

Algoritmi e Strutture Dati

Insiemi e dizionari Riassunto finale

Alberto Montresor

Università di Trento

2018/11/21

This work is licensed under a Creative Commons
Attribution-ShareAlike 4.0 International License.



Insiemi e dizionari

Insiemi

- Collezione di oggetti

Dizionari

- Associazioni chiave-valore

Implementazione

- Molte delle strutture dati viste finora
- Vantaggi e svantaggi

Insiemi realizzati con vettori booleani

Insieme

- Interi $1 \dots m$
- Collezione di m oggetti memorizzati in un vettore

Rappresentazione

- Vettore booleano di m elementi

Vantaggi

- Notevolmente semplice
- Efficiente verificare se un elemento appartiene all'insieme

Svantaggi

- Memoria occupata $O(m)$, indipendente dalle dimensioni effettive
- Alcune operazioni inefficienti – $O(m)$

Insiemi realizzati con vettori booleani

SET (vettore booleano)

boolean[] *V*

int *size*

int *dim*

SET Set(**int** *m*)

 SET *t* = **new** SET

t.size = 0

t.dim = *m*

t.V = [**false**] * *m*

return *t*

boolean contains(**int** *x*)

if $1 \leq x \leq \textit{dim}$ **then**

return *V*[*x*]

else

return **false**

int size()

return *size*

insert(**int** *x*)

if $1 \leq x \leq \textit{dim}$ **then**

if not *V*[*x*] **then**

size = *size* + 1

V[*x*] = **true**

remove(**int** *x*)

if $1 \leq x \leq \textit{dim}$ **then**

if *V*[*x*] **then**

size = *size* - 1

V[*x*] = **false**

Insiemi realizzati con vettori booleani

SET (vettore booleano)

SET union(SET A , SET B)

SET $C = \text{Set}(\max(A.\text{dim}, B.\text{dim}))$

for $i = 1$ **to** $A.\text{dim}$ **do**

if $A.\text{contains}(i)$ **then**

$C.\text{insert}(i)$

for $i = 1$ **to** $B.\text{dim}$ **do**

if $B.\text{contains}(i)$ **then**

$C.\text{insert}(i)$

SET intersection(SET A , SET B)

SET $C = \text{Set}(\min(A.\text{dim}, B.\text{dim}))$

for $i = 1$ **to** $\min(A.\text{dim}, B.\text{dim})$ **do**

if $A.\text{contains}(i)$ **and** $B.\text{contains}(i)$ **then**

$C.\text{insert}(i)$

SET difference(SET A , SET B)

SET $C = \text{Set}(A.\text{dim})$

for $i = 1$ **to** $A.\text{dim}$ **do**

if $A.\text{contains}(i)$ **and not** $B.\text{contains}(i)$

then

$C.\text{insert}(i)$

Java - class `java.util.BitSet`

Metodo	Operaz.
<code>void and(BitSet set)</code>	Union
<code>void or(BitSet set)</code>	Intersection
<code>int cardinality()</code>	Set size

Metodo	Operaz.
<code>void clear(int i)</code>	Remove
<code>void set(int i)</code>	Insert
<code>boolean get(int i)</code>	Contains

C++ STL

- `std::bitset` – Struttura dati bitset con dimensione fissata nel template al momento della compilazione.
- `std::vector<bool>` – Specializzazione di `std::vector` per ottimizzare la memorizzazione, dimensione dinamica.

Insiemi realizzati con liste / vettori non ordinati

Costo operazioni

- Operazioni di ricerca, inserimento e cancellazione: $O(n)$
- Operazioni di inserimento (assumendo assenza): $O(1)$
- Operazioni di unione, intersezione e differenza: $O(nm)$

SET difference(SET A , SET B)

SET C = Set()

foreach $s \in A$ **do**

if not $B.\text{contains}(s)$ **then**
 $C.\text{insert}(s)$

return C

Insiemi realizzati con liste / vettori ordinati

LIST intersection(LIST A , LIST B)

