

Problema del dizionario

Vogliamo in questa sezione discutere le
strutture dati o dizionari, ossia
array nel senso che "generalizzano" gli
come chiavi o punteggiatori, che
anche tipi diversi "da quello
intero" (il "non vanno out of
bound").

Supponiamo che U sia il nostro
universo delle chiavi (per esempio,
gli interi), e denotiamo con
 $S \subseteq U$ l'insieme delle chiavi
memorizzate.

Vogliamo definire le strutture dei
seguenti metodi:

- $\text{find}(x) \begin{cases} 1 & \text{se } x \in S \\ 0 & \text{altrimenti} \end{cases}$ } **STATICO**
 - $\text{insert}(x) \text{ — } S \leftarrow S \cup \{x\}$
 - $\text{remove}(x) \text{ — } S \leftarrow S \setminus \{x\}$
- DINAMICO**

Se V è anche ben ordinato allora possiamo implementare:

• $\text{pred}(x) \longrightarrow \max \{y \in S \mid y \leq x\}$

• $\text{succ}(x) \longrightarrow \min \{y \in S \mid y \geq x\}$

• $\text{intervallo}(a, b) \longrightarrow \{y \in S \mid a \leq y \leq b\}$

$\rightarrow x \in S \Rightarrow \text{pred}(x) = \text{succ}(x) = x$.

• $\text{count}(a, b) = \# \text{intervallo}(a, b)$.

Struttura	find	insert	delete
Array	$O(m)$	$O(1)$	$O(1)$
Array ord.	$O(\log(m))$	$O(m)$	$O(m)$
Linked list	$O(m)$	$O(1)$	$O(1)$ (con referim) $O(m)$ (con valore)
" ord.	$O(m)$	$O(m)$	"
Heap	$O(m)$	$O(\log(m))$ sift-up	$O(\log(m))$ sift-down
BST alt.h	$O(h)$	$O(h)$	$O(h)$
Hash	$O(1)$ medie	$O(1)$ medie	$O(1)$ medie

Studiosi questi hash, che ci
permettono per l'efficienza di
implem. // i dizionari
Hash

→ $[m] :=$ gruppo ins. di m el.

Una **FUNZIONE HASH** (o di hashing)
è una funz. $h: U \rightarrow [m]$
dove $|U| \gg m$.

→ l'idea è quella di **stendere**
 U su m insieme
mette più parole — è
vero che per Pigeonhole
vi saranno collisioni (ovvero
 h non è iniettiva), ma
cerchiamo di rendere la
scoperta di collisioni improbab.
nel caso medio.

Esempi classici MD5, SHA-1, SHA-256,
o basamenti $h(x) = x \% m$
per gli interi; un esempio
più moderno è quello
di prendere p primo in
 $[m+1, 2m]$ (teorema di Bertrand)
e usare $h(x) = (ax+b) \% p \% m$
con $a \in \mathbb{Z}_p^*$, $b \in \mathbb{Z}_p$.

UNIVERSAL
HASH

→ con quest'ultimo hash lo pred. di cellis. 1 $1/m$ (uniforme)

→ un hash si dice perfetto se non dà coll. $M = \underline{\underline{S}}$.

Per chiavi lunghe si usa la tecnica del **FOLDING**, si divide l'elem. potteroso e si effettua lo XOR del risultato dell'hashing sui singoli fattori.

Per le stringhe si divide in caratteri $a_0 \dots a_{n-1}$ e si usa $h(x) = \sum_{i=0}^{n-1} (a_i \cdot \sigma^{n-1-i}) \% m$ con $\sigma = |\Sigma|$ (l.g. $\sigma = 231$ per alfabeto l'ASCII).

Un altro hashing famoso è quello di Karp-Rabin, che annovera a $O(1)$ per "l'hash successivo".

