

Relazione del Progetto di Parallel Computing

Filtro di Bloom

Nome: Guido

Cognome: Ciardi

Matricola: 7090798

Email: guido.ciardi@stud.unifi.it

Per questo progetto sono state implementate una versione **sequenziale** ed una **parallela** inerenti alla creazione del filtro di Bloom.

Per la loro **implementazione** è stato utilizzato il **linguaggio Python** e, per la parallelizzazione, è stata utilizzata la libreria **Joblib**.

Descrizione del filtro di Bloom:

Un **filtro di Bloom** si definisce come un vettore di n bit che consente di **filtrare** un insieme di dati considerato. Il suo **obiettivo** è quello di filtrare l'insieme di dati considerato selezionando tra questi solo le tuple che soddisfano il criterio di selezione utilizzato.

Essendo un **metodo probabilistico** usato per testare se un elemento appartiene a un insieme o meno è possibile che si ottengano tra i risultati dei **falsi positivi**, ovvero degli elementi che non devono far parte dell'insieme filtrato dei dati ma che il filtro non è riuscito ad escludere.

Il filtro di Bloom si compone di:

1. Un **vettore di n bit**, inizialmente posti tutti a zero;
2. Una collezione di **funzioni hash** h_1, h_2, \dots, h_r , ciascuna delle quali trasforma i valori chiave in un intero compreso tra 0 ed $(n - 1)$;
3. Un **insieme S** di m valori dell'attributo chiave;

Dati questi componenti, il filtro cerca di selezionare tutti gli elementi del flusso di dati aventi valore del campo chiave nell'insieme S , e cerca di rifiutare la maggior parte degli elementi del flusso il cui valore di tale campo non è all'interno di S .

L'**idea** è quella di valutare le r funzioni hash sugli elementi dell'insieme S , andando successivamente a porre ad 1 il valore del bit del vettore prima citato, corrispondente al risultato dell'applicazione della funzione hash attualmente considerata.

Per verificare l'appartenenza di un elemento all'insieme S , si applicano le r funzioni hash a tale elemento e si verifica se i bit del vettore, nelle posizioni indicate dai valori risultanti, sono posti a 0 o ad 1.

- Se tutti i bit nelle posizioni corrispondenti nel vettore hanno valore 1, il dato viene fatto **passare** dal filtro;
- Se invece sono presenti uno o più bit posti a 0, il dato non può essere in S e si può **scartare**.

Filtro di Bloom con implementazione sequenziale:

Nella versione sequenziale del programma è stata definita inizialmente la dimensione **n** del vettore di bit usato per il filtraggio, il quale viene successivamente creato col nome di **“bloomBitVector”** ed inizializzato con tutti i suoi bit a 0 (con l’uso di un ciclo for).

In questo contesto sono state utilizzate come **dati** una serie di stringhe. Più in particolare è stato scelto di considerare l’insieme S come l’insieme delle parole più frequenti della lingua inglese, mentre come insieme da filtrare un insieme di altre parole inglesi, contenente però 14 delle parole presenti anche all’interno di S.

Entrambi gli insiemi di parole sono stati passati alle liste del programma `set_S` ed `extended_set` mediante i due rispettivi file **“most_uses_english_words.txt”** e **“parole_da_filtrare.txt”** all’interno dei quali erano contenuti.

Nel programma sono state poi definite 7 funzioni hash (da h_0 ad h_6), che prendono come parametro di input la parola su cui devono operare ed utilizzano logiche diverse nelle loro definizioni. Per tutte queste funzioni però è stato applicato al loro risultato il modulo n per far sì che fosse rispettato il vincolo per cui il loro risultato fosse compreso tra i valori 0 ed $(n-1)$.

Successivamente, per ogni parola presente all’interno dell’insieme S (**set_S**) iniziale, si pongono i rispettivi bit del filtro al valore 1 nelle posizioni calcolate dalle funzioni hash precedentemente definite mediante delle istruzioni con la forma illustrata di seguito:

$$\text{bloomBitVector}[h_i(\text{set_S}[i])] = 1$$

(dove nel codice al posto della “i” c’è il numero della corrispondente funzione hash).

In questo modo viene così conclusa quella che può essere definita come **fase di inizializzazione (e/o gestione) del filtro**.

Dopo questa prima fase si passa poi alla **fase di filtraggio** vero e proprio.

