

I. ANALIZA ALGORITMILOR

PRELIMINARII

Abu Ja`far Mohammed ibn Musa al-Khowarizmi (autor persan, sec. VIII-IX), a scris o carte de matematică cunoscută în traducere latină ca "Algorithmi de numero indorum", iar apoi ca "Liber algorithmi", unde "algorithm" provine de la "al-Khowarizmi", ceea ce literal înseamnă "din orașul Khowarizm". În prezent, acest oraș se numește Khiva și se află în Uzbekistan. Atât al-Khowarizmi, cât și alți matematicieni din Evul Mediu, înțelegeau prin *algoritm* o regulă pe baza căreia se efectuau calcule aritmetice. Astfel, în timpul lui Adam Riese (sec. XVI), algoritmi foloseau la: dublări, înjumătățiri, înmulțiri de numere. Alți algoritmi apar în lucrările lui Stifer ("Arithmetica integra", Nürnberg, 1544) și Cardano ("Ars magna sive de reguli algebraicis", Nürnberg, 1545). Chiar și Leibniz vorbește de "algoritmi de înmulțire". Termenul a rămas totuși multă vreme cu o întrebuințare destul de restrânsă, chiar și în domeniul matematicii.

Kronecker (în 1886) și Dedekind (în 1888) semnează actul de naștere al teoriei funcțiilor recursive. Conceptul de recursivitate devine indisolubil legat de cel de algoritm. Dar abia în deceniile al treilea și al patrulea ale secolului nostru, teoria recursivității și algoritmilor începe să se constituie ca atare, prin lucrările lui Skolem, Ackermann, Sudan, Gödel, Church, Kleene, Turing, Peter și alții.

1. ALGORITMI

Un algoritm constă într-o procedură de calcul bine definită care, într-o succesiune finită de pași, transformă o mulțime de valori (date de intrare) într-o altă mulțime de valori (date de ieșire sau rezultatul prelucrării).

Un algoritm constituie un mijloc de rezolvare a problemelor de calcul bine definite. Un enunț corect al unei astfel de probleme descrie, în termeni generali, relația intrare/ieșire corespunzătoare algoritmului problemei.

Conceptele de funcție calculabilă și algoritm formează obiectele principale de studiu în teoria calculabilității. Ele apar în diferite domenii ale informaticii și matematicii sub diverse grade de abstractizare.

Dezvoltarea unor științe implică, printre altele, construirea unor algoritmi care să rezolve o gamă cât mai variată de probleme.

Noțiunea de algoritm se întâlnește încă din clasele elementare relativ la cele patru operații fundamentale aplicate numerelor naturale reprezentate în baza zece. Mai târziu, se studiază algoritmul lui Euclid și algoritmi cu privire la rezolvarea unei ecuații, determinarea inversei unei matrice nesingulare etc.

Exemplu 1.1:

Să descriem un algoritm de înmulțire a două numere naturale x, y numit înmulțirea „La Russe”.

P1. Se consideră un tablou bidimensional cu două coloane. Pe prima linie se află x în prima coloană și y în a doua coloană.

P2. Numărul din prima coloană se împarte la 2 (împărțire întreagă), iar cel din a doua coloană se înmulțește cu 2.

P3. Pasul P2 se repetă până când în prima coloană se obține numărul natural 1.

P4. Se vor elimina liniile care conțin în prima coloană un număr par.

P5. Se vor însuma numerele naturale din coloana a doua care corespund liniilor rămase.

Suma obținută reprezintă produsul xy .

Fie numerele $x = 245$ și $y = 37$. Conform algoritmului descris se obține:

DEÎNMULȚIT	ÎNMULȚITOR
245	37
122**	74**
61	148
30**	296**
15	592
7	1184
3	2368
1	4736
xy	9065

Liniile marcate cu (**) se consideră eliminate.

Asupra algoritmului descris se pun o serie de întrebări.

a) **Condiția de oprire**

Răspunsul algoritmului constă în furnizarea produsului xy . Ultima linie din ciclu este cea care conține 1 în prima coloană. Odată cu ieșirea din ciclu, algoritmul se încheie.

b) Corectitudinea

Algoritmul calculează corect produsul celor două numere naturale. Se observă că dacă marcăm cu 1 liniile rămase și cu 0 pe cele eliminate, se obține numărul x scris în ordine inversă în binar ($245_{(10)} = 11110101_{(2)}$). De altfel acest algoritm calculează produsul în binar a două numere naturale.

c) Durata execuției (viteza de execuție)

Acest algoritm are aceeași viteză de execuție ca și algoritmul standard de înmulțire a două numere naturale (efectuând înmulțirea în binar). Există și alți algoritmi cu o viteză mai mare care calculează un astfel de produs.

d) Dimensiunea spațiului de memorie alocat etc.

Teoria mulțimilor constituie un instrument de descriere simplă și precisă a conceptului de funcție. Construirea unui model matematic pentru un astfel de concept este relativ ușor de realizat fiindcă avem de-a face cu un proces static. Descrierea conceptului de algoritm este mai dificilă deoarece construirea unui model matematic pentru un astfel de concept presupune un proces dinamic.

Au existat diverși matematicieni ca: E. Post, K. Gödel, A.M. Turing, A. Church etc. care, în mod independent, au dat câte o descriere a noțiunii de algoritm.