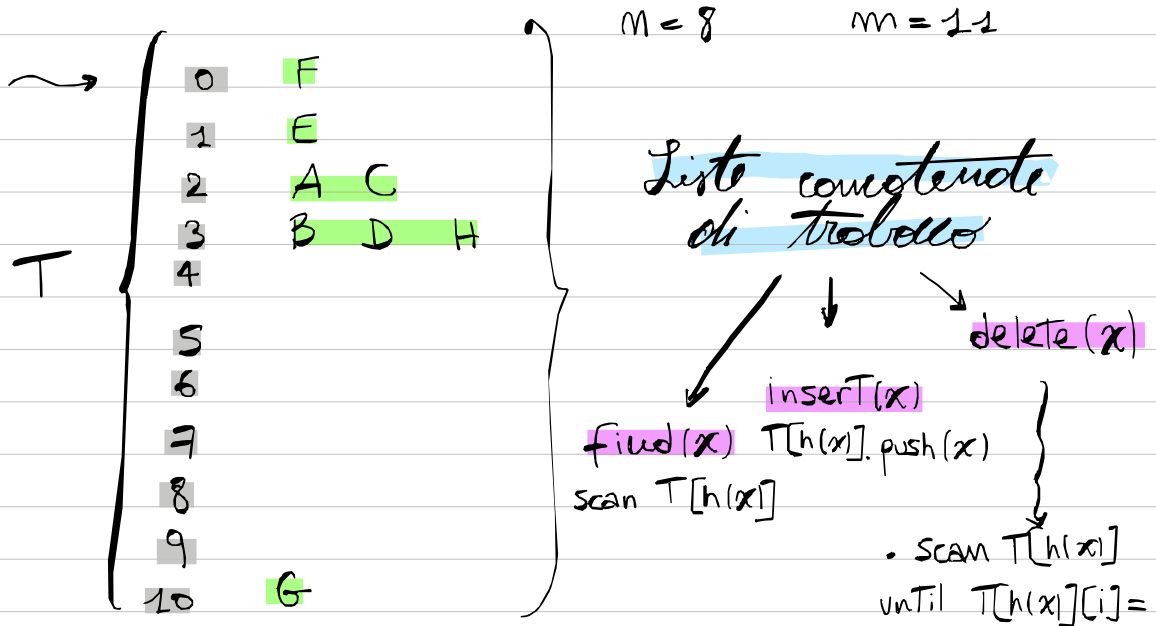
Se $h: U \rightarrow [m]$ e $|S| = m$, si definisce:

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

su pred. unif.

come **FATTORE DI CARICO**, indice median. quanti elem. ci sono in una lista - elem. delle liste di trabocco.

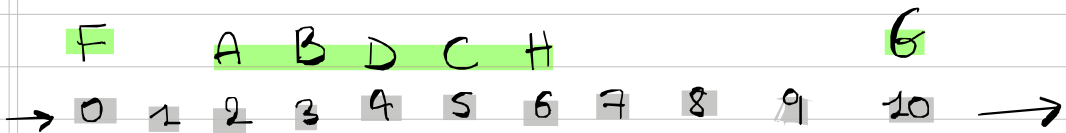
Esempio X E S B A D C F G E H \rightarrow
 $h(x)$ 3 2 3 2 0 10 1 3



Il costo medio delle op. è $O(1+d)$, mentre il costo pessimo è $O(M)$ (se gli elem. non vengono spinti unip. M T).

Un'alternativa alle liste di tubelle è l'**indirizzamento aperto**.

Esempio X E S B A D C F G H
 $h(x)$ 3 2 3 2 0 10 3
 $M=7$ $M=11$



CLUSTER: unione di alcune liste di nodi

indirizzamento aperto

- Inserimento: parte da $T[h(x)]$ e si muove nel primo posto libero in modo circolare verso destra
- Ricerca: parte da $T[h(x)]$ e scorre fino a trovare l'elem. o fino a quando trovi un posto libero (i.e. non c'è l'elem.)
- Concell: conellazione logica (come negli AVL)

L'inserimento in indirizzamento aperto è $O((1-\alpha)^{-2})$. Se $\tilde{T}(m, m)$ è il nr. prob. unif. medio di pos. esaminate per inserire una chiave, allora:

$$\begin{cases}
 \tilde{T}(m, m) = \frac{p}{m} (1 + \tilde{T}(m-1, m-1)) + \left(1 - \frac{p}{m}\right) \cdot 1 \\
 \tilde{T}(0, m) = 1 \quad \text{"nesso il resto"}
 \end{cases}$$

in port. $\tilde{T}(m, m) = 1 + \frac{m}{m} \tilde{T}(m-1, m-1)$, da cui:

$$\begin{aligned} \tilde{T}(m, m) &= 1 + \frac{m}{m} \left(1 + \frac{m-1}{m-1} \tilde{T}(m-2, m-2) \right) = \\ &= \dots = \\ &= 1 + \frac{m}{m} \left(1 + \frac{m-1}{m-1} \left(\dots \left(1 + \frac{1}{m-m+1} \cdot 1 \right) \dots \right) \right) \leq \\ &\leq \frac{m}{m-m} = (1-\alpha)^{-1} \rightarrow O((1-\alpha)^{-1}). \end{aligned}$$

Robin-Koop

→ tramite la tecnica del **ROLLING HASH** si possono comparare velocemente due stringhe per trovare una come sottostringa dell'altra.



stringa



sottostringa

→ si calcolano gli hash della sottostringa e della prima porzione della stringa.

