbf{0}$ ;
  {
    calcola  $\mathbf{d}_k \in \mathbb{R}^n$ ; /* direzione di discesa */
    calcola  $\alpha_k \in \mathbb{R}$ ; /* passo lungo  $\mathbf{d}_k$  */
     $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{d}_k$ ;
     $k := k + 1$ ;
  }
}

```

Si genera un punto iniziale $\mathbf{x}_0 \in \mathbb{R}^n$, si calcola il valore della funzione $f(\mathbf{x}_0)$ e del gradiente $\nabla f(\mathbf{x}_0)$. Se vale $\nabla f(\mathbf{x}_0) = \mathbf{0}$, allora $\mathbf{x}_0 \in PS$. Altrimenti, da \mathbf{x}_0 ci si sposta in cerca di un punto \mathbf{x}_1 , con un valore di $f(\mathbf{x}_1)$ migliore di quello di $f(\mathbf{x}_0)$. A questo scopo si sceglie una direzione di discesa \mathbf{d}_0 e lungo tale direzione ci si muove di un passo opportuno α_0 . Trovato \mathbf{x}_1 , se $\nabla f(\mathbf{x}_1) = \mathbf{0}$ allora $\mathbf{x}_1 \in PS$ e l'algoritmo si arresta, altrimenti si cerca un nuovo punto \mathbf{x}_2 e così via. L'iterazione generica è quindi semplicemente

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \tag{1}$$

Gli aspetti critici sono ovviamente la scelta di \mathbf{d}_k (direzione di discesa) e di α_k (passo lungo la direzione di discesa).

6.1 Determinazione del passo α_k

Affrontiamo per primo il problema di determinare in modo efficace il valore di α_k nell'ipotesi di conoscere la direzione di discesa \mathbf{d}_k . La scelta della direzione \mathbf{d}_k sarà l'argomento della Sezione 6.2.

Nel seguito, per semplificare la notazione ometteremo il pedice k , relativo all'indice dell'iterazione corrente nella relazione (1).

In primo luogo esprimiamo la variazione della funzione obiettivo f al variare di α :

$$\phi(\alpha) = f(\mathbf{x} + \alpha \mathbf{d}).$$

Poiché in un determinato punto \mathbf{x} e per una determinata direzione \mathbf{d} solo il valore α può variare con continuità, ci siamo ricondotti alla soluzione di un problema di PNL in una sola variabile. Dobbiamo cioè risolvere il problema

$$\min \phi(\alpha) = f(\mathbf{x} + \alpha \mathbf{d}), \quad \alpha > 0.$$

Iniziamo col calcolare la derivata della funzione $\phi(\alpha)$. Per prima cosa indichiamo con \mathbf{y} il punto incrementato in funzione di α :

$$\mathbf{y} = \mathbf{x} + \alpha \mathbf{d}.$$

Nell'ipotesi che $f(\mathbf{y})$ sia differenziabile con continuità, e che quindi il vettore gradiente $\nabla f(\mathbf{y})$ sia una funzione continua, possiamo ricavare la derivata di $\phi(\alpha) = f(\mathbf{y}(\alpha))$ come derivata di funzione composta:

$$\phi'(\alpha) = \frac{d\phi}{d\alpha} = \sum_{i=1}^n \frac{\partial f(\mathbf{y})}{\partial y_i} \frac{dy_i}{d\alpha}.$$

che si può riscrivere (grazie alla definizione di gradiente) come:

$$\phi'(\alpha) = \nabla f(\mathbf{y})^T \mathbf{d} = \nabla f(\mathbf{x} + \alpha \mathbf{d})^T \mathbf{d}.$$

$\phi'(\alpha)$ rappresenta il tasso di variazione del valore della funzione obiettivo man mano che ci allontaniamo dal punto corrente \mathbf{x} lungo la direzione \mathbf{d} . In particolare, per $\alpha = 0$, $\phi'(0)$ coincide con la derivata direzionale calcolata nel punto \mathbf{x} , $\nabla f(\mathbf{x})^T \mathbf{d}$ (cfr. Definizione 8 e Proposizione 7).

Abbiamo ora gli elementi formali che ci permetteranno di individuare il valore migliore per il passo α . Come spesso accade siamo però di fronte ad una scelta: da una lato vorremmo calcolare α in modo da ottenere una sostanziale riduzione del valore di f , ma allo stesso tempo non vorremmo spendere troppo tempo in questo calcolo. L'ideale sarebbe trovare il minimo globale della funzione $\phi(\alpha)$, per valori $\alpha > 0$, ma in generale è troppo costoso identificare tale valore. Anche trovare un minimo locale di $\phi(\alpha)$ con una certa accuratezza richiede generalmente troppe valutazioni della funzione obiettivo f e magari del gradiente $\nabla f(\mathbf{x})$. Le strategie più efficaci eseguono una ricerca monodimensionale *non esatta* per identificare una lunghezza del passo che permetta una adeguata riduzione di f col minimo sforzo computazionale.

Gli algoritmi di ricerca monodimensionali generano una sequenza di valori di α , fermandosi quando uno di tali valori soddisfa certe condizioni. Tali algoritmi identificano, in una prima fase, un intervallo entro il quale un adeguato valore del passo esiste, ed eseguono, in una successiva fase, una ricerca di tale valore dentro l'intervallo trovato.

Vediamo ora quali richieste deve soddisfare un buon algoritmo per la determinazione del passo α .

In primo luogo, osserviamo che la semplice diminuzione del valore di $f(\mathbf{x})$ ad ogni iterazione dell'algoritmo, non è sufficiente a garantire la convergenza del metodo di discesa ad un punto

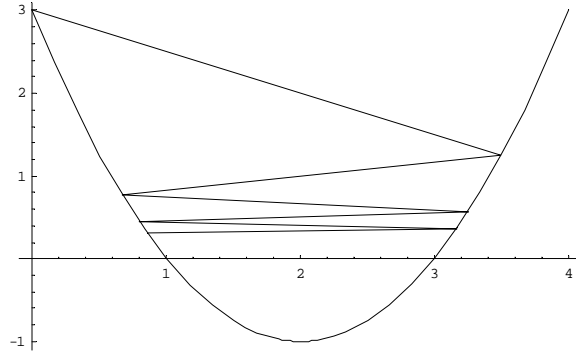


Figura 5: Convergenza a valori errati

stazionario.

Esempio. Consideriamo la funzione di una sola variabile $\phi(\alpha) = \alpha^2 - 4\alpha + 3$, convessa e avente come unico punto di minimo $\alpha^* = 2$, con valore $\phi(2) = -1$, e consideriamo il punto iniziale $\alpha_1 = 0$. Se la successione $\{\alpha_k\}$ di punti visitati dall'algoritmo viene generata dalla formula $\alpha_k = 2 + (-1)^k(1 + 1/k)$, otteniamo che ad ogni iterazione il valore della funzione obiettivo diminuisce: $\phi(\alpha_{k+1}) < \phi(\alpha_k)$, con $k = 1, 2, \dots$. Come si vede dalla Figura 5 la successione ha due punti limite, $\alpha' = 1$ e $\alpha'' = 3$, nei quali la funzione obiettivo vale 0. Dunque, la successione dei valori $\{\phi(\alpha_k)\}$ converge, ma al valore sbagliato. Il problema, in questo caso, è che la funzione ϕ diminuisce sì ad ogni iterazione, ma sempre di meno. Per evitare questo comportamento possiamo imporre una *condizione di sufficiente riduzione* del valore della funzione f .

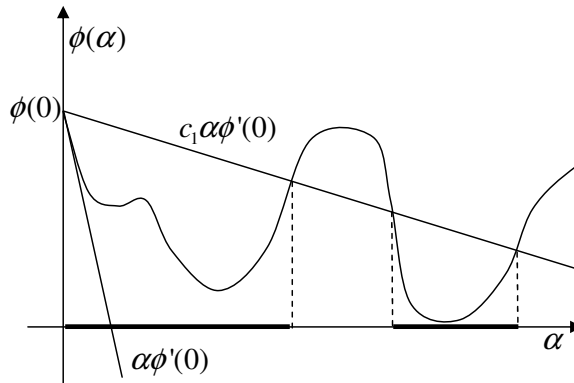


Figura 6: Condizione di Armijo

Condizione 1 *Condizione di sufficiente riduzione della f . Ad ogni iterazione dell'algoritmo, il punto incrementato deve soddisfare la seguente disuguaglianza:*

$$f(\mathbf{x} + \alpha \mathbf{d}) \leq f(\mathbf{x}) + c_1 \alpha \nabla f(\mathbf{x})^T \mathbf{d}, \quad (2)$$

con $c_1 \in (0, 1)$. In altre parole, la riduzione di f deve essere proporzionale sia alla lunghezza del passo α , sia alla derivata direzionale $\nabla f(\mathbf{x})^T \mathbf{d}$. La disuguaglianza (2) è detta *condizione di Armijo*. In termini di $\phi(\alpha)$ l'espressione diviene: $\phi(\alpha) \leq \phi(0) + \alpha c_1 \phi'(0)$. Come evidenzia

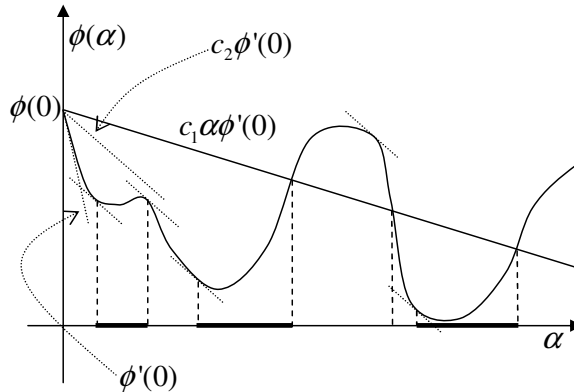


Figura 7: Condizioni di Wolfe

la Figura 6, consideriamo validi i valori di $\alpha > 0$ per i quali il grafico di $\phi(\alpha)$ sta al di sotto della retta di pendenza $c_1\phi'(0)$ (dunque negativa), passante per il punto $(0, \phi(0))$. La condizione di sufficiente riduzione da sola non basta ad assicurare che l'algoritmo faccia progressi ragionevoli in quanto essa è soddisfatta da tutti i valori di α sufficientemente piccoli. Per evitare passi troppo piccoli si aggiunge una seconda condizione.

Condizione 2 *Condizione di curvatura.* Ad ogni iterazione dell'algoritmo, il punto incrementato deve soddisfare la seguente disuguaglianza:

$$\nabla f(\mathbf{x} + \alpha \mathbf{d})^T \mathbf{d} \geq c_2 \nabla f(\mathbf{x})^T \mathbf{d} \quad (3)$$

con $c_2 \in (c_1, 1)$. In termini di $\phi(\alpha)$ l'espressione diviene: $\phi'(\alpha) \geq c_2\phi(0)$. La pendenza della funzione $\phi(\alpha)$ nei valori di α accettabili deve essere maggiore, di un fattore costante, alla pendenza che la funzione assume nel punto 0. Questo ha senso, poiché se la pendenza $\phi'(\alpha)$ è molto negativa abbiamo l'indicazione che possiamo ridurre f significativamente muovendoci ulteriormente nella direzione scelta. In sostanza, la condizione di curvatura vincola il passo α ad essere abbastanza lungo da percepire una significativa diminuzione (in valore assoluto) della derivata direzionale, fatto questo che indica un avvicinamento al minimo della funzione obiettivo. Si noti che anche i valori di α per i quali $\phi'(\alpha) > 0$ soddisfano la condizione di curvatura. Un valore tipicamente usato per c_1 è 10^{-4} , mentre il valore di c_2 dipende dal metodo usato per la scelta della direzione (0.9 se stiamo usando Newton o quasi-Newton e 0.1 se usiamo i metodi a gradiente coniugato). Le due condizioni appena introdotte sono note come *condizioni di Wolfe*. Come si vede dalla Figura 7, un valore del passo α può soddisfare le condizioni di Wolfe senza essere particolarmente vicino ad un minimo locale di $\phi(\alpha)$. Tali condizioni possono essere rinforzate per garantire che la lunghezza del passo stia almeno in un intorno di un minimo locale. Le *condizioni di Wolfe in senso forte* richiedono che α soddisfi, oltre alla condizione di Armijo, anche la disequazione:

$$|\nabla f(\mathbf{x} + \alpha \mathbf{d})^T \mathbf{d}| \leq c_2 |\nabla f(\mathbf{x})^T \mathbf{d}| \quad (4)$$

con $0 < c_1 < c_2 < 1$. In pratica si impedisce alla derivata $\phi'(\alpha)$ di essere troppo positiva. Le considerazioni appena introdotte vengono utilizzate all'interno di metodi iterativi, metodi ai quali dobbiamo ricorrere ogniqualvolta non si è in grado di calcolare analiticamente i valori di α per i quali $\phi'(\alpha) = 0$. Per nostra fortuna le condizioni sopra riportate valgono per un'ampia classe di funzioni.

Teorema 6 *Data una direzione di discesa \mathbf{d} per una funzione $f(\mathbf{x})$ di classe C^1 , se la funzione $\phi(\alpha)$ non è inferiormente limitata per $\alpha > 0$, allora esistono $0 < \alpha_1 < \alpha_2$ tali che per ogni $\alpha \in [\alpha_1, \alpha_2]$ valgono le condizioni di Wolfe anche in senso forte.*

Vediamo ora alcuni metodi per determinare il valore del passo α in modo da tener conto delle indicazioni fornite dalle condizioni di Wolfe.

6.1.1 Backtracking e metodo di Armijo

Questo metodo iterativo genera i valori di α in modo abbastanza accurato e converge con accettabile rapidità, pur facendo a meno della verifica puntuale della *condizione di curvatura* e utilizzando esplicitamente solo la condizione di sufficiente riduzione della funzione f . L'algoritmo è di tipo *backtracking*. Indicheremo con $\alpha_1, \alpha_2, \dots, \alpha_i, \dots$ i valori di α generati alle varie iterazioni, mentre α^* indica il valore restituito dall'algoritmo, che verrà quindi utilizzato come passo nell'iterazione (1).

Metodo di Armijo;
 {
 Scegli $\alpha_0 \in \mathbb{R}; \quad \alpha := \alpha_0;$
 While $f(\mathbf{x}_k + \alpha \mathbf{d}_k) > f(\mathbf{x}_k) + \alpha c_1 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k;$
 $\alpha := \sigma \alpha; \quad /* \text{backtracking} */$
 $\alpha^* := \alpha;$
 }

L'approccio backtracking consiste nello scegliere, inizialmente, un valore α_0 . Se già α_0 soddisfa la condizione di sufficiente riduzione, il procedimento termina e restituisce α_0 . Altrimenti, si moltiplica α_0 per un fattore di contrazione $0 < \sigma < 1/2$ e si prova il valore così generato. Il procedimento prosegue in questo modo fino a trovare un valore $\alpha^* = \sigma^i \alpha_0$ tale da soddisfare la condizione $\phi(\alpha^*) \leq \phi(0) + \alpha^* c_1 \phi'(0)$. L'idea è che il valore α^* restituito dal metodo, oltre a soddisfare tale condizione, non sia troppo piccolo, in quanto il valore trovato all'iterazione precedente, ossia $\sigma^{i-1} \alpha_0$, non era stato ritenuto soddisfacente, ossia era ancora troppo grande. La versione base, nella quale il metodo backtracking prevede di utilizzare ad ogni iterazione sempre lo stesso valore del fattore di contrazione σ , prende il nome di *metodo di Armijo*.

Teorema 7 *Dato un punto \mathbf{x} e una direzione di discesa \mathbf{d} , il metodo di Armijo determina in un numero finito di iterazioni un punto α^* tale che:*

$$f(\mathbf{x} + \alpha^* \mathbf{d}) < f(\mathbf{x})$$

Dimostrazione Per assurdo, se l'algoritmo non terminasse mai sarebbe sempre vera la relazione

$$\frac{f(\mathbf{x} + \sigma^j \alpha_0 \mathbf{d}) - f(\mathbf{x})}{\sigma^j \alpha_0} > c_1 \nabla f(\mathbf{x})^T \mathbf{d}$$

e al limite per $i \rightarrow \infty$ si avrebbe $\sigma^j \alpha_0 \rightarrow 0$, cioè $\nabla f(\mathbf{x})^T \mathbf{d} > c_1 \nabla f(\mathbf{x})^T \mathbf{d}$ che è impossibile per $0 < c_1 < 1$ poiché \mathbf{d} è direzione di discesa ($\nabla f(\mathbf{x})^T \mathbf{d} < 0$). \square

6.1.2 Metodo di ricerca esatto

Si tratta di calcolare analiticamente i valori α^* per i quali

$$\phi'(\alpha^*) = \nabla f(\mathbf{y})^T \mathbf{d} = \nabla f(\mathbf{x} + \alpha^* \mathbf{d})^T \mathbf{d} = 0.$$

Come si vede dalla Figura 8 è possibile dare una interpretazione geometrica alle condizioni

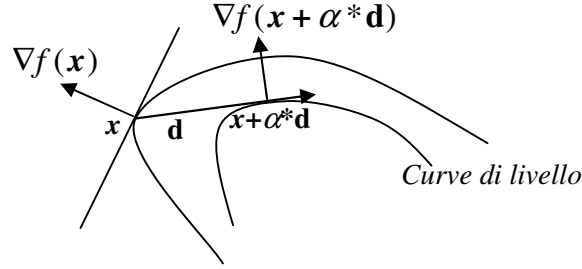


Figura 8: Scelta del passo ottima

analitiche che richiedono che un punto di minimo per la funzione $\phi(\alpha)$ sia un punto stazionario, quindi un punto in cui $\phi'(\alpha) = 0$. Consideriamo una direzione di discesa \mathbf{d} , che forma quindi con $\nabla f(\mathbf{x})$ un angolo ottuso, e spostiamoci in tale direzione. Il primo valore $\alpha^* > 0$ tale che $\phi'(\alpha^*) = 0$ individua un punto $\mathbf{x} + \alpha^* \mathbf{d}$, in cui l'angolo tra \mathbf{d} e $\nabla f(\mathbf{x} + \alpha^* \mathbf{d})$ è un angolo retto. Proseguendo ulteriormente nella direzione \mathbf{d} ci muoveremmo lungo una retta che giace nel piano tangente nel punto $\mathbf{x} + \alpha^* \mathbf{d}$ alla curva di livello passante per $\mathbf{x} + \alpha^* \mathbf{d}$ senza ottenere quindi alcun ulteriore miglioramento. Purtroppo solo raramente è possibile effettuare una ricerca esatta, essendo questa di solito onerosa dal punto di vista computazionale. Un caso però, merita un trattamento a parte: il caso di funzioni quadratiche (cfr. 4.3).

Nel caso di funzioni quadratiche con matrice Q simmetrica e definita positiva,

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x},$$

fissata una direzione di discesa \mathbf{d} , si ricava

$$\phi(\alpha) = f(\mathbf{x} + \alpha \mathbf{d}) = \frac{1}{2} (\mathbf{x} + \alpha \mathbf{d})^T Q (\mathbf{x} + \alpha \mathbf{d}) - \mathbf{b}^T (\mathbf{x} + \alpha \mathbf{d}).$$

Derivando $\phi(\alpha)$ rispetto ad α e uguagliando a zero si ottiene l'espressione

$$(Q(\mathbf{x} + \alpha \mathbf{d}) - \mathbf{b})^T \mathbf{d} = 0$$

dalla quale si ricava per α il valore ottimo

$$\alpha^* = -\frac{\mathbf{x}^T Q \mathbf{d} - \mathbf{b}^T \mathbf{d}}{\mathbf{d}^T Q \mathbf{d}} = -\frac{\nabla f(\mathbf{x})^T \mathbf{d}}{\mathbf{d}^T Q \mathbf{d}}. \quad (5)$$

Tale risultato ci tornerà utile quando per studiare le proprietà di convergenza di alcuni algoritmi li applicheremo a problemi di tipo quadratico.

6.1.3 Metodo di interpolazione

Questo metodo risulta efficace se la funzione $\phi(\alpha)$ non varia in maniera troppo accentuata da punto a punto. Esso può essere visto come un miglioramento del metodo di Armijo (cfr. Sezione 6.1.1) ed è parte integrante delle tecniche più sofisticate per la scelta di α .

Supponiamo che il valore iniziale α_0 sia dato. Se esso soddisfa la condizione di sufficiente riduzione, $\phi(\alpha_0) \leq \phi(0) + \alpha_0 c_1 \phi'(0)$ terminiamo la ricerca. Altrimenti, costruiamo una approssimazione quadratica di $\phi(\alpha)$ determinando i valori a, b , e c di un polinomio di secondo grado $\phi_q(\alpha) = a\alpha^2 + b\alpha + c$ in modo che soddisfi le condizioni $\phi_q(0) = \phi(0)$, $\phi'_q(0) = \phi'(0)$ e $\phi_q(\alpha_0) = \phi(\alpha_0)$.

Dalla prima relazione ricaviamo $c = \phi(0)$. Inoltre, da $\phi'_q(\alpha) = 2a\alpha + b$, e dalla seconda relazione ricaviamo $b = \phi'(0)$. Infine da $\phi_q(\alpha_0) = a\alpha_0^2 + \phi'(0)\alpha_0 + \phi(0)$ e dalla terza relazione ricaviamo $a = (\phi(\alpha_0) - \phi'(0)\alpha_0 - \phi(0))/\alpha_0^2$.

Tale parabola è convessa in quanto $a > 0$ (lo si dimostri per esercizio). Ora si cerca il minimo, α_1 della parabola, imponendo $\phi'_q(\alpha) = 2a\alpha + b = 0$. Otteniamo così il valore:

$$\alpha_1 = \frac{-\phi'(0)\alpha_0^2}{2(\phi(\alpha_0) - \phi'(0)\alpha_0 - \phi(0))}$$

A questo punto se α_1 soddisfa la condizione di sufficiente riduzione abbiamo terminato, altrimenti iteriamo il procedimento tenendo conto anche del nuovo valore appena calcolato, passando cioè ad un'approssimazione cubica. La determinazione del polinomio dipenderà dal passaggio della parabola

$$\phi_c(\alpha) = a\alpha^3 + b\alpha^2 + \alpha\phi'(0) + \phi(0)$$

per tre punti: $\phi_c(0) = \phi(0)$, $\phi'_c(0) = \phi'(0)$, $\phi_c(\alpha_0) = \phi(\alpha_0)$ e $\phi_c(\alpha_1) = \phi(\alpha_1)$. Si ricava un sistema lineare di due equazioni nelle due incognite a e b la cui soluzione è

$$\begin{bmatrix} a \\ b \end{bmatrix} = \frac{1}{\alpha_0^2 \alpha_1^2 (\alpha_1 - \alpha_0)} \begin{bmatrix} \alpha_0^2 & -\alpha_1^2 \\ -\alpha_0^3 & \alpha_1^3 \end{bmatrix} \begin{bmatrix} \phi(\alpha_1) - \phi(0) - \phi'(0)\alpha_1 \\ \phi(\alpha_0) - \phi(0) - \phi'(0)\alpha_0 \end{bmatrix}.$$

Derivando rispetto ad α ed uguagliando a zero si ricava che un minimo locale di $\phi_c(\alpha)$, a_2 , cade nell'intervallo $[0, \alpha_1]$ ed è dato da

$$\alpha_2 = \frac{-b + \sqrt{b^2 - 3a\phi'(0)}}{3a}.$$

Si verificano le condizioni di sufficiente riduzione per a_2 e se necessario si riapplica il metodo, escludendo il primo punto calcolato, tra quelli presenti, nell'approssimazione cubica. Come si può osservare questa procedura genera una sequenza decrescente di valori di α , tale che ogni α_i non è troppo più piccolo del suo predecessore α_{i-1} . Questa versione è pensata per effettuare il minor numero possibile di calcoli del vettore gradiente $\nabla f(\mathbf{x})$. L'interpolazione cubica è estremamente efficace e garantisce, di solito, un tasso di convergenza quadratico.

6.1.4 Inizializzazione del passo α_0

La scelta del valore al quale inizializzare α_0 dipende dal tipo di tecnica adottata per scegliere la direzione di discesa \mathbf{d} . Nei metodi di Newton e quasi-Newton, si dovrebbe sempre scegliere $\alpha_0 = 1$. Questo permette di usare sempre il passo unitario ogniqualvolta le condizioni di Wolfe

sono soddisfatte, garantendo il mantenimento dell'alto tasso di convergenza di questi metodi. Con le tecniche di gradiente o di gradiente coniugato, è invece importante usare l'informazione corrente per effettuare delle buone inizializzazioni.

Una strategia molto diffusa si basa sull'assunzione che non vari eccessivamente, fra una iterazione e la successiva, il valore della derivata direzionale. Quindi si sceglie α_0 in modo che $\alpha_0 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k = \alpha_{k-1} \nabla f(\mathbf{x}_{k-1})^T \mathbf{d}_{k-1}$:

$$\alpha_0 = \alpha_{k-1} \frac{\nabla f(\mathbf{x}_{k-1})^T \mathbf{d}_{k-1}}{\nabla f(\mathbf{x}_k)^T \mathbf{d}_k}.$$

Un'altra tecnica, approssima la funzione $\phi(\alpha)$ con una forma quadratica, imponendo l'uso dei tre valori: $f(\mathbf{x}_{k-1})$, $f(\mathbf{x}_k)$ e $\nabla f(\mathbf{x}_{k-1})^T \mathbf{d}_{k-1}$ e ricavando il minimo

$$\alpha_0 = \frac{2(f(\mathbf{x}_k) - f(\mathbf{x}_{k-1}))}{\phi'(0)}.$$

In questo caso, se si adotta il criterio di scegliere $\alpha_0 = \min\{1, 1.01\alpha_0\}$ si garantisce sempre la possibilità di scegliere il passo unitario, mantenendo il tasso di convergenza superlineare dei metodi di Newton e quasi-Newton.

6.1.5 Tecniche che non calcolano le derivate

Nei casi in cui non sia possibile, o troppo costoso, il calcolo del vettore gradiente, si possono adottare tecniche che non fanno uso delle informazioni del primo ordine. In questi casi è necessario individuare un intervallo di ricerca $[0, \alpha_0]$ ove cercare α^* , nel quale la funzione sia convessa o pseudo-convessa (convessa almeno nell'intervallo che si sta considerando). In questo caso si possono applicare tecniche che richiedono solo il calcolo della funzione obiettivo $f(\mathbf{x})$. Fra di esse ricordiamo:

- Sezione Aurea
- Fibonacci

e rimandiamo ad altri testi per un approfondimento.