Clasa tuturor algoritmilor, independent de modul de reprezentare a lor, posedă unele trăsături generale.

Pentru un algoritm trebuie precizat:

a) domeniul algoritmului;

b) descrierea propriu-zisă a algoritmului.

Domeniul algoritmului este o mulțime cel mult numărabilă de elemente asupra cărora acționează algoritmul respectiv în care se includ și rezultatele finale.

Descrierea propriu-zisă a algoritmului constă în instrucțiuni care se execută asupra elementelor de intrare prezente în domeniul algoritmului. Fiecare din aceste instrucțiuni pune în evidență un număr finit de operații elementare.

Iată descrierea algoritmului „La Russe” :

RUSSE(A, B)

1. **arrays** X, Y {inițializare}
2. $X[1] \leftarrow A; Y[1] \leftarrow B$
3. $i \leftarrow 1$ {se construiesc cele două coloane}
4. **while** $X[i] > 1$ **do**
5. $X[i+1] \leftarrow X[i] \text{ div } 2$ {div reprezintă împărțirea întreagă}
6. $Y[i+1] \leftarrow Y[i] + Y[i]$
7. $i \leftarrow i+1$
8. {adună numerele $Y[i]$ corespunzătoare numerelor $X[i]$ impare}
9. $prod \leftarrow 0$
10. **while** $i > 0$ **do**
11. **if** $X[i]$ este impar
12. **then** $prod \leftarrow prod + Y[i]$
13. $i \leftarrow i-1$
14. **return** $prod$

Acțiunea unui algoritm este subordonată:

a) unor reguli de operare asociate;

b) agentului de calcul.

Regulile de operare asociate unui algoritm solicită fiecărei instrucțiuni să indice locațiile unde se află datele de intrare cât și pe cele în care se depun rezultatele intermediare și finale.

Agentul de calcul poate fi uman, mecanic, electronic etc. și are rolul de a executa instrucțiunile prezente în algoritm.

O problemă rezolvabilă algoritmic (problemă pentru care există algoritm care să o rezolve) o gândim ca o funcție P total definită pe mulțimea tuturor datelor de intrare sau informațiilor inițiale DI (cel mult numărabilă), cu valori în mulțimea datelor (informațiilor) finale DF (cel mult numărabilă). O *instanță* a unei astfel de probleme constă într-un element dat $x \in DI$ care satisface condițiile impuse de enunțul problemei și este capabil să conducă la o soluție a acesteia. Dacă $P(x) \downarrow$ (există) și $P(x) = y$, se spune că y este *soluția* problemei P pentru intrarea x . În caz contrar, se spune că P nu are soluție pentru intrarea x .

În vederea prelucrării pe calculator, un astfel de algoritm este transpus într-un program scris într-un limbaj de programare. Acest program evaluează o funcție naturală. Dacă domeniul de definiție al funcției constituie *intrarea* pentru programul respectiv, atunci valorile funcției sunt de fapt rezultatele programului.

Exemplul 1.2.

Să considerăm problema identificării unei chei numerice date într-un tablou numeric unidimensional.

Datele de intrare sau intrarea algoritmului: dimensiunea n a tabloului, tabloul numeric $X = (x[1], \dots, x[n])$ și valoarea numerică y (cheia);

Datele de ieșire sau ieșirea algoritmului: poziția aceluiași element din tablou care coincide cu cheia dată y sau un mesaj de neidentificare în cazul în care această cheie nu coincide cu nici un element din tablou.

O instanță a problemei propuse constă într-un tablou numeric dat, fie el $X = (-11, 121, 1, 0, 7, -48)$ de dimensiune 6 și o cheie numerică, fie ea $y = 7$. Conform acestei instanțe algoritmul va furniza ca rezultat faptul că elementul $x[5]$ coincide cu cheia dată (soluție a problemei).

Un algoritm este **corect** dacă pentru orice instanță a sa el se încheie cu o ieșire corectă care este soluție a problemei de calcul rezolvată prin algoritmul respectiv. Unii algoritmi, deși sunt incorecți în sensul că nu conduc la nici o soluție în timp finit sau conduc la soluții eronate, ei pot fi utili atât timp cât erorile produse de ei pot fi controlate.

Comportarea aceluiași algoritm poate fi diferită în funcție de datele de intrare. De aceea se impune o mare atenție asupra acestora din urmă. O variabilă prezentă în intrarea unui algoritm poate identifica un tablou sau un obiect al unei date compuse și va fi tratată ca un pointer la elementele tabloului, respectiv la câmpurile (atributele) obiectului corespunzător.

Modul în care se definește **dimensiunea datelor de intrare** depinde de problema de calcul analizată. Ea poate fi exprimată prin:

a) numărul de elemente cuprinse în datele de intrare (de exemplu, dimensiunea unui tablou de numere întregi);

b) numărul total de biți din reprezentarea binară a datelor de intrare (vezi algoritmul de înmulțire a două numere mari întregi);

c) două numere naturale (de exemplu, pentru reprezentarea unui graf se solicită numărul de vârfuri și numărul de muchii ale sale);

d) valoarea numerică a intrării (pentru algoritmi numerici, de exemplu algoritmul „La Russe”).

