

Structuri de date arborescente

Arbori de compresie Huffman

În analiza algoritmilor pe care i-am studiat până acum prioritară a fost complexitatea timp. Acum ne punem problema elaborării unor algoritmi care să micșoreze spațiul necesar memorării unui fișier. Tehnicile de compresie sunt utile pentru fișiere text, în care anumite caractere apar cu o frecvență mai mare decât altele, pentru fișiere ce codifică imagini sau sunt reprezentări digitale ale sunetelor ori ale unor semnale analogice, ce pot conține numeroase motive care se repetă. Chiar dacă astăzi capacitatea dispozitivelor de memorare a crescut foarte mult, algoritmi de compresie a fișierelor rămân foarte importanți, datorită volumului tot mai mare de informații ce trebuie stocate. În plus, compresia este deosebit de utilă în comunicații, transmiterea informațiilor fiind mult mai costisitoare decât prelucrarea lor.

Una dintre cele mai răspândite tehnici de compresie a fișierelor text, care, în funcție de caracteristicile fișierului ce urmează a fi comprimat, conduce la reducerea spațiului de memorie necesar cu 20%-90%, a fost descoperită de D. Huffman în 1952 și poartă numele de *codificare (cod) Huffman*. În loc de a utiliza un cod în care fiecare caracter să fie reprezentat pe 7 sau pe 8 biți (lungime fixă), se utilizează un cod mai scurt pentru caracterele care sunt mai frecvente și coduri mai lungi pentru cele care apar mai rar.

Să presupunem că avem un fișier de 100.000 de caractere din alfabetul {a,b,c,d,e,f}, pe care dorim să-l memorăm cât mai compact. Dacă am folosi un cod de lungime fixă, pentru cele 6 caractere, ar fi necesari câte 3 biți. De exemplu, pentru codul:

	a	b	c	d	e	f
cod fix	000	001	010	011	100	101

ar fi necesari în total 300.000 biți.

Să presupunem acum că frecvențele cu care apar în text cele 6 caractere sunt:

	a	b	c	d	e	f
frecvență	45	13	12	16	9	5

Considerând următorul cod de lungime variabilă :

	a	b	c	d	e	f
cod variabil	0	101	100	111	1101	1100

ar fi necesari doar 224.000 biți (deci o reducere a spațiului de memorie cu aproximativ 25%).

Problema se reduce deci la a asocia fiecărui caracter un cod binar, în funcție de frecvență, astfel încât să fie posibilă decodificarea fișierului comprimat, fără ambiguități. De exemplu, dacă am fi codificat a cu 1001 și b cu 100101, când citim în fișierul comprimat secvența 1001 nu putem decide dacă este vorba de caracterul a sau de o parte a codului caracterului b.

Ideea de a folosi separatori între codurile caracterelor pentru a nu crea ambiguități ar conduce la mărirea dimensiunii codificării.

Pentru a evita ambiguitățile este necesar ca nici un cod de caracter să nu fie prefix al unui cod asociat al unui caracter (un astfel de cod se numește *cod prefix*).

Codul Huffman

D. Huffman a elaborat un algoritm *Greedy* care construiește un cod prefix optimal, denumit cod Huffman. Prima etapă în construcția codului Huffman este calcularea numărului de apariții ale fiecărui caracter în text. Există situații în care putem utiliza frecvențele standard de apariție a caracterelor, calculate în funcție de limbă sau de specificul textului.