6.2 Scelta della direzione \mathbf{d}_k

In questa sezione affrontiamo il problema di determinare la direzione \mathbf{d} nei metodi basati su ricerca lineare. Iniziamo proponendo una condizione la cui soddisfazione, assieme al rispetto delle condizioni introdotte per α , garantisce la convergenza degli algoritmi di discesa.

Condizione 3 *Una direzione di discesa \mathbf{d} soddisfa la condizione d'angolo se esiste $1 > \varepsilon > 0$ tale che:*

$$\nabla f(\mathbf{x})^T \mathbf{d} \leq -\varepsilon \|\nabla f(\mathbf{x})\| \cdot \|\mathbf{d}\| \quad (6)$$

La condizione impone di non scegliere mai direzioni troppo vicine al piano tangente alle curve di livello passanti per il punto corrente. L'interpretazione geometrica risulta più chiara se scriviamo la seguente espressione

$$\cos \theta_k = \frac{-\nabla f(\mathbf{x})_k^T \mathbf{d}_k}{\|\nabla f(\mathbf{x})_k\| \cdot \|\mathbf{d}_k\|} \quad (7)$$

dove θ_k misura l'angolo fra la direzione \mathbf{d}_k e la direzione dell'antigradiente. In tal modo possiamo riscrivere la condizione d'angolo come

$$\varepsilon \leq \cos \theta_k. \quad (8)$$

Si osservi che per $\mathbf{d} = -\nabla f(\mathbf{x})$ abbiamo $\cos \theta_k = 1$, quindi l'antigradiente certamente soddisfa la relazione (8). Nonostante la sua semplicità, tale ragionevole condizione è estremamente potente in quanto il suo soddisfacimento, assieme al soddisfacimento delle condizioni di Wolfe per il valore del passo α , è sufficiente a garantire la convergenza *globale* di qualsiasi algoritmo di discesa.

Teorema 8 *Sia data una funzione f limitata inferiormente in \mathbb{R}^n , di classe C^1 in un insieme aperto S contenente l'insieme di livello $L_f(\mathbf{x}_0) = \{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$, dove \mathbf{x}_0 è il punto di partenza di un algoritmo iterativo di tipo (1). Si supponga che ad ogni iterazione k la direzione \mathbf{d}_k soddisfi la condizione d'angolo ed il passo α_k soddisfi le condizioni di Wolfe, per opportuni valori dei parametri c_1, c_2 ed ε . Si supponga inoltre che il gradiente $\nabla f(\mathbf{x})$ sia Lipschitz continuo in S , cioè che esista una costante $L > 0$ tale che il gradiente $\nabla f(\mathbf{x})$ soddisfi la seguente relazione*

$$\|\nabla f(\mathbf{x}) - \nabla f(\tilde{\mathbf{x}})\| \leq L\|\mathbf{x} - \tilde{\mathbf{x}}\|, \text{ per tutte le coppie di punti } \mathbf{x}, \tilde{\mathbf{x}} \in S.$$

Allora la successione dei punti $\{\mathbf{x}_k\}$ è tale che il gradiente di $f(\mathbf{x})$ si annulla in un numero finito di passi o converge a zero asintoticamente:

$$\begin{aligned} &\text{esiste } k \text{ con } \nabla f(\mathbf{x}_k) = \mathbf{0} \\ &\text{oppure la successione } \{\nabla f(\mathbf{x}_k)\} \rightarrow \mathbf{0} \end{aligned}$$

Dimostrazione Dalla condizione di curvatura (3)

$$\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k \geq c_2 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k$$

ricaviamo

$$(\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) - \nabla f(\mathbf{x}_k))^T \mathbf{d}_k \geq (c_2 - 1) \nabla f(\mathbf{x}_k)^T \mathbf{d}_k$$

mentre la condizione di Lipschitz e la disuguaglianza di Cauchy-Schwarz² implicano che

$$(\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) - \nabla f(\mathbf{x}_k))^T \mathbf{d}_k \leq \alpha_k L \|\mathbf{d}_k\|^2.$$

Combinando queste due relazioni ricaviamo

$$\alpha_k \geq \frac{c_2 - 1}{L} \frac{\nabla f(\mathbf{x}_k)^T \mathbf{d}_k}{\|\mathbf{d}_k\|^2}.$$

Se sostituiamo il valore di α_k così ottenuto nella condizione di Armijo (2)

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + c_1 \alpha_k \nabla f(\mathbf{x}_k)^T \mathbf{d}_k,$$

²La disuguaglianza di Cauchy-Schwarz afferma che se \mathbf{x} e \mathbf{y} sono vettori di \mathbb{R}^n allora vale la relazione $(\mathbf{x}^T \mathbf{y})^2 \leq (\mathbf{x}^T \mathbf{x})(\mathbf{y}^T \mathbf{y})$

ricaviamo

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) - c_1 \frac{1 - c_2 (\nabla f(\mathbf{x}_k)^T \mathbf{d}_k)^2}{L \|\mathbf{d}_k\|^2}.$$

Usando la definizione (7) possiamo riscrivere questa relazione come

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) - c \cos^2 \theta_k \|\nabla f(\mathbf{x}_k)\|^2,$$

dove $c = c_1(1 - c_2)/L$. Sommando questa espressione su tutti gli indici minori o uguali a k ricaviamo

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_0) - c \sum_{j=0}^k \cos^2 \theta_j \|\nabla f(\mathbf{x}_j)\|^2.$$

Poiché la funzione $f(\mathbf{x})$ è inferiormente limitata, ne segue che il valore $f(\mathbf{x}_0) - f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)$ è limitato da una qualche costante positiva per tutti i k . Questo significa che vale la seguente relazione, detta *condizione di Zoutendijk*

$$\sum_{j=0}^{\infty} \cos^2 \theta_j \|\nabla f(\mathbf{x}_j)\|^2 < \infty. \quad (9)$$

La condizione di Zoutendijk implica che

$$\cos^2 \theta_j \|\nabla f(\mathbf{x}_j)\|^2 \rightarrow 0$$

Ora, poiché per ipotesi vale la condizione d'angolo

$$0 < \varepsilon \leq \cos \theta_k \quad \text{per tutti i valori di } k$$

segue immediatamente che

$$\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_j)\| = 0.$$

□

Si osservi che la condizione sulle linee di livello equivale a chiedere che la funzione non sia inferiormente illimitata ed è quindi ragionevole poiché altrimenti il problema non sarebbe ben posto. La condizione che il gradiente sia Lipschitz continuo è spesso richiesta per poter dimostrare la convergenza locale degli algoritmi ed è spesso soddisfatta in pratica. Vediamo ora, nel dettaglio, alcuni modi differenti di scegliere la direzione \mathbf{d} .

6.2.1 Metodo del gradiente

Nel metodo del gradiente la direzione, \mathbf{d} , coincide con l'antigradiente: $\mathbf{d} = -\nabla f(\mathbf{x})$. Tale direzione massimizza il decremento della funzione obiettivo. Nei metodi iterativi si ha quindi

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k). \quad (10)$$

Per ricavare una buona stima del valore del passo ottimo α_k possiamo approssimare $f(\mathbf{x})$ con una funzione quadratica, ottenuta fermandoci al secondo ordine nello sviluppo in serie di Taylor di f .

$$f(\mathbf{x}_k + \alpha \mathbf{d}_k) \approx \phi(\alpha) = f(\mathbf{x}_k) + \alpha \mathbf{d}_k^T \nabla f(\mathbf{x}_k) + \frac{1}{2} \alpha^2 \mathbf{d}_k^T H(\mathbf{x}_k) \mathbf{d}_k$$

Calcoliamo ora la derivata di $\phi(\alpha)$ rispetto ad α come fatto nella Sezione 6.1.2, ottenendo

$$\phi'(\alpha) = \mathbf{d}_k^T \nabla f(\mathbf{x}_k) + \alpha \mathbf{d}_k^T H(\mathbf{x}_k) \mathbf{d}_k$$

Se vale la condizione $\mathbf{d}_k^T H(\mathbf{x}_k) \mathbf{d}_k \neq 0$, allora possiamo imporre $\phi'(\alpha) = 0$ e ricavare

$$\alpha = -\frac{\mathbf{d}_k^T \nabla f(\mathbf{x}_k)}{\mathbf{d}_k^T H(\mathbf{x}_k) \mathbf{d}_k}$$

Metodo del gradiente;

```

{
  Scegli  $\mathbf{x}_0 \in \mathbb{R}^n$ ;  $k := 0$ ;
  While  $\nabla f(\mathbf{x}_k) \neq \emptyset$ ;
  {
     $\mathbf{d}_k := -\nabla f(\mathbf{x}_k)$  /* direzione di discesa */
    calcola  $\alpha_k \in \mathbb{R}$ ; /* passo lungo  $\mathbf{d}_k$  */
     $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{d}_k$ ;
     $k := k + 1$ ;
  }
}
```

Nell'ipotesi di ricerca lineare esatta ed approssimazione mediante una funzione quadratica, il passo iterativo (10) diviene:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{\nabla f(\mathbf{x}_k)^T \nabla f(\mathbf{x}_k)}{\nabla f(\mathbf{x}_k)^T H(\mathbf{x}_k) \nabla f(\mathbf{x}_k)} \nabla f(\mathbf{x}_k) \quad (11)$$

La direzione dell'antigradiente sembra la scelta più ovvia in un metodo di discesa: seguire ad ogni iterazione la direzione ed il passo di massimo decremento del valore della funzione obiettivo. Ma fino a che punto questa scelta è vincente? Per scoprirlo occorre analizzare il tasso di convergenza del metodo, e per farlo occorre scegliere una classe di problemi per la quale tale analisi sia praticabile.

6.2.2 Analisi del metodo del gradiente

Analizziamo le prestazioni del metodo del gradiente quando viene applicato a *problemi quadratici* con matrice Q definita positiva

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x}. \quad (12)$$

In questo caso il passo iterativo (11) diviene

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{(Q \mathbf{x}_k - \mathbf{b})^T (Q \mathbf{x}_k - \mathbf{b})}{(Q \mathbf{x}_k - \mathbf{b})^T Q (Q \mathbf{x}_k - \mathbf{b})} (Q \mathbf{x}_k - \mathbf{b})$$

Per questo tipo di problemi siamo in grado di ricavare analiticamente il valore della soluzione ottima, $\mathbf{x}^* = Q^{-1}\mathbf{b}$, e di conseguenza possiamo utilizzare questa informazione (di solito non disponibile, altrimenti non avremmo bisogno degli algoritmi iterativi...) per calcolare la forma esatta della riduzione dell'errore di approssimazione tra iterazioni successive. Per far ciò useremo la norma pesata

$$\|\mathbf{x}\|_Q^2 = \frac{1}{2}\mathbf{x}^T Q \mathbf{x}$$

e quindi la riduzione dell'errore calcolata sarà:

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q}{\|\mathbf{x}_k - \mathbf{x}^*\|_Q}$$

Come possiamo interpretare tale norma? Per rispondere dobbiamo effettuare alcune trasformazioni. Innanzitutto riscriviamo $f(\mathbf{x}^*)$.

$$f(\mathbf{x}^*) = \frac{1}{2}\mathbf{x}^{*T} Q \mathbf{x}^* - \mathbf{b}^T \mathbf{x}^* = \frac{1}{2}\mathbf{x}^{*T} Q \mathbf{x}^* - (Q \mathbf{x}^*)^T \mathbf{x}^* = -\frac{1}{2}\mathbf{x}^{*T} Q \mathbf{x}^*.$$

Ora misuriamo il quadrato della distanza fra il punto corrente \mathbf{x}_k ed il punto di ottimo con la nuova metrica:

$$\|\mathbf{x}_k - \mathbf{x}^*\|_Q^2 = \frac{1}{2}(\mathbf{x}_k - \mathbf{x}^*)^T Q (\mathbf{x}_k - \mathbf{x}^*) = \frac{1}{2}\mathbf{x}_k^T Q \mathbf{x}_k + \frac{1}{2}\mathbf{x}^{*T} Q \mathbf{x}^* - \mathbf{x}_k^T Q \mathbf{x}^* = f(\mathbf{x}_k) - f(\mathbf{x}^*). \quad (13)$$

dove, di nuovo, abbiamo fatto uso della relazione $Q \mathbf{x}^* = \mathbf{b}$ e del fatto che $\mathbf{x}_k^T \mathbf{b} = \mathbf{b}^T \mathbf{x}_k$. Questa norma misura quindi la distanza fra il valore corrente della funzione obiettivo ed il valore ottimo.

Con un po' di passaggi è possibile ricavare il seguente risultato intermedio.

Lemma 1 *Se l'algoritmo del gradiente viene applicato ad un problema di programmazione quadratica convessa la forma esatta della riduzione dell'errore di approssimazione tra iterazioni successive misurata con la norma pesata $\|\cdot\|_Q$ è*

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q}{\|\mathbf{x}_k - \mathbf{x}^*\|_Q} = \left\{ 1 - \frac{((Q \mathbf{x}_k - \mathbf{b})^T (Q \mathbf{x}_k - \mathbf{b}))^2}{((Q \mathbf{x}_k - \mathbf{b})^T Q (Q \mathbf{x}_k - \mathbf{b}))((Q \mathbf{x}_k - \mathbf{b})^T Q^{-1} (Q \mathbf{x}_k - \mathbf{b}))} \right\}^{\frac{1}{2}}$$

Dimostrazione Vedi Appendice \square

Abbiamo quindi una misura del decremento di f ad ogni iterazione. Il termine fra parentesi è però decisamente difficile da interpretare. Per nostra fortuna tale quantità si dimostra dominabile tramite una funzione degli autovalori di Q come ci dimostra il seguente teorema dovuto a Kantorovic.

Teorema 9 *Data una matrice Q definita positiva vale la seguente relazione*

$$\frac{(\mathbf{x}^T \mathbf{x})^2}{(\mathbf{x}^T Q \mathbf{x})(\mathbf{x}^T Q^{-1} \mathbf{x})} \geq \frac{4\lambda_m \lambda_M}{(\lambda_m + \lambda_M)^2}$$

valida per ogni $\mathbf{x} \in \mathbb{R}^n$ e dove λ_M e $\lambda_m > 0$ sono l'autovalore massimo e minimo, rispettivamente, di Q .

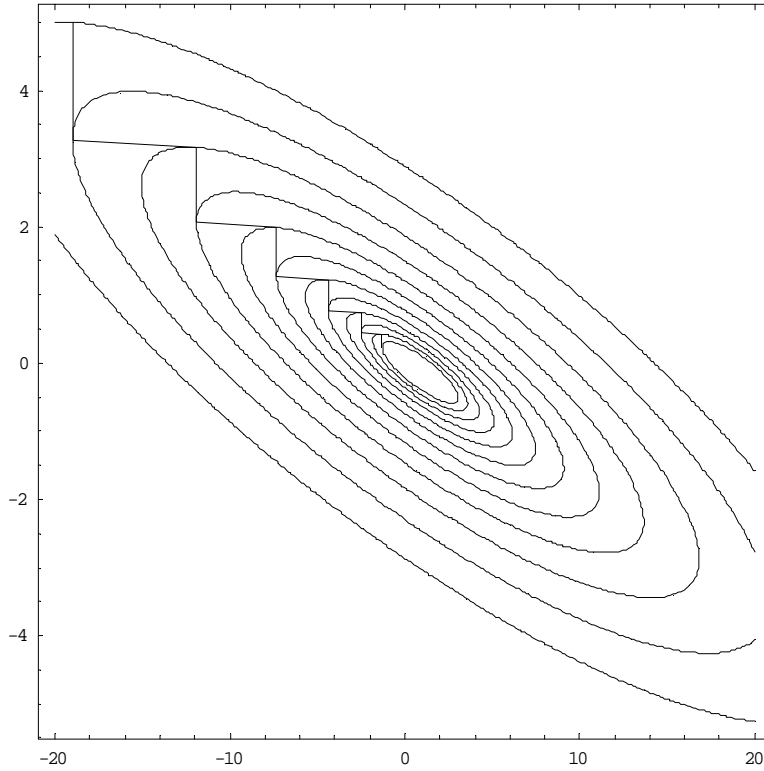


Figura 9: Andamento zigzagante del metodo del gradiente

Dimostrazione Vedi Appendice \square

Combinando i due precedenti risultati possiamo ricavare il seguente fondamentale risultato.

Teorema 10 *Il metodo del gradiente con scelta del passo ottimale, quando viene applicato al problema quadratico $\min \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x}$ dove Q è una matrice simmetrica definita positiva, produce una successione $\{\mathbf{x}_k\}$ tale che:*

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q}{\|\mathbf{x}_k - \mathbf{x}^*\|_Q} \leq \left(\frac{\lambda_M - \lambda_m}{\lambda_M + \lambda_m} \right) \quad (14)$$

dove λ_M e λ_m sono l'autovalore massimo e minimo, rispettivamente di Q .

Dimostrazione Denotando $\nabla f(\mathbf{x}_k) = (Q\mathbf{x}_k - \mathbf{b})$ con \mathbf{g}_k possiamo scrivere il risultato del Lemma 1 come

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q}{\|\mathbf{x}_k - \mathbf{x}^*\|_Q} = \left\{ 1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T Q \mathbf{g}_k)(\mathbf{g}_k^T Q^{-1} \mathbf{g}_k)} \right\}^{\frac{1}{2}}$$

e dal Teorema di Kantorovic ricaviamo

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q}{\|\mathbf{x}_k - \mathbf{x}^*\|_Q} \leq \left\{ 1 - \frac{4\lambda_m \lambda_M}{(\lambda_m + \lambda_M)^2} \right\}^{\frac{1}{2}} = \left(\frac{\lambda_M - \lambda_m}{\lambda_M + \lambda_m} \right)$$

\square

La disuguaglianza (14) mostra che l'algoritmo del gradiente converge con tasso di convergenza lineare.

Esempio La Figura 9 illustra l'andamento delle prime 12 iterazioni dell'algoritmo del gradiente applicato alla funzione quadratica (12) con $Q = \begin{pmatrix} 3 & 12 \\ 12 & 70 \end{pmatrix}$, definita positiva, con $\mathbf{b} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, a partire dal punto iniziale $\mathbf{x}_0 = \begin{pmatrix} -19 \\ 5 \end{pmatrix}$. Per via analitica si ricava $\mathbf{x}^* = \begin{pmatrix} 29/33 \\ -3/22 \end{pmatrix}$. In questo caso $\lambda_M/\lambda_m = 78.7297$.

In generale, al crescere del rapporto λ_M/λ_m (detto *numero condizione* di Q) le prestazioni dell'algoritmo degradano e l'andamento zigzagante presentato in Figura 9 si accentua. Per converso, come caso speciale, l'algoritmo del gradiente converge in un solo passo quando $\lambda_m = \lambda_M$, cioè quando tutti gli autovalori di Q sono uguali, e cioè quando Q è un multiplo della matrice identità. In questo caso la direzione punta all'ottimo e le curve di livello sono cerchi concentrici. Anche se la disuguaglianza (14) dà una stima per eccesso, essa risulta una accurata indicazione del comportamento dell'algoritmo per valori di $n > 2$.

Quanto detto vale quando l'algoritmo del gradiente viene applicato a funzioni quadratiche, ma cosa succede se lo applichiamo a funzioni qualsiasi? In pratica, il tasso di convergenza del metodo del gradiente è essenzialmente lo stesso per generiche funzioni obiettivo non lineari. Infatti, assumendo che ad ogni iterazione la lunghezza del passo α sia effettuata in modo ottimo, si può dimostrare il seguente teorema.

Teorema 11 *Si supponga $f : \mathbb{R}^n \rightarrow \mathbb{R}$ di classe C^2 e che le soluzioni generate dall'algoritmo del gradiente, con scelta del passo α ottima, convergano al punto \mathbf{x}^* , dove la matrice hessiana $H(\mathbf{x}^*)$ è definita positiva. Sia*

$$r \in \left(\frac{\lambda_M - \lambda_m}{\lambda_M + \lambda_m}, 1 \right)$$

dove λ_M e λ_m sono l'autovalore massimo e minimo, rispettivamente di $H(\mathbf{x}^)$. Allora esiste un \bar{k} tale che per tutti $i > \bar{k}$ vale*

$$\frac{f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)}{f(\mathbf{x}_k) - f(\mathbf{x}^*)} \leq r^2.$$

In generale, non possiamo aspettarci che il tasso di convergenza migliori se si adotta una scelta non ottima del passo α . Il Teorema 11 ci mostra che il metodo del gradiente può essere inaccettabilmente lento a convergere, anche quando la matrice hessiana è ragionevolmente ben condizionata (cioè con un valore del numero condizione non troppo grande).

Osserviamo ora, che generalmente il metodo del gradiente è globalmente convergente: è sufficiente che il passo venga scelto in modo da soddisfare le condizioni di Wolfe. Tuttavia la convergenza è estremamente lenta.

Quali sono le ragioni per delle prestazioni così deludenti?

Analizziamo innanzitutto il costo computazionale del metodo. Nel calcolare la direzione di

discesa esso richiede il calcolo del gradiente di $f(\mathbf{x})$, che ha un costo computazionale dell'ordine di $O(n^2)$. Nella scelta del passo α , se utilizziamo il metodo di Armijo ci limitiamo ad aggiungere un numero più o meno costante di valutazioni di $f(\mathbf{x})$ e di $\nabla f(\mathbf{x})$, mentre se utilizziamo l'approssimazione quadratica il costo aumenta sensibilmente: è necessario calcolare l'hessiana $H(\mathbf{x})$ con un costo computazionale dell'ordine di $O(n^3)$. In pratica, nessuno utilizza il metodo del gradiente con scelta del passo ottima, proprio per l'elevato costo computazionale. D'altronde è solo quando scegliamo il passo in modo ottimo che garantiamo la convergenza lineare! Quindi abbiamo convergenza lineare al costo dell'ordine di $O(n^3)$ per iterazione. La spiegazione per le prestazioni deludenti, non sta quindi nel basso costo computazionale del metodo, ma nell'uso non efficace delle informazioni elaborate: la matrice Hessiana, quando usata, non contribuisce alla scelta della direzione, ma solo del passo, inoltre il metodo non tiene alcuna traccia nelle scelte all'iterazione corrente della storia passata.

Diviene naturale chiedersi se non si possa far di meglio, nell'ipotesi di poter spendere comunque al più $O(n^3)$ passi di calcolo ad ogni iterazione. Il prossimo metodo risponde positivamente a questa domanda.

6.2.3 Metodo di Newton

Analogamente al metodo del gradiente, anche il metodo di Newton si basa sul concetto di minimizzare un'approssimazione quadratica della funzione f . In questo caso però, si desidera ricavare da tale approssimazione contemporaneamente il valore di α_k e quello di \mathbf{d}_k .

Sia $f(\mathbf{x})$ una funzione con Hessiana continua. Per valori sufficientemente piccoli della norma del vettore incremento $\mathbf{h}_k = \alpha_k \mathbf{d}_k$ è possibile scrivere:

$$f(\mathbf{x}_k + \mathbf{h}_k) \approx q(\mathbf{h}_k) = f(\mathbf{x}_k) + \mathbf{h}_k^T \nabla f(\mathbf{x}_k) + \frac{1}{2} \mathbf{h}_k^T H(\mathbf{x}_k) \mathbf{h}_k$$

Annullando il gradiente di $q(\mathbf{h}_k)$ si ricava

$$\nabla q(\mathbf{h}_k) = \nabla f(\mathbf{x}_k) + H(\mathbf{x}_k) \mathbf{h}_k = \mathbf{0},$$

da cui, se la matrice hessiana è non singolare, possiamo ottenere

$$\mathbf{h}_k = -H(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k).$$

Quindi quando la matrice hessiana nel punto \mathbf{x}_k è invertibile, ad esempio quando essa è definita positiva in quel punto, il vettore \mathbf{h}_k è ottimo unico.

Il metodo di Newton *puro* ha quindi la forma

$$\mathbf{x}_{k+1} = \mathbf{x}_k - H(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k) \quad (15)$$

Come si comporta il metodo di Newton? Iniziamo con l'osservare che quando esso viene applicato a problemi quadratici

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

e la matrice Q è definita positiva, il metodo trova l'ottimo in un passo. Infatti si ricava

$$\mathbf{x}_1 = \mathbf{x}_0 - H(\mathbf{x}_0)^{-1} \nabla f(\mathbf{x}_0) = \mathbf{x}_0 - Q^{-1}(Q\mathbf{x}_0 - \mathbf{b}) = \mathbf{x}_0 - \mathbf{x}_0 + Q^{-1}\mathbf{b} = Q^{-1}\mathbf{b}.$$

Se invece la matrice hessiana non è definita positiva, il metodo *non converge*.

<p>Metodo di Newton Puro;</p> <pre> { Scegli $\mathbf{x}_0 \in \mathbb{R}^n$; $k := 0$; While $\nabla f(\mathbf{x}_k) \neq \emptyset$; { $\mathbf{x}_{k+1} = \mathbf{x}_k - H(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$; $k := k + 1$; } }</pre>

Quanto osservato per i problemi quadratici ha importanti conseguenze:

le proprietà di convergenza del metodo sono fortemente legate alla definita positività della matrice hessiana. Quando la funzione $f(\mathbf{x})$ è di tipo generale la qualità della direzione dipende essenzialmente dal fatto che la matrice hessiana calcolata nel punto corrente \mathbf{x}_k sia definita positiva o meno.

Infatti, quando la matrice non è definita positiva in un punto, la direzione in quel punto può non essere una direzione di discesa.

Nei casi in cui invece la matrice hessiana è definita positiva, ed in particolare è tale in un intorno di un punto stazionario, vale il seguente fondamentale risultato:

Teorema 12 *Data una funzione f di classe C^2 , per la quale la matrice hessiana è definita positiva nell'intorno $I(\mathbf{x}^*, \varepsilon)$ di un punto stazionario \mathbf{x}^* e soddisfa la seguente condizione:*

$$\text{esiste } L > 0 : \|H(\mathbf{y}) - H(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\| \text{ per ogni } \mathbf{x}, \mathbf{y} \in I(\mathbf{x}^*, \varepsilon)$$

allora la successione $\{\mathbf{x}_k\}$ che si ottiene applicando il metodo di Newton puro è tale che, se il punto iniziale $\mathbf{x}_0 \in I(\mathbf{x}^, \varepsilon)$, la successione converge a \mathbf{x}^* con tasso di convergenza quadratico.*

Il miglioramento rispetto al metodo del gradiente è notevole in fatto di rapidità di convergenza. Ma tale miglioramento viene ad un costo, non solo computazionale: per convergere il metodo richiede di partire vicino ad un punto stazionario. Abbiamo quindi solo una proprietà di *convergenza locale* e non di convergenza globale. Inoltre, il teorema appena enunciato richiede qualcosa in più del fatto che la matrice hessiana sia definita positiva ad ogni iterazione. Infatti, la definita positività della matrice hessiana non basta a far sì che il metodo converga, poiché le direzioni di Newton potrebbero non soddisfare alla condizione d'angolo, potrebbero cioè essere quasi ortogonali al gradiente.

