

Esame di Programmazione su Architetture Parallele

Metodo del simplesso per la risoluzione della programmazione lineare

Belliato Riccardo (mat. 142652)

Simone Tomada

2022-08-04

Abstract

In questa relazione si propone una implementazione del metodo del simplesso a due fasi in CUDA per la risoluzione dei problemi di programmazione lineare in forma canonica.

Dopo una breve descrizione dell'algoritmo, seguirà la discussione su alcune scelte implementative.

Infine verranno valutate performance e scalabilità della soluzione proposta confrontando i tempi di esecuzione dell'algoritmo su istanze a dimensione crescente generate casualmente.

Contents

Introduzione al metodo del simplesso	1
Problemi di ottimizzazione lineare	1
Forma canonica e forma standard	2
Metodo del simplesso a due fasi	3
Scelte implementative e algoritmi utilizzati	6
Schema di parallelizzazione	6
Gestione della memoria	7
Ricerca del pivot e test di ottimalità	7
Eliminazione di Gauss	7
Aggiornamento della tabella	7
Risultati sperimentali	7

Introduzione al metodo del simplesso

Problemi di ottimizzazione lineare

Nell'ambito della Ricerca Operativa (Operations Research) uno dei principali argomenti è la cosiddetta **ottimizzazione lineare**, ossia lo studio di una classe di problemi del tipo:

$$\begin{aligned} &\min / \max c^T x \\ &\text{subject to} \\ &Ax \preceq b \end{aligned}$$

con $x \in \mathbb{R}^n$, $c \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^{1 \times m}$.

In altre parole, si vogliono trovare dei valori per le componenti del vettore x tali da massimizzare (o minimizzare) il valore di una *funzione lineare* (detta **funzione obiettivo**, mentre il vettore c è chiamato **vettore dei costi**), dati una serie di vincoli espressi nel sistema di disequazioni lineari $Ax \preceq b$ (dove A è detta **matrice dei vincoli** e b **vettore dei termini noti**).

Questa classe di problemi permette di modellare un gran numero di situazioni reali in molteplici ambiti (ottimizzazione dei costi, creazione di orari, *etc.*), oltre che alcuni problemi NP-hard come il *Knapsack* o il *Vertex Cover*.

Dal punto di vista dell'algebra lineare, un problema in n variabili non è altro che uno spazio \mathbb{R}^n , la funzione obiettivo è una retta nello spazio, mentre i vincoli definiscono un *poliedro* nello spazio.

Utilizzando i teoremi e le tecniche dell'algebra lineare è stato possibile creare degli algoritmi per risolvere i problemi di ottimizzazione, come il **simplex**, i quali si prestano molto bene ad essere parallelizzati (in quanto operano su matrici).

Problemi risolvibili, non risolvibili, illimitati

Dato un problema di ottimizzazione, questo può essere di tre tipi:

- avere una o più soluzioni ammissibili (*feasible*): ossia esistono uno o più vettori che moltiplicati per il vettore dei costi assegnano alla funzione obiettivo il valore massimo (o minimo possibile) e tutti i vincoli sono rispettati,
- non risolvibili (*infeasible*): se non esistono soluzioni
- illimitati (*unbounded*): se per una o più componenti della soluzione è possibile aumentarne (o ridurne) il valore all'infinito senza mai violare i vincoli

Forma canonica e forma standard

I problemi di massimizzazione possono essere convertiti in problemi di minimizzazione (e viceversa), così come è possibile manipolare le singole disequazioni. Queste operazioni servono a riportare i problemi in una forma che permetta di utilizzarli da parte dei solver. In particolare vengono utilizzate la forma *canonica* e la forma *standard*.