Fie $C = \{c_1, c_2, \dots, c_n\}$ mulțimea caracterelor dintr-un text, iar f_1, f_2, \dots, f_n , respectiv, numărul lor de apariții. Dacă l_i ar fi lungimea șirului ce codifică simbolul c_i , atunci lungimea totală a reprezentării ar fi :

$$L = \sum_{i=1}^n l_i \cdot f_i$$

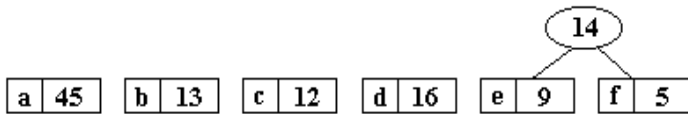
Scopul nostru este de a construi un cod prefix care să minimizeze această expresie. Pentru aceasta, construim un arbore binar complet în manieră *bottom-up* astfel :

1. Inițial, considerăm o partiție a mulțimii $C = \{ \{c_1, f_1\}, \{c_2, f_2\}, \dots, \{c_n, f_n\} \}$, reprezentată printr-o pădure de arbori formați dintr-un singur nod.
2. Pentru a obține arborele final, se execută $n-1$ operații de unificare. Unificarea a doi arbori A_1 și A_2 constă în obținerea unui arbore A , al cărui subarbore stâng este A_1 , subarbore drept A_2 , iar frecvența rădăcinii lui A este suma frecvențelor rădăcinilor celor doi arbori. La fiecare pas unificăm 2 arbori ale căror rădăcini au frecvențele cele mai mici.

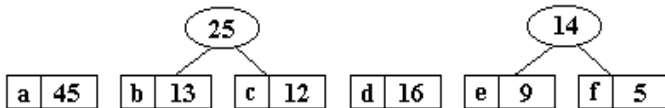
De exemplu, arborele Huffman asociat caracterelor $\{a, b, c, d, e, f\}$ cu frecvențele $\{45, 13, 12, 16, 9, 5\}$ se construiește pornind de la cele cinci noduri din figura 1:



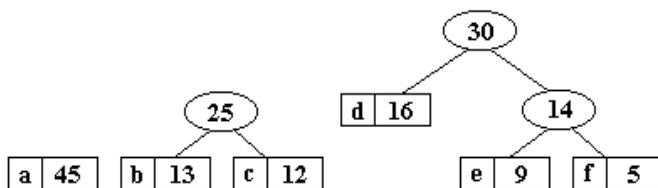
Pas 1: Unific arborii corespunzători lui e și f, deoarece au frecvențele cele mai mici:



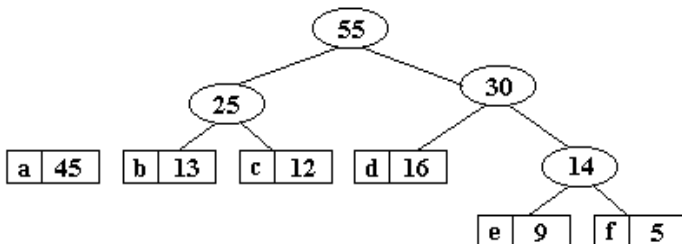
Pas 2: Unific arborii corespunzători lui b și c:



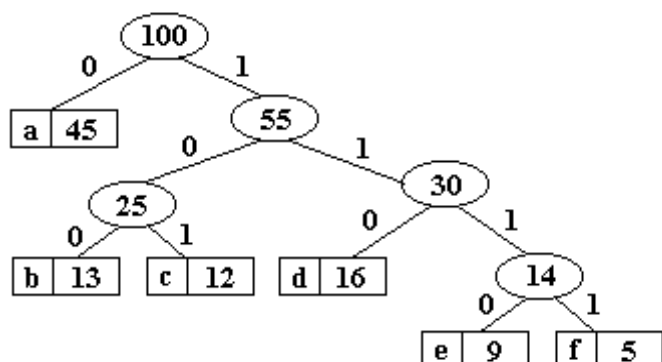
Pas 3: Unific arborele corespunzător lui d și arborele obținut la primul pas, cu rădăcina ce are frecvența 14:



Pas 4: Unific arborii ce au frecvențele 25 și 30.



Pas 5: Unificând ultimii doi arbori, obțin arborele Huffman asociat acestui set de caractere cu frecvențele specificate inițial.



Nodurile terminale vor conține un caracter și frecvența corespunzătoare caracterului; nodurile interioare conțin suma frecvențelor caracterelor corespunzătoare nodurilor terminale din subarbori.