Se invece valgono opportune condizioni sugli autovalori dell'inversa della matrice hessiana, condizioni che ne implicano la definita positività, allora la direzione di Newton soddisfa la condizione d'angolo.

Teorema 13 *Se $0 < m < \lambda_m H(\mathbf{x}_k)^{-1} \leq \lambda_M H(\mathbf{x}_k)^{-1} < M$ per ogni k (proprietà che implica che sia H che H^{-1} siano definite positive), la direzione di Newton soddisfa la condizione d'angolo con $\varepsilon = \frac{m}{M}$:*

$$\nabla f(\mathbf{x}_k)^T \mathbf{d}_k < -\frac{m}{M} \|\nabla f(\mathbf{x}_k)\| \cdot \|\mathbf{d}_k\|$$

Per riottenere una proprietà di *convergenza globale* è quindi necessario modificare il metodo.

Una prima modifica consiste nell'adattare dinamicamente la lunghezza del passo, che nel metodo puro è tenuto costantemente al valore $\alpha = 1$. Si tratta cioè di applicare ad ogni iterazione l'algoritmo di Armijo a partire da $\alpha_0 = 1$.

Una seconda modifica riguarda la scelta della direzione. Se la matrice hessiana non è definita positiva, si può perturbarla, perturbando di conseguenza la direzione \mathbf{d}_k . Delle diverse tecniche proposte in letteratura, nessuna è esente da controindicazioni. Un metodo semplice ed efficace consiste nel passare ad una matrice $(H(\mathbf{x}_k) + \gamma I)$ che, per tutti i valori di γ superiori ad una certa soglia (dipendente da $(H(\mathbf{x}_k))$), è definita positiva. Si risolve il sistema

$$(H(\mathbf{x}_k) + \gamma I)\mathbf{h}_k = -\nabla f(\mathbf{x}_k)$$

osservando che al crescere di γ , $\mathbf{h}_k \rightarrow -\nabla f(\mathbf{x}_k)$.

In alternativa, se non si desidera modificare la matrice hessiana, si può alternare l'utilizzo della direzione di Newton, con la direzione dell'antigradiente, scegliendo di volta in volta quale direzione adottare come esito di test sulla condizione d'angolo.

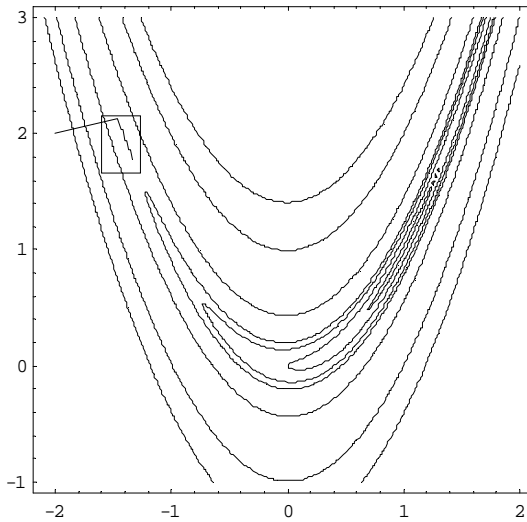
Dunque, se ad ogni iterazione l'Hessiana è definita positiva, e si garantisce il soddisfacimento della condizione d'angolo, il metodo di discesa che utilizza la direzione di Newton converge globalmente ad un punto stazionario. Peraltro, è possibile dimostrare che, da una certa iterazione k in poi, il passo $\alpha_k = 1$ soddisfa le condizioni del metodo di Armijo, e dunque si può di fatto usare il metodo di Newton puro, e la convergenza quadratica è ancora assicurata.

Riassumendo, se la matrice hessiana non è sempre definita positiva il metodo di discesa può risultare inapplicabile per due motivi: o perchè in alcuni punti la matrice hessiana risulta singolare, oppure perchè la direzione di Newton non è una direzione di discesa. Inoltre, può anche accadere che, pur risultando una direzione di discesa, la direzione di Newton e il gradiente siano talmente prossimi ad essere ortogonali che non conviene seguire la direzione di Newton. A questi inconvenienti è possibile ovviare in modo molto semplice, utilizzando la direzione dell'antigradiente quando non sia possibile o conveniente seguire la direzione di Newton.

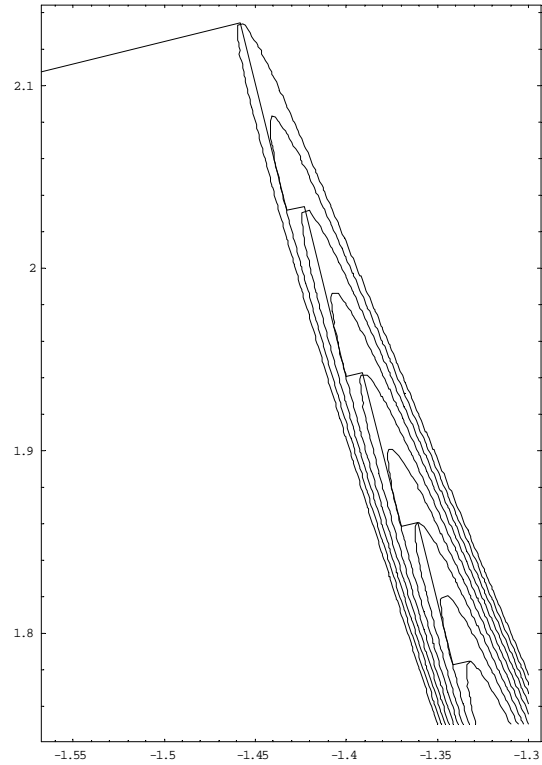

```

Metodo di Newton Modificato;
{
  Scegli  $\mathbf{x}_0 \in \mathbb{R}^n$ ;  $k := 0$ ;
  While  $\nabla f(\mathbf{x}_k) \neq \emptyset$ ;
  {
    if  $H(\mathbf{x}_k)$  è singolare then  $\mathbf{d}_k := -\nabla f(\mathbf{x}_k)$ ;
    else
    {
       $\mathbf{s} := -H(\mathbf{x}_k)^{-1}\nabla f(\mathbf{x}_k)$ ;
      if  $|\nabla f(\mathbf{x}_k)^T \mathbf{s}| < \varepsilon \|\nabla f(\mathbf{x}_k)\| \cdot \|\mathbf{s}\|$  then  $\mathbf{d}_k := -\nabla f(\mathbf{x}_k)$ ;
      else
        if  $\mathbf{s}$  è direzione di discesa then  $\mathbf{d}_k := \mathbf{s}$ ;
        else  $\mathbf{d}_k := -\mathbf{s}$ ;
      }
    }
    calcola  $\alpha_k \in \mathbb{R}$ ; /*passo lungo  $\mathbf{d}_k$ */
     $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ ;
     $k := k + 1$ ;
  }
}

```



(a) Curve di livello



(b) Dieci iterazioni di gradiente

Figura 10: Il metodo del gradiente applicato alla funzione di Rosenbrock.

Il metodo di Newton modificato in modo da considerare tutte queste possibilità è mostrato in figura, dove $\varepsilon > 0$ è una quantità sufficientemente piccola. In sostanza, come si può vedere il metodo sceglie come direzione di discesa quella di Newton, la sua opposta o l'antigradiente, e successivamente effettua una ricerca lineare con metodi standard, in modo da garantire la convergenza globale. L'algoritmo così ottenuto soddisfa la condizione d'angolo e dunque le proprietà di convergenza globale dell'algoritmo sono conservate.

Osservazione. Per trovare \mathbf{h} non si inverte $H(\mathbf{x})$, ma si risolve il sistema $H(\mathbf{x})\mathbf{h} = -\nabla f(\mathbf{x})$ mediante fattorizzazione di $H(\mathbf{x})$ come matrice in forma di prodotto LDL^T , con L diagonale e triangolare inferiore ed L^T triangolare superiore. Dalla fattorizzazione si può inoltre ricavare se la matrice hessiana è definita positiva.

Un ulteriore aspetto problematico del metodo di Newton, modificato o meno, è il richiedere la valutazione delle derivate seconde di $f(\mathbf{x})$ per il calcolo di $H(\mathbf{x}_k)$. In alternativa, si può ricorrere alle *differenze finite*. Si tratta di approssimare $H(\mathbf{x}_k)$ mediante $\frac{1}{2}(\overline{H}_k + \overline{H}_k^T)$ dove la colonna i -esima di \overline{H}_k è $\frac{1}{h}[\nabla f(\mathbf{x}_k + \mathbf{h}_i) - \nabla f(\mathbf{x}_k)]$, con $\mathbf{h}_i = (\dots, 0, h_i, 0, \dots)^T$, ed h_i è un opportuno incremento dell' i -esima variabile. Purtroppo la matrice risultante può non essere definita positiva (mentre questo serve), si devono calcolare n gradienti e risolvere un sistema lineare ad ogni passo, con un costo quindi dell'ordine di $O(n^3)$ per iterazione.

In sintesi il costo computazionale del metodo di Newton è $O(n^3)$ ad ogni iterazione ed è quindi proponibile solo per sistemi di dimensioni non superiori a qualche decina di migliaia di variabili.

6.2.4 Confronto fra i metodi del gradiente e di Newton

Esemplifichiamo le differenti prestazioni del metodo del gradiente e di quello di Newton su una funzione non convessa nota in letteratura come funzione di Rosenbrock

$$f(x, y) = 100(y - x^2)^2 + (1 - x)^2.$$

Tale funzione ammette un minimo assoluto in $(1, 1)$ di valore 0. Se applichiamo alla funzione di Rosenbrock il metodo del gradiente otteniamo l'andamento riportato in Figura 10, dove nella figura di destra viene riportato ingrandito, il riquadro, evidenziato a sinistra, relativo alle prime dieci iterazioni del metodo, a partire dal punto $(-2, 2)$. Tale andamento prosegue per centinaia di iterazioni prima di raggiungere il punto di minimo, con un percorso che segue l'andamento a *banana* evidenziato dalle curve di livello.

Se alla stessa funzione, e a partire dallo stesso punto, applichiamo l'algoritmo di Newton *puro*, con $\alpha = 1$, otteniamo l'andamento riportato in Figura 11, dove il punto di ottimo viene raggiunto alla quinta iterazione. Si evidenzia come alla seconda iterazione lo spostamento non porti ad un punto migliore rispetto a quello di partenza, con una direzione che è praticamente perpendicolare al gradiente ed un passo lunghissimo. Ciò nonostante l'algoritmo converge molto più rapidamente.

6.2.5 Metodi quasi Newton

L'esempio precedente ci invoglia a chederci se non si possano raggiungere le prestazioni del metodo di Newton ad un costo computazionale inferiore all' $O(n^3)$ richiesto ad ogni iterazione.

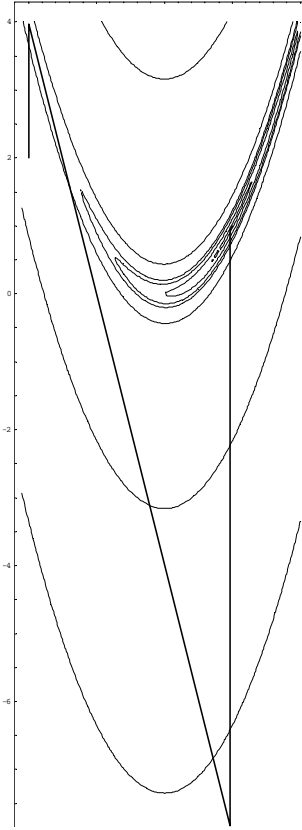


Figura 11: Il metodo di Newton applicato alla funzione di Rosenbrock.

Tale costo vale sia per il metodo puro che per le varie varianti di Newton modificato. In particolare, l'algoritmo del gradiente ha un costo computazionale dell'ordine di $O(n^2)$ ad iterazione, se si accetta di individuare la lunghezza del passo con metodi di ricerca lineare non esatti, come, ad esempio, quello di Armijo. In mancanza di alternative valide, il metodo del gradiente risulterebbe l'unico praticabile in presenza di funzioni di grandissime dimensioni. Per nostra fortuna, *esistono alternative valide*.

I metodi quasi Newton sono stati introdotti per ovviare all'eccessivo carico computazionale dei metodi di Newton puro e modificato, senza perdere le caratteristiche di convergenza dell'algoritmo. Lo scopo che questi metodi si prefiggono è di ridurre il costo computazionale ad ogni iterazione mediante l'approssimazione dell'inversa della matrice hessiana, $H(\mathbf{x}_k)^{-1}$ con una matrice definita positiva G_k , tale che la matrice al passo successivo, G_{k+1} , possa essere ricavata da G_k in $O(n^2)$ passi.

L'idea è quella

- (a) di ricavare la direzione come risultato dell'applicazione di un operatore lineare al vettore gradiente calcolato nel punto corrente:

$$\mathbf{d}_k := -G_k \nabla f(\mathbf{x}_k)$$

- (b) di calcolare il passo α_k mediante tecniche di ricerca monodimensionale,

(c) di calcolare il nuovo punto nel modo consueto

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

Per semplicità si inizializza G_0 con la matrice identità. Mentre per l'aggiornamento si opera nel seguente modo.

Come per il metodo di Newton puro, approssimiamo la funzione $f(\mathbf{x}_k + \mathbf{h}_k)$ con una espressione quadratica $q(\mathbf{h}_k)$, dove $\mathbf{h}_k = \alpha_k \mathbf{d}_k$:

$$f(\mathbf{x}_k + \mathbf{h}_k) \approx q(\mathbf{h}_k) = f(\mathbf{x}_k) + \mathbf{h}_k^T \nabla f(\mathbf{x}_k) + \frac{1}{2} \mathbf{h}_k^T H(\mathbf{x}_k) \mathbf{h}_k$$

da cui possiamo scrivere

$$\nabla q(\mathbf{h}_k) = \nabla f(\mathbf{x}_k) + H(\mathbf{x}_k) \mathbf{h}_k \approx \nabla f(\mathbf{x}_k + \mathbf{h}_k) = \nabla f(\mathbf{x}_{k+1}).$$

Se scriviamo $\mathbf{h}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ e definiamo $\mathbf{p}_k := \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$ possiamo utilizzare l'approssimazione quadratica per ricavare $\mathbf{p}_k \approx H(\mathbf{x}_k) \mathbf{h}_k$, ottenendo la seguente

Relazione guida

$$H(\mathbf{x}_k)^{-1} \mathbf{p}_k \approx \mathbf{h}_k$$

Quindi, dopo aver inizializzato la matrice G_0 con la matrice identità, imponiamo ad ogni iterazione k , che la matrice G_{k+1} soddisfi la relazione guida come uguaglianza:

$$G_{k+1} \mathbf{p}_k = \mathbf{h}_k \tag{16}$$

Il modo in cui viene soddisfatta l'equazione (16) differenzia le diverse varianti dei metodi quasi Newton.

```

Metodo quasi-Newton;
{
  Scegli  $\mathbf{x}_0 \in \mathbb{R}^n$ ;  $k := 0$ ;
   $G_0 := I$ ;
  While  $\nabla f(\mathbf{x}_k) \neq \emptyset$ ;
  {
     $\mathbf{d}_k := -G_k \nabla f(\mathbf{x}_k)$  /* direzione di discesa */
    calcola  $\alpha_k \in \mathbb{R}$ ; /* passo lungo  $\mathbf{d}_k$  */
     $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{d}_k$ ;
    calcola  $G_{k+1}$  a partire da  $G_k$ ;
     $k := k + 1$ ;
  }
}

```

6.2.6 Formula di aggiornamento di rango uno

Con la formula di aggiornamento di *rango uno* si impone

$$G_{k+1} \mathbf{p}_k = G_k \mathbf{p}_k + c \mathbf{u} \mathbf{u}^T \mathbf{p}_k = \mathbf{h}_k$$

cioè $G_{k+1} = G_k + c\mathbf{u}\mathbf{u}^T$.

Se scegliamo il coefficiente c in modo che valga la relazione $c\mathbf{u}^T\mathbf{p}_k = 1$, cioè

$$c := 1/(\mathbf{u}^T\mathbf{p}_k),$$

e imponiamo $\mathbf{u} = \mathbf{h}_k - G_k\mathbf{p}_k$, la formula di aggiornamento di G_{k+1} a partire da G_k diviene

$$G_{k+1} = G_k + \frac{(\mathbf{h}_k - G_k\mathbf{p}_k)(\mathbf{h}_k - G_k\mathbf{p}_k)^T}{(\mathbf{h}_k - G_k\mathbf{p}_k)^T\mathbf{p}_k}.$$

Come è facile verificare la matrice G_{k+1} può essere ricavata dalla matrice G_k con un costo computazionale dell'ordine di $O(n^2)$. Abbiamo quindi raggiunto il primo obiettivo: il costo computazionale rimane equivalente a quello del metodo del gradiente, almeno dal punto di vista asintotico. Vi è un secondo risultato positivo che si ottiene valutando le prestazioni di questo metodo quando esso viene applicato, al solito, alle funzioni quadratiche.

Proprietà 15 G_k converge a $H(\mathbf{x}_k)^{-1}$ per funzioni quadratiche anche se α_k non viene calcolato in modo ottimo.

Permangono però dei problemi:

- G_k non rimane definita positiva ad ogni iterazione e quindi non ci garantisce che la direzione scelta sia sempre una direzione di discesa.
- il valore $\mathbf{u}^T\mathbf{p}_k$ può diventare molto piccolo (\mathbf{u} può diventare prossimo al vettore nullo) introducendo problemi di instabilità numerica.

Per tali ragioni è stata introdotta una differente formula di aggiornamento.

6.2.7 Formula di aggiornamento di rango due (DFP)

Con la formula di aggiornamento di *rango due* si impone

$$G_{k+1} = G_k + c\mathbf{u}\mathbf{u}^T + d\mathbf{v}\mathbf{v}^T.$$

In questo caso, se si sceglie di imporre le relazioni

$$\mathbf{u} = \mathbf{h}_k, \quad c\mathbf{u}^T\mathbf{p}_k = 1, \quad \mathbf{v} = G_k\mathbf{p}_k \quad e \quad d\mathbf{v}^T\mathbf{p}_k = -1,$$

si ottiene la formula di aggiornamento seguente:

$$G_{k+1} = G_k + \frac{\mathbf{h}_k\mathbf{h}_k^T}{\mathbf{h}_k^T\mathbf{p}_k} - \frac{G_k\mathbf{p}_k\mathbf{p}_k^T G_k}{\mathbf{p}_k^T G_k\mathbf{p}_k}.$$

Questa formula di aggiornamento è alla base del cosiddetto *metodo DFP*, dai nomi degli autori Davidon, Fletcher e Powell, che gode di diverse proprietà.

Proprietà 16 Se $f(\mathbf{x})$ è una funzione qualsiasi il metodo DFP

- se G_0 è definita positiva (es. $G_1 = I$) genera matrici G_k definite positive, quindi ad ogni passo le direzioni \mathbf{d}_k sono di discesa;
- ha un costo computazionale dell'ordine $O(n^2)$ per iterazione;

- ha rapidità di convergenza superlineare.

Il costo computazionale rimane dell'ordine di quello del metodo del gradiente, inoltre le direzioni sono sempre direzioni di discesa e, importantissimo, ha una rapidità di convergenza notevolmente migliore di quella del metodo del gradiente, anche se non tale da uguagliare il metodo di Newton (quando quest'ultimo converge).

Proprietà 17 Se $f(\mathbf{x})$ è una funzione quadratica il metodo DFP con ricerca esatta del passo α_k

- termina in n iterazioni con $G_{n+1} = H(\mathbf{x}_n)^{-1}$;
- genera direzioni che soddisfano la seguente condizione³: $\mathbf{d}_i^T H \mathbf{d}_j = 0$, per ogni $i \neq j$;
- quando $G_0 = I$, genera gradienti coniugati (vedi nota a piè di pagina);
- $G_{k+1} \mathbf{p}_i = \mathbf{h}_i$ con $i = 1, \dots, k$ (la relazione (16) è soddisfatta in maniera retroattiva).

La precedente proprietà significa, in pratica, che quando questo metodo viene applicato alle funzioni quadratiche, esso converge in un numero di iterazioni al più pari alla dimensione dello spazio di ricerca, cioè in n passi.

Proprietà 18 Se $f(\mathbf{x})$ è una funzione convessa il metodo DFP con ricerca esatta del passo α_k converge.

Questo significa che la convergenza è garantita anche quando la matrice G_0 non è inizializzata alla matrice identità.

Ma si può fare di meglio.

6.2.8 Formula di aggiornamento di rango due inversa (BFGS)

In questa versione dei metodi *quasi Newton* si cerca una matrice B_k che approssimi direttamente la matrice Hessiana $H(\mathbf{x}_k)$ e non la sua inversa. Partiamo cioè dall'espressione $H(\mathbf{x}_k) \mathbf{h}_k \approx \mathbf{p}_k$ e imponiamo alla matrice B_k di soddisfarla come uguaglianza

$$B_{k+1} \mathbf{h}_k = \mathbf{p}_k \quad (17)$$

Cercando direttamente una formula di aggiornamento di *rango due* si impone

$$B_{k+1} = B_k + c \mathbf{u} \mathbf{u}^T + d \mathbf{v} \mathbf{v}^T$$

e con passaggi analoghi a quelli visti nella sezione (6.2.7) si perviene all'espressione

$$B_{k+1} = B_k + \frac{\mathbf{p}_k \mathbf{p}_k^T}{\mathbf{p}_k^T \mathbf{h}_k} - \frac{B_k \mathbf{h}_k \mathbf{h}_k^T B_k}{\mathbf{h}_k^T B_k \mathbf{h}_k}$$

dove vengono sostanzialmente invertiti i ruoli di \mathbf{h} e \mathbf{p} .

Naturalmente il metodo di Newton richiede di approssimare l'inversa di $H(\mathbf{x}_k)$, e dalla precedente espressione si può ricavare sia la G_k , che la relativa formula di aggiornamento

$$G_{k+1} = G_k + \left(1 + \frac{\mathbf{p}_k G_k \mathbf{p}_k^T}{\mathbf{h}_k^T \mathbf{p}_k} \right) \frac{\mathbf{h}_k \mathbf{h}_k^T}{\mathbf{h}_k^T \mathbf{p}_k} - \left(\frac{\mathbf{h}_k \mathbf{p}_k^T G_k + G_k \mathbf{p}_k \mathbf{h}_k^T}{\mathbf{h}_k^T \mathbf{p}_k} \right).$$

³Direzioni che soddisfano tali condizioni sono dette *direzioni coniugate* e vengono descritte nella sezione 6.3

Questa formula è nota come BFGS, dai nomi degli autori Broyden, Fletcher, Goldfarb e Shanno.

Il metodo che ne deriva è più robusto del metodo DFP rispetto agli errori di arrotondamento ed in particolare vale la seguente

Proprietà 19 *Data una funzione $f(\mathbf{x})$ il metodo BFGS converge globalmente purché la ricerca monodimensionale di α_k soddisfi la condizione di sufficiente riduzione della $f(\mathbf{x})$, relazione (2), e la condizione di Wolfe, relazione (3) (cfr. sezione 6.1).*

6.2.9 Famiglia di Broyden

Le formule di aggiornamento appena introdotte si possono vedere come particolarizzazioni di una famiglia di formule dette di *Broyden*

$$G_{k+1} = (1 - \phi)G_{k+1}^{DPF} + \phi G_{k+1}^{BFGS},$$

dove per le migliori proprietà e la stabilità del metodo si adotta $0 \leq \phi \leq 1$.

6.3 Metodi alle direzioni coniugate

Questi metodi sono stati sviluppati per risolvere efficacemente i modelli quadratici del tipo

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

Tali metodi convergono esattamente in un numero finito di iterazioni quando vengono applicati a funzioni quadratiche definite positive. Metodi che hanno tale proprietà di terminazione quadratica sono molto ricercati in quanto ci si aspetta che abbiano buone prestazioni anche per funzioni non quadratiche purché applicati nell'intorno di un minimo locale. Questo perché, dalle proprietà dello sviluppo in serie di Taylor, funzioni molto generali ben approssimano la forma quadratica nell'intorno di un minimo.

Due vettori \mathbf{u} e \mathbf{v} sono detti ortogonali se il loro prodotto scalare è nullo, cioè $\mathbf{u}^T \mathbf{v} = 0$. Similmente, due vettori sono detti *coniugati* rispetto ad una matrice simmetrica definita positiva Q se vale la relazione

$$\mathbf{u}^T Q \mathbf{v} = 0.$$

Il metodo si basa essenzialmente sulla seguente proprietà.

Teorema 14 *Siano \mathbf{u}_i , con $i = 1, \dots, n$, un insieme di vettori di \mathbb{R}^n mutuamente coniugati rispetto ad una matrice simmetrica definita positiva Q . Tali vettori formano una base per \mathbb{R}^n , cioè, per ogni vettore $\mathbf{x} \in \mathbb{R}^n$ vale*

$$\mathbf{x} = \sum_{i=1}^n \frac{\mathbf{u}_i(Q\mathbf{x})\mathbf{u}_i}{\mathbf{u}_i Q \mathbf{u}_i} = \sum_{i=1}^n \alpha_i \mathbf{u}_i, \quad \text{dove } \alpha_i = \frac{\mathbf{u}_i Q \mathbf{x}}{\mathbf{u}_i Q \mathbf{u}_i}$$

Useremo la precedente proprietà all'interno di un metodo iterativo

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

che, a partire dal vettore iniziale \mathbf{x}_0 , ci permetterà di determinare il minimo $\mathbf{x}^* = Q^{-1} \mathbf{b}$ della funzione quadratica $f(\mathbf{x})$ in un numero di passi minore o uguale ad n , senza la necessità di

invertire la matrice Q .

Il metodo usa come direzioni di discesa un insieme di n vettori \mathbf{d}_i , con $i = 0, \dots, n-1$, mutuamente coniugati rispetto a Q (che per ora assumeremo essere disponibili e successivamente vedremo come ricavare): $\mathbf{d}_i^T Q \mathbf{d}_j = 0$, per $i \neq j$.