LIST $C = \text{Set}()$

POS $p = A.\text{head}()$

POS $q = B.\text{head}()$

while not $A.\text{finished}(p)$ **and**
 not $B.\text{finished}(q)$ **do**

if $A.\text{read}(p) == B.\text{read}(q)$ **then**

$C.\text{insert}(C.\text{tail}(), A.\text{read}(p))$

$p = A.\text{next}(p)$

$q = B.\text{next}(q)$

else if $A.\text{read}(p) < B.\text{read}(q)$
 then

$p = A.\text{next}(p)$

else

$q = B.\text{next}(q)$

return C

Costo operazioni

- Ricerca:
 - $O(n)$ (liste)
 - $O(\log n)$ (vettori)
- Inserimento/cancellazione
 - $O(n)$
- Unione, intersezione e differenza:
 - $O(n)$

Insiemi – Strutture dati complesse

Alberi bilanciati

- Ricerca, inserimento, cancellazione: $O(\log n)$
- Iterazione: $O(n)$
- Con ordinamento
- Implementazioni:
 - Java TreeSet
 - Python OrderedSet
 - C++ STL set

Hash table

- Ricerca, inserimento, cancellazione: $O(1)$
- Iterazione: $O(m)$
- Senza ordinamento
- Implementazioni:
 - Java HashSet
 - Python set
 - C++ STL unordered_set

Insiemi – Riassunto

	contains lookup	insert	remove	min	foreach (Memoria)	Ordine
Vettore booleano	$O(1)$	$O(1)$	$O(1)$	$O(m)$	$O(m)$	Sì
Lista non ordinata	$O(n)$	$O(n)$	$O(n)$	$O(n)$	$O(n)$	No
Lista ordinata	$O(n)$	$O(n)$	$O(n)$	$O(1)$	$O(n)$	Sì
Vettore ordinato	$O(\log n)$	$O(n)$	$O(n)$	$O(1)$	$O(n)$	Sì
Alberi bilanciati	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(n)$	Sì
Hash (Mem. interna)	$O(1)$	$O(1)$	$O(1)$	$O(m)$	$O(m)$	No
Hash (Mem. esterna)	$O(1)$	$O(1)$	$O(1)$	$O(m+n)$	$O(m+n)$	No

$m \equiv$ dimensione del vettore o della tabella hash

Python – List

Operazione		Caso medio	Caso pessimo ammortizzato
<code>L.copy()</code>	Copy	$O(n)$	$O(n)$
<code>L.append(x)</code>	Append	$O(1)$	$O(1)$
<code>L.insert(i,x)</code>	Insert	$O(n)$	$O(n)$
<code>L.remove(x)</code>	Remove	$O(n)$	$O(n)$
<code>L[i]</code>	Index	$O(1)$	$O(1)$
<code>for x in L</code>	Iterator	$O(n)$	$O(n)$
<code>L[i:i+k]</code>	Slicing	$O(k)$	$O(k)$
<code>L.extend(s)</code>	Extend	$O(k)$	$O(k)$
<code>x in L</code>	Contains	$O(n)$	$O(n)$
<code>min(L), max(L)</code>	Min, Max	$O(n)$	$O(n)$
<code>len(L)</code>	Get length	$O(1)$	$O(1)$

$$n = \text{len}(L)$$

Python – Set

Operazione		Caso medio	Caso pessimo
<code>x in S</code>	Contains	$O(1)$	$O(n)$
<code>S.add(x)</code>	Insert	$O(1)$	$O(n)$
<code>S.remove(x)</code>	Remove	$O(1)$	$O(n)$
<code>S T</code>	Union	$O(n + m)$	$O(n \cdot m)$
<code>S&T</code>	Intersection	$O(\min(n, m))$	$O(n \cdot m)$
<code>S-T</code>	Difference	$O(n)$	$O(n \cdot m)$

$$n = \text{len}(S), m = \text{len}(T)$$

BitSet + Tabelle Hash = Bloom Filters

BitSet

Vantaggi

- 1 bit/oggetto

Svantaggi

- Elenco prefissato di oggetti

Tabelle Hash

Vantaggi

- Struttura dati dinamica

Svantaggi

- Alta occupazione di memoria

Bloom filters

Vantaggi

- Struttura dati dinamica
- Bassa occupazione di memoria (10 bit/oggetto)

Svantaggi

- Niente cancellazioni
- Risposta probabilistica
- No memorizzazione

Specifica

- `insert(k)`: Inserisce l'elemento x nel bloom filter
- `boolean contains(k)`
 - Se restituisce **false**, l'elemento x è sicuramente non presente nell'insieme
 - Se restituisce **true**, l'elemento x può essere presente oppure no (**falsi positivi**)

