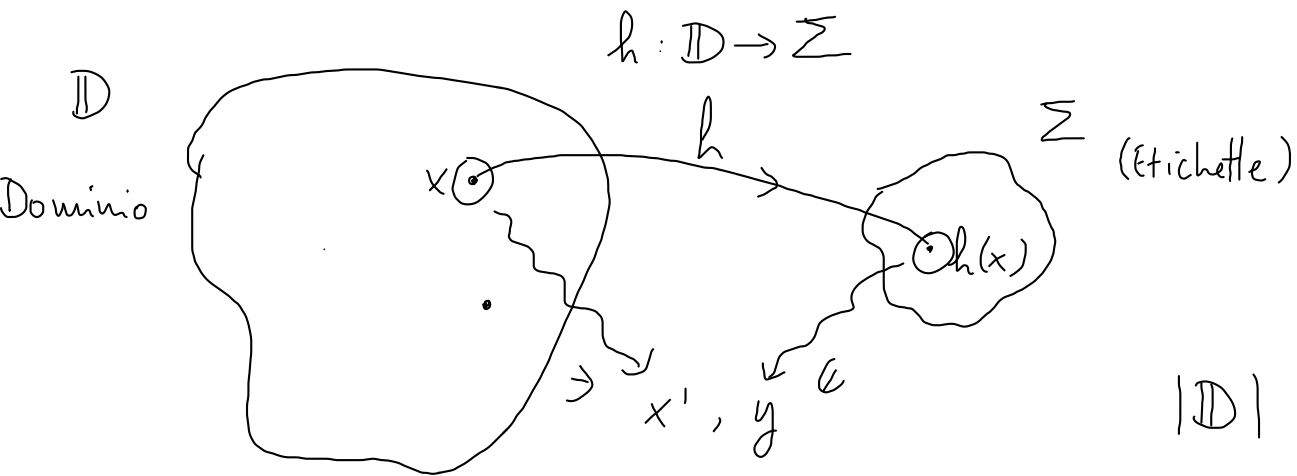


# FUNZIONI HASH



Ricevo  $x'$  e  $y$

Verifico se  $h(x') = y$

Se  $x' = x$ ,  $y = h(x) \Rightarrow y = h(x')$

Proprietà desiderabili: (informele)

1) Facile da calcolare

2) Se  $x \neq x'$ , allora  $h(x) \neq h(x')$

$$|\mathbb{D}| \gg |\Sigma|$$

Controlli di integrità

Codice fiscale



Problema del dizionario:

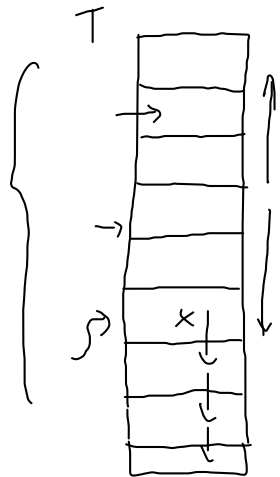
Struttura dati per memorizzare  $N$  elementi del dominio  $\mathbb{D}$

$(a_1, a_2, \dots, a_N)$  e supportare le seguenti operazioni:

- ricerca( $x$ ),  $x \in \mathbb{D}$ : esiste un elemento memorizzato  $a_i$  tale che  $a_i = x$ ?
- inserisci( $x$ ),  $x \in \mathbb{D}$ : aggiungi  $x$  agli elementi memorizzati
- cancella( $x$ ),  $x \in \mathbb{D}$ : rimuovi  $x$  dagli elementi memorizzati