Se osserviamo i passaggi del metodo iterativo possiamo ricavare che

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k = \mathbf{x}_{k-1} + \alpha_{k-1} \mathbf{d}_{k-1} + \alpha_k \mathbf{d}_k = \dots = \mathbf{x}_0 + \sum_{i=0}^k \alpha_i \mathbf{d}_i$$

Poiché i vettori \mathbf{d}_j formano una base per \mathbb{R}^n , noi sappiamo che ogni vettore \mathbf{x} si può esprimere come $\mathbf{x}_0 + \sum_{i=0}^{n-1} \alpha_i \mathbf{d}_i$. Ora, avendo a disposizione le direzioni coniugate \mathbf{d}_i , desideriamo scoprire quali valori devono assumere i moltiplicatori α_i affinché $\mathbf{x}_0 + \sum_{i=0}^{n-1} \alpha_i \mathbf{d}_i = \mathbf{x}^* = Q^{-1} \mathbf{b}$.

Poiché il punto di partenza \mathbf{x}_0 è dato e le \mathbf{d}_i sono assegnate, è possibile trasformare la funzione $f(\mathbf{x})$ in una funzione del vettore dei moltiplicatori incogniti α_i :

$$f(\mathbf{x}) = f\left(\mathbf{x}_0 + \sum_{i=0}^{n-1} \alpha_i \mathbf{d}_i\right) = F(\boldsymbol{\alpha}).$$

Abbiamo quindi

$$\begin{aligned} \min F(\boldsymbol{\alpha}) &= \min \frac{1}{2} \left(\mathbf{x}_0 + \sum_{i=0}^{n-1} \alpha_i \mathbf{d}_i\right)^T Q \left(\mathbf{x}_0 + \sum_{i=0}^{n-1} \alpha_i \mathbf{d}_i\right) - \mathbf{b}^T \left(\mathbf{x}_0 + \sum_{i=0}^{n-1} \alpha_i \mathbf{d}_i\right) = \\ &= \sum_{i=0}^{n-1} \min \left(\frac{1}{2} \left(\mathbf{x}_0 + \alpha_i \mathbf{d}_i\right)^T Q \left(\mathbf{x}_0 + \alpha_i \mathbf{d}_i\right) - \mathbf{b}^T \left(\mathbf{x}_0 + \alpha_i \mathbf{d}_i\right)\right) = \sum_{i=0}^{n-1} \min F_i(\alpha_i) \end{aligned}$$

dove si è sfruttato il fatto che $\mathbf{d}_i^T Q \mathbf{d}_j = 0$, per $i \neq j$. In questo modo, il problema originale è stato trasformato nella somma di n problemi monodimensionali indipendenti, nelle incognite $\alpha_0, \dots, \alpha_{n-1}$:

$$\min \left(\frac{1}{2} \alpha_i^2 \mathbf{d}_i^T Q \mathbf{d}_i + \alpha_i (Q \mathbf{x}_0)^T \mathbf{d}_i - \alpha_i \mathbf{b}^T \mathbf{d}_i\right),$$

dove in questo ultimo passaggio si è sfruttato il fatto che, poiché Q è simmetrica, $\mathbf{x}_i^T Q \mathbf{d}_i = \mathbf{d}_i^T Q \mathbf{x}_i$, $\mathbf{x}^T Q = (Q \mathbf{x})^T$, e si sono trascurati i termini indipendenti da α_i .

Risolvendo gli n problemi indipendenti, troveremo gli α_i^* ottimi e da questi ricaveremo

$$\mathbf{x}^* := \mathbf{x}_0 + \sum_{i=0}^{n-1} \alpha_i^* \mathbf{d}_i$$

Se deriviamo $F_i(\alpha_i)$ in funzione di α_i e eguagliamo a zero, otteniamo:

$$\alpha_i \mathbf{d}_i^T Q \mathbf{d}_i + (Q \mathbf{x}_0 - \mathbf{b})^T \mathbf{d}_i = 0$$

e quindi

$$\alpha_i^* = -\frac{(Q \mathbf{x}_0 - \mathbf{b})^T \mathbf{d}_i}{\mathbf{d}_i^T Q \mathbf{d}_i} = -\frac{\nabla f(\mathbf{x}_0)^T \mathbf{d}_i}{\mathbf{d}_i^T Q \mathbf{d}_i}$$

Riprendiamo ora il metodo iterativo

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k = \mathbf{x}_0 + \sum_{i=0}^k \alpha_i \mathbf{d}_i$$

e verifichiamo che, se utilizziamo i moltiplicatori ottimi α_i^* appena determinati, esso converge al passo n -esimo alla soluzione ottima.

Sostituendo a α_i i moltiplicatori ottimi α_i^* ricaviamo

$$\begin{aligned}\mathbf{x}_n &= \mathbf{x}_0 - \sum_{i=0}^{n-1} \frac{(Q\mathbf{x}_0 - \mathbf{b})^T \mathbf{d}_i}{\mathbf{d}_i^T Q \mathbf{d}_i} \mathbf{d}_i = \\ \mathbf{x}_n &= \mathbf{x}_0 - \sum_{i=0}^{n-1} \frac{\mathbf{d}_i^T (Q\mathbf{x}_0) \mathbf{d}_i}{\mathbf{d}_i^T Q \mathbf{d}_i} + \sum_{i=0}^{n-1} \frac{\mathbf{d}_i^T (Q(Q^{-1}\mathbf{b})) \mathbf{d}_i}{\mathbf{d}_i^T Q \mathbf{d}_i} \\ \mathbf{x}_n &= \mathbf{x}_0 - \mathbf{x}_0 + Q^{-1}\mathbf{b} = Q^{-1}\mathbf{b} = \mathbf{x}^*\end{aligned}$$

dove, nell'ultimo passaggio, abbiamo sfruttato il Teorema 14.

L'applicabilità del metodo dipende dal fatto che siano note n direzioni mutuamente coniugate. Il modo più diretto è calcolare gli autovettori di Q che godono la proprietà di essere sia mutuamente ortogonali che mutuamente coniugati rispetto a Q . Tuttavia la loro determinazione è altrettanto costosa quanto la soluzione del problema quadratico. Un modo alternativo è quello di modificare il processo di ortogonalizzazione di Gram-Schmidt in modo da generare direzioni coniugate invece che direzioni ortogonali. Questo approccio tuttavia richiede di memorizzare l'intero insieme di direzioni.

6.3.1 Il metodo di gradiente coniugato

Il metodo di *gradiente coniugato* è un metodo alle direzioni coniugate con la speciale proprietà di generare la direzione corrente \mathbf{d}_k sulla base della sola conoscenza (e memorizzazione) del vettore \mathbf{d}_{k-1} .

Nel richiamare brevemente il metodo ci svincoliamo anche dal suo utilizzo nell'ambito dei problemi di programmazione quadratica, e consideriamo la minimizzazione di funzioni $f(\mathbf{x})$ qualsiasi.

L'approccio proposto da Fletcher-Reeves (1964) propone le seguenti formule per ottenere direzioni coniugate:

$$\begin{aligned}\mathbf{d}_1 &= -\nabla f(\mathbf{x}_0) \\ \mathbf{d}_{k+1} &= -\nabla f(\mathbf{x}_k) + \beta_k \mathbf{d}_k \text{ e } \beta_k = \frac{\|\nabla f(\mathbf{x}_k)\|^2}{\|\nabla f(\mathbf{x}_{k-1})\|^2} \text{ per } k = 1, \dots, n-1\end{aligned}$$

Inoltre ad ogni passo è necessario eseguire una ricerca monodimensionale per determinare un minimo, o una sua approssimazione, per il passo α . Tuttavia, affinché le direzioni generate siano sempre di discesa è necessario che il passo α sia scelto ad ogni iterazione in modo da soddisfare le condizioni di Wolfe in senso forte (vedi formula (4)).

In alternativa, per β_k , si possono usare le formule di Polak-Ribière (1969):

$$\beta_k = \frac{(\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}))^T \nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_{k-1})\|^2}$$

che hanno la proprietà di ritornare alla direzione del gradiente ogniqualvolta la scelta della direzione coniugata non sia la scelta migliore, cioè quando $\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}) \approx 0$.

Si può dimostrare che per funzioni quadratiche i due metodi sono equivalenti e determinano la soluzione ottima se si adottano gli α_i^* ottimi.

Poichè non richiede operazioni matriciali questo approccio è adatto per problemi di grandi dimensioni.

7 Metodi di Trust-Region

Sia i metodi basati su ricerca lineare che i metodi di tipo *Trust-Region* determinano lo spostamento ad ogni iterazione con l'aiuto di un modello quadratico della funzione obiettivo. Essi usano però questo modello in modo diverso. Come abbiamo visto fino ad ora, i metodi basati su ricerca lineare, usano il modello quadratico per generare una direzione \mathbf{d}_k e poi si concentrano sulla scelta del passo α_k migliore. I metodi Trust-Region, invece, definiscono una regione (la Trust-Region, o Regione di Confidenza, appunto), di solito una (iper)sfera, intorno alla soluzione corrente, \mathbf{x}_k , nell'ipotesi che in tale regione il modello quadratico sia una rappresentazione adeguata della funzione obiettivo. Essi determinano simultaneamente la direzione e il passo che minimizzano il modello quadratico. Se la scelta del passo che deriva da questo approccio non soddisfa opportuni criteri, si rimane nel punto corrente, si riduce l'ampiezza della regione, in pratica il raggio dell'(iper)sfera, e si ripete il processo. Altrimenti, il nuovo punto corrente diventa la soluzione ottima del problema quadratico, e, se è il caso, si incrementa l'ampiezza della regione.

In generale, direzione e passo cambiano ogni volta che viene modificata l'ampiezza della Trust-Region. È quindi questa ampiezza, cioè il raggio dell'(iper)sfera, il punto chiave di questo metodo. Se la regione è troppo piccola, l'algoritmo perde l'opportunità di fare un passo lungo verso il minimo locale. Viceversa, se essa è troppo larga il modello quadratico potrebbe essere molto lontano dal rappresentare la funzione obiettivo nella Trust-Region, così da rendere necessaria una sua riduzione, perdendo una (o più) iterazione(i) senza, di fatto, muoversi. Dal punto di vista pratico si sceglie l'ampiezza della Trust-Region in base alle prestazioni che ha avuto l'algoritmo durante le iterazioni precedenti.

Nel seguito denoteremo con m_k la funzione modello adottata. Inoltre assumeremo che m_k sia un modello quadratico basato sullo sviluppo in serie di Taylor di $f(\mathbf{x})$ attorno al punto \mathbf{x}_k ,

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T H(\mathbf{x}_k) \mathbf{p}.$$

Se si denota con B_k una approssimazione della matrice hessiana mediante una qualche matrice simmetrica, il modello m_k viene definito come

$$m_k(\mathbf{p}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T B_k \mathbf{p},$$

La differenza fra $f(\mathbf{x}_k + \mathbf{p})$ e $m_k(\mathbf{p})$ è dell'ordine di $O(\|\mathbf{p}\|^2)$ che si riduce a $O(\|\mathbf{p}\|^3)$ quando B_k coincide con la matrice hessiana. In quest'ultimo caso il metodo prende il nome di Trust-Region Newton.

Ad ogni passo cerchiamo una soluzione al sottoproblema

$$\min_{\mathbf{p} \in \mathbb{R}^n} m_k(\mathbf{p}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T B_k \mathbf{p}, \text{ t.c. } \|\mathbf{p}\| \leq \Delta_k. \quad (18)$$

dove si denota con Δ_k l'ampiezza della Trust Region. Ci siamo spostati da un problema di minimo non vincolato ad un problema di ottimizzazione vincolata, dove sia la funzione obiettivo che i vincoli sono quadratici.

Quando la matrice B_k è definita positiva e vale la relazione $\|B_k^{-1} \nabla f(\mathbf{x}_k)\| \leq \Delta_k$, la soluzione del problema (18) è semplicemente il minimo non vincolato $\mathbf{p}^B = -B_k^{-1} \nabla f(\mathbf{x}_k)$. In questo

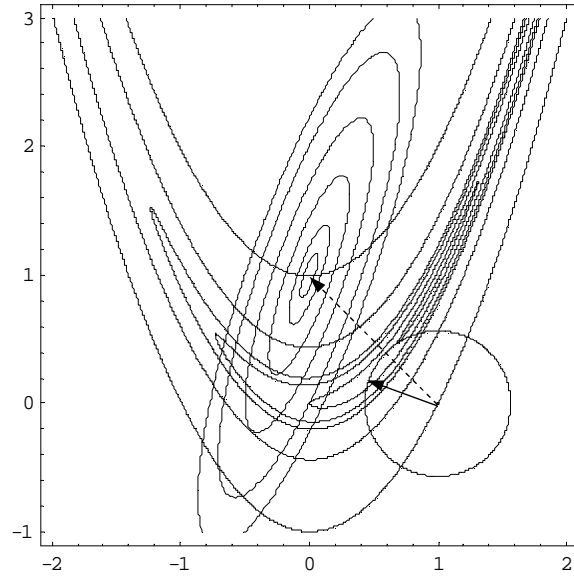


Figura 12: Esempio di modello quadratico per la funzione di Rosenbrock.

caso chiamiamo \mathbf{p}^B *full step*. In tutti gli altri casi occorre operare sul raggio Δ_k . Lo schema di un generico metodo Trust-Region è quindi definito da

- stabilire Δ_k
- risolvere il problema di minimo vincolato

In Figura 12 vediamo le curve di livello della funzione di Rosenbrock ed il relativo modello quadratico calcolato nel punto $(1, 0)$, $m_k(\mathbf{p}) = 601p_1^2 + 100p_2^2 - 400p_1p_2 + 400p_1 - 200p_2 + 100$. L'ottimo del modello quadratico non vincolato sta nel punto $(0, 1)$ indicato dalla freccia tratteggiata, mentre l'ottimo del modello con $\Delta = 0.4$ è indicato dalla freccia piena.

Come si è detto la scelta del raggio Δ_k della Trust Region è la chiave per un buon funzionamento del metodo. Tale scelta si basa su considerazioni attorno al rapporto fra la *riduzione attuale* del valore della funzione obiettivo determinata dal passo e la *riduzione predetta* dal modello

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})}{m_k(\mathbf{0}) - m_k(\mathbf{p}_k)}. \quad (19)$$

A seconda del segno di ρ_k possiamo avere vari casi. Se $\rho_k \leq 0$ il passo è peggiorante (o non migliorante) e viene rifiutato. Se è all'incirca pari a 1, m_k è una buona approssimazione della funzione obiettivo e quindi si può pensare di espandere la regione. Se ρ_k è positivo ma molto più piccolo di 1 la regione non viene modificata, mentre se il valore è prossimo a zero o negativo la regione viene ridotta, riducendo Δ_k .

```

Metodo Trust Region;
{
  Scegli  $\mathbf{x}_0 \in \mathbb{R}^n$ ;  $k := 0$ ;  $\bar{\Delta} > 0$ ,  $\Delta_0 \in (0, \bar{\Delta})$ ;  $\eta \in [0, \frac{1}{4}]$ ;
  While  $\nabla f(\mathbf{x}_k) \neq \emptyset$ ;
  {
     $\mathbf{p}_k := \operatorname{argmin}\{m_k(\mathbf{p}), \text{ t.c. } \|\mathbf{p}\| \leq \Delta_k\}$ ;
     $\rho_k := \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{p}_k)}{m_k(\mathbf{0}) - m_k(\mathbf{p}_k)}$ ;
    if  $\rho_k < \frac{1}{4}$  then
       $\Delta_{k+1} := \frac{1}{4}\Delta_k$ 
    else
      if  $\rho_k > \frac{3}{4}$  e  $\|\mathbf{p}_k\| = \Delta_k$  then
         $\Delta_{k+1} := \min\{2\Delta_k, \bar{\Delta}\}$ 
      else
         $\Delta_{k+1} := \Delta_k$ 
      if  $\rho_k > \eta$  then
         $\mathbf{x}_{k+1} := \mathbf{x}_k + \mathbf{p}_k$ 
      else
         $\mathbf{x}_{k+1} := \mathbf{x}_k$ 
       $k := k + 1$ ;
    }
  }
}

```

In questo algoritmo $\bar{\Delta}$ è un limite per eccesso sulla lunghezza dei passi e la regione viene allargata solo quando, nella soluzione ottima del problema (18), $\|\mathbf{p}_k\|$ raggiunge il bordo della regione.

Affinché il metodo Trust-Region possa funzionare in pratica è necessario risolvere efficientemente il sottoproblema (18), ed è su questo punto che concentriamo ora la nostra attenzione. Prima di procedere alla descrizione degli algoritmi vogliamo però caratterizzare in modo analitico le soluzioni del problema (18) mediante il seguente teorema, dove per semplicità adottiamo la notazione $\mathbf{g} = \nabla f(\mathbf{x})$ e sottintendiamo il pedice k , relativo all'iterazione corrente.

Teorema 15 *Il vettore \mathbf{p}^* è una soluzione globale ottima del problema di Trust-Region*

$$\min_{\mathbf{p} \in \mathbb{R}^n} m(\mathbf{p}) = f(\mathbf{x}) + \mathbf{g}^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T B \mathbf{p}, \text{ t.c. } \|\mathbf{p}\| \leq \Delta. \quad (20)$$

se e solo se \mathbf{p}^ è ammissibile ed esiste uno scalare $\lambda \geq 0$ tale che sono soddisfatte le seguenti condizioni*

$$(B + \lambda I) \mathbf{p}^* = -\mathbf{g}, \quad (21)$$

$$\lambda(\Delta - \|\mathbf{p}^*\|) = 0, \quad (22)$$

$$(B + \lambda I) \quad \text{è semidefinita positiva.} \quad (23)$$

Si osservi che la condizione (22) è una condizione di scarto complementare, ed implica che almeno una delle due quantità non negative, λ e $(\Delta - \|\mathbf{p}^*\|)$ deve essere zero. E quindi quando la soluzione giace strettamente dentro la regione di confidenza deve essere $\lambda = 0$ e

quindi $\mathbf{p}^* = -B^{-1}\mathbf{g}$, dalla (21), con B definita positiva dalla (23). Altrimenti si ha $\|\mathbf{p}^*\| = \Delta$ e λ può assumere un valore positivo. In questo caso, dalla (21) si ha

$$\lambda\mathbf{p}^* = -B\mathbf{p}^* - \mathbf{g} = -\nabla m(\mathbf{p}^*)$$

cioè la soluzione \mathbf{p}^* è collineare con l'antigradiente del modello quadratico.

7.1 Il punto di Cauchy

Come abbiamo visto, in particolare nel Teorema 8, la convergenza globale degli algoritmi basati su ricerca lineare viene garantita dal rispetto di condizioni non molto stringenti, sia sulla direzione \mathbf{d} , sia sul passo α . Qualcosa di simile accade con gli algoritmi di tipo Trust-Region. Infatti, per la convergenza globale di tali metodi è sufficiente trovare, ad ogni iterazione k , una soluzione approssimata, \mathbf{p}_k , del problema (18) che appartenga alla Trust-Region e garantisca una *sufficiente riduzione* del modello quadratico.

Tale riduzione viene quantificata in relazione al cosiddetto *punto di Cauchy*, che nel seguito verrà denotato come \mathbf{p}_k^C . Tale punto non è altro che il punto di minimo del modello quadratico lungo la direzione dell'antigradiente.

Sia

$$\mathbf{p}_k^S = -\frac{\Delta_k}{\|\nabla f(\mathbf{x}_k)\|} \nabla f(\mathbf{x}_k)$$

la direzione dell'antigradiente *normalizzata* in ragione dell'attuale ampiezza della Trust-Region. Cerchiamo ora il minimo τ_k del problema monodimensionale

$$\min_{\tau \geq 0} \mu_k(\tau) = \min_{\tau \geq 0} m_k(\tau\mathbf{p}_k^S) \text{ t.c. } \|\tau\mathbf{p}_k^S\| \leq \Delta_k$$

e definiamo $\mathbf{p}_k^C = \tau_k\mathbf{p}_k^S$.

Cominciamo col riscrivere esplicitamente il problema monodimensionale, dopo aver adottato la notazione \mathbf{g}_k al posto di $\nabla f(\mathbf{x}_k)$,

$$\mu_k(\tau) = f(\mathbf{x}_k) - \tau \frac{\Delta_k}{\|\mathbf{g}_k\|} \mathbf{g}_k^T \mathbf{g}_k + \frac{1}{2} \tau^2 \frac{\Delta_k^2}{\|\mathbf{g}_k\|^2} \mathbf{g}_k^T B_k \mathbf{g}_k.$$

Quando $\mathbf{g}_k^T B_k \mathbf{g}_k \leq 0$ e $\mathbf{g}_k \neq \mathbf{0}$ la funzione $\mu_k(\tau)$ decresce monotonicamente al crescere di τ e quindi τ_k è il più grande valore che soddisfa i limiti della Trust-Region, cioè $\tau_k = 1$.

Quando $\mathbf{g}_k^T B_k \mathbf{g}_k > 0$, $\mu_k(\tau)$ definisce un problema quadratico convesso, così il valore τ_k è il minimo fra il valore limite 1 e l'ottimo di tale problema. Tale ottimo si ricava derivando $\mu_k(\tau)$ rispetto a τ e imponendo l'uguaglianza a zero

$$\mu_k'(\tau) = -\frac{\Delta_k}{\|\mathbf{g}_k\|} \mathbf{g}_k^T \mathbf{g}_k + \tau \frac{\Delta_k^2}{\|\mathbf{g}_k\|^2} \mathbf{g}_k^T B_k \mathbf{g}_k = 0.$$

Quindi

$$\tau^* = \frac{\|\mathbf{g}_k\|^3}{\Delta_k \mathbf{g}_k^T B_k \mathbf{g}_k}.$$

Riassumendo abbiamo

$$\mathbf{p}_k^C = -\tau_k \frac{\Delta_k}{\|\mathbf{g}_k\|} \mathbf{g}_k$$

dove

$$\tau_k = \begin{cases} 1 & \text{se } \mathbf{g}_k^T B_k \mathbf{g}_k \leq 0 \\ \min\{\tau^*, 1\} & \text{altrimenti} \end{cases}$$

Ricavare il punto \mathbf{p}_k^C ha un costo computazionale dell'ordine $O(n^2)$ e, come detto, la convergenza dei metodi di Trust-Region è garantita quando ogni passo \mathbf{p}_k determina una riduzione del modello quadratico che è almeno qualche fissato multiplo positivo della riduzione fornita dal punto di Cauchy.

Dal punto di vista pratico, l'uso diretto dei punti di Cauchy nella soluzione del problema quadratico (18) comporta gli stessi inconvenienti dell'uso del metodo del gradiente all'interno dei metodi basati su ricerca lineare: convergenza lineare e prestazioni scadenti. È quindi necessario fare in modo che la matrice B_k rivesta un qualche ruolo nella determinazione della direzione oltre che del passo. Delle differenti varianti presenti in letteratura ci limitiamo ad introdurre un metodo che si può adottare quando la matrice B_k è definita positiva (ad esempio nei casi affrontati nella Sezione 8).

7.2 Il metodo *Dogleg*

Per semplicità anche nel seguito sottointenderemo il pedice k , relativo all'iterazione corrente, adotteremo ancora una volta la convenzione $\mathbf{g} = \nabla f(\mathbf{x})$, e faremo quindi riferimento al problema

$$\min_{\mathbf{p} \in \mathbb{R}^n} m(\mathbf{p}) = f(\mathbf{x}) + \mathbf{g}^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T B \mathbf{p}, \text{ t.c. } \|\mathbf{p}\| \leq \Delta. \quad (24)$$

Il metodo *Dogleg*, *a zampa di cane*, viene impiegato quando la matrice B è definita positiva. In tal caso la soluzione del problema quadratico non vincolato è $\mathbf{p}^B = -B^{-1}\mathbf{g}$, e tale soluzione viene scelta, cioè $\mathbf{p}^* = \mathbf{p}^B$, ogniqualvolta $\|\mathbf{p}^B\| \leq \Delta$, e quindi più frequentemente quando i valori di Δ sono grandi.

Quando la dimensione della regione Δ è piccola rispetto alla norma di \mathbf{p}^B , il vincolo $\|\mathbf{p}\| \leq \Delta$ garantisce che il termine quadratico possa essere trascurato nella soluzione di (20) e che si possa utilizzare la seguente approssimazione

$$\mathbf{p}^* \approx -\Delta \frac{\mathbf{g}}{\|\mathbf{g}\|}.$$

Per valori di Δ intermedi si deve procedere con scelte più sofisticate.

Il metodo *Dogleg* cerca una soluzione approssimata di \mathbf{p}^* come una concatenazione di due segmenti. Il primo segmento va dal punto corrente al punto che minimizza m lungo la direzione dell'antigradiente (vedi formula (11)), cioè

$$\mathbf{p}^U = -\frac{\mathbf{g}^T \mathbf{g}}{\mathbf{g}^T B \mathbf{g}} \mathbf{g}$$

mentre il secondo segmento unisce \mathbf{p}^U con \mathbf{p}^B . La traiettoria che ne deriva può venir descritta formalmente come $\tilde{\mathbf{p}}(\tau)$ per $\tau \in [0, 2]$, dove

$$\tilde{\mathbf{p}}(\tau) = \begin{cases} \tau \mathbf{p}^U, & \tau \in [0, 1] \\ \mathbf{p}^U + (\tau - 1)(\mathbf{p}^B - \mathbf{p}^U), & \tau \in [1, 2] \end{cases}$$

Il metodo *Dogleg* sceglie il valore di \mathbf{p} che minimizza il modello quadratico m lungo il cammino $\tilde{\mathbf{p}}(\tau)$, vincolato alla dimensione della regione. La soluzione di tale problema risulta semplice grazie alla seguente

Proprietà 20 Se la matrice B è definita positiva allora

- $\|\tilde{\mathbf{p}}(\tau)\|$ è monotona crescente
- $m(\tilde{\mathbf{p}}(\tau))$ è monotona decrescente

Tale proprietà comporta che il cammino $\tilde{\mathbf{p}}(\tau)$ intersechi il confine della regione in un solo punto nei casi in cui $\|\mathbf{p}^B\| \geq \Delta$, e non lo intersechi affatto altrimenti. Nel caso in cui $\|\mathbf{p}^B\| \geq \Delta$ si può calcolare il valore esatto di τ risolvendo l'equazione quadratica scalare

$$\|\mathbf{p}^U + (\tau - 1)(\mathbf{p}^B - \mathbf{p}^U)\| = \Delta$$

Dal punto di vista della convergenza, la proprietà 20 garantisce che il metodo *Dogleg* comporti ad ogni passo una riduzione che è almeno pari a quella garantita dal punto di Cauchy.