Pentru fiecare algoritm care rezolvă o anumită problemă P este necesar a se preciza modul de exprimare a dimensiunii datelor de intrare.

Dacă există un algoritm care rezolvă problema P nu înseamnă că el este unic.

De exemplu, există algoritmi ca: QuickSort (sortare rapidă), MergeSort (sortare prin fuziune), TreeSort (sortare prin arbori), sortare prin selecție și inserție etc. care sunt utilizați în același scop.

Prin urmare, apare necesitatea alegerii unui algoritm din clasa de algoritmi care rezolvă problema P care să corespundă unor cerințe. Algoritmul depinde de aplicație, implementare, mediu, frecvența utilizării etc. Compararea algoritmilor este un proces subtil care are în vedere mai multe aspecte. În continuare, vom căuta să analizăm câteva din aceste aspecte care joacă un rol important în elaborarea unui algoritm.

2. COMPLEXITATEA CALCULULUI

Teoria complexității constituie un domeniu foarte important al teoriei algoritmilor care investighează rezolvabilitatea individuală a problemelor sub anumite restricții de resurse.

Teoria complexității calculului cuprinde:

1) Un studiu analitic al complexității unor clase de probleme în raport cu diverse tipuri concrete de mașini matematice care le rezolvă;

2) Un studiu calitativ al complexității algoritmilor care rezolvă astfel de probleme.

Mai precis, **complexitatea calculului** se bazează pe legăturile existente între **problemele rezolvabile algoritmic** și **resursele de calcul**.

Înainte de a trece la analiza unui algoritm, trebuie să se cunoască o tehnologie de implementare a acestuia care va include un model pentru resursele sale și costurile corespunzătoare. Așadar, **resursele de calcul** sunt introduse prin definirea unui model de calcul care, de regulă, reprezintă o mașină. În acest sens pot fi utilizate: **mașini Turing de diverse tipuri, mașini cu acces aleatoriu (random-acces machine, RAM), circuite booleene** etc. Resursele tipice sunt: **timp, spațiu, benzi, capete de citire/scriere, nivele într-un circuit** etc.

Pentru tratarea celor două elemente fundamentale, **problemă și mașină**, considerăm trei nivele abstracte în complexitatea calculului.

a) Un nivel înalt care pune în evidență **complexitatea abstractă** a calculului, cunoscută sub numele de „**teoria Blum de complexitate a calculului**“, care este independentă atât de problemă cât și de mașină. O astfel de teorie este construită pe o enumerare efectivă a tuturor funcțiilor parțial recursive, o enumerare a așa numitelor funcții „cost-măsură“ și două axiome de bază (axiomele lui Blum).

b) Un nivel mediu care pune în evidență **complexitatea structurală**, independentă de problemă dar dependentă de mașină. Scopul unei astfel de complexități este de a studia capabilitatea unor modele de calcul concrete pentru resurse sau combinații de resurse, independent de problema propusă.

c) Un nivel inferior ce indică o *complexitate concretă* care depinde atât de mașină cât și de problemă. Scopul său este de a detecta o bună evaluare, cu un cost pe cât posibil minim, folosind o resursă (sau combinații de resurse) într-un model precizat pentru rezolvarea problemei.

În general, dificultatea unei probleme P se poate măsura prin resursele de calcul limitate asociate algoritmului secvențial care o rezolvă. Aceste resurse sunt *timpul și spațiul de memorie* necesare pentru execuția algoritmului respectiv. O astfel de măsură se poate evalua în diverse moduri. În termenii mașinii Turing, ca model de calcul, complexitatea temporală este stabilită de numărul de deplasări necesare realizării calculului respectiv, iar pentru un computer digital, de numărul de ciclări ale mașinii sau timpul real solicitat de mașină pentru acel calcul. În ceea ce privește complexitatea spațială, pentru mașina Turing, ea constă în numărul de celule (pătrate) ale benzii folosite în calcul, iar pentru computer, în numărul de bytes utilizați.

Prin urmare, ideea de *complexitate a unui algoritm* poate fi privită sub două aspecte:

a) *static*, prin structura sa.

b) *dinamic*, prin durata execuției sale.

Complexitatea spațială se exprimă prin dimensiunea spațiului de memorie necesar algoritmului, independent de comportarea pe care o are el pe parcursul execuției.

Complexitatea temporală se exprimă prin timpul necesar execuției sale. Axiomele lui Blum, de exemplu, pun în evidență măsuri de complexitate dinamică cu consecințe ca: teorema corelării măsurilor, teorema accelerării, teorema lacunei etc.

Există numeroase probleme de optimizare în rețele de transport, informatică, electronică, logistica construcțiilor, teoria comunicațiilor (în codurile corectoare de erori), criptografie etc. Pentru algoritmi de rezolvare a acestor probleme este obligatoriu un studiu al complexității lor. Complexitatea algoritmilor relativ la coduri, reprezentarea compactă a mesajelor în vederea transmiterii lor prin diverse canale, permite o analiză asupra eficacității metodelor de codificare/decodificare. Relativ la criptografie, studiul complexității implică elaborarea unor tehnici care să permită cifrarea mesajelor sub o formă care să le asigure, pe cât posibil, inviolabilitatea.