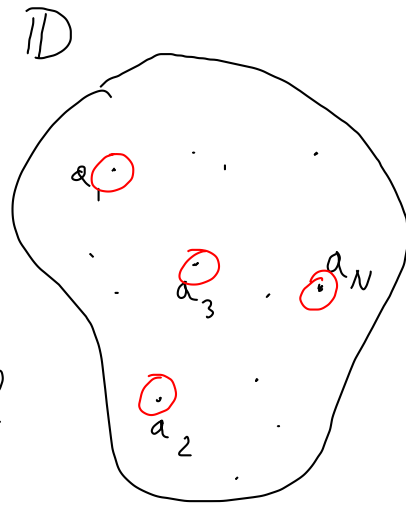
Misure di efficienza: tempo delle operazioni, spazio di memoria richiesto

Soluzione 1. Mantieni  $a_1, a_2, \dots, a_N$  in un array ordinato  $T$ :

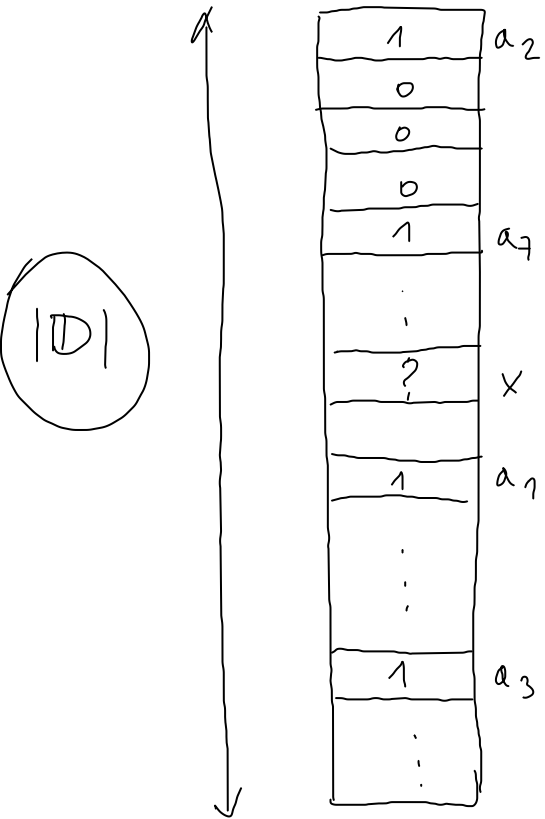


(Costruire  $T$  richiede tempo proporzionale a  $N \cdot \log N$ ) (tempo di preprocessing)

- ricerca( $x$ ): usa la ricerca binaria per determinare se  $x \in T$   
(tempo prop. a  $\log N$ )
- \* inserisci( $x$ ): aggiungi una riga, ricopiando le tabelle (tempo prop. a  $N$ )
- \* cancella( $x$ ): tempo prop. a  $N$



Soluzione 2 : Mantieni un array booleano  $B$  con un elemento per ogni  $z \in \mathbb{D}$  ponendo  $B[a_i] = 1$  per ogni  $i = 1, 2, \dots, n$  (e  $B[z] = 0$  per tutti gli altri  $z \in \mathbb{D}$ ).

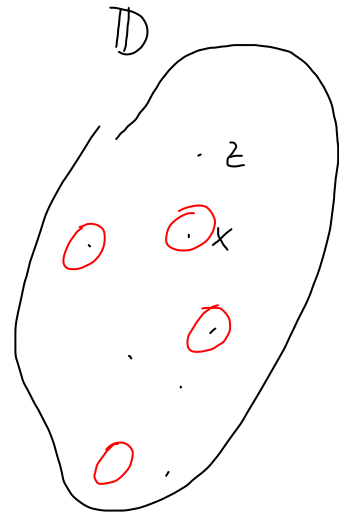


- ricerca( $x$ ) : verifica se  $B[x] = 1$  oppure no (tempo costante)
- inserisci( $x$ ) : imposta  $B[x] = 1$  (tempo " )
- cancella( $x$ ) : imposta  $B[x] = 0$  (tempo " )

Spazio richiesto è prop. a  $|\mathbb{D}| \gg N$

(anche il pre-processamento richiede tempo  $|\mathbb{D}|$ )

Domanda : esiste una struttura dati che supporti ricerca, inserimento e cancellazione in tempo costante ma in spazio prop. a  $N$ ?



