

6. HASHING E TABELLE HASH

Un'altra struttura dati di tipo INSIEME sono i DIZIONARI, usati per memorizzare INSIEMI DINAMICI di coppie $\langle \text{chiave}, \text{valore} \rangle$.

Le coppie sono indicizzate in base alle chiavi.

Le chiavi appartengono ad un UNIVERSO U , che è TOTALMENTE ORDINATO.

Supportano le seguenti operazioni:

- (1) INSERT ($\text{key}, \text{element}$)
- (2) DELETE (key)
- (3) SEARCH (key) $\rightarrow \text{element}$

Idealmente, vorremmo che il costo di ognuna di queste 3 operazioni sia $O(1)$; tuttavia, né con array non ordinati né ordinati, né con una lista e né con un Red-Black Tree riusciamo ad ottenere ciò.

Ma, può essere reso possibile tramite il HASHING.

• FUNZIONI HASH:

Una FUNZIONE HASH h mappa le chiavi dell'UNIVERSO U negli interi N .

$$h: U \mapsto N$$

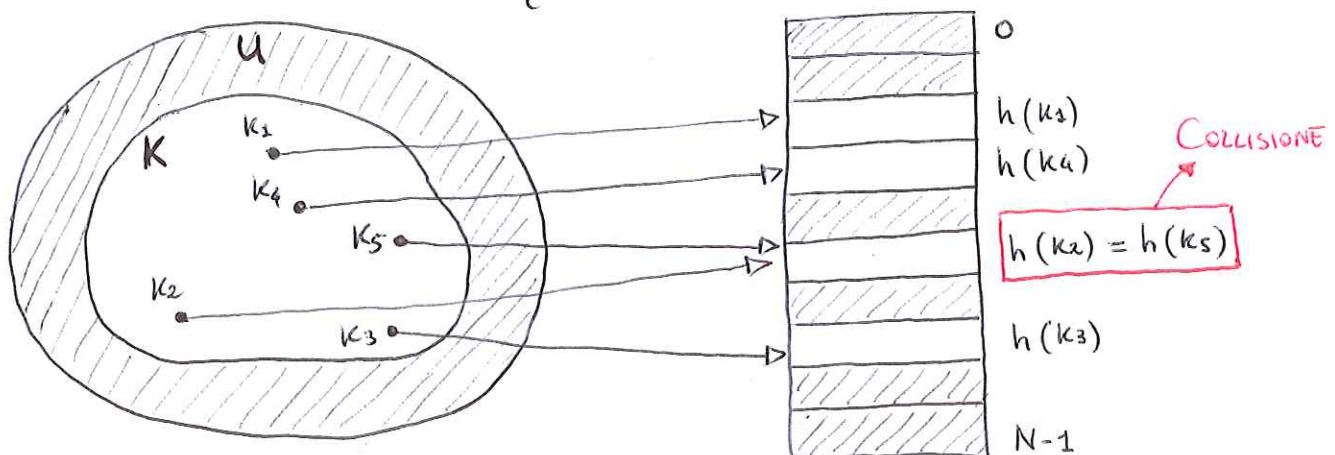
La coppia $\langle \text{chiave}, \text{valore} \rangle$ viene memorizzata in un array in POSIZIONE $h(\text{chiave})$.

In altre parole, sia A l'array $\Rightarrow A[h(\text{chiave})] \rightarrow \text{valore}$

Tale vettore viene denominato TABELLA HASH.

La tabella hash ha dimensione fissa e prefissata ad N . Le funzioni hash convertono le chiavi dell'universo U in un intero dell'INSIEME DELLE PSEUDOCIAVI $\{0, \dots, N-1\}$.

$$h: U \mapsto \{0, 1, \dots, N-1\}$$



COLLISIONI:

Si verificano quando 2 o più chiavi appartenenti ad U hanno lo stesso valore hash; cioè se: $K_i \neq K_j$, MA $h(K_i) = h(K_j)$.

• FUNZIONE HASH BANALE (IDENTITA'): TABELLE AD ACCESSO DIRETTO

Pensiamos come funzione hash la funzione IDENTITA': $h(K) = K$.

\Rightarrow La tabella di hash ha dimensione pari alla cardinalità dell'Universo delle chiavi: $N = |U|$.