Analiza unui algoritm conduce la stabilirea eficienței sale și implică, în principal, ideea de complexitate pentru că ea prevede resursele solicitate de algoritmul respectiv. Stabilirea complexității unui algoritm este necesară, pe de o parte, datorită faptului că mulți algoritmi sunt eliminați de practică deoarece necesită un timp de execuție și un spațiu de memorie foarte mari, iar pe de altă parte, se caută algoritmi cât mai performanți în vederea rezolvării cu calculatorul a unor probleme destul de complexe impuse de practică. Un calculator, oricât de performant ar fi, este în esență un automat finit, deci el are posibilități limitate. Performanța unui algoritm se stabilește conform unui criteriu care, de cele mai multe ori, îl constituie timpul de execuție al algoritmului respectiv.

Având în vedere faptul că pentru seturi de date de intrare diferite același algoritm folosește timpi de execuție diferiți este necesar a lua în considerație:

- *timpul în cazul cel mai favorabil* (durata minimă pentru execuția algoritmului);

- *timpul mediu* (raportul dintre suma timpilor necesari pentru toate seturile de date posibile și numărul acestor seturi);

- *timpul în cazul cel mai defavorabil* (durata maximă pentru execuția algoritmului).

Acesta din urmă oferă o margine superioară a timpului de execuție pentru toate intrările de o dimensiune fixă și reprezintă situația cea mai indicată pentru căutarea de informații într-o bază de date.

Să considerăm implementarea unui algoritm oarecare care să primească intrări de dimensiuni diferite. La o intrare de dimensiune n vom asocia o măsură $f(n)$ de folosire a resurselor sale care, de regulă, trebuie să se minimizeze. În mod normal, $f(n)$ va reprezenta timpul solicitat de algoritm pentru intrarea de dimensiune n în cazul cel mai defavorabil sau mediu (la alegerea noastră). Există o serie de factori care complică mai mult sau mai puțin calculul exact al lui $f(n)$.

De exemplu:

a) Timpul solicitat pentru execuția unei instrucțiuni individuale dintr-un program, cu aceeași instanță și pe aceeași mașină, poate varia.

b) Poate exista o mare variație de folosire a resurselor pentru intrări de o anumită dimensiune constantă n .

c) Cazul cel mai defavorabil poate fi mai rar întâlnit în practică, iar cazul mediu să fie nereprezentativ.

d) Unii algoritmi se pot executa mai bine pe anumite clase de intrări.

În acest sens se impun câteva sugestii, cum ar fi:

- Identificarea operațiilor abstracte folosite de algoritm. Pentru a obține o maximizare a independenței de mașină este necesară o analiză a acestor operații abstracte, pentru că resursele solicitate de o mare parte din ele vor fi neimportante. Unele din ele sunt utilizate o singură dată la inițializare.

- Un algoritm, în general, posedă cel puțin un ciclu. Acesta trebuie identificat pentru că instrucțiunile din interiorul său se execută mult mai des și ele vor juca un rol important în analiza timpului de execuție a algoritmului. Prin contorizarea acestor instrucțiuni se determină o margine superioară corespunzătoare pentru valoarea timpului de execuție în cazul cel mai defavorabil și, dacă este posibil, o configurare a cazului mediu. Această limită superioară este necesară pentru că, în mod practic, nu poate fi găsită o valoare exactă a timpului de execuție în cazul cel mai defavorabil.

- Timpul necesar unui algoritm pentru rezolvarea unei probleme cu o instanță dată, reprezintă numărul operațiilor elementare (primitive) efectuate de algoritm pentru acea instanță. Se acceptă o astfel de interpretare deoarece se presupune că execuția unei astfel de operații se produce într-un timp constant și nu depinde de mărimea operanzilor și nici de timpul de memorare a rezultatului. De aceea resursa timp asociată se analizează independent de sistemul de calcul pe care se implementează algoritmul respectiv.

- Deși progresele tehnologice actuale fac ca necesarul de memorie să treacă pe un plan secund, există totuși situații în care spațiul de memorie utilizat este folosit ca argument în stabilirea unor limite inferioare pentru timpul de execuție.

- Viteza de calcul ce caracterizează actuala generație de calculatoare are drept consecință faptul că problema timpului de execuție se pune, în general, pentru valori foarte mari ale dimensiunii intrării. Aceasta conduce la ideea unei analize a termenului dominant (cel care tinde cel mai repede la infinit) din formula de exprimare a numărului de operații elementare executate de către algoritm.

Dimensiunea informațiilor prelucrate în majoritatea problemelor pe care le întâlnim în practică nu rămâne constantă.

Pentru a înțelege noțiunea de ”procedeu mecanic de calcul“ au fost propuse numeroase modele de calcul formal. Conform celebrei teze a lui Church - tot ceea ce este intuitiv calculabil cu ajutorul unui algoritm este calculabil indiferent de modelul formal de calcul - aceste modele sunt echivalente. În mod clasic, în teoria complexității, se consideră drept model mașina Turing care prezintă dublul avantaj de a constitui un model de program pentru calculator, dar și un cronometru pentru timpul de execuție al său. Conform acestui model, complexitatea unui algoritm **A**, pentru orice dimensiune n a intrării, este o funcție $f(n)$ ce exprimă timpul maxim de execuție a algoritmului.

În descrierea algoritmilor, cel mai convenabil model este modelul RAM având în vedere și cele mai importante clase de complexitate studiate (timp și spațiu polinomial).