Def. (Famiglia di funzioni hash)

Dati  $\mathbb{D}$ ,  $\Sigma$ ,  $m > 1$ , una famiglia hash è un insieme  
 $\mathcal{H} = \{h_1, h_2, \dots, h_m\}$  con  $h_i : \mathbb{D} \rightarrow \Sigma$ , per ogni  $i = 1, 2, \dots, m$ .

Def. (Famiglia di funzioni hash  $\epsilon$ -universale)

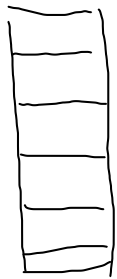
La famiglia  $\mathcal{H} = \{h_1, \dots, h_m\}$  di funzioni hash è  $\epsilon$ -universale ( $0 < \epsilon \leq 1$ )

se per ogni  $x \neq y$ ,  $x, y \in \mathbb{D}$  :

$$\Pr_i [h_i(x) = h_i(y)] \leq \epsilon$$

$\leftarrow$   $i$  è scelto uniformemente a caso in  $\{1, 2, \dots, m\}$

Soluzione 3 : simile alla soluzione 2, ma invece di mantenere un elemento dell'array  
per ogni  $z \in \mathbb{D}$ , ne mantengo uno per ogni possibile valore in  $\Sigma$ .

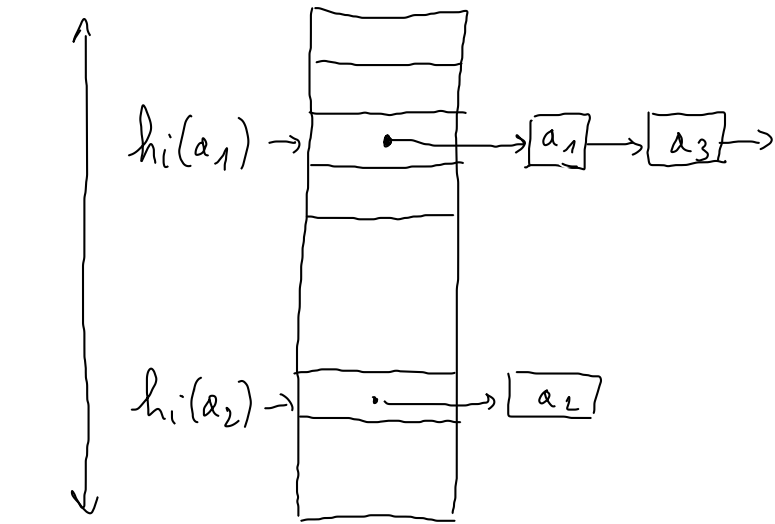


~~101~~  
121

Sia  $\mathcal{H} = \{h_1, \dots, h_m\}$  una famiglia hash  $\epsilon$ -universale.

Costruisco un array di liste collegate, una per ogni  $v \in \Sigma$ .

Tabella hash



Preprocessamento: Prendo  $h_i \in \mathcal{H}$  a caso.

Per ogni  $a_j$  ( $j = 1, 2, \dots, N$ )

aggiungi  $a_j$  alle lista collegata in posizione  $h_i(a_j)$

- ricerca( $x$ ): calcola  $h_i(x)$  e cerca  $x$  nelle lista collegata
- inserisci( $x$ ): calcola  $h_i(x)$  e aggiungi  $x$  alle lista collegata
- cancella( $x$ ): calcola  $h_i(x)$  e rimuovi  $x$  dalle lista collegata

Spazio complessivo: prop. a  $|\Sigma| + N$

(se  $|\Sigma| = O(N)$ , lo spazio totale è  $O(N)$ )

Costo inserimento: tempo costante (aggiungi  $x$  in fondo alle liste)

Costo cancellazione: uguale al tempo della ricerca

$h_i(a_2)$

$h_i(a_3) = h_i(a_1)$

$h_i(a_3)$

- Costo ricerca(x) : dipende dal numero di  $a_j$  tali che  $h_i(a_j) = h_i(x)$  -  
 ↑  
 "collisione"

Prop. Sia  $\mathcal{H} = \{h_1, \dots, h_m\}$  una famiglia hash  $\epsilon$ -universale.