Un problema (di massimizzazione) è in forma canonica se è nella forma

$$\begin{aligned} &\max c^T x \\ &\text{subject to} \\ &Ax \leq b \\ &x \geq 0 \end{aligned}$$

Considerando che qualsiasi disequazione nella forma $\alpha x \leq y$ può essere convertita in una equazione equivalente $\alpha x + \delta = b$ definiamo la forma standard di un problema in forma canonica

$$\begin{aligned} &\max d^T x \\ &\text{subject to} \\ &A'x = b \\ &x \geq 0 \end{aligned}$$

con $A' = (A|I) \in \mathbb{R}^{m \times (n+m)}$ e $d = (c|0) \in \mathbb{R}^{n+m}$

In altre parole aggiungiamo una nuova variabile al problema (detta variabile *slack*) per ogni disequazione.

La forma standard è quella che viene utilizzata dagli algoritmi di soluzione.

Metodo del simplesso a due fasi

Il metodo del simplesso è un'algoritmo per risolvere i problemi di ottimizzazione lineari.

Si basa sul teorema secondo il quale le soluzioni (dette *soluzioni di base*) di un qualsiasi problema in forma standard sono i *vertici* del poliedro costruito sui vincoli.

Ogni vertice è individuato da una o più *soluzioni di base ammissibili*, cioè tutti i valori delle variabili in base sono maggiori o uguali a 0, mentre il valore delle altre variabili è 0. Se in base sono presenti uno o più valori uguali a 0, si dice che la base è *degenere*.

Esistono diverse implementazioni del metodo del simplesso, per questo progetto verrà implementato il cosiddetto *simplesso a due fasi con il metodo del tableau*