8 Problemi ai Minimi Quadrati

All'inizio della Sezione 3 abbiamo introdotto come esempio, il problema di determinare i valori ottimi di un insieme di parametri di una funzione in modo da minimizzare la somma dei quadrati degli scarti fra i valori assunti dalla funzione ed i valori di un insieme di dati campionati. Riprendiamo l'esempio semplificandolo un po'.

Esempio Supponiamo di cercare una curva che si adatti ad un insieme di dati sperimentali, y_1, y_2, \dots, y_4 , campionati agli istanti t_1, t_2, \dots, t_4 . Il nostro modello ha la forma della funzione

$$\phi(t_i, \mathbf{x}) = e^{t_i x_1} + \cos(t_i x_2)$$

Il modello non lineare ha due parametri, x_1, x_2 . Desideriamo sceglierli in modo da far aderire il più possibile i valori $\phi(t_i, \mathbf{x})$ ai dati y_i . Nel seguito indichiamo con

$$r_i(\mathbf{x}) = \phi(t_i, \mathbf{x}) - y_i, \quad i = 1, \dots, m,$$

i cosiddetti residui, la somma dei cui quadrati intendiamo minimizzare.

Il problema è noto come *problema ai minimi quadrati* non lineare (least-squares problem) e la funzione obiettivo da minimizzare assume la forma

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^m r_i^2(\mathbf{x})$$

dove ogni residuo r_i è una funzione da \mathbb{R}^n in \mathbb{R} . Nel seguito si assumerà che $m \geq n$ anche se di solito la relazione è $m \gg n$. Si osservi che il calcolo della sola funzione obiettivo, per un dato valore delle variabili \mathbf{x} , può essere computazionalmente oneroso, anche se le variabili sono poche, dipendendo dal numero m di osservazioni, che di solito è molto elevato.

Per vedere come la forma speciale della funzione obiettivo $f(\mathbf{x})$ renda il problema di ottimizzazione ai minimi quadrati più semplice della maggior parte dei problemi di ottimizzazione non lineare, cominciamo col raggruppare i residui r_i in un *vettore dei residui* $\mathbf{r}(\mathbf{x})$ che vedremo come una funzione da \mathbb{R}^n in \mathbb{R}^m :

$$\mathbf{r}(\mathbf{x}) = (r_1(\mathbf{x}), r_2(\mathbf{x}), \dots, r_m(\mathbf{x}))^T.$$

Con questa notazione possiamo scrivere

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{r}(\mathbf{x})^T \mathbf{r}(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2.$$

Inoltre, le derivate di $f(\mathbf{x})$ si possono esprimere in forma compatta utilizzando lo Jacobiano

$$\mathbf{J}(\mathbf{x}) = \left[\frac{\partial r_i(\mathbf{x})}{\partial x_j} \right]_{\substack{i=1,2,\dots,m \\ j=1,2,\dots,n}} = \begin{bmatrix} \nabla r_1(\mathbf{x})^T \\ \nabla r_2(\mathbf{x})^T \\ \vdots \\ \nabla r_m(\mathbf{x})^T \end{bmatrix}$$

dove $\nabla r_i(\mathbf{x})$, con $i = 1, \dots, m$, denota il vettore gradiente dei residui $r_i(\mathbf{x})$.

Con questa notazione il vettore gradiente e la matrice hessiana di $f(\mathbf{x})$ si possono esprimere come segue:

$$\nabla f(\mathbf{x}) = \begin{pmatrix} \sum_{i=1}^m \frac{\partial r_i(\mathbf{x})}{\partial x_1} r_i(\mathbf{x}) \\ \vdots \\ \sum_{i=1}^m \frac{\partial r_i(\mathbf{x})}{\partial x_n} r_i(\mathbf{x}) \end{pmatrix} = \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x}), \quad (25)$$

$$H(\mathbf{x}) = \mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + \sum_{i=1}^m \nabla^2 r_i(\mathbf{x}) r_i(\mathbf{x}) \quad (26)$$

dove $\nabla^2 r_i(\mathbf{x})$ è la matrice hessiana dei residui.

Continua esempio. La matrice Jacobiana dell'esempio precedente risulta

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} t_1 e^{t_1 x_1} & -t_1 \sin(t_1 x_2) \\ t_2 e^{t_2 x_1} & -t_2 \sin(t_2 x_2) \\ t_3 e^{t_3 x_1} & -t_3 \sin(t_3 x_2) \\ t_4 e^{t_4 x_1} & -t_4 \sin(t_4 x_2) \end{bmatrix},$$

il vettore gradiente di $f(\mathbf{x})$ diviene

$$\nabla f(\mathbf{x}) = \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x}) = \begin{pmatrix} \sum_{i=1}^4 t_i e^{t_i x_1} r_i(\mathbf{x}) \\ -\sum_{i=1}^4 t_i \sin(t_i x_2) r_i(\mathbf{x}) \end{pmatrix},$$

il primo termine della matrice hessiana $H(\mathbf{x})$ diviene

$$\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) = \begin{pmatrix} \sum_{i=1}^4 t_i^2 e^{2t_i x_1} & -\sum_{i=1}^4 t_i^2 e^{t_i x_1} \sin(t_i x_2) \\ -\sum_{i=1}^4 t_i^2 e^{t_i x_1} \sin(t_i x_2) & \sum_{i=1}^4 t_i^2 \sin^2(t_i x_2) \end{pmatrix},$$

la matrice hessiana dei residui diviene

$$\nabla^2 r_i(\mathbf{x}) = \begin{pmatrix} t_i^2 e^{t_i x_1} & 0 \\ 0 & -t_i^2 \cos(t_i x_2) \end{pmatrix},$$

ed infine la matrice hessiana $H(\mathbf{x})$ risulta

$$H(\mathbf{x}) = \begin{pmatrix} \sum_{i=1}^4 t_i^2 e^{t_i x_1} (e^{t_i x_1} + r_i(\mathbf{x})) & -\sum_{i=1}^4 t_i^2 e^{t_i x_1} \sin(t_i x_2) \\ -\sum_{i=1}^4 t_i^2 e^{t_i x_1} \sin(t_i x_2) & \sum_{i=1}^4 t_i^2 (\sin^2(t_i x_2) - \cos(t_i x_2) r_i(\mathbf{x})) \end{pmatrix}.$$

In molte applicazioni le derivate parziali dei residui e quindi lo Jacobiano sono calcolabili in modo efficiente, per cui risulta utile applicare la formula (25). Inoltre è parimenti possibile calcolare il primo termine $\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x})$ della matrice hessiana in modo non dispendioso, non dovendo calcolare le derivate parziali seconde di $r_i(\mathbf{x})$. Questo aspetto è il punto distintivo dei problemi ai minimi quadrati ed è quello che viene sfruttato dagli algoritmi specializzati a risolvere tali problemi. L'idea infatti è che spesso il primo termine della matrice hessiana nell'espressione (26) è più rilevante del secondo, o poiché sono piccoli i residui $r_i(\mathbf{x})$, o poiché sono piccole le norme delle loro matrici hessiane $\nabla^2 r_i(\mathbf{x})$.

Nelle prossime due sezioni vedremo due metodi di risoluzione per problemi ai minimi quadrati non lineari, il primo appartenente alla famiglia dei metodi basati su ricerca lineare, mentre il secondo è l'antesignano dei metodi di Trust-Region.

8.1 Il metodo di Gauss-Newton

Il metodo di Gauss-Newton è un metodo basato su ricerca lineare

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$$

che si può interpretare come un metodo di Newton modificato. Per semplicità, nel seguito lasceremo cadere il pedice k relativo all'iterazione corrente. Per ricavare la direzione \mathbf{p} , invece di risolvere il sistema $H(\mathbf{x})\mathbf{p} = -\nabla f(\mathbf{x})$ si risolve un analogo sistema dove la matrice hessiana viene approssimata trascurando il secondo termine nell'espressione (26):

$$\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) \mathbf{p}^{GN} = -\mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x}). \quad (27)$$

Questo modello presenta una serie di vantaggi.

- Innanzitutto, in termini di tempo di calcolo, il risparmio dato dal non dover calcolare le derivate seconde dei residui può essere determinante in molte applicazioni, oltre al fatto che il termine $\mathbf{J}^T \mathbf{J}$ risulta semplice da calcolare una volta calcolato il termine $\mathbf{J}^T \mathbf{r}$.
- In secondo luogo, l'approssimazione introdotta spesso non è significativa, e questo comporta che il metodo di Gauss-Newton ha un tasso di convergenza simile a quello del metodo di Newton, specialmente in prossimità della soluzione ottima.
- Un terzo vantaggio è che, ogniqualvolta la matrice \mathbf{J} ha rango pieno e il gradiente $\nabla f(\mathbf{x})$ è non nullo, la direzione \mathbf{p}^{GN} è una direzione di discesa per $f(\mathbf{x})$ e quindi adatta ad essere utilizzata nei metodi basati su ricerca lineare. Per riconoscere tale proprietà si ricava

$$(\mathbf{p}^{GN})^T \nabla f(\mathbf{x}) = (\mathbf{p}^{GN})^T \mathbf{J}^T \mathbf{r} = -(\mathbf{p}^{GN})^T \mathbf{J}^T \mathbf{J} \mathbf{p}^{GN} = -\|\mathbf{J} \mathbf{p}^{GN}\|^2 \leq 0.$$

Tale disuguaglianza è stretta, a meno di non essere in un punto stazionario dove si annulla il gradiente $\mathbf{J}^T \mathbf{r} = \mathbf{0}$ e quindi $\mathbf{J} \mathbf{p}^{GN} = \mathbf{0}$.

- Un ulteriore vantaggio deriva dal fatto che \mathbf{p}^{GN} è la soluzione del modello ai minimi quadrati lineare

$$\min_{\mathbf{p} \in \mathbb{R}^n} f(\mathbf{p}) = \frac{1}{2} \|\mathbf{J} \mathbf{p} + \mathbf{r}\|^2.$$

per il quale sono disponibili algoritmi specializzati molto efficienti basati principalmente su particolari fattorizzazioni della matrice $\mathbf{J}^T \mathbf{J}$.

Le implementazioni più diffuse del metodo Gauss-Newton eseguono una ricerca lineare per la determinazione del passo α lungo la direzione \mathbf{p}^{GN} , e richiedono che il valore di α soddisfi le condizioni di Armijo (vedi equazione (2)) e di Wolfe (vedi equazione (3)).

Il metodo di Gauss-Newton converge globalmente qualora si possa garantire che lo Jacobiano $\mathbf{J}(\mathbf{x})$ soddisfi la seguente relazione

$$\|\mathbf{J}(\mathbf{x})\mathbf{z}\| \geq \gamma\|\mathbf{z}\| \quad (28)$$

in tutti i punti \mathbf{x} nell'intorno dell'insieme di livello $\{\mathbf{x}|f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$, dove $\gamma > 0$ e \mathbf{x}_0 è la soluzione di partenza dell'algoritmo.

Nei casi in cui non vale la relazione (28) si dice che la matrice \mathbf{J} è *rank-deficient* e come conseguenza la matrice di coefficienti $\mathbf{J}^T\mathbf{J}$ risulta singolare. Il sistema (27) ammette ancora soluzione, in effetti ne ammette infinite, e questo provoca difficoltà di convergenza. In questi casi risulta più efficace il metodo descritto nella prossima sezione.

8.2 Il metodo di Levenberg-Marquardt

Mentre il metodo di Gauss-Newton è un metodo basato su ricerca lineare, il metodo di Levenberg-Marquardt è un metodo di tipo Trust-Region. Anche in questo caso la matrice hessiana viene approssimata limitandosi a considerare il primo termine dell'espressione (26) ed il metodo che si ricava risulta efficace anche in presenza di matrici \mathbf{J} *rank-deficient*.

Il problema che viene risolto ad ogni iterazione è

$$\min_{\mathbf{p} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{J}\mathbf{p} + \mathbf{r}\|^2, \text{ t.c. } \|\mathbf{p}\| \leq \Delta. \quad (29)$$

dove ancora una volta si è ommesso il pedice k relativo all'iterazione corrente. Il relativo modello quadratico risulta

$$\min_{\mathbf{p} \in \mathbb{R}^n} m(\mathbf{p}) = \frac{1}{2} \|\mathbf{r}\|^2 + \mathbf{p}^T \mathbf{J}^T \mathbf{r} + \frac{1}{2} \mathbf{p}^T \mathbf{J}^T \mathbf{J} \mathbf{p}. \quad (30)$$

Il Teorema 15 ci permette di caratterizzare la soluzione del problema (29) nel seguente modo.

- Quando la soluzione \mathbf{p}^{GN} dell'equazione (27) è interna alla regione di confidenza ($\|\mathbf{p}^{GN}\| < \Delta$) questa è anche la soluzione del problema (29).
- Altrimenti esiste un valore $\lambda > 0$ tale che la soluzione \mathbf{p}^{LM} del problema (29) sta nel bordo della regione di confidenza ($\|\mathbf{p}^{LM}\| = \Delta$) e vale la relazione

$$(\mathbf{J}^T \mathbf{J} + \lambda \mathbf{I}) \mathbf{p}^{LM} = -\mathbf{J}^T \mathbf{r}.$$

I metodi più efficaci cercano di identificare iterativamente il valore λ che permette di soddisfare quest'ultima uguaglianza. Essi si basano su particolari fattorizzazioni della matrice $\mathbf{J}^T \mathbf{J} + \lambda \mathbf{I}$. Il metodo di Levenberg-Marquardt converge sotto le lasche ipotesi di convergenza dei metodi di Trust-Region.

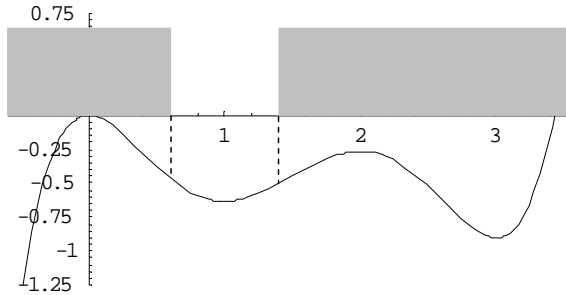


Figura 13: Esempio di funzione obiettivo convessa solo nell'insieme ammissibile.

9 Ottimizzazione vincolata

Consideriamo ora il generico problema di ottimizzazione vincolata

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ g_j(\mathbf{x}) \leq & 0 \quad j = 1, \dots, k; \\ h_j(\mathbf{x}) = & 0 \quad j = 1, \dots, h \end{aligned}$$

con $\mathbf{x} \in \mathbb{R}^n$.

Inizieremo con lo studiare le condizioni analitiche di ottimalità. Nel caso dei problemi non vincolati, tutti i minimi locali soddisfano le condizioni necessarie di ottimalità e, almeno in linea di principio, i minimi locali possono venir cercati all'interno dell'insieme dei punti stazionari. Nel caso dei problemi vincolati, invece, non sempre si possono ricavare tutti i minimi locali anche quando risulta possibile imporre il soddisfacimento delle condizioni analitiche.

Nel seguito analizzeremo in primo luogo la regione ammissibile e le relative proprietà, in seguito distingueremo i problemi vincolati analizzando separatamente il caso con soli vincoli di uguaglianza, quello con soli vincoli di disuguaglianza ed infine il caso generale.

Inoltre tratteremo a parte il caso di vincoli lineari ed i problemi quadratici.

Osserviamo innanzitutto che la presenza di vincoli può sia rendere più semplice un problema non vincolato difficile, sia rendere difficile un problema che in assenza di vincoli sarebbe semplice.

In Figura 13 vediamo l'andamento della funzione monodimensionale

$$f(x) = \frac{1}{5}x^5 - \frac{3}{2}x^4 + \frac{11}{3}x^3 - 3x^2.$$

Tale funzione è non convessa in \mathbb{R} , mentre è convessa in opportuni intervalli. Ad esempio, nell'intervallo evidenziato $X = [0.6, 1.4]$. In particolare, qualsiasi tecnica di ottimizzazione monodimensionale a partire da un punto $x_0 \in X$ determinerebbe rapidamente l'ottimo globale del problema.

Al contrario invece, in Figura 14 vediamo le curve di livello e la regione ammissibile, in chiaro, del problema

$$\begin{aligned} \min \quad & f(x, y) = (x - 1)^2 + (y + 1)^2 \\ g_1(x, y) = & 1 + \frac{1}{4} \sin(8x) - y \leq 0; \\ g_2(x, y) = & -y \leq 0. \end{aligned}$$

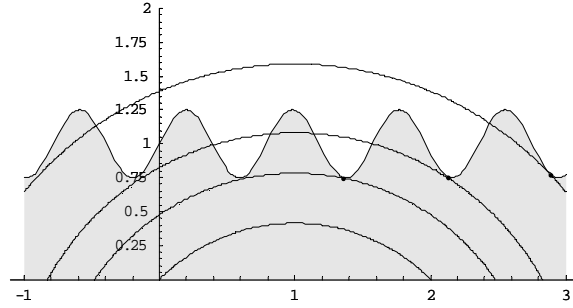


Figura 14: Esempio di funzione obiettivo convessa e infiniti minimi locali.

La funzione obiettivo è convessa ed ammette un unico punto stazionario, $(1, -1)$, che sarebbe un minimo globale per il problema non vincolato. Il problema vincolato ammette invece infiniti minimi locali, tre dei quali evidenziati in figura.

9.1 Condizioni analitiche: vincoli di uguaglianza

Analizziamo in primo luogo un problema vincolato da soli vincoli di uguaglianza

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ h_j(\mathbf{x}) &= 0 \quad j = 1, \dots, h < n \end{aligned}$$

Nel 1760 Lagrange trasformò questo problema vincolato in un problema non vincolato mediante l'introduzione dei cosiddetti moltiplicatori di Lagrange, λ_j , con $j = 1, \dots, h$ nella formulazione della cosiddetta *funzione Lagrangiana*

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{j=1}^h \lambda_j h_j(\mathbf{x}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}).$$

Le condizioni necessarie per l'esistenza di un minimo del problema vincolato con vincoli di uguaglianza possono essere date in termini della funzione Lagrangiana e dei moltiplicatori di Lagrange.

Teorema 16 *Sono date una funzione $f(\mathbf{x})$ ed h funzioni $h_j(\mathbf{x})$, con $j = 1, \dots, h$ di classe C^1 . Condizioni necessarie, nell'ipotesi che i vettori gradienti delle funzioni h_j calcolati nel punto \mathbf{x}^* siano tra loro linearmente indipendenti, affinché \mathbf{x}^* sia un minimo locale del problema con vincoli di uguaglianza è che esista $\boldsymbol{\lambda}^*$ tale che $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ sia un punto stazionario della funzione $L(\mathbf{x}, \boldsymbol{\lambda})$, cioè:*

$$\frac{\partial L}{\partial x_i} = \frac{\partial f(\mathbf{x}^*)}{\partial x_i} + \sum_{j=1}^h \lambda_j^* \frac{\partial h_j(\mathbf{x}^*)}{\partial x_i} = 0, \quad i = 1, 2, \dots, n \quad (31)$$

$$\frac{\partial L}{\partial \lambda_j} = h_j(\mathbf{x}^*) = 0, \quad j = 1, 2, \dots, h \quad (32)$$

Le condizioni rappresentano un sistema di $n + h$ equazioni in $n + h$ incognite. Il secondo gruppo di h condizioni coincide con la richiesta che i vincoli di uguaglianza siano soddisfatti nel punto di ottimo. Il primo gruppo di n condizioni coincide con $\nabla f(\mathbf{x}^*) + J(\mathbf{x}^*)^T \boldsymbol{\lambda}^* = \mathbf{0}$, ovvero

$$-\nabla f(\mathbf{x}^*) = \sum_{j=1}^h \lambda_j^* \nabla h_j(\mathbf{x}^*)$$

che dal punto di vista geometrico esprime la richiesta che l'antigradiente della funzione obiettivo si possa ottenere come combinazione lineare dei gradienti dei vincoli di uguaglianza.

Esempio. Consideriamo il problema

$$\begin{aligned} \min \quad & f(x, y) = (x - 2)^2 + (y - 2)^2 \\ & h_1(x, y) = 1 - x^2 - y^2 = 0. \end{aligned}$$

Il punto di ottimo è $(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2})$ e in tale punto l'antigradiente $-\nabla f(x, y) = -(2(x-2), 2(y-2))^T$ vale $(4 - \sqrt{2}, 4 - \sqrt{2})^T$. Nello stesso punto il gradiente di h_1 , $\nabla h(x, y) = (-2x, -2y)^T$, vale $(-\sqrt{2}, -\sqrt{2})$. Come illustra la Figura 15 di destra i due vettori sono collineari e $\lambda^* = \frac{4-\sqrt{2}}{\sqrt{2}}$.

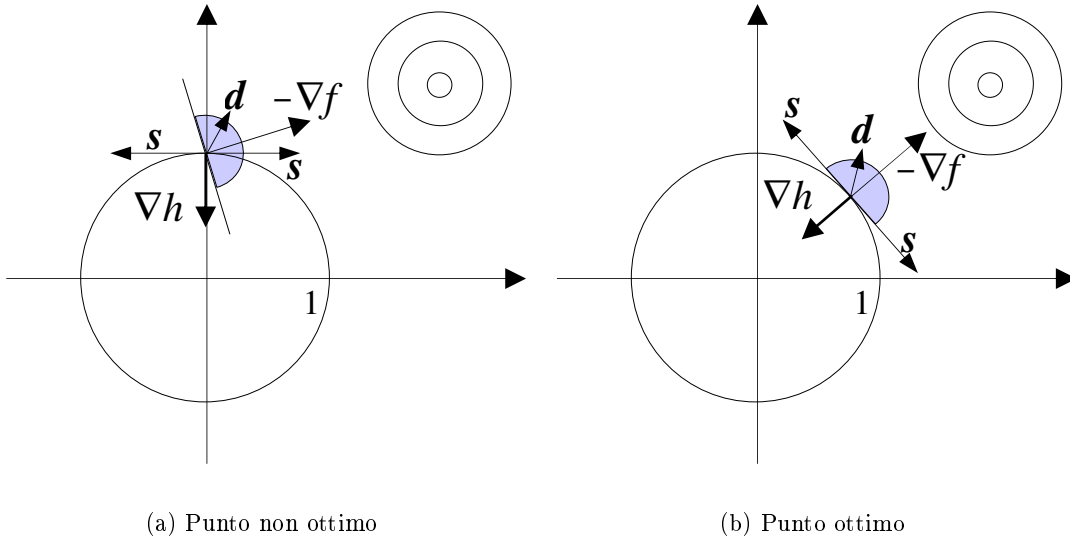


Figura 15: Condizione di ottimalità per problemi con vincoli di uguaglianza.

Ma che cosa caratterizza tale punto geometricamente? Osserviamo nella Figura 15 di sinistra il punto non ottimale $P_1 = (0, 1)$. In tale punto sono evidenziate le direzioni \mathbf{s} ortogonali al vettore $\nabla h(x, y)$, cioè tali che $\nabla h(x, y)^T \mathbf{s} = 0$. L'insieme di queste direzioni definisce un iperpiano (in questo caso una retta), diciamo F , che rappresenta l'approssimazione al primo ordine della funzione $h(x, y)$ in P_1 . In P_1 è inoltre evidenziato, mediante un semicerchio ombreggiato, il sottospazio, diciamo D , definito da tutte le direzioni di discesa \mathbf{d} , che formano cioè con l'antigradiente $-\nabla f(x, y)$ un angolo θ tale che $\cos \theta > 0$, quindi tali che

$$\nabla f(x, y)^T \mathbf{d} < 0.$$

Il punto P_1 non è ottimo in quanto in tale punto esistono direzioni che appartengono sia ad F che a D , seguendo le quali, almeno per un tratto infinitesimo, si migliora la funzione obiettivo e si continua a soddisfare il vincolo di uguaglianza.

Nella Figura 15 di destra il punto è ottimo proprio perché in tale punto non esistono direzioni che appartengono a $F \cap D$: le direzioni miglioranti non appartengono alla linearizzazione del vincolo di uguaglianza. Tale condizione si ha quando i vettori $-\nabla f(x, y)$ e $\nabla h(x, y)$ sono

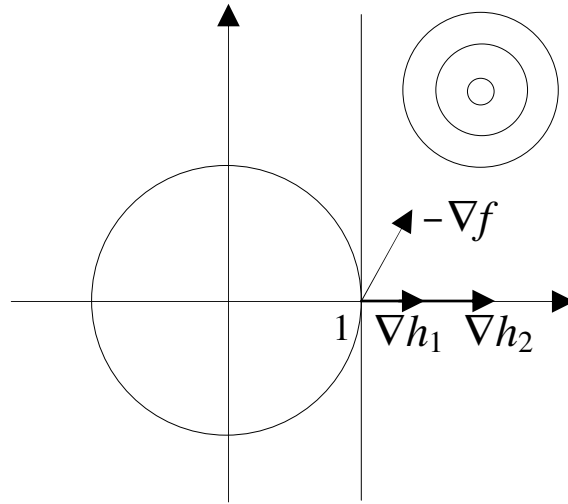


Figura 16: Esempio di gradienti di vincoli di uguaglianza linearmente dipendenti.

collineari.

Il Teorema 16 fa però menzione dell'ulteriore richiesta che:

[...] i vettori gradienti delle funzioni h_j calcolati nel punto \mathbf{x}^* siano tra loro linearmente indipendenti [...]

Vediamo con un esempio il significato geometrico di tale richiesta.

Esempio. Consideriamo il problema

$$\begin{aligned} \min \quad & f(x, y) = (x - 2)^2 + (y - 2)^2 \\ & h_1(x, y) = x^2 + y^2 - 1 = 0; \\ & h_2(x, y) = x - 1 = 0. \end{aligned}$$

Il punto di ottimo è l'unico punto ammissibile $(1, 0)$. In tale punto i vettori $\nabla h_1(x, y)$ e $\nabla h_2(x, y)$ sono $(2, 0)^T$ e $(1, 0)^T$, rispettivamente, sono collineari e quindi linearmente dipendenti. Nello stesso punto l'antigradiente della funzione obiettivo, $-\nabla f(x, y)$, vale $(2, 4)^T$ e per nessun valore di λ_1 e λ_2 può venir soddisfatto il sistema di equazioni

$$\begin{pmatrix} 2 \\ 4 \end{pmatrix} = \lambda_1 \begin{pmatrix} 2 \\ 0 \end{pmatrix} + \lambda_2 \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

D'altro canto se la funzione obiettivo fosse stata $f(x, y) = (x - 2)^2 + y^2$ il sistema avrebbe avuto soluzione anche in presenza di vincoli i cui gradienti non sono linearmente indipendenti. Questo perché l'antigradiente della funzione obiettivo sarebbe stato comunque generabile mediante una combinazione lineare del sottoinsieme dei gradienti dei soli vincoli linearmente indipendenti (in questo caso uno solo).

