

Proiectarea algoritmilor

Lucrare de laborator nr. 10

Paradigma *greedy*

Interclasarea optimală

Cuprins

1	Arbori binari ponderați pe frontieră	1
1.1	Descriere	1
1.2	Algoritm pentru construirea unui arbore cu lungimea externă ponderată minimă . . .	2
1.3	Implementarea algoritmului pentru construirea unui arbore cu lungimea externă ponderată minimă	2
2	Interclasarea optimală	3
2.1	Descriere	3
2.2	Algoritm	3
3	Sarcini de lucru și barem de notare	6

1 Arbori binari ponderați pe frontieră

1.1 Descriere

Considerăm arbori binari cu proprietatea că orice vârf are 0 sau 2 succesori și vârfurile de pe frontieră au ca informații (etichete, ponderi) numere, notate cu $info(v)$. Convenim să numim acești arbori ca fiind *ponderați pe frontieră*. Pentru un vârf v din arborele t notăm cu d_v lungimea drumului de la rădăcina lui t la vârfurile v . *Lungimea externă ponderată* a arborelui t este:

$$LEP(t) = \sum_{v \text{ pe frontiera lui } t} d_v \cdot info(v)$$

Modificăm acești arbori etichetând vârfurile interne cu numere ce reprezintă suma etichetelor din cele două vârfuri fii. Pentru orice vârf intern v avem $info(v) = info(v_1) + info(v_2)$, unde v_1, v_2 sunt fiii lui v (Figura 1).

Lema 1.1 Fie t un arbore binar ponderat pe frontieră. Atunci:

$$LEP(t) = \sum_{v \text{ intern în } t} info(v)$$

Lema 1.2 Fie t un arbore din $\mathcal{T}(x)$ cu LEP minimă și v_1, v_2 două vârfuri pe frontiera lui t . Dacă $info(v_1) < info(v_2)$ atunci $d_{v_1} \geq d_{v_2}$.

Lema 1.3 Presupunem $x_0 \leq x_1 \leq \dots \leq x_{n-1}$. Există un arbore în $\mathcal{T}(x)$ cu LEP minimă și în care vârfurile etichetate cu x_0 și x_1 (vârfurile sunt situate pe frontieră) sunt frați.

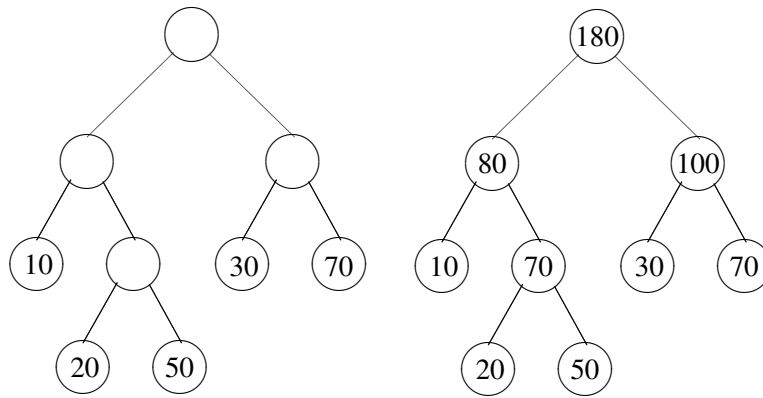


Figura 1: Arbore ponderat pe frontieră, înainte și după modificare

1.2 Algoritm pentru construirea unui arbore cu lungimea externă ponderată minimă

Ideea algoritmului rezultă direct din Lema 1.3. Presupunem $x_0 \leq x_1 \leq \dots \leq x_{n-1}$. Știm că există un arbore optim t în care x_0 și x_1 sunt memorate în vârfuri frate. Tatăl celor două vârfuri va memora $x_0 + x_1$. Prin ștergerea celor două vârfuri ce memorează x_0 și x_1 se obține un arbore t' . Fie $t1'$ un arbore optim pentru secvența $y = (x_0 + x_1, x_2, \dots, x_{n-1})$ și $t1$ arborele obținut din $t1'$ prin "agățarea" a două vârfuri cu informațiile x_0 și x_1 de vârful ce memorează $x_0 + x_1$. Avem $LEP(t1') \leq LEP(t')$ ce implică $LEP(t1) = LEP(t1') + x_0 + x_1 \leq LEP(t') + x_0 + x_1 = LEP(t)$. Cum t este optim, rezultă $LEP(t1) = LEP(t)$ și de aici t' este optim pentru secvența y . Considerăm în loc de secvențe de numere secvențe de arbori.