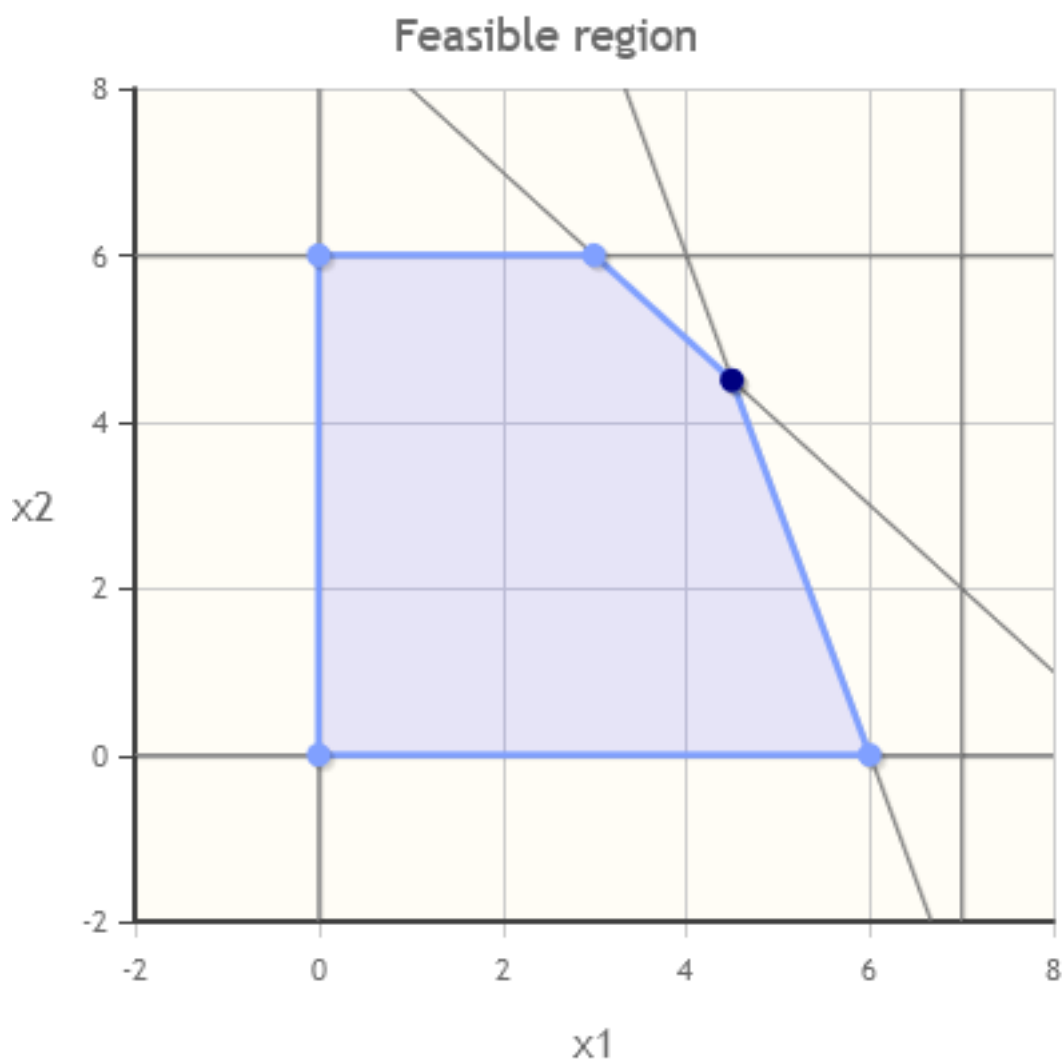


Figure 1: Esempio di problema con due variabili. Il vertice evidenziato in blu è la soluzione

Cos'è il tableau

Il tableau di un simplesso è la struttura dati su cui viene eseguito l'algoritmo del simplesso.

Si tratta di una matrice composta nel seguente modo:

	x_1	x_2	s_1	s_2	s_3	RHS	
Vettore dei costi	-3	-2	0	0	0	0	Valore funzione obiettivo
	2	1	1	0	0	10	
	1	1	0	1	0	8	
	1	0	0	0	1	4	
	Matrice dei vincoli					Termini noti	

Nella prima riga della matrice (R_0) è presente il vettore dei costi (coefficienti della funzione obiettivo) con i segni invertiti. La base viene memorizzata in un vettore a parte.

Algoritmo

1. Controllo se tutti i coefficienti del vettore dei costi sono maggiori o uguali a 0. Se sì, allora il mio problema è risolto e procedo a ricavare la soluzione, altrimenti proseguo
2. Scansiono il vettore dei costi alla ricerca della *posizione* con l'elemento minimo del vettore (C_p). Questa sarà la variabile da far entrare in base¹
3. Costruisco il *vettore degli indicatori* dividendo il vettore dei termini noti con il vettore colonna della variabile che entra in base
4. Trovo la posizione dell'elemento minimo nel vettore degli indicatori (R_p). Questa operazione serve a individuare quale variabile deve uscire dalla base.

Infatti, essendo il tableau una matrice di rango pieno, la base sarà sempre composta da tante colonne quanti sono i vincoli e la scelta della variabile uscente dipenderà unicamente dal valore dell'indicatore (al valore dell'indicatore più basso corrisponde al maggiore incremento possibile della funzione obiettivo al passo successivo). Per questo motivo andando a selezionare l' i -esima riga del tableau si andrà a sostituire l' i -esima variabile contenuta nel vettore della base con la variabile entrante.

5. Individuo il *pivot* $p = R_p \cap C_p$.
6. Siccome una variabile può essere in base solo se nel suo vettore colonna è presente un unico 1, mentre tutti gli altri componenti sono uguali a 0, procedo ad aggiornare il tableau tramite operazioni di Gauss:

- $R_p = \frac{1}{p} R_p$
- $R_i = R_i - \frac{C_p[i]}{p} R_p \quad \forall i = 1, \dots, n \wedge i \neq p$

Con $C_p[i]$ l'elemento a riga i del vettore colonna C_p .

7. Ritorno al punto 1.
8. Una volta trovato l'ottimo ricavo la soluzione nel seguente modo:
 - Se una variabile è l' i -esima variabile di base leggo il valore di quest'ultima dall' i -esima riga del vettore dei termini noti
 - altrimenti il valore di quella variabile è pari a 0

Inoltre leggo il valore ottimo della funzione obiettivo dal primo elemento dell'ultima colonna

¹Si noti che è possibile unire queste prime due fasi in una sola: una volta trovato l'elemento minimo è sufficiente verificare che sia < 0 , se sì allora proseguo, altrimenti termino il ciclo (se il minimo è ≥ 0 , allora anche tutti gli altri elementi del vettore lo sono)

Complessità Considerando che l'algoritmo non fa altro che esplorare le soluzioni di base, l'algoritmo esplora al massimo $O\binom{n+m}{m}$ (con n numero di variabili e m il numero di vincoli).