Există câteva funcții tipice de exprimare a timpului de execuție a unui algoritm. Putem obține, în mod obișnuit, o astfel de funcție printr-o relație de recurență.

Cel mai adesea se va întâlni funcția de forma

$f(n) = cg(n) + tm$ unde $c > 0$ este o constantă, n reprezintă dimensiunea intrării, iar tm este un ”termen mic” care este semnificativ doar pentru valori mici ale lui n sau pentru algoritmi sofisticăți.

Funcția $g(n)$ poate fi:

- constantă (algoritm se execută în timp constant);
- $\log n$ (algoritm se execută în timp log);
- n^m ($m = 0, 1, 2, \dots$) (algoritm se execută în timp polinomial);

Cazuri particulare:

a) n (algoritm se execută în timp liniar);

b) $n \log n$ (algoritm se execută în timp liniar log);

c) n^2 (algoritm se execută în timp pătratic);

d) n^3 (algoritm se execută în timp cubic);

• $n!$ (algoritm se execută în timp factorial);

• k^n ($k > 1$, constantă) (algoritm se execută în timp exponențial).

Exemplul 2.1.

Să considerăm un algoritm care efectuează operații de punere/extragere a unei date de 32 biți efectuate pe o stivă cu elemente întregi. Astfel de operații vor necesita un timp constant. Un algoritm care se execută într-un timp constant mic se zice că este ”**perfect**”.

Exemplul 2.2.

Să presupunem un algoritm a cărui intrare constă dintr-un șir de n caractere. El poate prelucra intrarea sa caracter cu caracter, solicitând același timp pentru fiecare astfel de caracter. În acest caz, $f(n) = n$. Un astfel de algoritm care se execută în timp liniar se zice că este ”**excelent**”.

Exemplul 2.3.

Un algoritm poate cicla direct intrarea sa formată din numere întregi sub forma

$$f(n) = n + f(n - 1).$$

Iterând această relație se obține $f(n) = n + f(n - 1) = n + ((n-1) + f(n - 2)) = \dots$

$= n + (n - 1) + (n - 2) + \dots + 2 + f(1) = n^2/2 + n/2 + k$, unde k este o constantă. Pentru un n suficient de mare, k și $n/2$ sunt mici față de n^2 . Un astfel de algoritm se execută într-un timp pătratic.

Un algoritm care solicită o prelucrare a celor n^2 perechi de caractere ce corespund unei intrări care constă dintr-un șir de n caractere, solicită tot un timp pătratic de execuție.

Pentru intrări de dimensiuni mari algoritmi care se execută în timp pătratic vor evolua mai lent.

Exemplul 2.4.

Un algoritm poate prelucra, la fiecare pas, jumătate din intrarea sa. Așadar,

$$f(n) = f(n/2) + 1$$

(cu aproximare în cazul în care n nu este o putere naturală a lui 2).

Fie $n = 2^x$. Atunci, $f(n) = f(2^x) = 1 + f(2^{x-1}) = 1 + (1 + f(2^{x-2})) = \dots = x + f(2^0) = x + k$, unde k este o constantă.

În acest caz, $f(2^x)$ este aproximativ x , iar $f(n)$ este aproximativ $\log_2 n$. Un astfel de algoritm se execută în timp

log. Acești algoritmi au o creștere foarte lentă și sunt foarte buni în practică.

Este foarte cunoscut algoritmul "Divide et Impera" folosit în Quicksort, Mergesort etc. Pentru acest algoritm, $f(n) = n + 2f(n/2)$.

Pentru $n = 2^x$, folosind un mic artificiu, se obține:

$$f(2^x)/2^x = 2^x/2^x + 2 * f(2^{x-1})/2^x = 1 + f(2^{x-1})/2^{x-1} = 1 + (1 + f(2^{x-2})/2^{x-2}) = \dots = x + f(2^0)/2^0 = x + k, \text{ unde } k \text{ este o constantă.}$$

Rezultă că $f(2^x) = 2^x(x + k) = 2^x x + tm$. Prin urmare, $f(n) = n \log_2 n + tm$. În concluzie, algoritmul respectiv se execută în timp **liniar log**.

Exemplul 2.5.

Să presupunem că intrarea unui algoritm constă într-o mulțime de n numere întregi. Se cere să se găsească o submulțime a acestei mulțimi cu proprietatea că suma elementelor sale este nulă.

Printr-o cercetare exhaustivă, în cazul cel mai defavorabil, se vor verifica toate cele 2^n submulțimi posibile. Un astfel de algoritm se va executa într-un timp exponențial.

Exemplul 2.6.

Un algoritm, prin care se cere obținerea tuturor anagramelor corespunzătoare unui cuvânt format din n caractere la intrare, se va executa într-un timp factorial, pentru că există $n!$ astfel de posibilități.

Funcția $n!$ crește aproximativ cu aceeași viteză ca n^n , dar mai repede decât 2^n .

În 1964, Cobham a introdus clasa P a **problemelor rezolvabile algoritmic în timp polinomial**, adică o problemă ce se rezolvă printr-un algoritm A care pentru orice număr natural n , funcția sa de complexitate este mărginită superior de un polinom cunoscut $p(n)$ în variabila n , adică $f(n) \leq kp(n)$, unde k este o constantă nenulă. În această clasă se întâlnesc așa numitele **probleme "facile"**. Se mai spune despre un algoritm de complexitate polinomială că reprezintă rezultatul unei cunoașteri detaliate a problemei pe care o rezolvă.