Arborele Huffman obținut va permite asocierea unei codificări binare fiecărui caracter. Caracterele fiind frunze în arborele obținut, se va asocia pentru fiecare deplasare la stânga pe drumul de la rădăcină la nodul terminal corespunzător caracterului un 0, iar pentru fiecare deplasare la dreapta un 1.

Obținem codurile :

	a	b	c	d	e	f
cod	0	100	101	110	1110	1111

Observații

- Caracterele care apar cel mai frecvent sunt mai aproape de rădăcină și astfel lungimea codificării va avea un număr mai mic de biți.
- La fiecare pas am selectat cele mai mici două frecvențe, pentru a unifica arborii corespunzători. Pentru extragerea eficientă a minimumului vom folosi un min-heap. Astfel timpul de execuție pentru operațiile de extragere minim, inserare și ștergere va fi, în cazul cel mai defavorabil, de $O(\log n)$.

Algoritm de construcție a arborelui Huffman

Pas 1. Inițializare :

- fiecare caracter reprezintă un arbore format dintr-un singur nod;
- organizăm caracterele ca un min-heap, în funcție de frecvențele de apariție;

Pas 2. Se repetă de n-1 ori :

- extrage succesiv X și Y, două elemente din heap
- unifică arborii X și Y :
 - creează Z un nou nod ce va fi rădăcina arborelui
 - $Z \rightarrow st := X$
 - $Z \rightarrow dr := Y$
 - $Z \rightarrow frecv := X \rightarrow frecv + Y \rightarrow frecv$
- inserează Z în heap;

Pas 3. Singurul nod rămas în heap este rădăcina arborelui Huffman. Se generează codurile caracterelor, parcurgând arborele Huffman.

Analiza complexității

Inițializarea heapului este liniară. Pasul 2 se repetă de n-1 ori și presupune două operații de extragere a minimumului dintr-un heap și de inserare a unui element în heap, care au timpul de execuție de $O(\log n)$. Deci complexitatea algoritmului de construcție a arborelui Huffman este de $O(n \log n)$.

Corectitudinea algoritmului

Trebuie să demonstrăm că algoritmul lui Huffman produce un cod prefix optimal.

Lema 1.

Fie $x, y \in C$, două caractere care au frecvența cea mai mică. Atunci există un cod prefix optimal în care x și y au codificări de aceeași lungime și care diferă doar prin ultimul bit.

Demonstrație:

Fie T arborele binar asociat unui cod prefix optimal oarecare și fie a, b două noduri de pe nivelul maxim, care sunt fiii aceluiași nod. Putem presupune fără a restrânge generalitatea că $f[a] \leq f[b]$ și $f[x] \leq f[y]$. Cum x și y au cele mai mici frecvențe rezultă că $f[x] \leq f[a]$ și $f[y] \leq f[b]$. Transformăm arborele T în T' , schimbând pe a cu x :

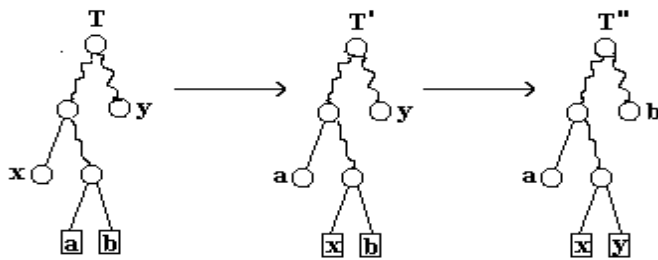


Fig. 7.