Siamo in presenza di una delle particolarità dell'ottimizzazione vincolata. L'imposizione delle condizioni analitiche di ottimalità non garantisce di individuare tutti i punti di minimo locale a meno che i vincoli che definiscono la regione ammissibile non soddisfino alcune peculiari

condizioni. Poiché il problema nasce dal fatto che non sempre lo spazio definito dall'approssimazione lineare delle funzioni h_j (e come vedremo in seguito, g_j) costituisce una buona approssimazione locale di tali funzioni, occorre identificare in quali casi ciò accade. Le condizioni in esame sono dette condizioni di *qualificazione dei vincoli* e devono valere in un punto di ottimo \mathbf{x}^* , sia per i vincoli di uguaglianza sia per i vincoli di disuguaglianza che sono soddisfatti come uguaglianza in \mathbf{x}^* .

Condizione 4 Qualificazione dei vincoli *In un punto \mathbf{x}^* diciamo che sono soddisfatte le condizioni di qualificazione dei vincoli se esiste un vettore \mathbf{h} tale che $\nabla g_j(\mathbf{x}^*)^T \mathbf{h} < 0$, in corrispondenza di tutti gli indici j tali che $g_j(\mathbf{x}^*) = 0$, $\nabla h_j(\mathbf{x}^*)^T \mathbf{h} = 0$ con $j = 1, 2, \dots, h$ ed i vettori $\nabla h_j(\mathbf{x}^*)$ con $j = 1, 2, \dots, h$ sono linearmente indipendenti.*

Le condizioni di qualificazione dei vincoli sono sempre soddisfatte se:

- i gradienti dei vincoli di uguaglianza e dei vincoli attivi in \mathbf{x}^* sono fra loro linearmente indipendenti, ovvero se lo Jacobiano dei vincoli di uguaglianza e dei vincoli attivi ha rango massimo in \mathbf{x}^* ;
- se tutti i vincoli sono lineari e
- se tutti i vincoli sono convessi e la regione ammissibile contiene almeno un punto interno.

Definizione 15 *Si dice punto regolare un punto \mathbf{x}^* che soddisfa le condizioni di qualificazione dei vincoli.*

Richiamiamo infine che nel caso di funzioni convesse le condizioni introdotte nel Teorema 16 diventano *condizioni sufficienti*

Teorema 17 *Sono date una funzione convessa $f(\mathbf{x})$ ed h funzioni convesse $h_j(\mathbf{x})$, con $j = 1, \dots, h$ di classe C^1 . Condizioni sufficienti, nell'ipotesi che la matrice Jacobiana delle funzioni $h_j(\mathbf{x})$, calcolate nel punto \mathbf{x}^* , sia di rango h , affinché \mathbf{x}^* sia un minimo locale del problema con vincoli di uguaglianza è che esista $\boldsymbol{\lambda}^*$ tale che $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ sia un punto stazionario della funzione $L(\mathbf{x}, \boldsymbol{\lambda})$.*

9.1.1 Funzione obiettivo quadratica e vincoli di uguaglianza lineari

Vediamo ora come si applica il metodo lagrangiano alla minimizzazione di una funzione obiettivo quadratica *definita positiva*, soggetta a vincoli di uguaglianza lineari

$$\begin{aligned} \min \quad & f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x} \\ \text{t.c.} \quad & A \mathbf{x} = \mathbf{d} \end{aligned}$$

Per ipotesi Q è definita positiva ed A è una matrice ($h \times n$) di rango pieno con $h < n$. Poiché i vincoli sono lineari essi soddisfano le *condizioni di qualificazione*. In questo caso la funzione lagrangiana è

$$L(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x} + \boldsymbol{\lambda}^T (\mathbf{d} - A \mathbf{x})$$

e le condizioni necessarie per l'esistenza di un minimo vincolato in \mathbf{x}^* è che esista un vettore $\boldsymbol{\lambda}^*$ tale che:

$$\begin{aligned}\nabla_{\mathbf{x}}L(\mathbf{x}^*, \boldsymbol{\lambda}^*) &= Q\mathbf{x}^* - \mathbf{b} - A^T\boldsymbol{\lambda}^* = \mathbf{0} \\ \nabla_{\boldsymbol{\lambda}}L(\mathbf{x}^*, \boldsymbol{\lambda}^*) &= A\mathbf{x}^* - \mathbf{d} = \mathbf{0}\end{aligned}$$

Tali condizioni possono essere riscritte nella forma:

$$\begin{bmatrix} Q & -A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x}^* \\ \boldsymbol{\lambda}^* \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ \mathbf{d} \end{bmatrix}$$

la cui soluzione è

$$\begin{bmatrix} \mathbf{x}^* \\ \boldsymbol{\lambda}^* \end{bmatrix} = \begin{bmatrix} Q & -A^T \\ A & 0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{b} \\ \mathbf{d} \end{bmatrix}$$

9.1.2 Da vincoli di disuguaglianza a vincoli di uguaglianza

Riprendiamo ora il problema generale

$$\begin{aligned}\min \quad & f(\mathbf{x}) \\ g_j(\mathbf{x}) & \leq 0 \quad j = 1, \dots, k; \\ h_j(\mathbf{x}) & = 0 \quad j = 1, \dots, h\end{aligned}$$

Una prima tecnica per la soluzione del problema generale lo riconduce ad un problema con soli vincoli di uguaglianza mediante l'introduzione di variabili ausiliarie. Ciascun vincolo di disuguaglianza viene trasformato in un vincolo di uguaglianza mediante l'aggiunta di una variabile ausiliaria. Si passa da $g_i(\mathbf{x}) \leq 0$ a $g_i(\mathbf{x}) + \theta_i^2 = 0$. A differenza dell'analoga trasformazione che permette di passare dalla forma generale alla forma standard nell'ambito della programmazione lineare, in questo caso le variabili ausiliarie vengono elevate al quadrato. Elevare al quadrato le variabili ausiliarie elimina la necessità di introdurre le condizioni di non negatività $\theta_i \geq 0$, cioè altre disuguaglianze. Il problema diventa quindi:

$$\begin{aligned}\min \quad & f(\mathbf{x}) \\ g_j(\mathbf{x}) + \theta_j^2 & = 0 \quad j = 1, \dots, k; \\ h_j(\mathbf{x}) & = 0 \quad j = 1, \dots, h\end{aligned}$$

a cui corrisponde il modello lagrangiano:

$$L(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^k \lambda_j (g_j(\mathbf{x}) + \theta_j^2) + \sum_{j=1}^h \mu_j h_j(\mathbf{x})$$

Le condizioni necessarie perché un punto \mathbf{x} sia di minimo sono

$$\begin{aligned}\frac{\partial L}{\partial x_i} &= \frac{\partial f(\mathbf{x})}{\partial x_i} + \sum_{j=1}^k \lambda_j \frac{\partial g_j(\mathbf{x})}{\partial x_i} + \sum_{j=1}^h \mu_j \frac{\partial h_j(\mathbf{x})}{\partial x_i} = 0, \quad i = 1, 2, \dots, n \\ \frac{\partial L}{\partial \theta_j} &= 2\lambda_j \theta_j = 0, \quad j = 1, 2, \dots, k \\ \frac{\partial L}{\partial \lambda_j} &= g_j(\mathbf{x}) + \theta_j^2 = 0, \quad j = 1, 2, \dots, k \\ \frac{\partial L}{\partial \mu_j} &= h_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, h\end{aligned}$$

Tali condizioni rappresentano un sistema di $n + 2k + h$ equazioni in $n + 2k + h$ incognite, la cui (eventuale) soluzione individua i punti candidati ad essere soluzione del problema di ottimizzazione.

Si osservi che le k relazioni $2\lambda_j\theta_j = 0$, con $j = 1, 2, \dots, k$, sono relazioni di scarto complementare, in quanto impongono che il moltiplicatore λ_j sia nullo ogniqualvolta il vincolo $g_j(\mathbf{x}) \leq 0$ è soddisfatto come disuguaglianza stretta, e che il vincolo sia soddisfatto come uguaglianza ogniqualvolta λ_j è diverso da zero.

9.2 Il caso generale: le condizioni KKT

Partendo dal problema generale

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ g_j(\mathbf{x}) \leq 0 \quad & j = 1, \dots, k; \\ h_j(\mathbf{x}) = 0 \quad & j = 1, \dots, h \end{aligned} \quad (33)$$

e dal relativo modello lagrangiano:

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^k \lambda_j g_j(\mathbf{x}) + \sum_{j=1}^h \mu_j h_j(\mathbf{x})$$

Karush (1939) e Kuhn e Tucker (1951) derivarono indipendentemente alcune condizioni necessarie affinché un punto \mathbf{x} sia un punto di minimo locale. Tali condizioni sono ora note come *condizioni KKT*.

Teorema 18 *Sia dato un problema in forma generale (33), dove le funzioni $f(\mathbf{x})$, $g_j(\mathbf{x})$ e $h_j(\mathbf{x})$ sono tutte di classe C^1 . Se \mathbf{x}^* è un minimo locale e in \mathbf{x}^* valgono le condizioni di qualificazione dei vincoli di uguaglianza e di quelli di disuguaglianza attivi, allora esistono moltiplicatori di Lagrange $\boldsymbol{\lambda}^*$ e $\boldsymbol{\mu}^*$, tali che le seguenti condizioni sono soddisfatte,*

$$\begin{aligned} \frac{\partial f(\mathbf{x}^*)}{\partial x_i} + \sum_{j=1}^k \lambda_j^* \frac{\partial g_j(\mathbf{x}^*)}{\partial x_i} + \sum_{j=1}^h \mu_j^* \frac{\partial h_j(\mathbf{x}^*)}{\partial x_i} &= 0, & i = 1, \dots, n \\ g_j(\mathbf{x}^*) &\leq 0, & j = 1, \dots, k \\ \lambda_j^* g_j(\mathbf{x}^*) &= 0, & j = 1, \dots, k \\ h_j(\mathbf{x}^*) &= 0, & j = 1, \dots, h \\ \lambda_j^* &\geq 0, & j = 1, \dots, k \end{aligned}$$

Si osservi che le k relazioni $\lambda_j^* g_j(\mathbf{x}^*) = 0$, con $j = 1, 2, \dots, k$, sono relazioni di scarto complementare, in quanto impongono che $\lambda_j^* = 0$ ogniqualvolta vale $g_j(\mathbf{x}^*) > 0$, e che valga $g_j(\mathbf{x}^*) = 0$, cioè il vincolo sia attivo ogniqualvolta vale $\lambda_j > 0$. Indichiamo con $I \subseteq \{1, 2, \dots, k\}$ il sottoinsieme degli indici $1, 2, \dots, k$, che corrispondono ai vincoli di disuguaglianza attivi. Tenendo a mente queste relazioni il primo gruppo di n condizioni coincide con

$$-\nabla f(\mathbf{x}^*) = \sum_{j \in I} \lambda_j^* \nabla g_j(\mathbf{x}^*) + \sum_{j=1}^h \mu_j^* \nabla h_j(\mathbf{x}^*) \quad (34)$$

che esprime la

richiesta che nel punto di ottimo \mathbf{x}^* l'antigradiente della funzione obiettivo si possa ottenere come combinazione lineare non negativa dei gradienti dei vincoli di disuguaglianza attivi e come combinazione lineare dei vincoli di uguaglianza.

Abbiamo già incontrato questa relazione nell'ambito della programmazione lineare, ma mentre nel contesto della PL siamo sempre in grado di calcolare i valori delle variabili duali, qui è vero solo se sono soddisfatte le condizioni di qualificazione dei vincoli.

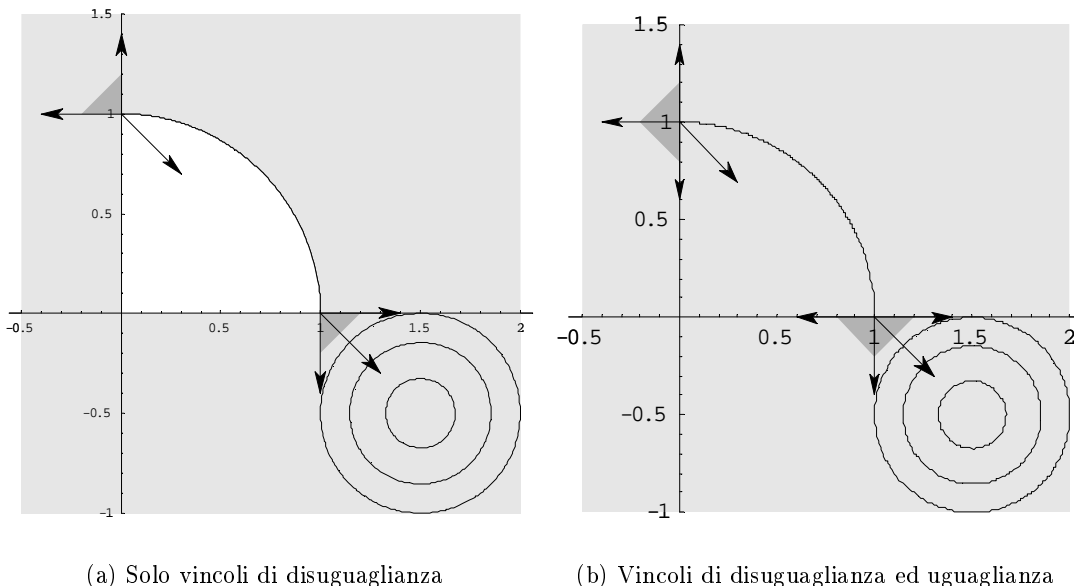


Figura 17: Condizioni di ottimalità con vincoli di disuguaglianza.

Esempio. Consideriamo il problema di Figura 17 a sinistra

$$\begin{aligned}
 \min \quad & f(x, y) = (x - 1.5)^2 + (y + 0.5)^2 \\
 & g_1(x, y) = -x \leq 0; \\
 & g_2(x, y) = -y \leq 0; \\
 & g_3(x, y) = x^2 + y^2 - 1 \leq 0.
 \end{aligned}$$

Il punto di ottimo è $(1, 0)$. In tale punto sono attivi i vincoli g_2 e g_3 , i vettori $\nabla g_2(x, y)$ e $\nabla g_3(x, y)$ sono $(0, -1)^T$ e $(2, 0)^T$, rispettivamente, e sono linearmente indipendenti. Nello stesso punto l'antigradiente della funzione obiettivo, $\nabla f(x, y)$, vale $(1, -1)^T$ ed il vettore λ^* vale quindi $(0, 1, 1/2)^T$. Come si vede l'antigradiente giace all'interno del cono generato dalle combinazioni lineari non negative dei gradienti dei vincoli attivi nel punto.

Consideriamo ora il punto *non* ottimo $(0, 1)$. In tale punto sono attivi i vincoli g_1 e g_3 , i vettori $\nabla g_1(x, y)$ e $\nabla g_3(x, y)$ sono $(-1, 0)^T$ e $(0, 2)^T$, rispettivamente, e sono linearmente indipendenti. Nello stesso punto l'antigradiente della funzione obiettivo, $\nabla f(x, y)$, vale $(3, -3)^T$ e l'unico modo di generarlo mediante combinazione lineare dei gradienti dei vincoli attivi nel punto è per mezzo del vettore $\lambda^T = (-3, 0, -3/2)^T$ non ammissibile. Come si vede l'antigradiente giace all'esterno del cono generato dalle combinazioni lineari *non negative* dei gradienti dei vincoli attivi nel punto.

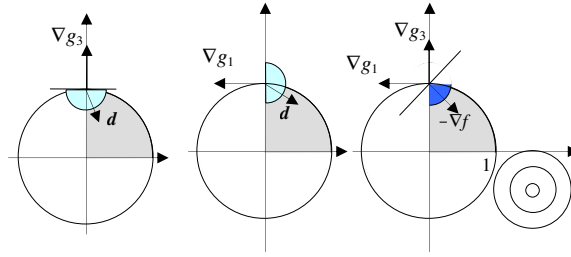


Figura 18: Condizioni di non ottimalità con vincoli di disuguaglianza.

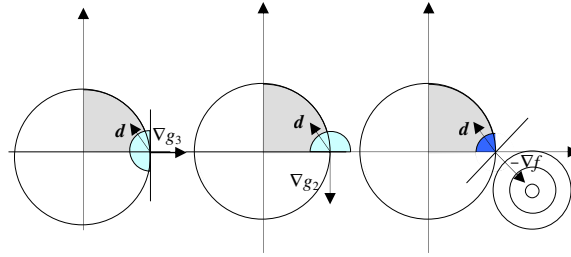


Figura 19: Condizioni di ottimalità con vincoli di disuguaglianza.

Consideriamo ora il problema di Figura 17 a destra

$$\begin{aligned} \min \quad & f(x, y) = (x - 1.5)^2 + (y + 0.5)^2 \\ & g_1(x, y) = -x \leq 0; \\ & g_2(x, y) = -y \leq 0; \\ & h_1(x, y) = x^2 + y^2 - 1 = 0; \end{aligned}$$

Il punto di ottimo è ancora $(1, 0)$ e valgono le considerazioni fatte per il problema precedente, salvo che ora risultano accettabili anche valori negativi del moltiplicatore μ_1 (che prende il posto di λ_3), visto che il relativo vincolo è di uguaglianza.

Esempio continua. In Figura 18 vediamo l'interpretazione geometrica delle condizioni di ottimalità in presenza di vincoli di disuguaglianza. Qui abbiamo rappresentato in grigio la regione ammissibile. Consideriamo il punto *non* ottimo $(0, 1)$. Poiché i vincoli di disuguaglianza sono posti in forma di \leq , i loro gradienti puntano verso l'esterno della regione ammissibile e le direzioni \mathbf{d} tali che $\nabla g_j(\mathbf{x})^T \mathbf{d} \leq 0$ sono direzioni seguendo le quali, almeno per un infinitesimo spostamento, si rimane all'interno della regione ammissibile. Tali direzioni sono evidenziate in figura mediante un semicerchio colorato per entrambi i vincoli attivi $g_1(x, y)$ e $g_3(x, y)$. Le direzioni ammissibili di spostamento sono quindi date dall'intersezione dei diversi semispazi, intersezione evidenziata nella figura di destra, ed un punto non è ottimo se tale intersezione contiene delle direzioni di discesa, come nel caso evidenziato in figura.

In Figura 19 (anche qui abbiamo rappresentato in grigio la regione ammissibile) vediamo come nel caso del punto di ottimo $(1, 0)$ invece, l'intersezione dei semispazi che definiscono direzioni ammissibili per i diversi vincoli attivi nel punto, non contenga direzioni di discesa.

Definiamo ora formalmente il concetto di *direzione ammissibile* appena introdotto.

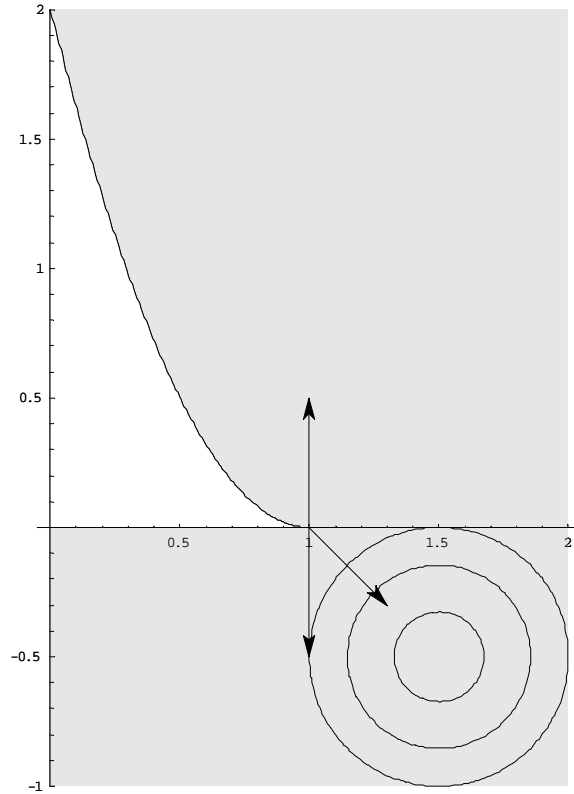


Figura 20: Esempio di gradienti di vincoli attivi linearmente dipendenti.

Definizione 16 Dato un punto ammissibile \mathbf{x} e l'insieme $I \subseteq \{1, 2, \dots, k\}$ degli indici che corrispondono a vincoli di disuguaglianza attivi in \mathbf{x} , chiamiamo insieme delle direzioni ammissibili, l'insieme

$$F(\mathbf{x}) = \{\mathbf{d} \mid \nabla h_j(\mathbf{x})^T \mathbf{d} = 0, j = 1, \dots, h; \nabla g_j(\mathbf{x})^T \mathbf{d} \leq 0, j \in I\}.$$

Esempio. Consideriamo ora il problema di Figura 20

$$\begin{aligned} \min \quad & f(x, y) = (x - 1.5)^2 + (y + 0.5)^2 \\ & g_1(x, y) = -2(x - 1)^3 + y \leq 0; \\ & g_2(x, y) = -y \leq 0. \end{aligned}$$

Il punto di ottimo è $\mathbf{x}^* = (1, 0)$. In tale punto sono attivi entrambi i vincoli g_1 e g_2 , i vettori $\nabla g_1(x, y)$ e $\nabla g_2(x, y)$ sono $(0, 1)^T$ e $(0, -1)^T$, rispettivamente, e sono linearmente *dipendenti*. In questo esempio $F(\mathbf{x}^*) = \{(d, 0)^T \mid d \in \mathbb{R}\}$. In \mathbf{x}^* l'antigradiente della funzione obiettivo, $\nabla f(x, y)$, vale $(1, -1)^T$ e non può esistere alcun vettore $\boldsymbol{\lambda}^*$.

Dal punto di vista operativo, la determinazione di un punto di ottimo per mezzo delle condizioni analitiche nel caso di un problema in forma generale (33) diventa un problema di tipo *combinatorio*. Da un punto di vista puramente teorico si tratta di generare tanti sistemi di equazioni non lineari quanti sono i sottoinsiemi di vincoli di disuguaglianza, dove per ciascun sistema, il relativo insieme di vincoli di disuguaglianza è imposto *attivo* (cioè soddisfatto come

uguaglianza). La maggior parte di tali sistemi non avrà soluzioni ammissibili nello spazio delle variabili \mathbf{x} (coinvolgendo vincoli g_j o h_j che non hanno punti in comune nelle rispettive frontiere), mentre altri non saranno ammissibili nello spazio dei moltiplicatori λ_j , poichè alcuni di essi risulteranno negativi. Naturalmente ci riferiamo a casi ideali, poichè come abbiamo visto trattando i problemi non vincolati, quasi mai è possibile determinare i punti stazionari imponendo la soluzione dei sistemi di equazioni derivanti dall'imposizione delle condizioni analitiche. Nel caso dei problemi vincolati c'è l'ulteriore complicazione derivante dall'esistenza di punti non regolari che, come evidenziato dagli esempi nelle Figure 16 e 20, possono essere punti di ottimo globale non ricavabili dalla soluzione di alcuno dei possibili sistemi di equazioni non lineari.

Richiamiamo infine che nel caso di funzioni convesse le condizioni KKT sono sufficienti

Teorema 19 *Sia dato un problema in forma generale (33), dove le funzioni $f(\mathbf{x})$, $g_j(\mathbf{x})$ e $h_j(\mathbf{x})$ sono tutte di classe C^1 . Se le funzioni $f(\mathbf{x})$, $g_j(\mathbf{x})$ e $h_j(\mathbf{x})$ sono convesse allora le condizioni KKT sono condizioni sufficienti.*

9.3 Condizioni di ottimalità del secondo ordine

Nel definire le condizioni che un punto di minimo vincolato deve soddisfare abbiamo usato fino ad ora solo condizioni del primo ordine. Le informazioni che ci provengono dall'uso del secondo termine nello sviluppo in serie di Taylor, permettono di trattare lo status delle direzioni ammissibili che sono ortogonali alla direzione del gradiente della funzione obiettivo. Iniziamo definendo l'insieme delle direzioni *critiche*.

Definizione 17 *Dato un punto ammissibile \mathbf{x}^* e due vettori di moltiplicatori $\boldsymbol{\lambda}^*$ e $\boldsymbol{\mu}^*$ che soddisfano le condizioni KKT, con $I \subseteq \{1, 2, \dots, k\}$ a denotare l'insieme degli indici che corrispondono a vincoli di disuguaglianza attivi in \mathbf{x}^* , chiamiamo insieme delle direzioni critiche, l'insieme*

$$C(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \{\mathbf{d} \in F(\mathbf{x}^*) \mid \nabla g_j(\mathbf{x}^*)^T \mathbf{d} = 0, j \in I, \text{ con } \lambda_j^* > 0\}.$$

Per definizione valgono quindi le relazioni

$$\begin{aligned} \lambda_j^* \nabla g_j(\mathbf{x}^*)^T \mathbf{d} &= 0, \quad \forall \mathbf{d} \in C(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*), \quad j \in I, \\ \mu_j^* \nabla h_j(\mathbf{x}^*)^T \mathbf{d} &= 0, \quad \forall \mathbf{d} \in C(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*), \quad j = \dots, h. \end{aligned}$$

Dalla relazione 34 possiamo perciò ricavare

$$-\nabla f(\mathbf{x}^*)^T \mathbf{d} = \sum_{j \in I} \lambda_j^* \nabla g_j(\mathbf{x}^*)^T \mathbf{d} + \sum_{j=1}^h \mu_j^* \nabla h_j(\mathbf{x}^*)^T \mathbf{d} = 0 \quad \forall \mathbf{d} \in C(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*).$$

Le direzioni critiche sono perciò ortogonali alla direzione del gradiente della funzione obiettivo.

In Figura 21 (la regione ammissibile è in grigio) vediamo il problema

$$\begin{aligned} \min \quad & f(x, y) = (x - 1.5)^2 + y^2 \\ & g_1(x, y) = -x \leq 0; \\ & g_2(x, y) = -y \leq 0; \\ & g_3(x, y) = x^2 + y^2 - 1 \leq 0; \end{aligned}$$

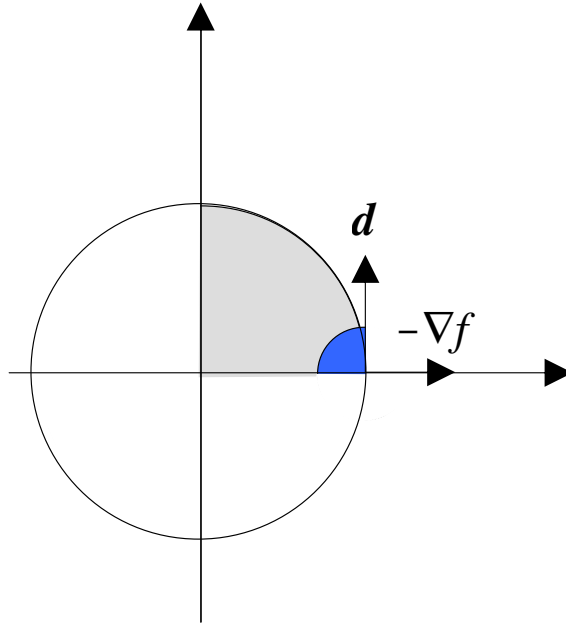


Figura 21: Direzioni critiche.