Notatii: $t(x_i)$ = arborele format dintr-un singur vârf etichetat cu x_i și $rad(t)$ = rădăcina arborelui t .

Premise: Inițial se consideră n arbori cu un singur vârf, care memorează numerele $x_i, i = 0, \dots, n-1$.

```

procedure lep(x, n)
  1:  B ← {t(x0), ..., t(xn-1)}
  2:  while (#B > 1) do
  3:    alege t1, t2 din B cu info(rad(t1)), info(rad(t2)) minime
  4:    construiește arborele t în care subarborii rădăcinii
  5:      sunt t1, t2 și info(rad(t)) = info(rad(t1)) + info(rad(t2))
  6:    B ← (B \ {t1, t2}) ∪ {t}
end

```

1.3 Implementarea algoritmului pentru construirea unui arbore cu lungimea externă ponderată minimă

- Dacă mulțimea B este implementată printr-o listă liniară, atunci în cazul cel mai nefavorabil operația 3 este are timpul de execuție $O(n)$, iar operația 6 are timpul de execuție $O(1)$.
- Dacă mulțimea B este implementată printr-o listă liniară ordonată, atunci în cazul cel mai nefavorabil operația 3 are timpul de execuție $O(1)$, iar operația 6 are timpul de execuție $O(n)$.
- Dacă mulțimea B este implementată printr-un *heap*, atunci în cazul cel mai nefavorabil operația 3 are timpul de execuție $O(\log n)$, iar operația 6 are timpul de execuție $O(\log n)$.

Concluzie: *heapul* este alegerea cea mai bună pentru implementarea mulțimii B .

2 Interclasarea optimală

2.1 Descriere

Se consideră m secvențe sortate a_0, \dots, a_{m-1} care conțin n_0, \dots, n_{m-1} , respectiv, elemente dintr-o mulțime total ordonată. Interclasarea celor m secvențe constă în execuția repetată a următorului proces:

Se extrag din mulțime două secvențe și se pune în locul lor secvența obținută prin interclasarea acestora.

Procesul se repetă până când se obține o singură secvență sortată cu cele $n_0 + \dots + n_{m-1}$ elemente.

Problema constă în determinarea unei alegeri pentru care numărul total de transferuri de elemente să fie minim.

Un exemplu este dat de *sortarea externă*

- Presupunem că avem de sortat un volum mare de date ce nu poate fi încărcat în memoria internă.
- Se partiționează colecția de date în mai multe secvențe ce pot fi ordonate cu unul dintre algoritmi de sortare internă.
- Secvențele sortate sunt memorate în fișiere pe suport extern.
- Sortarea întregii colecții se face prin interclasarea fișierelor ce memorează secvențele sortate.

Considerăm problema interclasării a două secvențe sortate:

Fie date două secvențe sortate $x = (x_0, \dots, x_{p-1})$ și $y = (y_0, \dots, y_{q-1})$ ce conțin elemente dintr-o mulțime total ordonată. Să se construiască o secvență sortată $z = (z_0, \dots, z_{p+q-1})$ care să conțină cele $p + q$ elemente ce apar în x și y .

Utilizăm notația $z = \text{merge}(x, y)$ pentru a nota faptul că z este rezultatul interclasării secvențelor x și y . Numărul de comparații executate de algoritmul $\text{merge}(x, y)$ este cel mult $p + q - 1$, iar numărul de elemente transferate este $p + q$.

Revenim la problema interclasării a m secvențe. Considerăm un exemplu: Fie $m = 5, n_0 = 20, n_1 = 60, n_2 = 70, n_3 = 40, n_4 = 30$. Un mod de alegere a secvențelor pentru interclasare este următorul:

$$\begin{aligned} b_0 &= \text{merge}(a_0, a_1) \\ b_1 &= \text{merge}(b_0, a_2) \\ b_2 &= \text{merge}(a_3, a_4) \\ b &= \text{merge}(b_1, b_2) \end{aligned} \quad (\text{vezi Figura 7.a})$$

Numărul de transferuri al acestei soluții este $(20 + 60) + (80 + 70) + (40 + 30) + (150 + 70) = 80 + 150 + 70 + 220 = 520$.

Există alegeri mai bune? Răspunsul este afirmativ!!

