

# Curs cu tema „Hashing”

## Grup țintă: elevii componenți ai centrului InfO(1)

*prof. Dana Lica*

Multe aplicații necesită o structură de date dinamică, care să permită executarea eficientă a operațiilor *Insert*, *Delete* și *Search*. O astfel de structură de date poartă numele de *dicționar*.

Tabelele *hash* sunt o variantă de implementare a unui dicționar, foarte ușor de implementat și frecvent utilizată în practică. În cazul cel mai defavorabil, o tabelă *hash* se poate comporta ca o listă simplu înlănțuită (deci operațiile specificate se execută în timp liniar). Dar în practică, utilizând *hashing* cele 3 operații se execută foarte eficient.

Să considerăm  $A$  o mulțime cu  $n$  valori care trebuie stocate în structura de date (aceste valori sunt denumite chei) și  $H$  un vector cu  $m$  componente (denumit tabelă *hash*).

Vom considera o funcție *hash* prin care se asociază fiecărei chei un număr întreg din intervalul  $[0, m-1]$ .

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

Funcția *hash* are proprietatea că valorile ei sunt uniform distribuite în intervalul specificat. Aceste valori vor fi utilizate ca indici în vectorul  $H$ . Mai exact, cheia  $x$  va fi plasată în tabela *hash* pe poziția  $h(x)$ .

### Funcții hash

O funcție *hash* bună trebuie să distribuie uniform cheile în tabela *hash*. Din păcate această condiție este greu de îndeplinit, pentru că nu cunoaștem structura cheilor.

În general, funcțiile *hash* sunt definite pe mulțimea numerelor naturale (cheile sunt considerate numere naturale). În cazul în care cheile nu sunt numere naturale, ele pot fi transformate în numere naturale.

De exemplu, pentru o cheie șir de caractere, putem considera că șirul este un număr natural scris în baza 128 (pentru codul ASCII) sau 256 (pentru codul ASCII extins).

## Metoda 1: Restul împărțirii la m

$$h(x) = x \% m$$

În acest caz este indicat să evităm ca m să fie de forma  $2^p$  (pentru că atunci  $h(x)$  ar avea ca valoare ultimii p biți ai lui x). Exceptând cazul în care știm că ultimii toate secvențele binare de lungime p apar cu aceeași probabilitate ca ultimi p biți ai cheii, este de preferat să utilizăm o funcție în care se utilizează toți biții cheii.

Un număr prim apropiat de o putere a lui 2 este adesea o bună alegere pentru m.

### Observație

Putem considera orice funcție *hash* de forma  $h(x) = (ax+b) \% m$ , cu  $a \neq 0$ .

Exemplu de funcție *hash* pentru chei șir de caractere:

```
int hashfunction(char *s)
{ int i;
  for (i=0; *s; s++) i = 131*i + *s;
  return( i % m ); }
```

## Metoda 2: Înmulțirea

Se înmulțește cheia x cu un număr subunitar  $0 < y < 1$ , considerăm partea fracționară a lui  $x*y$ , o înmulțim cu m și reținem doar partea întreagă a rezultatului.

$$h(x) = [m * (\text{frac}(x*y))] ]$$

Un avantaj al acestei metode este că nu există valori "rele" pentru m, prin urmare de obicei m poate fi o putere a lui 2, pentru a implementa eficient funcția hash.

## Implementarea celor trei operații în C++ folosind containerul <vector> din STL

```
#include <cstdio>
#include <vector>
#define Mod 666013

using namespace std;

vector <int> H[Mod];
vector <int> :: iterator it;
int i, n, op, x, ind;

vector <int> :: iterator find_v(int x){
    for(it= H[ind].begin(); it!=H[ind].end(); it++)
```

```

        if (*it==x) return it;
    return it;
}

void insert_v(int x){
    if (it == H[ind].end()) H[ind].push_back(x);
}

void erase_v(int x){
    if (it!=H[ind].end()) H[ind].erase(it);
}

void write_v (int x){
    if(it==H[ind].end()) printf("0\n"); else printf("1\n");
}

int main(){
    freopen ("hashuri.in","r", stdin);
    freopen("hashuri.out","w", stdout);

    scanf("%d\n", &n);

    for (int i=1; i<=n; i++){
        scanf("%d %d\n", &op, &x);

        ind = x % Mod;
        it = find_v(x);

        if(op==1) insert_v(x);
        if(op==2) erase_v(x);
        if(op==3) write_v(x);
    }

    return 0;
}

```

## Aplicația pe șiruri de caractere: Algoritmul Rabin-Karp

**ENUNȚ:** Se dau două șiruri  $A$  și  $B$  formate din litere mici și mari ale alfabetului englez și din cifre. Se cere găsirea tuturor aparițiilor șirului  $A$  ca subsecvență a șirului  $B$ .

```

1 function RabinKarp(string s[1..n], string sub[1..m])
2     hsub := hash(sub[1..m]); hs := hash(s[1..m])
3     for i from 1 to n-m+1
4         if hs = hsub
5             if s[i..i+m-1] = sub // compararea O(m)
6                 return i
7             hs := hash(s[i+1..i+m])
8     return not found

```

Consider șirul S= ABDCABAC de lungime N  
Și subșirul ABA de lungime M

Consider două funcții de dispersie pe care dacă le aplic asupra subșirului și obțin valorile:

```
H1=(A*101^2+B*101^1+C) % 666013 (unde A,B,C sunt inlocuite cu 65, 66,67 adica cu ASCII)
H2=(A*109^2+B*109^1+C) % 10003
for(i=1;i<=m) {
h1=(h1*101%666013+(int)sub[i])%666013; h2=h2*109%10003+(int)sub[i]%10002; }
(bazele le-am considerat nr. prime mai mari decât ,Z' respectiv 101 și 109)
```

Se traversează liniar șirul și se numără de câte ori subșirul curent la traversare obține aceleași două valori pe cele două funcții de dispersie. Dacă da, îl consider corect și-l număr.  
Dacă sunt pe poziția I în șirul S inițial, atunci pentru a calcula valoarea pentru subșirul care se termină pe poziția I, calculez în O(1) astfel:

Pe șirurile noastre luate ca exemplu,

s=ABDCABAC de lungime N  
și  
subșir ABA de lungime M

după prima șevantă verificată din șir ABD, cu funcția H1 se obține valoarea

$$V1 = (A \cdot 101^2 + B \cdot 101 + D) \% 666013$$

Din valoarea curentă V1 „iese”  $A \cdot 101^2$  și „intră” C. Calcul realizat în O(1):

Notez

$$K = 666013 \text{ și } Q = 10003$$
$$P = 101^{(M-1)} \text{ și } R = 109^{(M-1)};$$

Atunci algoritmul se poate scrie astfel:

```
for (i = M; i < N; i++)
{
V1 = ( (V1 - (S[ i - M]*P) % K + K) * 101 + S[i]) % K;
V2 = ( (V2 - (S[ i - M]*R) % Q + Q) * 109 + S[i]) % Q;
if (V1==H1 && V2==H2) nr++;
}
```

### Aplicații propuse pentru exercițiu:

<http://infoarena.ro/problema/eqs>

<http://infoarena.ro/problema/strmatch>