Questa è stata definita nella funzione **“matchingCheck”**, che prende come parametri di input la stringa (parola) da considerare ed il vettore di bit (**bloomBitVector**) da usare per il filtraggio. Nel caso in cui il bit del vettore che si trova alla posizione calcolata applicando la funzione hash attualmente considerata sulla stringa abbia valore 0, si passa all’iterazione successiva del ciclo (e si considera dunque la parola successiva dell’insieme). Se invece tutte le funzioni hash calcolate sulla parola attuale portano in posizioni del vettore di bit che hanno valore pari ad 1, verrà stampata tale stringa come risultato attuale.

Questo procedimento consente dunque di trovare tutte le **parole comuni** sia al primo che al secondo insieme e le riproduce in output come illustrato nell'esempio di esecuzione sotto riportato.

Esempio di esecuzione:

```
Corrispondenze trovate:  
ground  
mind  
wonder  
hot  
come  
did  
my  
sound  
no  
most  
number  
who  
over  
know
```

Filtro di Bloom con implementazione parallela:

L'**implementazione parallela** risulta dal punto di vista strutturale la stessa di quella sequenziale, con l'aggiunta però dell'utilizzo della libreria **Joblib di Python**.

Inizialmente è stata aggiunta una variabile **"numThread"**, nella quale viene definito il numero di threads che si vogliono usare per parallelizzare il codice, mentre la libreria Joblib è stata utilizzata in **3 parti** del programma:

1. Il **primo punto** in cui questa si incontra nel codice è all'interno della **funzione hash h5**. La logica di questa funzione sfrutta la codifica Unicode dei caratteri della stringa passata come parametro in input e memorizza in una lista **"result"** i loro valori. Successivamente si calcola la corrispondente **matrice di Vandermonde di ordine n** usando la funzione **"vander"** del modulo **numpy** di Python ed effettua la somma dei valori interni a tale matrice, restituendola come risultato.

Questa funzione può essere parallelizzata usando la libreria Joblib per la costruzione della lista result. In questo caso la funzione su cui si parallelizza è la funzione “**ord()**” che, usando il for parallelizzato mediante “**Parallel**”, viene applicata a tutti i caratteri della parola attualmente considerata.

2. Il **secondo punto** riguarda la **gestione parallela dell’array di bit** mediante l’uso delle funzioni hash precedentemente definite. Mentre nel codice sequenziale ciascuna funzione hash era applicata internamente ad un ciclo e veniva eseguita in maniera sequenziale, è stato pensato in questo caso di parallelizzare creando un **ciclo parallelo per ognuna** delle 7 funzioni hash a disposizione. Ciascun ciclo parallelo applica una funzione hash ad ogni parola dell’insieme considerato, e restituisce come output un vettore (lungo come quello di input) contenente i risultati delle varie funzioni hash. Questa strategia consente di parallelizzare l’inizializzazione dell’array di bit sfruttando il **multiprocessing**.
3. La **terza parte** del programma su cui si può agire per parallelizzare è quella relativa al **filtraggio** vero e proprio dei dati. In questo caso si definisce un ciclo parallelo che consente di considerare ogni parola del set da filtrare (“**extended_set**”), applicando a ciascuna la funzione di filtraggio “**matchingCheck**”. Siccome in questa avviene il controllo del valore dei bit di controllo per la parola considerata, parallelizzando preventivamente sul ciclo che considererà le varie parole da filtrare si riesce a velocizzare l’esecuzione del programma.

Confronto delle prestazioni:

Il confronto delle prestazioni tra la versione sequenziale e quella parallela è avvenuto basandosi sui **tempi di esecuzione** dei due programmi e/o di alcune parti al loro interno.

Inizialmente consideriamo una normale esecuzione sia nel caso sequenziale che in quello parallelo agendo al massimo delle prestazioni, ovvero utilizzando **8 threads** durante l'esecuzione (del caso parallelo) ed utilizzando in entrambi i casi una dimensione del vettore di bit pari ad **n = 30000**.

In questa situazione il **programma sequenziale** mostra in output i seguenti risultati:

```
Tempo di esecuzione h5:  9.037000894546509
Tempo di inizializzazione:  9.03799843788147
Corrispondenze trovate:
ground
mind
wonder
hot
come
did
my
sound
no
most
number
who
over
know
Tempo di matching:  0.5094091892242432
Tempo complessivo di esecuzione:  9.5484299659729
```

Nell'output, come si può notare dall'immagine precedente, vengono illustrati:

- Il tempo relativo alla fase di inizializzazione (e dunque anche gestione) del filtro;
- Le corrispondenze (parole) trovate del secondo insieme (fatte passare, dunque, dal filtro);
- Il tempo complessivo di esecuzione del programma.