2.2 Algoritm

Unei alegeri i se poate atașa un arbore binar în modul următor:

- informațiile din vârfuri sunt lungimi de secvențe;
- vârfurile de pe frontieră corespund secvențelor inițiale a_0, \dots, a_{m-1} ;
- vârfurile interne corespund secvențelor intermediare.

Se observă ușor că aceștia sunt arbori ponderați pe frontieră și numărul de transferuri de elemente corespunzător unei alegeri este egală cu LEP a arborelui asociat. Așadar, alegerea optimă corespunde arborelui cu LEP minimă.

Pentru exemplul anterior ($m = 5, n_0 = 20, n_1 = 60, n_2 = 70, n_3 = 40, n_4 = 30.$), soluția optimă dată de algoritmul greedy este:

$$\begin{aligned} b_0 &= \text{merge}(a_0, a_4) \\ b_1 &= \text{merge}(a_3, b_0) \\ b_2 &= \text{merge}(a_1, a_2) \\ b &= \text{merge}(b_1, b_2) \end{aligned} \quad (\text{vezi Figura 7.b})$$

Numărul de comparații este $50 + 90 + 130 + 220 = 490$.

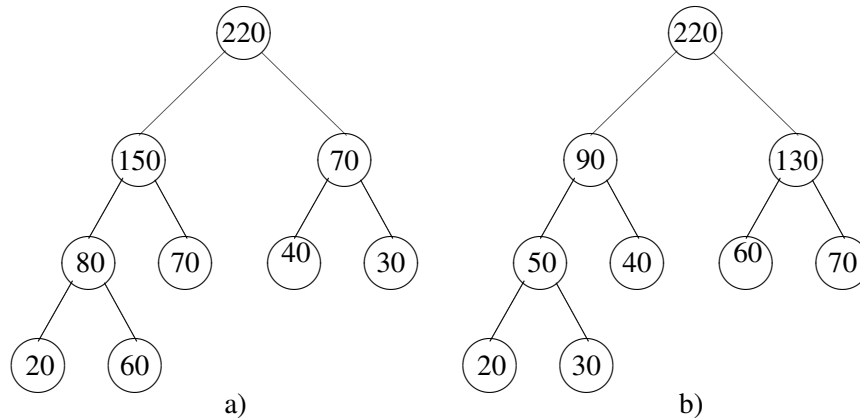


Figura 2: Arborii asociați celor două soluții de întrecere

Presupunem că intrarea este memorată într-un tablou T de structuri cu două câmpuri:

1. $T[i].secv$ conține adresa secvenței sortate a_i ;
2. $T[i].n$ conține numărul de elemente din secvența a_i .

Algoritmul care urmează utilizează reprezentarea arborilor prin tablouri de structuri.

Notăm cu H tabloul ce reprezintă arborele de interclasare. Tabloul H conține structuri formate din trei câmpuri:

1. $H[i].elt$;
2. $H[i].fst$ = indicele fiului de pe partea stângă;
3. $H[i].fdr$ = indicele fiului de pe partea dreaptă.

Semnificația câmpului $H[i].elt$ este următoarea:

1. dacă i este nod intern, atunci $H[i].elt$ reprezintă informația calculată din nod;
2. dacă i este pe frontieră (corespunde unui mesaj), atunci $H[i].elt$ este adresa din T a secvenței corespunzătoare.

Notăm cu $val(i)$ funcția care întoarce informația din nodul i , calculată ca mai sus.

Tabloul H , care în final va memora arborele de interclasare optimă, va memora pe parcursul construcției acestuia colecțiile intermediare de arbori.

În timpul execuției algoritmului de construcție a arborelui, H este compus din trei părți (Figura 4):

Partea I: un min-heap care va conține rădăcinile arborilor din colecție;