Notăm cu $C(T)$ costul arborelui T :

$$C(T) = \sum_{c \in C} f(c) \cdot \text{niv}_T(c)$$

$$C(T) - C(T') = f(x) \cdot \text{niv}_T(x) + f(a) \cdot \text{niv}_T(a) - f(x) \cdot \text{niv}_{T'}(x) - f(a) \cdot \text{niv}_{T'}(a) \Leftrightarrow$$

$$C(T) - C(T') = (f(a) - f(x)) \cdot (\text{niv}_T(a) - \text{niv}_T(x)) \geq 0, \text{ pentru că } \begin{cases} \text{niv}_{T'}(x) = \text{niv}_T(a) \\ \text{niv}_{T'}(a) = \text{niv}_T(x) \end{cases}$$

În mod analog, construim arborele T'' , schimbând pe x cu y . Obținem:

$$C(T') - C(T'') = (f(b) - f(y)) \cdot (\text{niv}_T(b) - \text{niv}_T(y)) \geq 0.$$

Deci $C(T) \geq C(T'')$, dar cum T era optimal deducem că T'' este arbore optimal și în plus x și y sunt noduri terminale situate pe nivelul maxim, fii ai aceluiași nod.

Observație

Din lema deducem că pentru construcția arborelui optimal putem începe prin a selecta după o strategie Greedy caracterele cu cele mai mici frecvențe.

Lema 2.

Fie T un arbore binar complet corespunzător unui cod prefix optimal pentru alfabetul C . Dacă x și y sunt două noduri, fii ai aceluiași nod z în T , atunci arborele T' obținut prin eliminarea nodurilor x și y este arborele binar complet asociat unui cod prefix optimal pentru alfabetul

$$C' = (C - \{x, y\}) \cup \{z\}, f[z] \text{ fiind } f[x] + f[y].$$

Demonstrație:

$$\begin{aligned} \forall c \in C - \{x, y\}, \text{ niv}_{T'}(c) &= \text{niv}_T(c) \\ \text{niv}_{T'}(x) &= \text{niv}_{T'}(y) = \text{niv}_T(z) + 1 \\ C(T) &= \sum_{c \in C} f(c) \cdot \text{niv}_T(c) = C(T') + (f(x) + f(y)) \\ &= (f(x) + f(y)) \cdot (\text{niv}_T(z) + 1) = f(z) \cdot \text{niv}_T(z) + f(x) + f(y) = \\ &= f(z) \cdot \text{niv}_{T'}(z) + (f(x) + f(y)). \end{aligned}$$

Deci,

$$C(T) = \sum_{c \in C} f(c) \cdot \text{niv}_T(c) = C(T') + (f(x) + f(y))$$

Dacă presupunem prin reducere la absurd că arborele T' nu este optimal pentru alfabetul $C' = (C - \{x, y\}) \cup \{z\} \Rightarrow \exists T''$, ale cărui frunze sunt caractere din C' , astfel încât $C(T'') < C(T)$.

Cum z este nod terminal în T'' , putem construi un arbore T_1 pentru alfabetul C , atârând pe x și y ca fii ai lui z .

$C(T_1) = C(T'') + f(x) + f(y) < C(T)$, ceea ce contrazice ipoteza că T este arbore optimal pentru C . Deci, T' este un arbore optimal pentru C' .

Corectitudinea algoritmului lui Huffman este o consecință imediată a celor două leme.

Observație

Metoda de compresie bazată pe arbori Huffman statici are o serie de dezavantaje:

1. Fișierul de compactat trebuie parcurs de două ori: o dată pentru a calcula numărul de apariții ale caracterelor în text, a doua oară pentru a comprima efectiv fișierul. Deci metoda nu poate fi folosită pentru transmisii de date pentru care nu este posibilă reluarea.

2. Pentru o dezarhivare ulterioară a fișierului, așa cum este și firesc, este necesară memorarea arborelui Huffman utilizat pentru comprimare, ceea ce conduce la mărirea dimensiunii codului generat.

3. Arborele Huffman generat este static, metoda nefiind capabilă să se adapteze la variații locale ale frecvențelor caracterelor.

Aceste dezavantaje au fost în mare parte eliminate în metodele care folosesc arbori Huffman dinamici, ideea fiind ca la fiecare nouă codificare, arborele Huffman să se reorganizeze astfel încât caracterul respectiv să aibă eventual un cod mai scurt.