Despre un algoritm care nu este de complexitate polinomială se zice că se comportă "**exponențial**". În aceeași clasă a algoritmilor exponențiali sunt încadrați și algoritmi de complexitate $n \log n$ deși nu corespund în sensul strict matematic al noțiunii respective. Având în vedere viteza de creștere sau ordinul de creștere a timpului de execuție în raport cu n a funcțiilor exponențiale față de cele polinomiale, în general, algoritmi exponențiali sunt inutilizabili în practică. Despre un algoritm de complexitate exponențială se spune că este o variantă a unei enumerări totale a căilor de identificare a soluțiilor problemei respective.

Există unii algoritmi cu comportare exponențială care, pentru valori relativ mici ale lui n , sunt mai eficienți decât cei alternativi cu comportare polinomială. De asemenea, unii algoritmi exponențiali se comportă acceptabil pentru unele probleme particulare (vezi metoda simplex de rezolvare a problemelor de programare liniară).

În general, despre o problemă pentru care s-a dovedit că nu există un algoritm de complexitate polinomială pentru rezolvarea ei, adică nu este **tractabilă**, se spune că este **intractabilă**. Dificultatea unei astfel de probleme trebuie remarcată prin faptul că

timpul necesar pentru găsirea unei soluții este de ordin exponențial, adică funcția nu se poate exprima printr-o expresie care să poată fi mărginită superior de un polinom în variabila n .

Dacă se are în vedere ca model, mașina de calcul în paralel (capabilă de a efectua simultan un număr relativ mare de calcule independente), **clasa problemelor dificile** (cele care nu sunt facile) se poate împărți în două subclase:

a) subclasa problemelor pentru care nu există algoritmi nedeterminiști de complexitate polinomială pentru rezolvarea lor;

b) subclasa problemelor pentru care există algoritmi nedeterminiști de complexitate polinomială pentru rezolvarea lor (**NP**).

Trebuie observat faptul că majoritatea problemelor care par a fi dificile admit algoritmi nedeterminiști de complexitate polinomială pentru rezolvarea lor. Noțiunea de "**algoritm determinist**" trebuie înțeleasă în sensul că structura sa nu permite o alegere între mai multe căi posibile în vederea obținerii rezultatului dorit. Conceptul de "**algoritm nedeterminist**" trebuie înțeles în sensul că el se poate afla în mai multe stări independente care nu se pot alege după anumite criterii și nici genera aleator. Un algoritm se numește **nedeterminist polinomial** dacă complexitatea calculelor, efectuate pe orice cale a arborelui ce descrie ramificările procesului de căutare a soluției, este polinomială.

Având în vedere aceste interpretări, conceptul de **algoritm determinist polinomial** este evident.

Clasa problemelor "**dificile în sensul cel mai tare**" este formată din acele probleme pentru care s-a demonstrat că nu există algoritmi pentru rezolvarea lor (de exemplu, problema a zecea a lui Hilbert cu privire la rezolvabilitatea în numere întregi a ecuațiilor polinomiale).

3. COMPLEXITATE ASIMPTOTICĂ

O simplă caracterizare a eficienței unui algoritm constă în viteza de creștere a timpului său de execuție care oferă posibilitatea de a compara performanțele relative ale unor algoritmi alternativi.

În general, este destul de dificil de a determina expresia exactă ce definește timpul de execuție al unui algoritm în funcție de dimensiunea problemei, de aceea se stabilesc anumite limite între care el poate varia. Complexitatea se exprimă în funcție de această dimensiune.

Când datele de intrare au dimensiuni suficient de mari, astfel încât numai viteza de creștere a timpului de execuție să fie relevantă, vom studia eficiența *asimptotică* a algoritmilor. Mai exact, se va urmări viteza de execuție a mai multor algoritmi cu același obiectiv, sau, altfel zis, viteza de creștere a funcțiilor a căror expresie reprezintă timpul de execuție a unor astfel de algoritmi. În practică există un interes deosebit asupra modului de creștere a timpului de execuție a unui algoritm o dată cu creșterea dimensiunii datelor de intrare până la cazul limită în care această dimensiune crește la infinit. Se folosește, în acest sens, conceptul de *timp asimptotic de execuție a unui algoritm*. Pentru exprimarea unui astfel de timp se va folosi un limbaj riguros care cuprinde următoarele simboluri asimptotice: Θ , O , o , Ω , ω . Se vor folosi funcții al căror domeniu de definiție îl reprezintă mulțimea numerelor naturale N . Notăția se poate extinde, în unele situații, la domeniul numerelor reale sau restrânge la o submulțime a mulțimii numerelor naturale.

Fie \mathfrak{R} mulțimea funcțiilor nenegative care au ca domeniul de definiție mulțimea numerelor naturale $N = \{0, 1, 2, \dots\}$.

Fie funcțiile $f(n), g(n)$ $f(n) \in \mathfrak{R}, g(n) \in \mathfrak{R}$.