dove l'insieme delle direzioni critiche è $C(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \{(1, d)^T \mid d \geq 0\}$.

Vediamo ora delle condizioni necessarie del secondo ordine, definite nei termini delle direzioni critiche.

Teorema 20 *Sia dato un problema in forma generale (33), dove le funzioni $f(\mathbf{x})$, $g_j(\mathbf{x})$ e $h_j(\mathbf{x})$ sono tutte di classe C^2 . Se \mathbf{x}^* è un minimo locale e in \mathbf{x}^* valgono le condizioni di qualificazione dei vincoli di uguaglianza e di quelli di disuguaglianza attivi, e $\boldsymbol{\lambda}^*$ e $\boldsymbol{\mu}^*$ sono vettori di moltiplicatori che soddisfano le condizioni KKT, allora vale*

$$\mathbf{d}^T \nabla_{\mathbf{x}, \mathbf{x}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \mathbf{d} \geq 0 \quad \forall \mathbf{d} \in C(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*).$$

Nel teorema appena enunciato la matrice hessiana della funzione lagrangiana è data da

$$\nabla_{\mathbf{x}, \mathbf{x}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = H(\mathbf{x}^*) + \sum_{j \in I} \lambda_j^* \nabla^2 g_j(\mathbf{x}^*) + \sum_{j=1}^h \mu_j^* \nabla^2 h_j(\mathbf{x}^*).$$

In pratica si richiede che la matrice hessiana della funzione lagrangiana sia semi definita positiva nell'insieme delle direzioni critiche. Se tale matrice è definita positiva in questo insieme le condizioni risultano sufficienti, e non è più necessario soddisfare le condizioni di qualificazione dei vincoli.

Teorema 21 *Sia dato un problema in forma generale (33), dove le funzioni $f(\mathbf{x})$, $g_j(\mathbf{x})$ e $h_j(\mathbf{x})$ sono tutte di classe C^2 . Se \mathbf{x}^* è un punto ammissibile e $\boldsymbol{\lambda}^*$ e $\boldsymbol{\mu}^*$ sono vettori di moltiplicatori che soddisfano le condizioni KKT, e vale*

$$\mathbf{d}^T \nabla_{\mathbf{x}, \mathbf{x}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \mathbf{d} > 0 \quad \forall \mathbf{d} \in C(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$$

allora \mathbf{x}^ è un minimo locale in senso stretto del problema (33).*

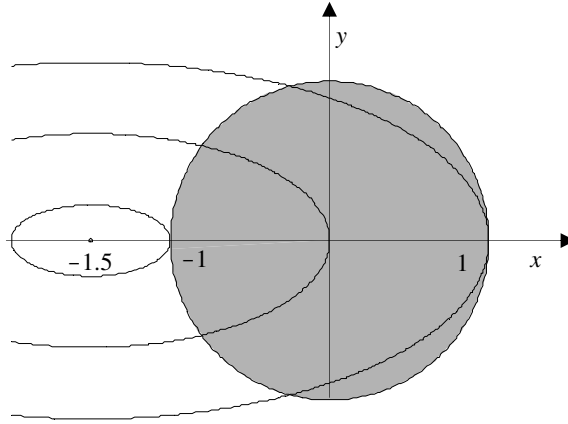


Figura 22: Esempio sulle condizioni del secondo ordine.

Esempio In Figura 22 (la regione ammissibile è in chiaro) vediamo il problema

$$\begin{aligned} \min \quad & f(x, y) = 2(x + 1.5)^2 + 10y^2 \\ & g_1(x, y) = 1 - x^2 - y^2 \leq 0; \end{aligned}$$

che ammette un punto di minimo globale in $(-1.5, 0)^T$ dove g_1 non è attivo e $\lambda^* = 0$, ed un punto di minimo locale in senso stretto $\tilde{\mathbf{x}} = (1, 0)$ dove invece $g_1(1, 0) = 0$. Verifichiamo questa seconda affermazione. In $\tilde{\mathbf{x}}$ valgono le condizioni KKT con

$$\begin{pmatrix} 4(x + 1.5) \\ 20y \end{pmatrix} - \lambda_1 \begin{pmatrix} 2x \\ 2y \end{pmatrix} = \begin{pmatrix} 10 \\ 0 \end{pmatrix} - \lambda_1 \begin{pmatrix} 2 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

con $\lambda_1^* = 5$. La matrice hessiana della funzione lagrangiana risulta

$$\nabla_{\mathbf{x}\mathbf{x}}^2 L(\tilde{\mathbf{x}}, \lambda_1^*) = \begin{pmatrix} 4 & 0 \\ 0 & 20 \end{pmatrix} - \lambda_1^* \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} = \begin{pmatrix} 4 - 2\lambda_1^* & 0 \\ 0 & 20 - 2\lambda_1^* \end{pmatrix} = \begin{pmatrix} -6 & 0 \\ 0 & 10 \end{pmatrix}$$

Il gradiente del solo vincolo attivo in $\tilde{\mathbf{x}}$ è $\nabla g_1(\tilde{\mathbf{x}}) = (2, 0)^T$ e l'insieme delle direzioni critiche ad esso ortogonali è $C(\tilde{\mathbf{x}}, \lambda_1^*) = \{(0, d)^T \mid d \in \mathbb{R}\}$. Quindi abbiamo,

$$\mathbf{d}^T \nabla_{\mathbf{x}\mathbf{x}}^2 L(\tilde{\mathbf{x}}, \lambda^*) \mathbf{d} = \begin{pmatrix} 0 \\ d \end{pmatrix}^T \begin{pmatrix} -6 & 0 \\ 0 & 10 \end{pmatrix} = \begin{pmatrix} 0 \\ d \end{pmatrix} = 10d^2 > 0.$$

Quindi nel punto $\tilde{\mathbf{x}}$ sono soddisfatte le condizioni sufficienti del secondo ordine ed esso è un punto di minimo locale in senso stretto.

9.4 Punti di sella e dualità

Richiamiamo qui brevemente alcuni risultati relativi alla teoria della dualità per problemi di programmazione non lineare. Considereremo solo il caso di problemi con soli vincoli di disuguaglianza, che nel seguito indicheremo come problema primale (P), e nel quale sia la funzione obiettivo che i vincoli sono funzioni convesse.

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ & g_i(\mathbf{x}) \leq 0 \quad i = 1, \dots, k; \end{aligned} \tag{35}$$

Il relativo modello lagrangiano è:

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{j=1}^k \lambda_j g_j(\mathbf{x})$$

Innanzitutto introduciamo la nozione di punto di sella della funzione lagrangiana.

Definizione 18 *Chiamiamo punto di sella della Lagrangiana un punto $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ tale che:*

$$L(\mathbf{x}^*, \boldsymbol{\lambda}) \leq L(\mathbf{x}^*, \boldsymbol{\lambda}^*) \leq L(\mathbf{x}, \boldsymbol{\lambda}^*) \quad \text{per ogni } \mathbf{x} \text{ e per ogni } \boldsymbol{\lambda} \geq \mathbf{0}$$

In pratica in un punto di sella, fissato $\boldsymbol{\lambda}^*$, il punto \mathbf{x}^* è un minimo della funzione $L(\mathbf{x}, \boldsymbol{\lambda}^*)$, mentre, fissato \mathbf{x}^* , il punto $\boldsymbol{\lambda}^*$ è un massimo della funzione $L(\mathbf{x}^*, \boldsymbol{\lambda})$.

Osserviamo ora che uno degli aspetti problematici delle condizioni KKT è che, da un lato sono generalmente solo condizioni necessarie, e dell'altro esse si applicano solo se le funzioni f , e g_j sono di classe C^1 , sono cioè differenziabili con continuità. Tali richieste possono essere rimosse formulando le condizioni di KKT nei termini di proprietà di punto di sella della funzione lagrangiana.

Teorema 22 *Un punto $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ con $\boldsymbol{\lambda}^* \geq \mathbf{0}$ è un punto di sella della funzione lagrangiana del problema primale se e solo se valgono le seguenti condizioni:*

$$\begin{aligned} \mathbf{x}^* &= \arg \min L(\mathbf{x}, \boldsymbol{\lambda}^*) \\ g_j(\mathbf{x}^*) &\leq 0 \quad j = 1, \dots, k \\ \lambda_j g_j(\mathbf{x}^*) &= 0 \quad j = 1, \dots, k \end{aligned}$$

Teorema 23 *Se il punto $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ con $\boldsymbol{\lambda}^* \geq \mathbf{0}$ è un punto di sella della funzione lagrangiana del problema primale allora \mathbf{x}^* è la soluzione ottima del problema primale.*

Il principale vantaggio di questi due teoremi è che forniscono condizioni necessarie per problemi di ottimizzazione che non sono né convessi né differenziabili. Qualunque tecnica di ricerca può essere usata per minimizzare $L(\mathbf{x}, \boldsymbol{\lambda}^*)$ al variare di \mathbf{x} . Naturalmente rimane il problema di conoscere il vettore dei moltiplicatori ottimi $\boldsymbol{\lambda}^*$. Dal punto di vista pratico si possono ottenere stime di $\boldsymbol{\lambda}^*$ usando tecniche iterative o risolvendo il cosiddetto problema duale. A partire dalla funzione lagrangiana introduciamo ora il modello duale del problema primale (P). In primo luogo definiamo la funzione duale

$$h(\boldsymbol{\lambda}) = \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda})$$

Va notato che, per un dato vettore di moltiplicatori $\boldsymbol{\lambda}$ la soluzione ottima $\mathbf{x}^*(\boldsymbol{\lambda})$ del problema $\min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda})$ non necessariamente soddisfa $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$, e addirittura un minimo $\mathbf{x}^*(\boldsymbol{\lambda})$ potrebbe non esistere per ogni valore di $\boldsymbol{\lambda}$.

Definiamo quindi l'insieme, D , dei vettori di moltiplicatori $\boldsymbol{\lambda}$ per i quali un minimo $\mathbf{x}^*(\boldsymbol{\lambda})$ esiste:

$$D = \{ \boldsymbol{\lambda} \mid \exists h(\boldsymbol{\lambda}), \boldsymbol{\lambda} \geq \mathbf{0} \}$$

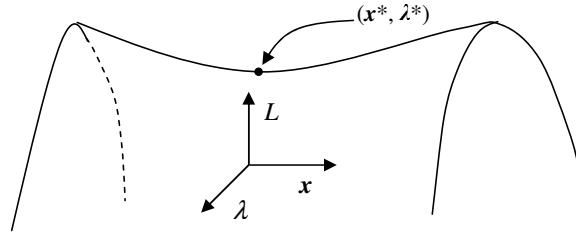


Figura 23: Punto di sella della funzione Lagrangiana

Definizione 19 *Il problema*

$$\max_{\lambda \in D} h(\lambda) = \max_{\lambda \in D} \left\{ \min_x L(\mathbf{x}, \lambda) \right\}$$

è detto duale del problema (P).

In analogia a quanto visto nell'ambito della programmazione lineare, che, per inciso, diventa una caso particolare della teoria qui presentata, è possibile ricavare relazioni di dualità debole e forte fra il problema primale ed il suo duale.

Teorema 24 *La funzione duale $h(\lambda)$ soddisfa la relazione $h(\lambda) \leq f(\mathbf{x})$ per tutti i punti \mathbf{x} che soddisfano $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$, e per tutti i $\lambda \in D$.*

Dimostrazione

$$h(\lambda) = \min_x L(\mathbf{x}, \lambda), \quad \lambda \in D \quad (36)$$

$$h(\lambda) \leq \min_x L(\mathbf{x}, \lambda), \quad \lambda \in D, \quad \mathbf{g}(\mathbf{x}) \leq \mathbf{0} \quad (37)$$

$$h(\lambda) \leq f(\mathbf{x}) + \sum_{j=1}^k \lambda_j g_j(\mathbf{x}), \quad \lambda \in D, \quad \mathbf{g}(\mathbf{x}) \leq \mathbf{0} \quad (38)$$

$$h(\lambda) \leq f(\mathbf{x}), \quad \lambda \in D, \quad \mathbf{g}(\mathbf{x}) \leq \mathbf{0} \quad (39)$$

□

Quindi la funzione duale restituisce un limite inferiore al valore della soluzione ottima del problema primale. Il limite inferiore massimo corrisponde al valore della soluzione ottima del problema duale.

Teorema 25 *Il punto $(\mathbf{x}^*, \lambda^*)$ con $\lambda^* \geq \mathbf{0}$ è un punto di sella della funzione lagrangiana del problema primale se e solo se:*

- \mathbf{x}^* è una soluzione del problema primale
- λ^* è una soluzione del problema duale
- $f(\mathbf{x}^*) = h(\lambda^*)$.

La Figura 23 illustra visivamente le condizioni di ottimalità descritte nel teorema 25.

Le implicazioni algoritmiche del teorema di dualità sono che il problema primale può essere risolto eseguendo i seguenti passi:

1. Risolvi il problema di ottimizzazione non vincolato duale $\max_{\lambda \in D} h(\lambda)$ per ricavare $\lambda^* \geq \mathbf{0}$.

2. Noto $\boldsymbol{\lambda}^* \geq \mathbf{0}$, risolvi il problema primale $\min_x L(\boldsymbol{x}^*, \boldsymbol{\lambda}^*)$ per ricavare $\boldsymbol{x}^*(\boldsymbol{\lambda}^*)$.
3. Verifica se $\boldsymbol{x}^*, \boldsymbol{\lambda}^*$ soddisfa le condizioni di KKT.

9.5 Programmazione quadratica con vincoli di disuguaglianza lineari

Un problema con soli vincoli di disuguaglianza di estremo interesse è il caso di una funzione quadratica definita positiva, soggetta a vincoli lineari di disuguaglianza

$$\begin{aligned} \min \quad & f(\boldsymbol{x}) = \frac{1}{2} \boldsymbol{x}^T Q \boldsymbol{x} - \boldsymbol{b}^T \boldsymbol{x} \\ \text{t.c.} \quad & A \boldsymbol{x} \leq \boldsymbol{d} \end{aligned}$$

Rappresenta la generalizzazione della programmazione quadratica con vincoli lineari. La soluzione ottima, \boldsymbol{x}^* , potrebbe essere all'interno o sul confine della regione ammissibile. Se \boldsymbol{x}^* si trova all'interno, allora nessun vincolo è attivo, e $\boldsymbol{x}^* = \boldsymbol{x}_0 = -Q^{-1}\boldsymbol{b}$.

Altrimenti, almeno un vincolo lineare è attivo nella soluzione ottima. Se l'insieme di vincoli attivi in \boldsymbol{x}^* è noto, allora il problema è enormemente semplificato. Supponiamo che tale insieme sia noto, allora posso rappresentare questo insieme *attivo* come: $A' \boldsymbol{x} = \boldsymbol{d}'$ e applicare la teoria di Lagrange (cfr. Sezione 9.1.1).

Nel risolvere problemi di programmazione quadratica, il lavoro più oneroso è rappresentato dall'identificazione dell'insieme dei vincoli attivi. Un metodo immediato, ma applicabile quando il numero di vincoli non è eccessivo, è quello di Theil e Van de Panne. Si tratta di identificare l'insieme dei vincoli attivi S partendo dall'insieme vuoto, $S = \emptyset$, risolvendo il problema nell'insieme S e verificando se la soluzione ottima viola qualche vincolo. Se non si violano vincoli il procedimento si arresta altrimenti si generano tutti gli insiemi S_1, S_2, \dots, S_r formati da un solo vincolo. Si risolvono gli r problemi risultanti e si verifica se esistono soluzioni che non violano alcun vincolo. In caso affermativo la migliore di esse è la soluzione ottima, altrimenti si generano tutti gli insiemi S_1, S_2, \dots, S_q formati da coppie, terne, ecc. di vincoli e si itera il processo.

E' importante richiamare che, come messo in evidenza nella Sezione 9.1.1, la tecnica restituisce non solo il valore della soluzione ottima \boldsymbol{x}^* ma anche quella dei corrispondenti moltiplicatori.

In generale l'importanza dei problemi quadratici con vincoli di disuguaglianza lineari risiede nel fatto che esso risulta il sotto problema chiave da risolvere all'interno delle moderne tecniche di risoluzione dei problemi in forma generale (cfr. Sezione 9.10).

Poiché molto raramente è possibile utilizzare direttamente le condizioni di ottimalità per individuare in modo rapido un punto stazionario o, meglio ancora, un punto di minimo, sono stati introdotti numerosi algoritmi di tipo iterativo per approssimare le soluzioni del generico problema di ottimizzazione vincolata.

9.6 Metodi con funzione di penalità

Un primo tipo di approccio, si basa sull'idea di ricondurre la soluzione di un problema vincolato a quella di un problema non vincolato. Tale approccio è di tipo sequenziale, cioè è basato sulla soluzione di una successione di problemi non vincolati, in modo tale che le soluzioni ottime convergano a quella del problema vincolato. In tali problemi si fa uso di una funzione continua

$p(\mathbf{x})$, detta di penalità, tale che $p(\mathbf{x})$ è nulla nei punti \mathbf{x} che rispettano i vincoli e maggiore di zero altrove. Consideriamo per primo un problema con soli vincoli di uguaglianza

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ h_j(\mathbf{x}) &= 0 \quad j = 1, \dots, h; \end{aligned}$$

In tali modelli una tipica scelta per la funzione di penalità è

$$p(\mathbf{x}) = \sum_{j=1}^h h_j^2(\mathbf{x})$$

Tale scelta garantisce che la funzione di penalità regolare ed in particolare sia differenziabile nei punti in cui i vincoli sono soddisfatti. Il modello con funzione di penalità diviene

$$\min q(\mathbf{x}) = f(\mathbf{x}) + \alpha \sum_{j=1}^h h_j^2(\mathbf{x}).$$

Intuitivamente, maggiore è il valore del parametro α e maggiore è la probabilità che la soluzione ottima del modello penalizzato soddisfi i vincoli di uguaglianza. La tecnica prevede quindi di risolvere una successione di problemi penalizzati per valori crescenti di α , ottenendo così una successione di punti che convergono alla soluzione ottima del problema vincolato.

In questo caso, le condizioni necessarie del primo e del secondo ordine affinché \mathbf{x}^* sia un punto di minimo del problema non vincolato sono

$$\nabla q(\mathbf{x}^*) = \nabla f(\mathbf{x}^*) + 2\alpha \sum_{j=1}^h h_j(\mathbf{x}^*) \nabla h_j(\mathbf{x}^*) = \mathbf{0},$$

e che la matrice hessiana di $q(\mathbf{x})$ in \mathbf{x}^*

$$\nabla^2 q(\mathbf{x}^*) = \nabla^2 f(\mathbf{x}^*) + 2\alpha \sum_{j=1}^h (h_j(\mathbf{x}^*) \nabla^2 h_j(\mathbf{x}^*) + \nabla h_j(\mathbf{x}^*) \nabla h_j(\mathbf{x}^*)^T)$$

sia semi definita positiva.

Si può dimostrare che facendo crescere α a infinito, la successione degli ottimi del problema penalizzato, $\mathbf{x}^*(\alpha)$, tende ad un minimo locale del problema vincolato, inoltre si ha

$$\lim_{\alpha \rightarrow \infty} 2\alpha h_j(\mathbf{x}^*(\alpha)) = \lambda_j^*$$

dove λ_j^* è il valore ottimo del moltiplicatore di Lagrange associato all' j -esimo vincolo.

Dalla condizione di ottimalità del secondo ordine si può osservare che la matrice hessiana della funzione obiettivo del problema penalizzato è data dalla somma di due parti. La prima parte è

$$\nabla^2 f(\mathbf{x}^*) + 2\alpha \sum_{j=1}^h h_j(\mathbf{x}^*) \nabla^2 h_j(\mathbf{x}^*)$$

che, per quanto appena detto, facendo crescere α a infinito, tende alla forma

$$\nabla^2 f(\mathbf{x}^*) + \sum_{j=1}^h \lambda_j^* \nabla^2 h_j(\mathbf{x}^*)$$

cioè alla matrice Hessiana della funzione Lagrangiana nel punto di ottimo.

La seconda parte è

$$\sum_{j=1}^h 2\alpha \nabla h_j(\mathbf{x}^*) \nabla h_j(\mathbf{x}^*)^T$$

che al tendere di α a infinito diventa illimitata in norma. La conseguenza di questo fatto è che, sebbene da un punto di vista teorico il metodo converga, da un punto di vista pratico l'Hessiana della funzione obiettivo penalizzata diviene sempre più malcondizionata man mano che ci si avvicina al punto ottimo \mathbf{x}^* .

Questa difficoltà può essere ovviata usando funzioni di penalità diverse, che non richiedano di far tendere α a infinito, ma in genere questo porta a perdere la differenziabilità della funzione $q(\mathbf{x})$, introducendo difficoltà di tipo diverso.

9.7 Metodi di barriera

Consideriamo ora un problema con soli vincoli di disuguaglianza

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ g_j(\mathbf{x}) \leq 0 \quad & j = 1, \dots, k; \end{aligned}$$

Dividiamo la regione ammissibile del problema nell'insieme frontiera $S_f := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{g}(\mathbf{x}) = \mathbf{0}\}$ e nell'insieme dei punti interni $S_{int} := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{g}(\mathbf{x}) < \mathbf{0}\}$. I metodi a barriera sono applicabili nell'ipotesi che l'insieme S_{int} non sia vuoto. In tali metodi si fa uso di una *funzione barriera* $v(\mathbf{x})$ per l'insieme ammissibile del problema che è continua in S_{int} , e tale che $v(\mathbf{x}) \rightarrow \infty$ quando $\mathbf{x} \rightarrow S_f$. Il modello con funzione a barriera diviene quindi

$$\min b(\mathbf{x}) = f(\mathbf{x}) + \alpha v(\mathbf{x}).$$

E' un modello non vincolato dove si è creato un effetto barriera che impedisce a un punto che si trovi in S_{int} di uscire dalla regione ammissibile. L'effetto barriera cresce al crescere del parametro α .

Contrariamente al metodo che usa le funzioni di penalità, nel metodo a barriera si lavora con punti che sono sempre ammissibili. Si può dimostrare che sotto ipotesi abbastanza blande, per valori decrescenti del parametro α la successione delle soluzioni ottime dei problemi non vincolati converge a un minimo locale del problema vincolato.

La funzione barriera più usata è la funzione logaritmica

$$v(\mathbf{x}) = - \sum_{i=1}^k \log(-g_i(\mathbf{x}))$$

Anche in questo caso il problema principale sta nel malcondizionamento della Hessiana della funzione $b(\mathbf{x})$ al decrescere di α . Un'ulteriore difficoltà è che questi metodi richiedono che il punto di partenza \mathbf{x}_0 sia ammissibile e questo non è semplice facile da ottenere.

Recentemente l'interesse per i metodi a barriera si è risvegliato grazie al fatto che si sono rivelati utili nella risoluzione di particolari problemi lineari.

9.8 Metodo del gradiente proiettivo

Questo metodo è dovuto a Rosen (1960, 1961). Consideriamo un problema con vincoli di uguaglianza lineari

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ t.c. \quad & A\mathbf{x} = \mathbf{b} \end{aligned}$$

Si parte da una soluzione iniziale \mathbf{x}' ammissibile, $A\mathbf{x}' = \mathbf{b}$, e si cerca $\mathbf{x} = \mathbf{x}' + \alpha \mathbf{d}$ (nuova soluzione migliorata). La direzione \mathbf{d} deve soddisfare condizioni di normalizzazione, $\|\mathbf{d}\| = 1$,

deve conservare l'ammissibilità, $A(\mathbf{x}' + \alpha \mathbf{d}) - \mathbf{b} = \mathbf{0}$, e soprattutto deve garantire il miglior decremento del valore della funzione obiettivo in \mathbf{x}' . La condizione di normalizzazione si può scrivere come $1 - \mathbf{d}^T \mathbf{d} = 0$, la condizione di ammissibilità si traduce in $\lambda A \mathbf{d} = \mathbf{0}$, cioè $A \mathbf{d} = \mathbf{0}$, mentre la richiesta di una direzione di massima discesa di $f(\mathbf{x})$ equivale a trovare una direzione \mathbf{d} tale che sia minimizzata la derivata direzionale $\nabla f(\mathbf{x}')^T \mathbf{d}$ in \mathbf{x}' . Da queste condizioni ricaviamo un problema di ottimizzazione vincolata in cui le nuove variabili sono \mathbf{d} .

$$\begin{aligned} \min \quad & \nabla f(\mathbf{x}')^T \mathbf{d} \\ \text{t.c.} \quad & 1 - \mathbf{d}^T \mathbf{d} = 0 \\ & A \mathbf{d} = \mathbf{0} \end{aligned}$$

Applicando la tecnica lagrangiana per un problema soggetto a vincoli di uguaglianza ricaviamo:

$$L(\mathbf{d}, \boldsymbol{\lambda}, \lambda_0) = \nabla f(\mathbf{x}')^T \mathbf{d} + \boldsymbol{\lambda}^T A \mathbf{d} + \lambda_0(1 - \mathbf{d}^T \mathbf{d})$$

Le condizioni necessarie sono

$$\begin{aligned} \nabla_{\mathbf{d}} L &= \nabla f(\mathbf{x}') + \boldsymbol{\lambda}^T A - 2\lambda_0 \mathbf{d} = \mathbf{0} \\ \nabla_{\boldsymbol{\lambda}} L &= A \mathbf{d} = \mathbf{0} \\ \nabla_{\lambda_0} L &= (1 - \mathbf{d}^T \mathbf{d}) = 0 \end{aligned}$$

Dalla prima condizione ricaviamo:

$$\mathbf{d} = \frac{1}{2\lambda_0} (\nabla f(\mathbf{x}') + \boldsymbol{\lambda}^T A).$$

Sostituendo questo risultato nella terza condizione ricaviamo:

$$1 = \frac{1}{4\lambda_0^2} (\nabla f(\mathbf{x}') + \boldsymbol{\lambda}^T A)^T (\nabla f(\mathbf{x}') + \boldsymbol{\lambda}^T A)$$

ed infine:

$$\lambda_0 = \pm \frac{1}{2} \|\nabla f(\mathbf{x}') + \boldsymbol{\lambda}^T A\|$$

Sostituendo l'espressione per λ_0 nell'equazione che definisce la direzione \mathbf{d} si ottiene

$$\mathbf{d} = \pm \frac{(\nabla f(\mathbf{x}') + \boldsymbol{\lambda}^T A)}{\|\nabla f(\mathbf{x}') + \boldsymbol{\lambda}^T A\|}$$

Fra i due versi (+/-), si sceglie quello negativo che individua una direzione di decrescita di $\nabla f(\mathbf{x}')^T \mathbf{d}$. Questo garantisce inoltre che la matrice hessiana della funzione lagrangiana rispetto a \mathbf{d} , $\nabla_{\mathbf{d}}^2 L = -2\lambda_0$, sia definita positiva, assicurandoci che la funzione di $\nabla f(\mathbf{x}')^T \mathbf{d}$ assume un valore minimo rispetto a \mathbf{d} .

