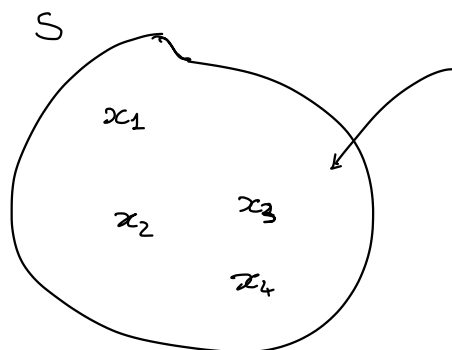


Algoritmi e Strutture Dati (26/10/2021)

* Hash Tables

collezione di oggetti

x $\left\{ \begin{array}{l} x.key \\ \text{dati satellite} \end{array} \right.$



Insert (S, x)
Delete (S, x)
Search (S, key)

efficiente:

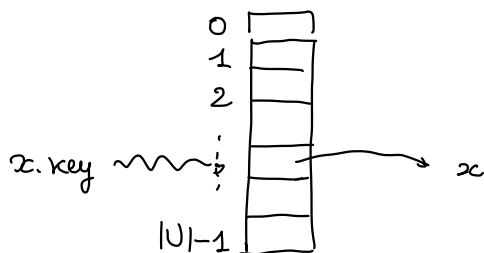
→ costo medio $\Theta(1)$
→ " peggiore $\Theta(n)$

* Indirizzamento diretto

insieme (universo) di chiavi $U = \{0, 1, \dots, |U|-1\}$

idea: collezione come un array $T[0..|U|-1]$

un elemento x è inserito in $T[x.key]$



$T[j] = \begin{cases} \text{l'elemento } x \text{ sulla collezione} \\ \text{t.c. } x.key = j & (\text{se c'è}) \\ \text{nil} & \text{altrimenti} \end{cases}$

problemi:

- ① non è possibile avere oggetti con la stessa chiave
- ② ok se il numero $|U|$ è "piccolo"

Insert (T, x)
 $T[x.key] = x$ } $\Theta(1)$

Search (T, k)
return $T[k]$ } $\Theta(1)$

Delete (T, x)
 $T[x.key] = \text{nil}$ } $\Theta(1)$

problema :

U = stringhe di 8 caratteri

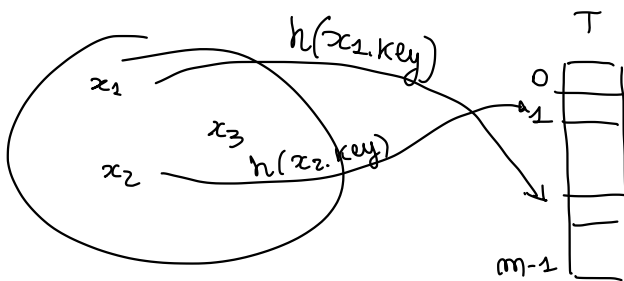
1 car. = 8 bit 2^8 caratteri

$$\underbrace{2^8 \times \dots \times 2^8}_8 \text{ stringhe} = 2^{64} \text{ stringhe} \\ \cong 16 \cdot 10^8$$

* Hash Tables

obiettivo : usare quantità di memoria proporzionale al # di elementi da memorizzare

Idea : uso tabella T di dimensione $m \ll |U|$



$$h: U \rightarrow [0..m-1]$$

funzione di hashing

* Collisione : k_1 e k_2 t.c. $h(k_1) = h(k_2)$

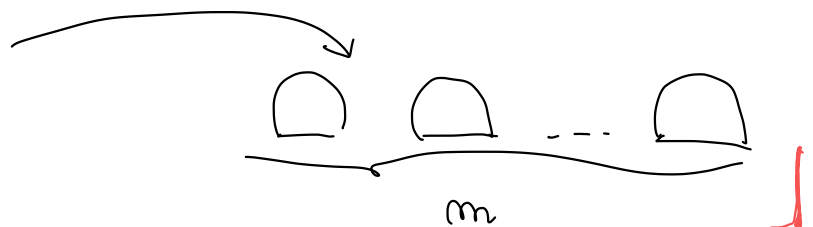
inevitabile se $m \ll |U|$

indico con n = num. di elementi memorizzati
se $n > m$ } \rightarrow collisione

PIGEON HOLE PRINCIPLE

Dati A, B finiti $|A| > |B|$

ogni funzione ogni $f: A \rightarrow B$ non iniettivo.