Bloom filter

Trade-off fra occupazione di memoria e probabilità di falso positivo

- Sia ϵ la probabilità di falso positivo
- I bloom filter richiedono $1.44 \log_2(1/\epsilon)$ bit per elemento inserito

ϵ	Bit
10^{-1}	4.78
10^{-2}	9.57
10^{-3}	14.35
10^{-4}	19.13

Applicazioni dei Bloom Filter

Chrome Safe Browsing

- Chrome contiene un database delle URL associate a siti con malware, costantemente aggiornato
- Fino al 2012, memorizzato con un Bloom Filter
- Chrome verifica l'appartenenza di ogni URL al database
 - Se la risposta è **false**, non appartiene
 - Se la risposta è **true**, potrebbe appartenere e viene fatta una verifica tramite un servizio centralizzato di Google

Qualche dato (da prendere cum grano salis)

- Nel 2011, 650k URL memorizzati in 1.94MB
- 25 bit per URL, $\epsilon \approx 10^{-5}$

Applicazioni dei Bloom Filter

Ogni qual volta una verifica locale permette di evitare un'operazione più costosa, quali operazioni di I/O e comunicazioni di rete

They say...

- **Facebook** uses bloom filters for typeahead search, to fetch friends and friends of friends to a user typed query.
- **Apache HBase** uses bloom filter to boost read speed by filtering out unnecessary disk reads of HFile blocks which do not contain a particular row or column.
- **Transactional Memory** (TM) has recently applied Bloom filters to detect memory access conflicts among threads.
- When you log into **Yahoo** mail, the browser page requests a bloom filter representing your contact list

Implementazione

- Un vettore booleano A di m bit, inizializzato a **false**
- k funzioni hash $h_1, h_2, \dots, h_k : U \rightarrow [0, m - 1]$

k_1

k_2

k_3

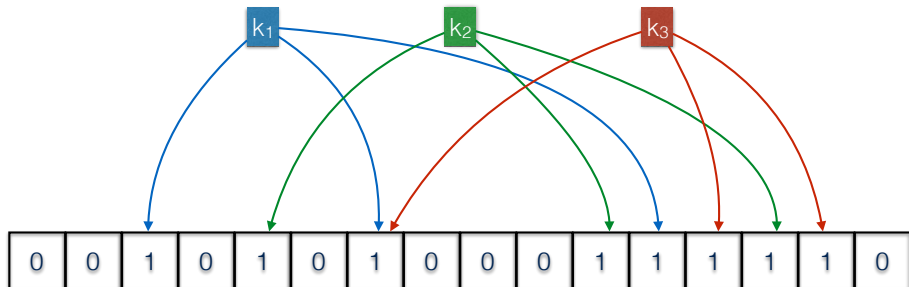


Implementazione

`insert(k)`

for $i = 1$ **to** k **do**

$A[h_i(k)] = \text{true}$



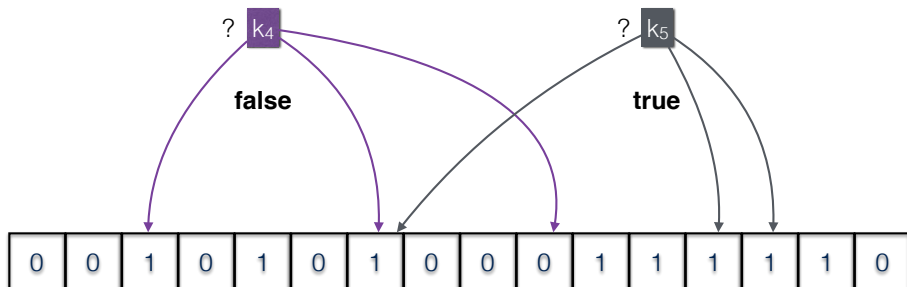
Implementazione

boolean contains(k)

for $i = 1$ **to** k **do**

if $A[h_i(k)] == \text{false}$ **then return false**

return true



Qualche formula (senza dimostrazione)

- Dati n oggetti, m bit, k funzioni hash, la probabilità di un falso positivo è pari a:

$$\epsilon = \left(1 - e^{-kn/m}\right)^k$$

- Dati n oggetti e m bit, il valore ottimale per k è pari a

$$k = \frac{m}{n} \ln 2$$

- Dati n oggetti e una probabilità di falsi positivi ϵ , il numero di bit m richiesti è pari a:

$$m = -\frac{n \ln \epsilon}{(\ln 2)^2}$$