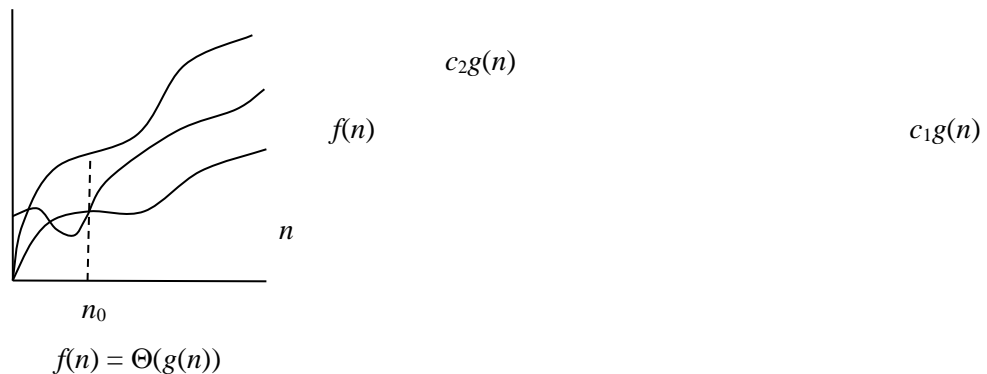
3.1. Θ - notația

Definiția 3.1.1. $f(n) = \Theta(g(n))$ ($n \rightarrow \infty$) dacă există constantele $c_1 > 0, c_2 > 0$ și $n_0 \in N$, astfel încât $\forall n \geq n_0, 0 \leq c_1 g(n) \leq f(n) \leq c_2 g(n)$.

Notația $n \rightarrow \infty$ indică, o dată în plus, faptul că relația asimptotică are loc pentru valori suficient de mari ale lui n . Într-o scriere asimptotică, chiar dacă nu este prezentă această notație cu scopul de a simplifica scrierea, ea este subînțeleasă.

Din punct de vedere matematic este corect a folosi notația $f(n) \in \Theta(g(n))$, pentru că $\Theta(g(n)) = \{f(n) / \text{există constantele } c_1 > 0, c_2 > 0 \text{ și } n_0 \in N$

, astfel încât $\forall n \geq n_0, 0 \leq c_1 g(n) \leq f(n) \leq c_2 g(n)$.



În literatura de specialitate, se folosește, printr-un abuz de utilizare a semnelui de egalitate, notația $f(n) = \Theta(g(n))$ ($n \rightarrow \infty$). Se acceptă ambele notații și pentru celelalte simboluri care vor urma. Se spune despre funcțiile $f(n)$ și $g(n)$ că au aproape aceeași viteză de creștere având în vedere constantele pozitive multiplicative c_1, c_2 sau că ele sunt funcții egale până la un factor constant. Vom spune că funcția $g(n)$ constituie o *margine asimptotică tare sau strânsă* (inferioară și superioară) pentru funcția $f(n)$.

Din definiția mulțimii $\Theta(g(n))$ se observă că este necesar ca elementele sale să fie funcții nenegative pentru valori ale lui n suficient de mari. O funcție asimptotic pozitivă este strict pozitivă pentru orice valoare a lui n suficient de mare ($n \geq n_0$).

În mod intuitiv, pe baza definiției simbolului Θ , într-o expresie polinomială ce reprezintă timpul de execuție a unui algoritm se ia în calcul doar termenul dominant (de ordin maxim) și se neglijează toți termenii de ordin inferior cât și coeficientul termenului dominant.

Fie $f(n) = (3/7)n^2 - 4n$. Pentru a arăta că $f(n) = \Theta(n^2)$ trebuie să determinăm constantele pozitive n_0, c_1, c_2 astfel încât

$0 \leq c_1 n^2 \leq (3/7)n^2 - 4n \leq c_2 n^2$. Pentru $\forall n \geq 1, c_2 \geq 3/7$ se observă că are loc inegalitatea $3/7 - 4/n \leq c_2$, iar pentru $\forall n \geq 1, c_1 \leq 1/35$ are loc inegalitatea $3/7 - 4/n \geq c_1$. În concluzie, există constantele

$c_1 = 1/35$, $c_2 = 3/7$ și $n_0 = 10$ astfel încât să aibă loc dubla inegalitate, adică $f(n) = \Theta(n^2)$. Adesea, dacă i se dă constantei c_1 o valoare ceva mai mică decât coeficientul termenului dominant, iar lui c_2 o valoare puțin mai mare decât același coeficient, definiția simbolului Θ ar putea fi îndeplinită. Coeficientul termenului dominant poate fi ignorat pentru că el ar modifica cele două constante c_1 și c_2 doar cu un factor constant.

Să observăm că aceste constante, pentru o aceeași funcție, nu sunt unice și ele depind de funcția respectivă.

Prin relația $0 \leq f(n) \leq cg(n)$, se înțelege că funcția $f(n)$ are ordinul de creștere al funcției $g(n)$ pentru valori ale lui n suficient de mari, până la un anumit factor constant c . Logaritmând această relație, se obține $\log(f(n)) \leq c_1 + \log(g(n))$ ($c_1 = \log c$, este o constantă). Logaritmând și relația $0 \leq c'g(n) \leq f(n)$, se obține $c_2 + \log(g(n)) \leq \log(f(n))$ ($c_2 = \log c'$ este o constantă). Se observă că $f(n) = \Theta(g(n))$ ($n \rightarrow \infty$) dacă și numai dacă există o constantă c astfel încât pentru un n suficient de mare, $\log(f(n))$ și $\log(g(n))$ să difere prin cel mult o astfel de constantă multiplicativă.