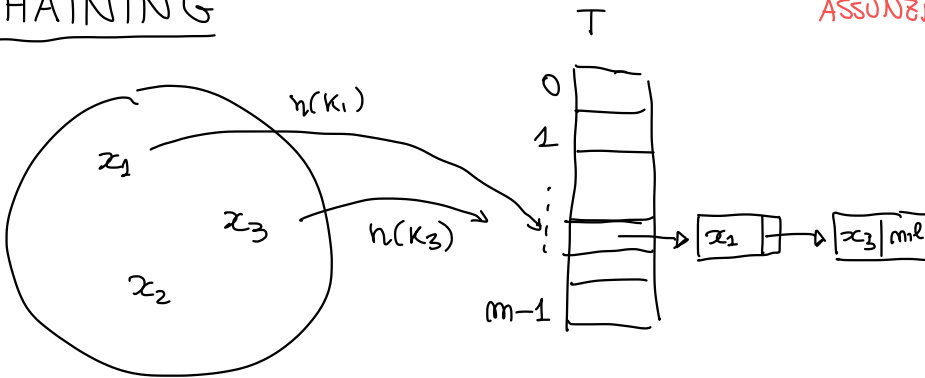


2 possibili soluzioni

Ⓐ chaining

Ⓑ open addressing

Ⓐ CHAINING



ASSUNZIONE: h (funzione di hash)
si calcola in $O(1)$

$T[0, \dots, m-1]$ contiene liste

$T[j] =$ lista degli elementi x t.c. $h(x.key) = j$

Operazioni

Insert (T, x)

inserisce x nella lista $T[h(x.key)]$ (in testa) } Ⓐ(1)

Delete (T, x)

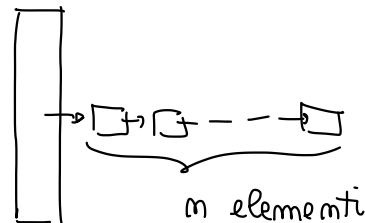
cancella x dalla lista $T[h(x.key)]$



Search (T, k)

cerca un elemento con chiave k nella lista $T[h(k)]$

caso peggiore Ⓐ(m)



CASO MEDIO ?

m = # celle tabella

n = # elementi inseriti

$$\alpha = \frac{n}{m}$$

↑
fattore di carico

< 1

$= 1$

> 1

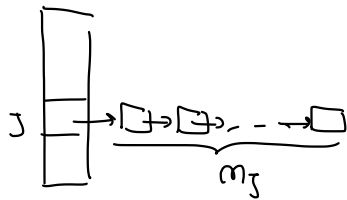
↳ qualità della funzione hash

distribuzione dell'input

ipotesi : HASHING UNIFORME SEMPLICE

ogni elemento x in input ha la stessa probabilità $\frac{1}{m}$ di essere "indirizzato" in una qualunque delle m celle

* indicato con m_j = lunghezza lista in $T[j]$



$$E[m_j] = \sum_{i=1}^m \frac{1}{m} \cdot 1 = \frac{n}{m} = \alpha$$

* costo medio search

(• a •) Ricerca di una chiave non presente k

→ calcolo $h(k) = j$

→ accesso alla lista $T[j]$

→ scorrimento della lista $T[j]$ fino in fondo (m_j elementi) m_j

$$\Theta(1 + m_j)$$

media $\Theta(1 + \alpha)$

(• b •) Ricerca di una chiave presente

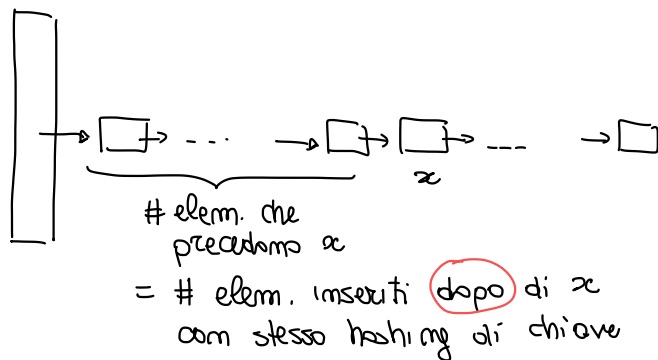
→ calcolo $h(k) = j$

→ accesso alla lista $T[j]$

→ scorrimento della lista $T[j]$ fino a trovare elem. con chiave k

$$1 + \frac{m_j}{2} \quad \xrightarrow{\text{media}} \quad 1 + \frac{\alpha}{2}$$