Se combaciano, si comparano le due stringhe. Si parte poi da successive porzioni, aggiornando i hash. Se combaciano, si comparano le stringhe, altrimenti si ripete.

Spesso come hash per $C_0 C_1 C_2 \dots C_{k-1} \dots$
uso: (RABIN-KARP FINGERPRINT)

$$C_0 \sigma^{k-1} + C_1 \sigma^{k-2} + \dots + C_{k-1} \pmod{p}$$

dove $\sigma = |\Sigma|$ e Σ è l'alfabeto usato.

Per aggiornarlo è sufficiente computare:

$$H = \sigma(H - C_0 \sigma^{k-1}) + C_k \pmod{p}$$

calcolato con la **post-exp.**

Un'altra alternativa è il **cuckoo hashing**,
che sfrutta due funzioni hash
invece che una sola:

Cuckoo hashing

- $O(1)$
 - **delete** (caso pessimo)
 - **find** (pessimo)
- $O(1)$ ammortizzato.
 - **insert** (medio)

• **delete**: elimina X in $T[h_1(X)]$ e $T[h_2(X)]$

• **find**: controlla $T[h_1(X)]$ e $T[h_2(X)]$

• **insert**: se libero, mette X in $T[h_2(X)]$; **alt**ram, mette X in $T[h_2(X)]$ se libero.
Se entrambi sono occupati inizia il seguente ciclo:

- $i = h_2(X)$ [o $h_2(X)'$]
- while $T[i] \neq \text{vuoto}$:
 - $\text{tmp} = T[i]$
 - $T[i] = X$
 - $X = \text{tmp}$
 - $i \leftarrow \{h_2(X), h_2(X)'\} \setminus \{i\}$

|| Se il ciclo viene eseguito più di n volte, si effettua il **REHASHING**.

REHASHING: si ricalcola nuove h_2 e h_2' , si riempie da capo T .

Analisi dell'incremento

Si possono verificare tre situazioni:

(a.) la cella trovata è vuota ($O(1)$)

(b.) la cella non è vuota e il ciclo termina con costo di lunghezza l ($O(1+l)$)

(c.) il ciclo non termina \rightarrow **REHASHING**

Per (b.) è possibile stimare la lunghezza media \tilde{l} di un cammino supponendo $m > 2mC$ con $C > 2$.

In tal caso — per inclusione —
 $P(h_2(x) \rightsquigarrow J) \leq \frac{1}{C^l} \frac{1}{m}$

"si occorrono da $h_2(x)$ a J con cammino lungo J "

prob. di scegliere $h_2(x)$

$\binom{m}{m}$ è la prob. a ogni passo di trovare la cella occupata e si usa $m > 2mC$ per la stima...

Quindi:

$$\tilde{l} \leq \sum_{l \geq 1} l \cdot \frac{1}{C^l m} \leq \frac{1}{m} \frac{C}{(C-1)^2}$$

è una serie derivata... $\rightarrow = O(1)$

Quindi in media ci vuole $O(1 + \tilde{l}) = O(1)$.

Nel caso (c), possiamo stimare la probabilità di fare un rehashing:

$$\begin{aligned}
 P(\text{rehashing}) &\leq P(\text{esiste un ciclo}) \leq \\
 &\leq \sum_{i=1}^m P(i \rightsquigarrow i) \leq \\
 P(UX_i) &\leq \sum_{i=1}^m \sum_{l \geq 1} P(i \xrightarrow{l} i) \leq \\
 &\leq \sum_{i=1}^m P(X_i) \leq \sum_{i=1}^m \sum_{l \geq 1} \frac{1}{c^l} \frac{1}{m} \stackrel{(*)}{=} \\
 &= \frac{1}{m} \sum_{i=1}^m \frac{1}{c-1} = \frac{1}{c-1} = p
 \end{aligned}$$

\leadsto quindi $P(\neq \text{rehashing}) \leq p^t$

$$\begin{aligned}
 \# \text{ medio di rehashing} &= \sum_{t \geq 1} t p^{t-1} \stackrel{(**)}{=} \\
 &= \frac{p}{(1-p)^2} = o(1)
 \end{aligned}$$

$$* \sum_{l \geq 1} \frac{1}{c^l} = \frac{c}{c-1} - 1 = \frac{1}{c-1}$$

$$\begin{aligned}
 ** \sum_{t \geq 1} t p^t &= p \sum_{t \geq 1} t p^{t-1} = p \left(\frac{1}{1-p} \right)' = \\
 &= \frac{p}{(1-p)^2} \quad (\text{dovché } c > 2, p < 1 \dots)
 \end{aligned}$$