

Tabele de dispersie

Ne punem problema menținerii unui set finit de date, fiecare dintre ele identificată printr-o cheie unică, și notăm cu K mulțimea cheilor din acest set, numită *mulțimea cheilor actuale*. Prin natura lor, fie că sunt întregi de lungime fixă, dar mare, fie că sunt stringuri de caractere, cheile actuale K fac parte dintr-o mulțime U , de cardinalitate mult mai mare decât K , numită *universul cheilor*, lucru pe care îl notăm $|K| \ll |U|$. Pe această mulțime dorim să facem operații de inserare și ștergere, și, de asemenea, foarte frecvent, operații de căutare. Ideal ar fi ca operația de căutare să aibă timp $O(1)$ sau apropiat de el, în orice caz timp constant (și mic) în raport cu cheile din U . Structura de date care ne permite accesul în timp $O(1)$ la orice componentă este vectorul. Dacă ne-am putea permite menținerea în memorie a unui vector T a cărui mulțime de indici să fie chiar U , atunci pentru orice cheie actuală $k \in K$, în locația $T[k]$ am putea menține data identificată de cheia k . Un asemenea T se numește **tabel cu adresare directă**. Din păcate, pentru multe probleme concrete nu ne putem permite această soluție deoarece ar conduce la o risipă de spațiu. Dăm mai jos câteva exemple.

1. Vrem să menținem înregistrări cu date despre cei 68 de angajați ai unei firme. Fiecare angajat este identificat cu ajutorul unui număr de cod compus din 4 cifre. Cardinalul mulțimii cheilor actuale este $|K| = 68$, dar universul cheilor U este format din mulțimea tuturor întregilor cu 4 cifre în scriere zecimală, deci $|U| = 10^4 = 10000$. Un tabel cu adresare directă ar însemna să menținem în memorie un vector cu 10000 de componente, dintre care doar 68 ar fi efectiv folosite.
2. Vrem să menținem înregistrări cu date despre populația României. Fiecare locuitor este identificat cu ajutorul unor chei unice, codul numeric personal, un întreg compus din 12 cifre. Cardinalul mulțimii cheilor actuale este aproximativ 22 de milioane, $|K| = 22 \cdot 10^6$. Universul cheilor va fi mulțimea tuturor întregilor de 12 cifre, deci va avea dimensiunea $|U| = 10^{12}$.
3. Anuarul telefonic al unei localități menține informații despre abonații identificați cu ajutorul numelui și prenumelui, deci un șir de maximum 20 de caractere deci va avea dimensiunea $|U| = 26^{20}$.

În toate situațiile descrise mai sus este nerealistă folosirea unui tabel cu adresare directă, adică a unui tabel T indexat chiar după U . Trebuie să găsim o mulțime de indici pentru T , de cardinal m , mult mai mic decât $|U|$, eventual apropiat de cardinalul lui K , și o metodă de a dispersa cheile din K pe mulțimea de indici $\{0, 1, \dots, m-1\}$. Cu alte cuvinte, avem nevoie de o funcție, definită pe universul cheilor și cu valori în mulțimea indicilor T :

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

cu ajutorul căreia să găsim locul fiecărei date din mulțimea noastră: data identificată printr-o cheie k , $k \in K$, să se afle în locația $T[h(k)]$ a tabelului T .

O asemenea funcție se numește **funcție de dispersie**. Pentru a fi o funcție bună de dispersie, funcția h trebuie să îndeplinească anumite cerințe, și anume:

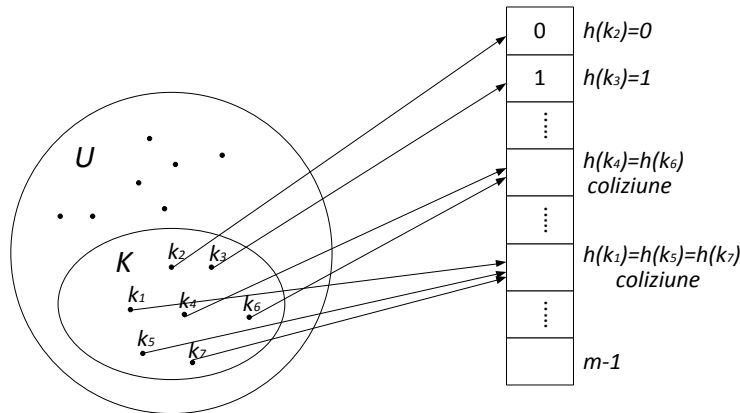
1. Să se poată calcula rapid. Timpul de calcul al lui $h(k)$ ne va da timpul de acces la componenta $T[h(k)]$.
2. Codomeniul ei să fie cât mai mic. Este dezirabil ca m să fie cât mai apropiat de $|K|$, în orice caz diferența dintre ele să fie acceptabilă ca dimensiune.
3. h să fie cât mai aproape de o funcție surjectivă, adică să avem cât mai puține locații goale. Această cerință se leagă de precedenta.
4. h să fie cât mai aproape de o funcție injectivă pe K . Cu alte cuvinte, pentru chei $k_1, k_2 \in K$, $k_1 \neq k_2$ să avem $h(k_1) \neq h(k_2)$.

Numim *coliziune* a cheilor k_1 și k_2 situația în care $h(k_1) = h(k_2)$. Cerința 4 se exprimă și sub forma „ne dorim să avem cât mai puține coliziuni”. Din păcate coliziunile nu pot fi evitate complet și trebuie găsite metode de rezolvare a lor.

Pentru a folosi un tabel de dispersie programatorul trebuie să rezolve două probleme: să aleagă o funcție de dispersie și să opteze pentru o metodă de rezolvare a coliziunilor. Aceasta din urmă se împarte în două mari clase :

- rezolvarea coliziunilor prin înlănțuire (când locația $T[h(k)]$ conține o listă înlănțuită a tuturor înregistrărilor cu chei ce au colizionat cu k) și
- rezolvarea coliziunilor prin adresare directă, când se folosesc metode mai sofisticate de căutare a unei locații libere în T pentru o cheie care colizionează.

Ne vom ocupa pe rând de cele 2 tipuri de probleme.



Funcții de dispersie

Ne ocupăm în acest paragraf de prezentarea câtorva tehnici euristice de creare a unor funcții de dispersie bune. Am exprimat informal în paragraful precedent câteva cerințe pentru o asemenea funcție

$$h : U \rightarrow \{0, 1, \dots, m-1\}.$$

Vom lucra la **ipoteza hashingului simplu uniform (HSU)**, care, informal, spune că orice cheie este distribuită prin funcția h cu probabilitate egală în oricare din locațiile $\{0, \dots, m-1\}$. Mai precis, pentru oricare $h \in K$ probabilitatea ca $h(k) = i$ este $\frac{1}{m}$ pentru orice $i \in [0 \dots m-1]$, și este independentă de restul cheilor. Dacă P este probabilitatea de distribuție a cheilor din K în U , adică dacă $P(k)$ = probabilitatea de a extrage cheia k din U , atunci ipoteza HSU se formulează:

$$\sum_{\{k|h(k)=i\}} P(k) = \frac{1}{m}, \forall i = 0, 1, \dots, m-1$$

Deoarece P este de obicei necunoscută, nu este în general posibil să verificăm că ipoteza HSU este satisfăcută.

În cele ce urmează vom interpreta cheile și numerele naturale, adică vom considera că $U < N$. Atunci când cheile sunt stringuri, metoda folosită pentru a le converti la numere naturale este următoarea: se înlocuiește fiecare caracter cu codul său ASCII și se calculează valoarea șirului rezultat ca întreg în baza 128.

De exemplu: AB se scrie $6566_{128} = 65 \cdot 128 + 66 = 8386$ ABC se scrie $656667_{128} = 65 \cdot 128^2 + 66 \cdot 128 + 67 = 1073605$

Metoda diviziuni

Se numesc funcții de dispersie obținute prin metoda diviziunii funcțiile de forma

$$h : U \rightarrow \{0, \dots, m-1\}, h(k) = k \bmod m.$$

Ele sunt printre cele mai comune pentru că sunt ușor de calculat.

Alegerea unei astfel de funcții revine practic la alegerea lui m , dimensiunea tabelului de dispersie. Deci, ea va fi guvernată, pe de o parte, de dimensiunile lui K și de modul în care rezolvăm coliziunile. De exemplu, dacă ne propunem să le rezolvăm prin înălțuire și să menținem în medie 3 date într-o locație, m ar putea fi apropiat de $\frac{|K|}{3}$. Dacă însă vrem ca toate elementele din K să-și găsească locul în tabel, atunci am putea alege un m apropiat de $\frac{3|K|}{2}$, acceptând deci o risipă de spațiu egală cu $\frac{|K|}{2}$.

