

Hierarchical hash tree

Luigi Foscari

Si vuole creare un dizionario che permette le operazioni di inserimento, cancellazione e ricerca. I valori da inserire nella struttura appartengono ad un insieme \mathcal{X} per cui esistono k feature map $f_i : \mathcal{X} \mapsto \mathcal{H} \quad i = 1, \dots, k$ ed esistono t funzioni di hash $h_{i,b} : \mathcal{H} \mapsto \{0, \dots, b-1\} \quad i = 1, \dots, t$.

La struttura è organizzata come un albero, in cui ad ogni nodo corrisponde una lista di collisione o una tabella di dimensione b , una delle mappe di feature e una funzione di hash, nelle celle della tabella sono presenti puntatori a nodi figli. La struttura è inizializzata con un singolo nodo composto da solo un bucket.

La procedura di inserimento per un elemento $x \in \mathcal{X}$ in un nodo y avviene nel seguente modo

- Se y è una lista di collisione x è aggiunto in testa.
- Se y è una tabella con feature map f e funzione di hash h , si ripete l'inserimento ricorsivamente nel nodo puntato dalla posizione $h(f(x))$ della tabella.

Se un bucket in un nodo t raggiunge un limite massimo¹ viene **scisso**, l'operazione di scissione controlla tutti gli elementi nel bucket e trova la funzione di hash che distribuisce più equamente i valori; scelta una funzione h_b e una feature map f , viene poi creata una tabella di taglia b e i valori del bucket sono inseriti in questa tabella.

L'intera struttura dati può fare a meno delle feature map se si utilizzano funzioni di hash che hanno come dominio \mathcal{U} , ma questa opzione è meno interessante perché l'estrazione di feature permette di operare in modo più espressivo sui dati. Esempio pratico: l'universo contiene informazioni su libri, una feature che si è scelto di estrarre è la data di pubblicazione, si possono quindi cercare tutti i testi usciti nel medesimo anno in tempo efficiente.

Senza scendere nei dettagli la struttura dati dell'albero è espressa ricorsivamente in questo modo

```
type 'a tree =  
| Leaf of 'a list  
| Node of 'a tree array
```

¹Bisogna fare ulteriori considerazioni su questo valore

Dato un bucket di elementi bisogna trovare la funzione di hash tra quelle a disposizione che ottiene i valori più sparsi, è possibile utilizzare l'indice di Gini normalizzato, indichiamo con x i valori nel campione e con ϕ le frequenze relative:

$$\{x_1, \dots, x_n\} \rightarrow \{\phi_1, \dots, \phi_k\} \quad I = \left(1 - \sum_{i=0}^k \phi_i\right) \frac{k}{k-1}$$

Un algoritmo ingenuo che utilizza questo indice è il seguente, assumiamo per semplicità che ad ogni mappa di feature sia associata una e una sola funzione di hash

Algorithm 1: Algoritmo per la scelta della funzione di hash

Input: B bucket, f_1, \dots, f_k feature maps, h_1, \dots, h_k funzioni di hash
Output: Una feature map

```

1 Sia  $M$  una matrice  $k \times |B|$ 
2 for  $i = 0, \dots, k-1$  do
3   for  $j = 0, \dots, |B|-1$  do
4      $M[i][j] \leftarrow h_i(f_i(B_j))$ 
5   end
6 end
7 Sia  $G$  un array di taglia  $k$ 
8 for  $i = 0, \dots, k-1$  do
9    $G[i] \leftarrow$  indice di Gini normalizzato della  $i$ -esima riga di  $M$ 
10 end
11  $r \leftarrow$  indice dell'elemento più grande di  $G$ 
12 return  $f_r$ 
```

Sarebbe anche possibile definire le feature map in modo che abbiamo ognuna un codominio distinto, ma questo renderebbe necessario definire almeno una funzione di hash per ognuna di essere, rendendo più complicato il procedimento di scelta delle feature.

Si può rendere l'algoritmo più efficiente? Anche nel caso in cui si rimuova il rilassamento.