Per cui nei casi peggiori il simplesso può raggiungere complessità esponenziali. Tuttavia nel corso degli anni si è visto EMPIRICAMENTE come questa complessità venga raggiunta in pochissimi casi specifici, anzi l'algoritmo si è rivelato molto efficiente su problemi “reali”, dove in media raggiunge una complessità nell'ordine di $O(mn)$, anche a confronto con altri algoritmi di soluzione (polinomiali per costruzione) come il metodo dei punti interni.

Da queste considerazioni si può concludere anche come la complessità del simplesso dipenda molto dall'istanza del problema fornito in input.

Fase 1

In questa fase viene risolto un problema ausiliario

$$\begin{aligned} \max & 0x - 1y \\ \text{subject to} & \\ & Ax + y = b \\ & x \geq 0, y \geq 0 \end{aligned}$$

In altre parole si tratta di inserire m nuove variabili al problema (una per vincolo). Queste sono dette *variabili artificiali*.

Questa fase ha tre scopi:

- capire se il problema iniziale è infeasible: in questo caso il valore finale della funzione obiettivo è negativo
- capire se il problema ha una o più soluzioni degeneri: in questo caso la base ottenuta alla fine di questa fase contiene una o più variabili artificiali
- se non si presenta nessuno di questi due casi, fornire in input alla fase successiva una base ammissibile di partenza

Ad esempio per il problema

$$\begin{aligned} \max & -x_1 + x_2 \\ \text{subject to} & \\ & x_1 + x_2 \leq 1 \\ & 3x_1 + 2x_2 \leq 6 \end{aligned}$$

Il corrispondente problema ausiliario sarà

$$\begin{aligned} \max & 0x_1 + 0x_2 + 0s_1 + 0s_2 - a_1 - a_2 \\ \text{subject to} & \\ & x_1 + x_2 + s_1 + a_1 = 1 \\ & 3x_1 + 2x_2 + s_2 + a_2 = 6 \end{aligned}$$

Si avrà quindi il seguente tableau

x_1	x_2	s_1	s_2	a_1	a_2	RHS
0	0	0	0	1	1	0
1	1	1	0	1	0	1
3	2	0	1	0	1	6

Con base di partenza $B = \{a_1, a_2\}$

Si noti che nonostante a_1 e a_2 siano in base, la funzione obiettivo è espressa in funzione di esse (le variabili di base devono avere come coefficiente della funzione obiettivo 0, altrimenti l'algoritmo non è corretto).

Per risolvere questo problema si sottrae alla prima riga tutte le altre (eliminazione di Gauss):

$$R_0 = R_0 - R_i \quad \forall i = 1, \dots, m$$

x_1	x_2	s_1	s_2	a_1	a_2	RHS
-4	-3	-1	-1	0	0	-7
1	1	1	0	1	0	1
3	2	0	1	0	1	6

In questo modo la funzione obiettivo non è più espressa in termini delle variabili di base e si può proseguire a risolvere il tableau.

Alla fine di questa fase si otterrà una base ammissibile di partenza per la fase successiva (a meno di problema degenerare o infeasible).

Fase 2

Terminata la prima fase si eliminano le variabili artificiali dal tableau e si sostituisce la funzione obiettivo con quella del problema originale.

Durante questa operazione può capitare lo stesso problema visto in fase 1, ossia la funzione obiettivo è espressa in funzione di qualche variabile di base. Anche in questo caso la soluzione è quella di procedere a una serie di eliminazioni di Gauss:

$$R_0 = R_0 - \frac{R_0[B[i]]}{R_i[B[i]]} R_i \quad \forall i = 1, \dots, m$$

Con $R_0[B[i]]$ e $R_i[B[i]]$ rispettivamente il valore della prima riga e della i -esima riga nella colonna della i -esima variabile di base.