più precisi: cerca x con chiave $x.key = k$



siano $x_1 \dots x_m$ gli elementi inseriti

tempo medio per cercare un elemento presente è

$$\begin{aligned}
 & \frac{1}{m} \sum_{i=1}^m \underbrace{\text{tempo atteso per trovare } x_i}_{\substack{\text{valore atteso del \# elem. inseriti} \\ \text{dopo di } x_i \text{ con stesso hashing} \\ x_{i+1} \dots x_m}} \\
 &= \frac{1}{m} \sum_{i=1}^m \sum_{j=i+1}^m \underbrace{\text{prob } h(x_i.key) = h(x_j.key)}_{1/m} \\
 &= \frac{1}{m} \sum_{i=1}^m \underbrace{\sum_{j=i+1}^m 1/m}_{m-i \text{ termini}} = \frac{1}{m} \sum_{i=1}^m \frac{m-i}{m} = \frac{1}{m} \sum_{h=0}^{m-1} h \\
 & \qquad \qquad \qquad \frac{m(m-1)}{2} \\
 &= \frac{m(m-1)}{2m} = \frac{m}{2m} - \frac{1 \cdot m}{2m \cdot m} = \frac{1}{2} - \frac{1}{2m}
 \end{aligned}$$

$$\underbrace{(1 + \frac{1}{2})}_{(1)} = (1 + \frac{1}{2m})$$

IN CONCLUSIONE: costo della ricerca (in media)

- elemento assente (1)
- elemento presente $(1 + \frac{1}{2m})$

(1) ←

NOTA: se assumiamo

$$m \leq C \cdot n$$

$$\alpha = \frac{n}{m} \leq C$$

(C costante)

* Come definire (buone) funzioni di hash?

NOTA: se le chiavi fossero numeri reali $K \in [0, 1)$

$$h(K) = \lfloor m K \rfloor$$

* Metodo della divisione

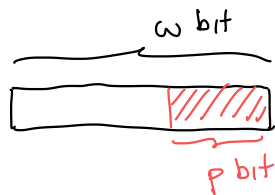
chiavi $U = \{0, 1, \dots, |U|-1\}$

$$h: U \rightarrow \{0, 1, \dots, m-1\} \quad m < |U|$$

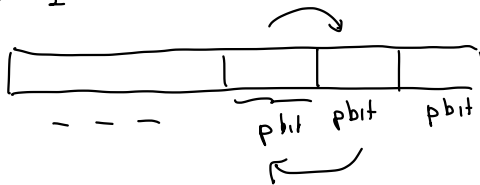
$$h(K) = K \bmod m$$

→ scelta di m è critica

→ $m = 2^p$



→ $m = 2^p - 1$



$h(K) = h(K')$ se K si ottiene "permutando" K'

numero divisibile per 9 se \sum cifre divisibile per 9

$$\text{numero in base } b \quad a_1 \dots a_n \bmod (b-1) = \left(\sum_{i=1}^n a_i \right) \bmod (b-1)$$

ESERCIZIO

* Metodo della moltiplicazione

se $x \in [0, 1)$ $h(x) = \lfloor m x \rfloor$

ma la nostra chiave è $K \in \{0, 1, \dots, |U|-1\}$

dato K

→ fisso costante $A \in (0, 1)$ ($0 < A < 1$)

$$\rightarrow K \cdot A = \text{-----}, \text{-----}$$
$$\boxed{K \cdot A \bmod 1}$$
$$(K \cdot A - \lfloor K \cdot A \rfloor)$$

$$h(K) = \lfloor m (K \cdot A \bmod 1) \rfloor$$

PRO :

→ m mem outio

→ A " " ($A = \frac{\sqrt{5}-1}{2}$)

CONTRO : ~~costo~~ ?

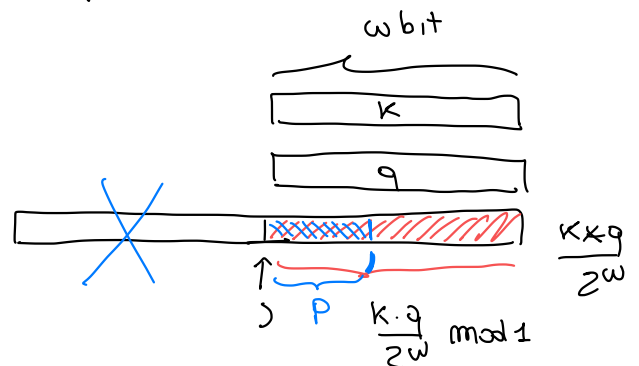
(NO)

se $m = 2^p$

$$A = \frac{q}{2^w}$$

$0 < q < 2^w$
 $w = \# \text{ bit parola}$

$$K \cdot A = K \cdot \frac{q}{2^w}$$



$$h(K) = \lfloor m (K \cdot A \bmod 1) \rfloor$$

$$= \lfloor 2^p \underbrace{\frac{K \cdot q}{2^w} \bmod 1} \rfloor$$

= p bit più significativi
della parola meno significative di $K \times q$