Dar acestea nu sunt singurele considerente care guvernează alegerea lui m . Important este ca cheile să fie bine dispersate pe mulțimea $\{0, 1, \dots, m-1\}$ și să avem cât mai puține coliziuni. Vom evita valori pentru m de forma lui 2, deoarece, dacă $m = 2$, atunci $h(k) = k \bmod 2^p$ reprezintă doar ultimii p biți

din scrierea binară a lui k . Am pierdut în felul acesta informația conținută în ceilalți biți, lucru nerecomandabil: dacă nu avem indicații contrare, trebuie să facem ca $h(k)$ să depindă de toți biții lui k . Numerele apropiate de puteri ale lui 2 sunt și ele de evitat. Un exemplu de rezultat care ne conduce la această decizie este următorul: dacă $m = 2^p - 1$ și cheile k sunt șiruri de caractere în baza 2^p , atunci 2 șiruri care diferă doar printr-o transpoziție a două caractere adiacente vor avea aceeași valoare prin h .

Vom evita și valori pentru m de forma puterilor lui 10, căci dacă am alege un asemenea m , atunci $h(k)$ nu ar fi decât de o parte a caracterelor ce apar în scrierea lui k în baza 10. Valori bune pentru această metodă sunt m , un număr prim și cât mai departe de puteri ale lui 2. Evident, ne putem propune să indexăm tabelul după mulțimea de indici $\{1, 2, \dots, m\}$, în caz în care funcțiile din această clasă vor fi de forma $h(k) = k \bmod m + 1$.

Să considerăm exemplul firmei cu 68 de angajați, identificați cu chei întregi cu 4 cifre în scrierea în baza 10. Să presupunem că alegem pentru tabel o dimensiune apropiată de $\frac{3|K|}{2} = \frac{3}{2} \cdot 68$. O bună alegere pentru m ar fi $m = 97$, deoarece este prim, apropiat de $\frac{3}{2} \cdot 68$ și destul de departe de puteri ale lui 2.

Să considerăm cheile $k_1 = 3205$, $k_2 = 7148$ și $k_3 = 2345$.

Cu funcția de dispersie dată de metoda diviziunii

$$h(k) = k \bmod 97$$

Obținem pentru aceste chei următoarele adrese din tabel: $h(k_1) = 3205 \bmod 97 = 4$, $h(k_2) = 7148 \bmod 97 = 67$, $h(k_3) = 2345 \bmod 97 = 17$

Metoda multiplicării

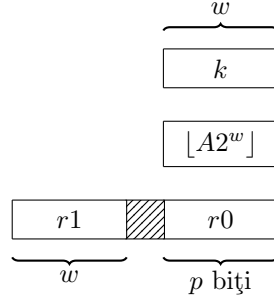
Fie A o constantă pozitivă subunitară, $0 < A < 1$. Să considerăm funcția dată de

$$h(k) = \lfloor m(kA \bmod 1) \rfloor$$

$kA \bmod 1$ reprezintă partea fracționară a lui kA , $kA - [kA]$, care se înmulțește cu m , iar rezultatul se trunchiază la cel mai apropiat întreg. Observăm că valoarea lui m nu e critică pentru această metodă, iar puterile lui 2 sunt alegeri bune pentru m , pentru că funcția se va calcula ușor. Fie $m = 2^p$. Atunci

$$h(k) = \lfloor 2^p(kA - [kA]) \rfloor = \lfloor kA2^p - 2^p[kA] \rfloor$$

Fie w lungimea unui cuvânt mașina în biți și presupunem că k se poate reprezenta pe w biți de forma



$$K \lfloor A2^w \rfloor = r_1 2^w + r_0$$

Deoarece $kA2^p = (kA2^w) 2^{p-w}$ înseamnă că $h(k) = (r_1 2^w + r_0) 2^{p-w}$.