La chiave K diventa l'indice dell'array del corrispondente valore. In questo modo si evitano le COLLISIONI, ma per dimensioni già mediamente grandi di U la TABELLA DI HASH esplode \rightarrow NON EFFICIENTE!

Idealmente vorremmo funzioni hash senza collisioni, cioè perfette.

Una funzione HASH è PERFETTA se è INIETTIVA:

$$\forall K_1, K_2 \in U : K_1 \neq K_2 \Rightarrow h(K_1) \neq h(K_2)$$

Purtroppo però, l'UNIVERSO DELLE CHIAVI è spesso molto VASTO, SPARSO e SCONOSCIUTO, ed è impraticabile ottenere una funzione hash perfetta. Quindi, bisogna cercare di LIMITARE LA PROBABILITA' DI COLLISIONE.

• UNIFORMITA' SEMPLICE:

Si cerca di avere una DISTRIBUZIONE UNIFORME delle chiavi negli indici $[0, \dots, N-1]$.

* Principio di UNIFORMITA' SEMPLICE:

Sia $P(K)$ la probabilità che la chiave K venga inserita nella tabella

Sia $Q(i)$ la probabilità che una chiave finisca nell'elemento i -esimo della tabella

$$\Rightarrow Q(i) = \sum_{K \in U : h(K)=i} P(K)$$

Si ha UNIFORMITA' SEMPLICE se: $\forall i \in [0, \dots, N-1] : Q(i) = \frac{1}{N}$

! Molto difficile da realizzare, poiché deve essere nota a priori la distribuzione di probabilità P .

\rightarrow STRADA POCO PRATICABILE!

COME REALIZZARE UNA FUNZIONE DI HASH

Le chiavi possono essere di qualsiasi tipo \rightarrow OCCORRE INTERPRETARE la loro RAPPRESENTAZIONE come un NUMERO (di solito binario)

Per esempio, come interpretare le STRINGHE:

- $\text{ord}(c)$: valore BINARIO del char c , in qualche codifica (ASCII, UNICODE, ...)
- $\text{bin}(k)$: RAPPRESENTAZIONE BINARIA della chiave k , concatenando i valori binari ($\text{ord}(c)$) di ogni char della stringa
- $\text{int}(b)$: valore numerico DECIMALE associato al numero binario b
- $\text{int}(k) = \text{int}(\text{bin}(k))$
- FUNZIONE HASH: ESTRAZIONE

L'ESTRAZIONE consiste nel prendere una certa parte di bit (di solito i meno significativi) della rappresentazione binaria della chiave come risultato della funzione hash:

sia $N = 2^p$, sia b un sottoinsieme di p bit presi da $\text{bin}(k)$

$$\Rightarrow \boxed{h(k) = \text{int}(b)}$$

Esempio: $N = 2^{16} = 65536$; sia $k = \text{"Alessandro"}$

Prendiamo i 16 bit meno significativi della RAPP. BINARIA di "Alessandro"

$$\rightarrow 01110010 \ 01101111 \Rightarrow h(\text{"Alessandro"}) = \text{int}(01110010 \ 01101111) = 29295$$

* PROBLEMA: probabilità di COLLISIONE molto elevata!

FUNZIONE HASH: XOR

$N = 2^p$; $h(k) = \text{int}(b)$, dove b è lo XOR effettuato a gruppi di p bit di $\text{bin}(k)$. Può essere necessario il PADDING, ovvero riempire con degli zeri gli spazi vuoti.

* MA LA PROBABILITÀ DI COLLISIONE è ancora TROPPO ALTA!

Per esempio: $h(\text{"accettare"}) = h(\text{"concedere"}) = 29696$

• FUNZIONE HASH: METODO DELLA DIVISIONE

Si prende N DISPARI, preferibilmente NUMERO PRIMO e distante da potenze di 2.

$$h(k) = \text{int}(b) \bmod N \rightarrow \text{Resto della divisione per } N$$

ESEMPIO: $N = 571$

$$\Rightarrow h(\text{"rome"}) = 416$$

$$h(\text{"romano"}) = 523$$

$$h(\text{"orme"}) = 318$$

$$h(\text{"seme"}) = 322$$

NON VANNO BENE

- $N = 2^p \rightarrow$ significherebbe utilizzare solo i p bit meno significativi
- $N = 2^p - 1 \rightarrow$ permutazioni di sottostringhe di lunghezza 2^p hanno lo stesso valore di hash.