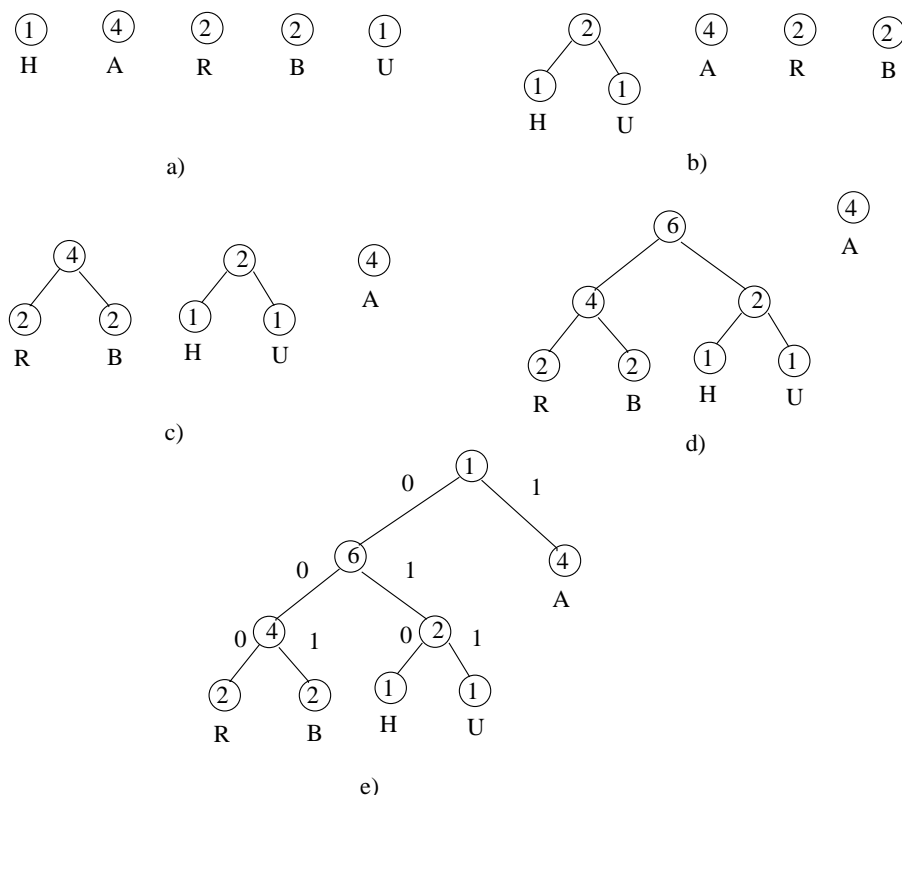


Figura 3: Construcția arborelui Huffman pentru HARABABURA

heap-ul rădăcinilor	noduri care nu sunt rădăcini	zonă vidă
---------------------	------------------------------	-----------

Figura 4: Organizarea tabloului H

Partea a II-a: conține nodurile care nu sunt rădăcini;

Partea a III-a: zonă vidă în care se poate extinde partea din mijloc.

Un pas al algoritmului de construcție ce realizează selecția *greedy* presupune parcurgerea următoarelor etape:

1. Mutarea rădăcinii cu informația cea mai mică pe prima poziție liberă din zona a treia, să zicem k . Aceasta este realizată de următoarele operații:

- a) copierea rădăcinii de pe prima poziție din heap pe poziția k :

$$H[k] \leftarrow H[1]$$

$$k \leftarrow k + 1$$

- b) mutarea ultimului element din heap pe prima poziție:

$$H[1] \leftarrow H[m]$$

$$m \leftarrow m - 1$$

- c) refacerea min-heapului.

2. Copierea rădăcinii cu informația cea mai mică pe prima poziție liberă din zona a treia, fără a o elimina din min-heap:

$$H[k] \leftarrow H[1]$$
$$k \leftarrow k + 1$$

3. Construirea noii rădăcini și memorarea acesteia pe prima poziție în min-heap (în locul celei copiate anterior).
4. Refacerea min-heapului.

Algoritmul rezultat are timpul de execuție $O(n \log n)$.

3 Sarcini de lucru și barem de notare

Sarcini de lucru:

1. Scrieți o funcție C/C++ care implementează un algoritm de construcție a arborelui de interclasare optimală a unei mulțimi de secvențe sortate.
2. Date fiind n secvențe sortate, scrieți un program care să determine o alegere optimă în cazul interclasării a n secvențe sortate

Barem de notare:

1. Implementarea algoritmului de construcție a arborelui de interclasare optimală: 6p
2. Determinarea unei alegeri optime în cazul interclasării a n secvențe sortate: 3p
3. Baza: 1p

Temă suplimentară:

1. Considerăm un graf $G = (V, E)$. Scrieți un program C/C++ pentru determinarea drumurilor minime între un vârf sursă și celelalte vârfuri.

Bibliografie

- [1] Lucanu, D. și Craus, M., *Proiectarea algoritmilor*, Editura Polirom, 2008.
- [2] Moret, B.M.E. și Shapiro, H.D., *Algorithms from P to NP: Design and Efficiency*, The Benjamin/Cummings Publishing Company, Inc., 1991.