Rimane da ricavare $\boldsymbol{\lambda}$. Dalla seconda condizione si ottiene:

$$A(\nabla f(\mathbf{x}') + \boldsymbol{\lambda}^T A) = \mathbf{0}$$

Così se $\mathbf{d} \neq \mathbf{0}$, si ricava:

$$A A^T \boldsymbol{\lambda} = -A \nabla f(\mathbf{x}')$$

con soluzione:

$$\boldsymbol{\lambda} = -(A A^T)^{-1} A \nabla f(\mathbf{x}')$$

La direzione \mathbf{d} , detta direzione del *gradiente proiettivo*, è così data da:

$$\mathbf{d} = -\frac{(I - A^T(AA^T)^{-1}A) \nabla f(\mathbf{x}')}{\|(I - A^T(AA^T)^{-1}A) \nabla f(\mathbf{x}')\|}$$

L'espressione trovata si può interpretare nel seguente modo. La direzione di massima decrescita per la funzione $f(\mathbf{x})$ in \mathbf{x}' è la direzione dell'antigradiente, $-\nabla f(\mathbf{x}')$. La direzione trovata è la direzione dell'antigradiente **proiettata** nell'iperpiano $A\mathbf{x} = \mathbf{b}$. La matrice $P = (I - A^T(AA^T)^{-1}A)$ è detta perciò *matrice di proiezione*.

Nella pratica, $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha\mathbf{d}$, spesso si usa la direzione non normalizzata $\mathbf{d} = -P\nabla f(\mathbf{x}')$, e si determina il passo α con i metodi già visti in precedenza (es. Armijo, Fibonacci, Sezione Aurea, ecc.).

Consideriamo ora un problema con vincoli di uguaglianza non necessariamente lineari

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ & h_j(\mathbf{x}) = 0, \quad j = 1, \dots, h \end{aligned}$$

Si linearizzano i vincoli nell'intorno della soluzione ammissibile corrente \mathbf{x}' , mediante lo sviluppo in serie di Taylor arrestato al primo ordine, ricavando

$$h_j(\mathbf{x}) = h_j(\mathbf{x}') + \nabla h_j(\mathbf{x}')^T(\mathbf{x} - \mathbf{x}'),$$

che permette quindi la seguente approssimazione lineare (tenuto conto che $h_j(\mathbf{x}') = 0$) nell'intorno di \mathbf{x}'

$$\nabla h_j(\mathbf{x}')^T \mathbf{x} - \nabla h_j(\mathbf{x}')^T \mathbf{x}' = 0, \quad j = 1, \dots, h.$$

Ora, ponendo $A = \left[\frac{\partial h(\mathbf{x}')}{\partial \mathbf{x}} \right]^T$, e $\mathbf{b} = \left[\frac{\partial h(\mathbf{x}')}{\partial \mathbf{x}} \right]^T \mathbf{x}'$, il problema linearizzato in \mathbf{x}' diventa

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{t.c.} \quad & A\mathbf{x} = \mathbf{b} \end{aligned}$$

Si applica la tecnica della matrice di proiezione $P(\mathbf{x}') = (I - A^T(AA^T)^{-1}A)$, che ora però dipende da \mathbf{x}' attraverso la matrice A , e si usa la direzione $\mathbf{d} = -P(\mathbf{x}')\nabla f(\mathbf{x}')$.

Poiché in generale, dato $\mathbf{x}_k = \mathbf{x}'$, per ogni scelta del passo $\alpha > 0$, il nuovo punto $\mathbf{x}'' = \mathbf{x}_k + \alpha\mathbf{d}$, non soddisfa necessariamente i vincoli, $\mathbf{h}(\mathbf{x}'') \neq \mathbf{0}$, occorre apportare un passo correttivo, $\mathbf{x}'' \rightarrow \mathbf{x}_{k+1}$. Tale passo è calcolato in modo che la sua proiezione in \mathbf{x}_{k+1} sia nulla, cioè

$$P(\mathbf{x}_k)(\mathbf{x}_{k+1} - \mathbf{x}'') = \mathbf{0},$$

il che equivale a chiedere che la correzione sia ortogonale alla direzione di discesa calcolata in \mathbf{x}_k , e inoltre si vuole valga $\mathbf{h}(\mathbf{x}_{k+1}) = \mathbf{0}$.

Imponendo

$$(I - A^T(AA^T)^{-1}A)(\mathbf{x}_{k+1} - \mathbf{x}'') = \mathbf{0},$$

con A calcolata in \mathbf{x}_k , si ricava

$$\mathbf{x}_{k+1} - A^T(AA^T)^{-1}A\mathbf{x}_{k+1} = \mathbf{x}'' - A^T(AA^T)^{-1}A\mathbf{x}'',$$

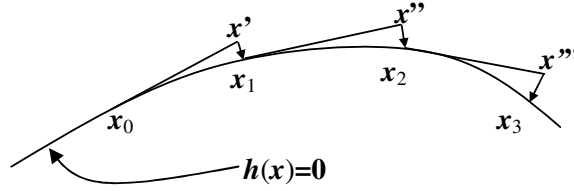


Figura 24: Passo di correzione nel gradiente proiettivo

cioè

$$\mathbf{x}_{k+1} - A^T(AA^T)^{-1}h(\mathbf{x}_{k+1}) \approx \mathbf{x}'' - A^T(AA^T)^{-1}h(\mathbf{x}''),$$

per l'approssimazione $h(\mathbf{x}) \approx A\mathbf{x} - \mathbf{b}$, che, imponendo $h(\mathbf{x}_{k+1}) = \mathbf{0}$, si riduce a

$$\mathbf{x}_{k+1} \approx \mathbf{x}'' - A^T(AA^T)^{-1}h(\mathbf{x}'').$$

La Figura 24 illustra la dinamica dell'aggiornamento mediante il passo di correzione. Il passo di correzione viene ripetuto fino a quando $h(\mathbf{x}_{k+1})$ è sufficientemente piccolo, mentre l'intero algoritmo viene fatto arrestare quando $P(\mathbf{x}')\nabla f(\mathbf{x}') \approx \mathbf{0}$.

9.9 Metodo dei lagrangiani aumentati

Questo metodo combina il classico metodo della funzione lagrangiana (cfr. Sezione 9.2) con la tecnica della funzione di penalità (cfr. Sezione 9.6). Come abbiamo visto nell'approccio lagrangiano il punto di minimo del problema vincolato coincide con un punto stazionario $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ della funzione lagrangiana che, in generale, è difficile trovare analiticamente. D'altro canto, nell'approccio con funzioni di penalità il minimo della funzione di penalità approssima il minimo vincolato ma al crescere dell'accuratezza richiesta cresce il malcondizionamento della matrice hessiana della funzione di penalità.

Nel *metodo dei moltiplicatori* (Bertsekas 1976) i due approcci vengono combinati per dare un problema non vincolato e non malcondizionato. L'idea è quella di approssimare i moltiplicatori di Lagrange.

Consideriamo il caso di un problema con soli vincoli di uguaglianza.

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ & h_j(\mathbf{x}) = 0, \quad j = 1, \dots, h \end{aligned}$$

e introduciamo la corrispondente funzione dei **langrangiani aumentati**:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \rho) = f(\mathbf{x}) + \sum_{j=1}^h \lambda_j h_j(\mathbf{x}) + \rho \sum_{j=1}^h h_j^2(\mathbf{x})$$

Con tutti i moltiplicatori $\lambda_j = 0$ la funzione \mathcal{L} si riduce alla funzione di penalità, mentre se sono noti i valori ottimi λ_j^* è possibile dimostrare (Fletcher 1987) che per ogni valore $\rho > 0$ la minimizzazione di $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \rho)$ rispetto a \mathbf{x} fornisce la soluzione ottima \mathbf{x}^* del problema.

Se il vettore $\boldsymbol{\lambda}^k$ è una buona approssimazione di $\boldsymbol{\lambda}^*$, allora è possibile approssimare l'ottimo attraverso la minimizzazione non vincolata di $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^k, \rho)$ senza richiedere valori di ρ eccessivamente grandi.

Il valore di ρ deve essere abbastanza grande da garantire che la funzione $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^k, \rho)$ abbia un minimo locale rispetto a \mathbf{x} e non solo un punto di stazionarietà.

Per comprendere la tecnica dei lagrangiani aumentati è necessario confrontare le condizioni di stazionarietà di L ed \mathcal{L} in \mathbf{x}^* .

Per \mathcal{L} :

$$\frac{\partial \mathcal{L}}{\partial x_i} = \frac{\partial f}{\partial x_i} + \sum_{j=1}^h (\lambda_j^k + 2\rho h_j) \frac{\partial h_j}{\partial x_i} = 0, \quad i = 1, \dots, n.$$

Per L :

$$\frac{\partial L}{\partial x_i} = \frac{\partial f}{\partial x_i} + \sum_{j=1}^h \lambda_j^k \frac{\partial h_j}{\partial x_i} = 0, \quad i = 1, \dots, n.$$

Il confronto indica che al tendere del punto di minimo di \mathcal{L} a \mathbf{x}^* , allora:

$$\lambda_j^k + 2\rho h_j \rightarrow \lambda_j^*$$

Questa osservazione ha suggerito il seguente schema per approssimare $\boldsymbol{\lambda}^*$. Dato un vettore di moltiplicatori approssimato $\boldsymbol{\lambda}^k$, con $k = 1, 2, \dots$, e un valore $\rho > 0$, si risolve $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^k, \rho)$ in funzione di \mathbf{x} , con un qualsiasi approccio iterativo per l'ottimizzazione non vincolata, ricavando \mathbf{x}_k^* . I valori della nuova stima di $\boldsymbol{\lambda}^*$ sono calcolati come

$$\lambda_j^{k+1} := \lambda_j^k + 2\rho h_j(\mathbf{x}_k^*)$$

In letteratura sono stati proposti anche schemi per variare iterativamente anche il parametro ρ .

9.10 SQP (Sequential Quadratic Programming)

Il metodo SQP si basa sull'applicazione del metodo di Newton per determinare il punto $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ dalle condizioni di KKT del problema di ottimizzazione vincolato. Si può dimostrare (Bazaraa et al.1993) che la determinazione del passo di Newton è equivalente alla soluzione di un problema di programmazione quadratica. Consideriamo il problema generale

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ g_i(\mathbf{x}) \quad & \leq 0 \quad i = 1, \dots, k; \\ h_j(\mathbf{x}) \quad & = 0 \quad j = 1, \dots, h \end{aligned} \tag{40}$$

ed il relativo modello lagrangiano:

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^k \lambda_j g_j(\mathbf{x}) + \sum_{j=1}^h \mu_j h_j(\mathbf{x})$$

E' data una stima $(\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k)$, con $\boldsymbol{\lambda}_k \geq 0$, $k = 1, 2, \dots$, del valore della soluzione e dei relativi moltiplicatori lagrangiani ottimi e si è ricavata la matrice hessiana della funzione lagrangiana

$$\nabla^2 L(\mathbf{x}_k) = H(\mathbf{x}_k) + \sum_{j=1}^k \lambda_j^k \nabla^2 g_j(\mathbf{x}_k) + \sum_{j=1}^h \mu_j^k \nabla^2 h_j(\mathbf{x}_k).$$

Si può dimostrare che il valore del passo di Newton \mathbf{d} che fornisce \mathbf{x}_{k+1} a partire da \mathbf{x}_k ,

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \mathbf{d}_k$$

è dato dalla soluzione del seguente problema quadratico con vincoli lineari di uguaglianza e disuguaglianza:

$$\begin{aligned} \min \phi(\mathbf{d}) &= f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 L(\mathbf{x}_k) \mathbf{d} \\ \mathbf{g}(\mathbf{x}_k) + \left[\frac{\partial \mathbf{g}(\mathbf{x}_k)}{\partial \mathbf{x}} \right]^T \mathbf{d} &\leq \mathbf{0}, \\ \mathbf{h}(\mathbf{x}_k) + \left[\frac{\partial \mathbf{h}(\mathbf{x}_k)}{\partial \mathbf{x}} \right]^T \mathbf{d} &= \mathbf{0} \end{aligned}$$

Si osservi che la soluzione del problema quadratico non solo restituisce \mathbf{d} e quindi \mathbf{x}_{k+1} , ma anche i vettori di moltiplicatori lagrangiani $\boldsymbol{\lambda}_{k+1}$ e $\boldsymbol{\mu}_{k+1}$ (cfr. Sezione 9.5). In questo modo è possibile immediatamente formulare il modello quadratico alla successiva iterazione.

La condizione di arresto è l'approssimarsi di \mathbf{d} al vettore nullo.

In pratica il modello può non convergere se si parte da punti lontani dai minimi. Di conseguenza si preferisce adottare un modello del tipo

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$$

nel quale il valore del passo α_k si ricava minimizzando una cosiddetta *funzione di merito* mediante tecniche di ricerca monodimensionale.

Una funzione di merito frequentemente adottata è

$$f_m(\mathbf{x}_k) = f(\mathbf{x}_k) + \gamma \left(\sum_{j=1}^k \max\{0, g_j(\mathbf{x}_k)\} + \sum_{j=1}^h |h_j(\mathbf{x}_k)| \right)$$

dove $\gamma = \max\{\lambda_1, \lambda_2, \dots, \lambda_k, |\mu_1|, |\mu_2|, \dots, |\mu_h|\}$. Poiché tale funzione non è differenziabile si adottano come tecniche di ottimizzazione monodimensionale metodi stabili quali la Sezione Aurea o Fibonacci.

10 Appendice

Definizione 20 Dato un punto $\mathbf{x} \in \mathbb{R}^n$ e una funzione $f(\mathbf{x}) : X \rightarrow \mathbb{R}^n$ di classe $C^1(\mathbf{x})$, definiamo vettore gradiente di $f(\mathbf{x})$ in \mathbf{x} (o semplicemente gradiente), il vettore:

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_2} \quad \cdots \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \right]^T$$

Definizione 21 Dato un punto $\mathbf{x} \in \mathbb{R}^n$ e una funzione $f(\mathbf{x}) : X \rightarrow \mathbb{R}^n$ di classe $C^2(\mathbf{x})$, definiamo matrice Hessiana di $f(\mathbf{x})$ in \mathbf{x} , la matrice:

$$\nabla^2 f(\mathbf{x}) = H(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix}^T$$

Definizione 22 Una funzione $f(\mathbf{x}) : X \rightarrow \mathbb{R}^n$, con $X \subseteq \mathbb{R}^n$, si dice lipschitziana con costante di Lipschitz L se $|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|$, $\forall \mathbf{x}, \mathbf{y} \in X$

Definizione 23 Chiamiamo curva di livello di valore C , di una funzione $f(\mathbf{x}) : \mathbb{R} \rightarrow \mathbb{R}^n$, il luogo dei punti

$$\{\mathbf{x} : f(\mathbf{x}) = C\}.$$

Proprietà 21 Dato un punto $\mathbf{x} \in \mathbb{R}^n$ e una funzione $f(\mathbf{x}) : \mathbb{R} \rightarrow \mathbb{R}^n$ di classe $C^1(\mathbf{x})$ il vettore gradiente calcolato in \mathbf{x} è ortogonale al piano tangente in \mathbf{x} alla curva di livello passante per \mathbf{x} .

Definizione 24 Data una funzione $f : X \rightarrow \mathbb{R}$ con $X \subseteq \mathbb{R}^n$, un punto $\mathbf{x}^* = [\mathbf{y}_0, \mathbf{z}_0]^T \in X$ è un punto di sella della funzione f se esiste $\varepsilon > 0$ tale che $\forall \mathbf{y} \in I(\mathbf{y}_0, \varepsilon)$ e $\forall \mathbf{z} \in I(\mathbf{z}_0, \varepsilon)$ vale la relazione

$$f(\mathbf{y}, \mathbf{z}_0) \leq f(\mathbf{y}_0, \mathbf{z}_0) \leq f(\mathbf{y}_0, \mathbf{z})$$

Lemma 1 Se l'algoritmo del gradiente viene applicato ad un problema di programmazione quadratica convesso la forma esatta della riduzione dell'errore di approssimazione tra iterazioni successive misurata con la norma pesata $\|\cdot\|_Q$ è

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q}{\|\mathbf{x}_k - \mathbf{x}^*\|_Q} = \left\{ 1 - \frac{((Q\mathbf{x}_k - \mathbf{b})^T(Q\mathbf{x}_k - \mathbf{b}))^2}{((Q\mathbf{x}_k - \mathbf{b})^T Q(Q\mathbf{x}_k - \mathbf{b}))((Q\mathbf{x}_k - \mathbf{b})^T Q^{-1}(Q\mathbf{x}_k - \mathbf{b}))} \right\}^{\frac{1}{2}}$$

Dimostrazione La riduzione dell'errore di approssimazione tra iterazioni successive del metodo del gradiente applicato a funzioni quadratiche può essere riscritta nel seguente modo

$$\frac{E_{k+1}}{E_k} = \left\{ 1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T Q \mathbf{g}_k)(\mathbf{g}_k^T Q^{-1} \mathbf{g}_k)} \right\}$$

dove abbiamo adottato la seguente notazione:

$$\mathbf{g}_k = \nabla f(\mathbf{x}_k), \quad E_k = \|\mathbf{x}_k - \mathbf{x}^*\|_Q^2, \quad E_{k+1} = \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q^2$$

e abbiamo utilizzato la relazione $\mathbf{g}_k = Q\mathbf{x}_k - \mathbf{b} = Q\mathbf{x}_k - Q\mathbf{x}^* = Q(\mathbf{x}_k - \mathbf{x}^*)$.
Iniziamo a calcolare

$$E_{k+1} = \frac{1}{2}(\mathbf{x}_k - \alpha\mathbf{g}_k - \mathbf{x}^*)^T Q(\mathbf{x}_k - \alpha\mathbf{g}_k - \mathbf{x}^*) = E_k + \frac{1}{2}\alpha^2 \mathbf{g}_k^T Q \mathbf{g}_k - \alpha \mathbf{g}_k^T (Q\mathbf{x}_k - Q\mathbf{x}^*)$$

Sostituendo ad α il valore del passo ottimo, $\frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T Q \mathbf{g}_k}$ si ricava

$$E_{k+1} = E_k + \frac{1}{2} \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{(\mathbf{g}_k^T Q \mathbf{g}_k)^2} \mathbf{g}_k^T Q \mathbf{g}_k - \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T Q \mathbf{g}_k} \mathbf{g}_k^T \mathbf{g}_k = E_k - \frac{1}{2} \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T Q \mathbf{g}_k}.$$

Ricaviamo ora

$$E_k = \frac{1}{2}(\mathbf{x}_k - \mathbf{x}^*)^T Q(\mathbf{x}_k - \mathbf{x}^*) = \frac{1}{2} \mathbf{g}_k^T Q^{-1} \mathbf{g}_k$$

dove abbiamo di nuovo fatto uso di $\mathbf{g}_k = Q(\mathbf{x}_k - \mathbf{x}^*)$ e dell'identità $QQ^{-1} = I$.
Possiamo ora raccogliere E_k a fattor comune

$$E_{k+1} = E_k - \frac{1}{2} \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T Q \mathbf{g}_k} \frac{E_k}{E_k} = E_k \left\{ 1 - \frac{(\mathbf{g}_k^T \mathbf{g}_k)^2}{\mathbf{g}_k^T Q \mathbf{g}_k \mathbf{g}_k^T Q^{-1} \mathbf{g}_k} \right\}$$

da cui segue l'espressione cercata. \square

Definizione 25 Una matrice T è ortogonale se valgono le relazioni

$$TT^T = T^T T = I.$$

In altre parole, l'inversa di una matrice ortogonale è la sua trasposta. Vale inoltre $\det T = \det T^T = \pm 1$.

Proprietà 22 Sia data una matrice Q simmetrica e definita positiva. Q può essere riscritta mediante la seguente decomposizione spettrale

$$Q = \sum_{i=1}^n \lambda_i t_i t_i^T,$$

dove $\lambda_1, \lambda_2, \dots, \lambda_n$ sono gli n autovalori di Q e t_1, t_2, \dots, t_n sono i corrispondenti n autovettori. In termini matriciali tale decomposizione può essere riscritta come

$$Q = T \Lambda T^T$$

dove

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n), \quad T = [t_1 | t_2 | \dots | t_n]$$

e la matrice T è una matrice ortogonale. In altre parole, la matrice Q può essere diagonalizzata nel seguente modo

$$T^{-1}QT = \Lambda.$$

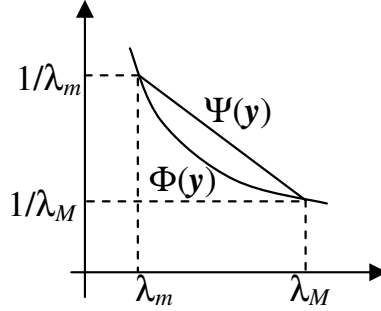


Figura 25: Andamento dei grafici di $\Phi(\mathbf{y})$ e $\Psi(\mathbf{y})$.

Teorema 9 Data una matrice Q definita positiva vale la seguente relazione

$$\frac{(\mathbf{x}^T \mathbf{x})^2}{(\mathbf{x}^T Q \mathbf{x})(\mathbf{x}^T Q^{-1} \mathbf{x})} \geq \frac{4\lambda_m \lambda_M}{(\lambda_m + \lambda_M)^2}$$

valida per ogni $\mathbf{x} \in \mathbb{R}^n$ e dove λ_M e $\lambda_m > 0$ sono l'autovalore massimo e minimo, rispettivamente, di Q .

Dimostrazione Si consideri la decomposizione spettrale di Q , $Q = T\Lambda T^T$ e si effettui la trasformazione $T^{-1}\mathbf{x} = \mathbf{z}$, cioè $\mathbf{x} = T\mathbf{z}$. Possiamo riscrivere

$$\frac{(\mathbf{x}^T \mathbf{x})^2}{(\mathbf{x}^T Q \mathbf{x})(\mathbf{x}^T Q^{-1} \mathbf{x})} = \frac{(z^T T^T T z)^2}{(z^T T^{-1} Q T z)(z^T T^{-1} Q^{-1} T z)}$$

che, sfruttando le proprietà delle matrici ortogonali si riduce a

$$\frac{(z^T z)^2}{(z^T \Lambda z)(z^T \Lambda^{-1} z)} = \frac{(\sum_{i=1}^n z_i^2)^2}{(\sum_{i=1}^n \lambda_i z_i^2)(\sum_{i=1}^n z_i^2 / \lambda_i)}$$

Dividendo per $(\sum_{i=1}^n z_i^2)^2$ otteniamo

$$\frac{1}{(\sum_{i=1}^n \lambda_i (z_i^2 / \sum_{i=1}^n z_i^2))(\sum_{i=1}^n (z_i^2 / \sum_{i=1}^n z_i^2) / \lambda_i)}$$

Effettuiamo ora un cambiamento di variabili definendo $y_i = \frac{z_i^2}{\sum_{i=1}^n z_i^2}$, con $0 \leq y_i \leq 1$ e $\sum_{i=1}^n y_i = 1$, e riscriviamo l'espressione come

$$\frac{1}{(\sum_{i=1}^n \lambda_i y_i)(\sum_{i=1}^n y_i / \lambda_i)} = \frac{1 / \sum_{i=1}^n \lambda_i y_i}{\sum_{i=1}^n y_i / \lambda_i} = \frac{\Phi(\mathbf{y})}{\Psi(\mathbf{y})}$$

Le funzioni $\Phi(\mathbf{y})$ e $\Psi(\mathbf{y})$ sono ottenute da combinazioni lineari di λ_i e $1/\lambda_i$, e ogni λ_i e $1/\lambda_i$ può essere generato come combinazione lineare convessa dei valori y_i . In Figura 25 si vede l'andamento delle funzioni $\Phi(\mathbf{y})$ e $\Psi(\mathbf{y})$ per fissati valori degli autovalori λ_i . La funzione $\Phi(\mathbf{y})$ può essere minorata dall'iperbole passante per i punti $(\lambda_m, 1/\lambda_m)$ e $(\lambda_M, 1/\lambda_M)$, mentre la funzione $\Psi(\mathbf{y})$ può essere maggiorata dal segmento di retta compreso fra gli stessi punti. Possiamo quindi scrivere

$$\frac{\Phi(\mathbf{y})}{\Psi(\mathbf{y})} \geq \min_{\lambda_m \leq \lambda \leq \lambda_M} F(\lambda) = \min_{\lambda_m \leq \lambda \leq \lambda_M} \frac{1/\lambda}{(\lambda_m + \lambda_M - \lambda)/\lambda_m \lambda_M}$$

Derivando ed uguagliando a zero $F(\lambda)$ si perviene all'espressione

$$\frac{1}{\lambda_m \lambda_M \lambda} = \frac{\lambda_m + \lambda_M - \lambda}{\lambda_m \lambda_M \lambda^2} \quad \text{da cui si ricava} \quad \lambda^* = \frac{\lambda_m + \lambda_M}{2},$$

che sostituito in $F(\lambda)$ ci dà $F(\lambda^*) = \frac{4\lambda_m \lambda_M}{(\lambda_m + \lambda_M)^2}$, il che dimostra la tesi. \square

Riferimenti bibliografici

- [1] Francesco Maffioli. *Elementi di Programmazione matematica*. Volume secondo. Ed. Masson, 1991.
- [2] Jan A. Snyman. *Practical Mathematical Optimization*. Ed. Springer, 2005.
- [3] J. Nocedal, S.J. Wright. *Numerical Optimization*. Ed. Springer, 2006.