COME GESTIRE LE COLLISIONI

Cosa fare in caso di COLLISIONE?

- (1) In caso di INSERIMENTO \rightarrow trovare una posizione alternativa
- (2) In caso di RICERCA \rightarrow se non troviamo la chiave cercata bisogna cercarla altrove

Comunque sia, queste operazioni dovrebbero avere nel caso medio un costo $O(1)$

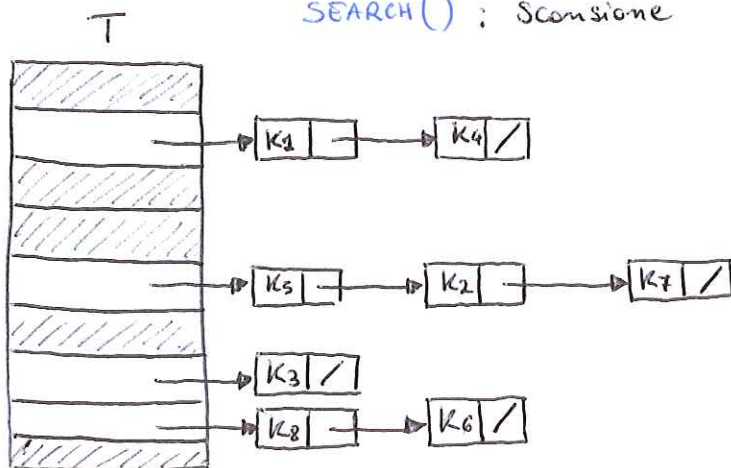
• LISTE DI TRABOCCO:

Le chiavi con lo stesso valore hash vengono organizzate in una LISTA SINGOLARMENTE COLLEGATA (oppure ARRAY DINAMICO).

Le operazioni fondamentali sono: INSERT(): inserimento in testa

DELETE(): Sconsione

SEARCH(): Sconsione



ANALISI COMPLESSITA': Liste di Trabocco

Consideriamo il FATTORE DI CARICO $\alpha = \frac{m}{N}$, dove m è il numero di elementi inseriti nelle tabelle hash. Ci dice la percentuale di riempimento della tabella.

Assumiamo che il COSTO DI CALCOLO di $h(\cdot)$ sia $\Theta(1)$.

• CASO PESSIMO: tutte le chiavi in un'unica lista di trabocco

INSERT(): $\Theta(1)$, poiché in testa

DELETE()/SEARCH(): $\Theta(m)$

• CASO MEDIO: Se assumiamo l'hashing con DISTRIBUZIONE UNIFORME delle chiavi, la lunghezza di una lista è $\alpha = \frac{m}{N}$

INSERT(): $\Theta(1)$

DELETE()/SEARCH(): $\begin{cases} \text{se con SUCCESSO: } \Theta(1) + \frac{\alpha}{2} \\ \text{se con INSUCCESSO: } \Theta(1) + \alpha \end{cases}$

\Rightarrow Se $m = O(N) \Rightarrow \boxed{\alpha = O(1)}$

Con un appropriato dimensionamento della Tabella Hash, è possibile avere nel caso medio costo $O(1)$ per tutte le operazioni.

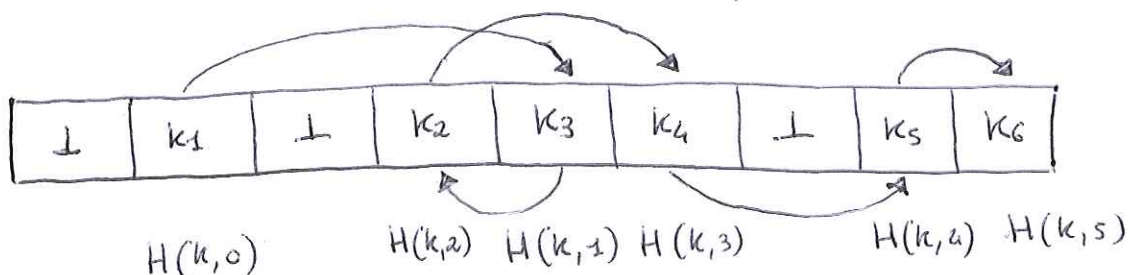
Intuitivamente, di media dev'esserci 1 elemento per lista

GESTIONE DELLE COLLISIONI: INDIRIZZAMENTO APERTO

Se si vuole evitare l'utilizzo di STRUTTURE DATI collegate come le liste, le collisioni devono essere RISOLTE ALL'INTERNO DELLA TABELLA HASH.