Analizzando invece l'**output del programma parallelo**, si osservano i seguenti risultati:

```

Tempo di esecuzione h5: 3.281970500946045
Tempo di inizializzazione: 4.377262353897095
Corrispondenze trovate:
mind
ground
hot
my
no
wonder
did
come
most
sound
who
over
know
number
Tempo di matching: 0.2814290523529053
Tempo complessivo di esecuzione: 4.65869140625

```

Guardando i tempi di inizializzazione (e gestione) del filtro notiamo già un'ottima riduzione dei tempi di esecuzione, passando da circa 9.038 secondi a circa 4.377 secondi:

TEMPI DI ESECUZIONE	
CASO SEQUENZIALE	CASO PARALLELO
9.038	4.377

Calcolando dunque lo **speedup** risultante si ottiene:

$$Speedup = \frac{T_S}{T_P} = \frac{9.038}{4.377} = 2.0649$$

Quindi questo risultato ci dice che con i parametri specificati la versione parallela è **due volte** più veloce rispetto a quella sequenziale.

Una parte in cui i tempi si accorciano molto è data dalla funzione h5. In questa, infatti, i tempi passano da 9.037 a 3.282. In questo caso lo speedup risulta essere:

$$Speedup = \frac{T_S}{T_P} = \frac{9.037}{3.282} = 2.7535$$

Ovvero si arriva ad ottenere la funzione parallelizzata che è quasi 3 volte più veloce rispetto a quella sequenziale.

Per quanto riguarda invece il tempo di matching si passa da 0.5904 a 0.2814. In questo caso si ha:

$$Speedup = \frac{T_S}{T_P} = \frac{0.5904}{0.2814} = 2.0981$$

Se invece si passa all'analisi del tempo di esecuzione complessivo del programma, mentre il programma sequenziale viene eseguito in 9.5484 secondi, quello parallelo riduce i suoi tempi di esecuzione a soli 4.6587 secondi, generando dunque il seguente speedup:

$$Speedup = \frac{T_S}{T_P} = \frac{9.5484}{4.6487} = 2.05400$$

Questi tempi sono tutti relativi all'esecuzione definita con un 30000 bit per il filtro.

Ponendoci in un caso ancora più particolare, aumentando ad esempio questo numero di bit a 500000, otteniamo i seguenti risultati:

```
Tempo di esecuzione h5: 157.98341393470764
Tempo di inizializzazione: 158.0486147403717
Corrispondenze trovate:
hot
come
did
my
sound
no
most
number
who
over
know
Tempo di matching: 5.371013402938843
Tempo complessivo di esecuzione: 163.42062401771545
```

Eseguendo il programma parallelo si ottiene invece:

```
Tempo di esecuzione h5: 60.8713583946228
Tempo di inizializzazione: 62.0406231880188
Corrispondenze trovate:
my
hot
come
sound
did
no
most
number
who
over
know
Tempo di matching: 4.906532526016235
Tempo complessivo di esecuzione: 66.94812560081482
```

Come prima possiamo confrontarne i tempi di esecuzione:

Per il tempo di inizializzazione si ottiene:

$$Speedup = \frac{T_S}{T_P} = \frac{158.049}{62.041} = 2.5475$$

Per il tempo di esecuzione relativo ad h5:

$$Speedup = \frac{T_S}{T_P} = \frac{157.983}{60.871} = 2.5954$$

Per il tempo di matching:

$$Speedup = \frac{T_S}{T_P} = \frac{5.371}{4.907} = 1.0946$$

E, infine, per il tempo complessivo di esecuzione del programma si ottiene:

$$Speedup = \frac{T_S}{T_P} = \frac{163.421}{66.948} = 2.4410$$

Il quale risulta aumentato rispetto al caso testato prima in cui n era uguale a 30000.

Adesso, una prova interessante da fare, è invece quella in cui si analizzano i tempi di esecuzione partendo dal caso sequenziale ed aumentando via via il numero di threads utilizzati nella versione parallela. Queste esecuzioni è stato deciso di testarle con $n = 500000$:

Tempi di esecuzione con $n = 500000$	
Sequenziale	163.420624
nThreads = 2	96.676440
nThreads = 4	74.713442
nThreads = 6	69.186874
nThreads = 8	66.948126
nThreads = 10	74.268388

Da questa tabella si nota che i tempi di esecuzione diventano sempre più piccoli fino al raggiungimento del **valore minimo** nel caso in cui il loro numero sia pari al numero di **processori logici** della macchina utilizzata (che nel caso in questione è, infatti, 8). Continuando ad aumentare il numero di threads è visibile invece un **peggioramento** dei tempi di esecuzione, i quali tendono a diventare poi via via più lunghi.