Knuth recomandă pentru A numărul de aur, $\frac{\sqrt{5}+1}{2}$, a cărui valoare aproximată cu 7 zecimale este 1,6180339. Revenind la exemplul nostru, să alegem pentru această metodă $A = 0,6180339^1$ și $m = 2^6 = 64$. Atunci, pentru cele 3 chei din exemplu și funcția $h(k) = \lfloor m(kA \bmod 1) \rfloor$, avem: $h(k_1) = h(3205) = 51$ deoarece partea fracțională a lui $3205 \cdot A$ este 0,7986, care, înmulțit cu 64 ne dă 51,1107, a cărui parte întreagă este 51. Analog, obținem $h(k_2) = h(7148) = 45$ și $h(k_3) = h(2345) = 18$.

Metoda împachetării

Metoda împachetării sau folding generează funcții de dispersie în felul următor. Fiecare cheie k , despre care am făcut presupunerea că este întreg, este partiționată în bucăți de lungimi egale (eventual cu excepția ultimei), k_1, k_2, \dots, k_r . Dacă k este un întreg scris în baza 10 bucățile vor fi șiruri de caractere în baza 10, de aceeași lungime fixă, l , pe care din nou le putem considera întregi. O funcție de dispersie

$$h : U \rightarrow [0 \dots 10^2 - 1]$$

este dată de formula:

$$h(k) = (k_1 + k_2 + \dots + k_r) \bmod 10^l$$

ceea ce înseamnă că se adună cele r bucăți, ca întregi, și din sumă se păstrează doar l cifre, cele mai puțin semnificative. Alte funcții de dispersie se obțin inversând ordinea cifrelor în toate bucățile pare sau în toate bucățile impare. De exemplu:

$$h(k) = (\overline{k_1} + k_2 + \overline{k_3} + \dots) \bmod 10^l$$

unde k este întregul obținut prin inversarea ordinii caracterelor. Pe exemplul nostru metoda împachetării se aplică în felul următor: spargem fiecare cheie de

¹Pentru calcularea valorii $h(k)$ nu păstrăm decât partea fracționară a produsului kA . Prin urmare, este suficient să alegem A ca fiind partea fracționară a numărului $\frac{\sqrt{5}+1}{2}$, adică 0,6180339, care este chiar $\frac{\sqrt{5}+1}{2}$.

4 caractere în 2 bucăți de câte 2 caractere, pe care le adunăm și oțitem, dac e cazul, cifra cea mai nesemnificativ pentru ca rezultatul s aib lungimea tot de 2 caractere. Avem deci $h(k_1) = h(3205) = 32 + 5 = 37$ $h(k_2) = h(7148) = (71 + 48) \bmod 100 = 19$ $h(k_3) = h(2345) = 23 + 45 = 68$ Pentru funcția $h'(k) = k_1 + \overline{k_2}$, în care a doua funcție o inversm în oglind înainte de adunare, obținem $h'(3205) = 32 + \overline{05} = 32 + 50 = 82$ $h'(7148) = 71 + \overline{48} = (71 + 84) \bmod 100 = 55$ $h'(2345) = 23 + \overline{45} = 23 + 54 = 77$

Metoda ptratului

Presupunem c l este lungimea fixat a lui $h(k)$ ca întreg în baza 10. Fie

$$h : U \rightarrow [0 \dots 10^l], h(k) = (k^2 \text{div} 10^{c_1}) \bmod 10^{c_2}.$$

Se ridic la ptrat, iar din întregul lung astfel obținut se elimin c_1 cifre dintre cele mai puțin semnificative și c_2 cifre dintre cele mai semnificative, pstrndu-se exact l cifre din „centrul” lui k^2 . Cu alte cuvinte c_1 și l sunt fixate, iar $c_{\text{@}}$ se obține din formula:

$$l = \text{lung}(k^2) - c_1 - c_2.$$

S aplicm aceast metod exemplului nostru. O cheie de 4 cifre va avea ptratul de 7 sau 8 cifre, dintre care s eliminm $c_1 = 3$ dintre cele din dreapta, s pstrm $l = 2$ cifre, și s eliminm și restul de cifre „centrale”, adic a 4-a și a 5-a din stnga. $k_1 = 3205 : k_1^2 = 3205^2 = 10272025$ și $h(k_1) = 72$ $k_2 = 7148 : k_2^2 = 7148^2 = 51093904$ și $h(k_2) = 93$ $k_3 = 2345 : k_3^2 = 2345^2 = 5499025$ și $h(k_3) = 99$