Si utilizza una FUNZIONE HASH ESPANSA $H(k, i)$ che calcola la posizione dell'elemento di chiave k all' i -esimo tentativo.

Per fare ciò si genera una SEQUENZA DI ISPEZIONE delle tabelle hash; tuttavia, la tabella può andare in overflow per $\alpha = \frac{m}{N} = 1$.



SEQUENZE DI ISPEZIONE

Determiniamo le sequenze con cui vengono "ispezionate" le posizioni successive delle tabelle hash:

- La funzione $H(k, i)$ DEVE PERMETTERE di ISPEZIONARE TUTTE LE POSIZIONI possibili della tabella hash;

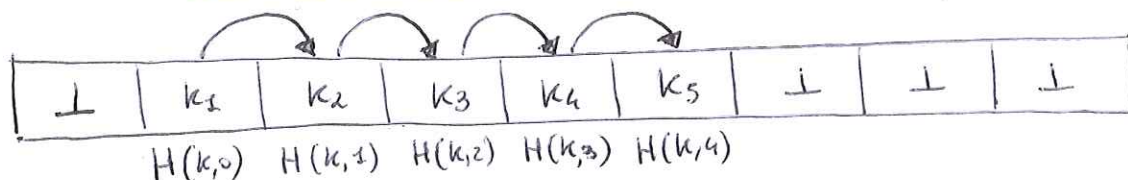
Ci sono tecniche differenti:

- ▷ SCANSIONE LINEARE;
- ▷ SCANSIONE QUADRATICA;
- ▷ HASHING DOPPIO

• SCANSIONE LINEARE:

Si cerca nelle posizioni adiacenti ad $h(k)$ linearmente:

$$H(k, i) = h(k) + c \cdot i \mod N$$



Si ha il problema dell' AGGLOMERAZIONE PRIMARIA (primary clustering):

Si generano LUNGHE SOTTOSEQUENZE OCCUPATE ed i tempi medi per le operazioni che richiedono ispezione crescono molto.

• SCANSIONE QUADRATICA:

Si cerca quadraticamente nelle posizioni adiacenti ad $h(k)$:

$$H(k, i) = h(k) + c \cdot i^2 \mod N$$

Sono possibili al più N sequenze di ispezione.

Si ha il problema dell' AGGLOMERAZIONE SECONDARIA (secondary clustering):

se $h(k_1) = h(k_2)$, le sequenze di ispezione per le 2 chiavi sono le stesse.

• HASHING DOPPIO:

Una funzione hash per determinare la posizione iniziale e un'altra funzione hash per generare la sequenza d'ispezione:

$$H(k, i) = h_1(k) + i \cdot h_2(k) \mod N$$

Sono possibili al più N^2 sequenze di ispezione.

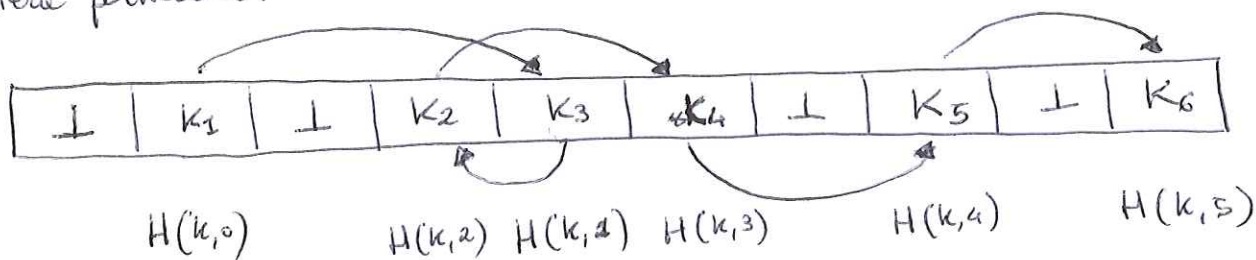
BISOGNA SCEGLIERE $h_2(\cdot)$ in modo da consentire UN'ISPEZIONE COMPLETA:

- Si ottiene solo se $h_2(k)$ ed N sono COPRIMI:

- (1) $N = 2^p$ e $h_2(k)$ genera solo numeri dispari
- (2) N PRIMO e $h_2: U \mapsto \{0, \dots, N-1\}$

CANCELLAZIONE NELL'INDIRIZZAMENTO APERTO:

Nel caso di indirizzamento aperto, la cancellazione deve essere eseguita in maniera particolare.