Si noti come questo passo sia una generalizzazione di quanto già visto in precedenza per la fase 1.

A questo punto si può procedere con l'esecuzione dell'algoritmo del simplesso, con l'aggiunta di un passaggio tra il passo 2. e il passo 3. per controllare se il problema è unbounded: se tutti i coefficienti del vettore colonna che sta per entrare in base sono ≤ 0 allora il problema è unbounded su quella variabile, altrimenti si può proseguire con la ricerca del pivot.

Scelte implementative e algoritmi utilizzati

Schema di parallelizzazione

Per introdurre parallelismo si è deciso di parallelizzare i singoli passi che lo compongono, mantenendo invece sequenziale la struttura dell'algoritmo in sé.

Così facendo non si prevede una riduzione della complessità teorica dell'algoritmo, ma ottenere un miglioramento empirico sui tempi necessari a risolvere modelli anche molto complessi (in quanto si va a ridurre la complessità delle singole operazioni sui vettori e sulle matrici).

Gestione della memoria

Implementando gli algoritmi in CUDA, merita particolare risalto il modo in cui viene organizzata la memoria:

- i dati del problema originale sono memorizzati in memoria host, mentre il tableau viene costruito e memorizzato nella memoria globale della scheda video
- i dati del problema originale vengono memorizzati utilizzando memoria *page-locked*, questo per permettere i trasferimenti di memoria in parallelo con i CUDA Stream (e quindi ottimizzare la fase di costruzione del tableau in global memory)
- il vettore della base è memorizzato come memoria *mapped page-locked*, in modo da potervi accedere sia lato host che lato GPU in qualsiasi momento senza dover gestire i trasferimenti
- siccome l'algoritmo opera spesso su singole righe o colonne è fondamentale memorizzare i dati in memoria globale in modo da minimizzare il numero di accessi strided o disallineati. Per questo motivi si è deciso il seguente schema di memorizzazione (visibile in figura 2):
 - la matrice del tableau viene allocata con *cudaMallocPitch* ed è memorizzata in maniera *lineare* per **colonne**. In altre parole l'implementazione lavorerà sul tableau trasposto. Questo perché l'algoritmo tende ad accedere più alle singole colonne che alle singole righe. Infatti così facendo sarà necessario implementare solo un kernel che accede alla memoria in maniera strided (laddove la classica memorizzazione per righe avrebbe richiesto tre),
 - il vettore dei termini noti non viene memorizzato in fondo al tableau, ma nella prima riga della matrice. In questo modo nel passaggio dalla fase 1 alla fase 2 non è necessario spostare la colonna nella matrice, ma è sufficiente troncare il numero di righe della matrice da una variabile,
 - il vettore dei costi non viene memorizzato nella matrice del tableau, ma in un vettore a parte e il valore della funzione obiettivo è nel primo elemento di questo vettore.

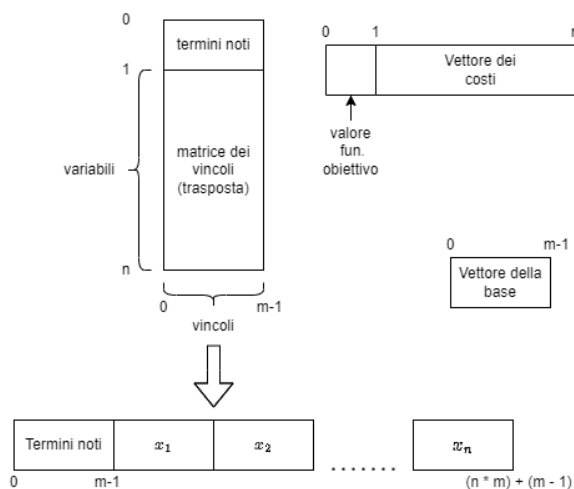


Figure 2: Schema di memorizzazione in memoria globale

Ricerca del pivot e test di ottimalità

Eliminazione di Gauss

Aggiornamento della tabella

Risultati sperimentali