Allora, per ogni  $x, a_1, a_2, \dots, a_N \in \mathbb{D}$  distinti,

$$\mathbb{E}_i \left[ \underbrace{|\{a_j : h_i(a_j) = h_i(x)\}|}_{\# \text{ di collisioni}} \right] \leq \epsilon \cdot N$$

scelta  
aleatoria  
di  $i$

Dim. Fissiamo  $j, 1 \leq j \leq N$ . Per definizione di famiglia hash  $\epsilon$ -universale,

$$\begin{aligned} \Pr_i [h_i(x) = h_i(a_j)] &\leq \epsilon. \quad \text{Ma } \mathbb{E}_i [\# \text{ collisioni}] = \mathbb{E}_i \left[ \sum_{j=1}^N \begin{cases} 1 & \text{se } h_i(a_j) = h_i(x) \\ 0 & \text{se } h_i(a_j) \neq h_i(x) \end{cases} \right] \\ &= \sum_{j=1}^N \Pr_i [a_j \text{ collida con } x] \leq \sum_{j=1}^N \epsilon = \epsilon N. \quad \text{QED.} \end{aligned}$$

Conclusione: se esiste una famiglia hash con dominio  $\mathbb{D}$  e codominio  $\Sigma$   $\varepsilon$ -universale, se  $\varepsilon = O(1/N)$  e  $|\Sigma| = O(N)$ , allora esiste una struttura dati (randomizzata) che dati  $N$  elementi  $a_1, a_2, \dots, a_N \in \mathbb{D}$  supporta ricerca, inserimento e cancellazione in tempo atteso costante, e in spazio proporzionale a  $N$ .  $\leftarrow (|\Sigma| + N)$

Data  $\mathcal{H} = \{h_1, \dots, h_n\}$ , possiamo considerare il codice associato  $\subseteq \Sigma^n$   
 $C_{\mathcal{H}} : \mathbb{D} \rightarrow \Sigma^n$  definito da  $C_{\mathcal{H}}(x) = (h_1(x), h_2(x), \dots, h_n(x))$   
 $\begin{matrix} \cap & \cap & \cap \\ \Sigma & \Sigma & \Sigma \end{matrix}$

Viceversa, dato un codice  $C(n, k)$  sull'alfabeto  $\Sigma$ , possiamo considerare

$\mathcal{H}_C = \{h_1, \dots, h_n\}$  dove  $h_i : \sum_x^k \rightarrow \Sigma$  definita da  $h_i(x) \triangleq [C(x)]_i$  per ogni  $x \in \sum^k$   
 $\begin{matrix} \downarrow \\ \text{parola di} \\ \text{codice che} \\ \text{codifica } x \end{matrix}$   $\begin{matrix} \uparrow \\ \text{prendo la componente} \\ i\text{-esima} \end{matrix}$

Prop. Se  $\mathcal{H}$  è  $\varepsilon$ -universale, allora  $C_{\mathcal{H}}$  ha distanza minima  $d_{\min} \gg (1-\varepsilon)n$ .

Viceversa, se  $C$  è un codice  $C(n, k, \delta n)$  (per qualche  $\delta \in (0, 1)$ ) allora la famiglia  $\mathcal{H}_C$  è una famiglia  $(1-\delta)$ -universale.

~~Dim~~

-> Progettare una buona famiglia di funzioni hash significa progettare un buon codice di canale.