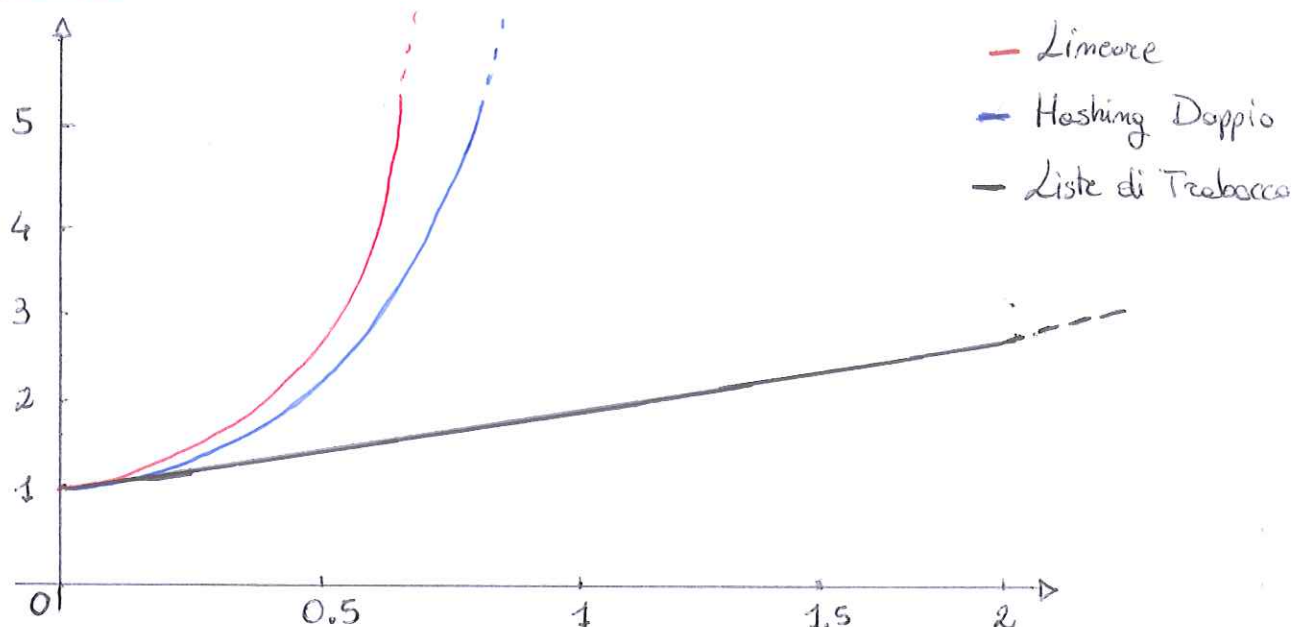


Se facessi DELETE (K_4) \Rightarrow Al posto di K_4 vorrei \perp

Poi, se eseguiessi SEARCH (K_6) \Rightarrow La sequenza di ispezione mi porterebbe prima nella cella che conteneva in precedenza K_4 e, trovando \perp , verrebbe concluso che K_6 non c'è.

- Si introduce l'elemento speciale DEL, per indicare che un elemento è stato cancellato. In caso di ispezione, DEL viene considerato come un elemento qualsiasi che però non deve mai essere restituto.

COMPLESSITA':



Le TABELLE HASH più efficienti sono quelle con l'utilizzo di liste collegate per la gestione delle collisioni, cioè le LISTE DI TRABOCCO !

TABELLA HASH ADATTATIVA

A prescindere delle strategie per gestire le collisioni, è chiaro che è conveniente non far aumentare il FATTORE DI CARICO $\alpha (= \frac{m}{N})$; l'unico modo per farlo è \rightarrow AUMENTARE N .

Anziché utilizzare come Tabella Hash un semplice array, si può far uso di un ARRAY DINAMICO:

Sopra una certa soglia del fattore di carico, si RADDOPPIANO le dimensioni della tabella hash e si REINSERISCONO TUTTE LE CHIAVI. Nel caso di indirizzamento aperto, si possono rimuovere anche tutti i DEL.

\rightarrow COSTO: $O(1)$ AMMORTIZZATO