Se spune că implementarea unui algoritm se execută în timp **log**, respectiv în timp **liniar**, dacă funcția sa de complexitate $f(n)$ este din $\Theta(\log(n))$, respectiv $\Theta(n)$.

3.2. O - notația

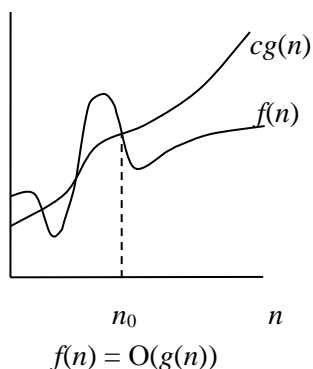
Definiția 3.2.1 $f(n) = O(g(n))$ ($n \rightarrow \infty$) dacă $\exists c > 0, \exists n_0 \in N$ astfel încât $\forall n \geq n_0, 0 \leq f(n) \leq c g(n)$.

Printr-o astfel de notație se pune în evidență faptul că funcția $f(n)$ crește strict mai lent sau cel mult la fel ca funcția $g(n)$.

Se mai spune că funcția $g(n)$ domină funcția $f(n)$ sau că funcția $f(n)$ este dominată de funcția $g(n)$. Cu alte cuvinte, prin notația $f(n) = O(g(n))$ ($n \rightarrow \infty$) înțelegem faptul că un multiplu constant al lui $g(n)$ constituie o margine asimptotică superioară a lui $f(n)$. Așadar

$$O(g(n)) = \{f(n) \in \mathfrak{R} / \exists c > 0, \exists n_0 \in N$$

astfel încât $0 \leq f(n) \leq g(n), \forall n \geq n_0\}$.



În literatura de specialitate există o distincție clară între o margine asimptotică superioară și o margine asimptotică tare (strânsă) pe care o oferă Θ . Notația O oferă o margine asimptotică superioară care nu este neapărat o margine asimptotică tare.

Exemplul 3.2.1.

$5n^2 = O(n^2)$ constituie o margine asimptotică superioară tare.

$5n = O(n^2)$ este o margine asimptotică superioară.

Din faptul că $f(n) = \Theta(g(n))$ rezultă că $f(n) = O(g(n))$ pentru că notația Θ este mai puternică decât notația O , adică $\Theta(g(n)) \subseteq O(g(n))$.

O funcție liniară $f(n) = an + b = \Theta(n)$ ($a > 0$). În schimb, folosind notația O , este adevărată și relația $f(n) = O(n^2)$. Într-adevăr trebuie să determinăm cele două constante strict pozitive n_0 și c astfel încât să aibă loc relația $0 \leq an + b \leq cn^2$ pentru orice $n \geq n_0$. Dacă dreapta intersectează parabola trebuie ca ecuația atașată să aibă rădăcini reale. De aceea, vom considera expresia $a^2 + 4bc$ de forma $(a + 2b)^2$. Printr-un calcul elementar se observă că pentru orice $n \geq n_0 = 1$ și $c = a + |b|$ ($c > 0$) relația este adevărată. Dacă dreapta nu intersectează parabola, relația cerută este evident adevărată.

Printr-o afirmație de genul „ **timpul de execuție al unui algoritm A este $O(n^k)$** “ se înțelege că timpul de execuție în cazul cel mai defavorabil (exprimat printr-o funcție de n) este $O(n^k)$, adică oricare ar fi datele de intrare de dimensiune n , pentru fiecare valoare a lui n , timpul de execuție respectiv este $O(n^k)$, unde k este o constantă întregă.

Propoziția 3.2.1. O condiție suficientă ca o funcție $g(n)$ să dea o comportare asimptotică de tip O pentru o funcție $f(n)$, adică $f(n) = O(g(n))$, este să existe o constantă $n_0 \in N$ astfel încât pentru orice $n \geq n_0$ să fie îndeplinite următoarele condiții:

a) $g(n) > 0$;

b) $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)}$ să existe și să fie finită.

Demonstrație: Fie $l = \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)}$. Aceasta înseamnă că $\forall \varepsilon > 0, \exists n_0(\varepsilon)$ astfel încât

$$\left| \frac{f(n)}{g(n)} - l \right| < \varepsilon, \forall n \geq n_0(\varepsilon). \text{ Adică } -\varepsilon < \frac{f(n)}{g(n)} - l < \varepsilon \Rightarrow l - \varepsilon < \frac{f(n)}{g(n)} < l + \varepsilon. \text{ Alegând}$$

$$\varepsilon < 1 \Rightarrow \frac{f(n)}{g(n)} < l + 1, \forall n > n_0(\varepsilon). \text{ Există două constante pozitive } n_0(\varepsilon), n_1^0 \geq l + 1 \text{ astfel încât}$$

$$\forall n > n_0(\varepsilon), g(n) > 0, 0 \leq f(n) < n_1^0 * g(n). \text{ Așadar } f(n) = O(g(n)).$$

Această condiție nu este și necesară. Justificați această afirmație printr-un contraexemplu.

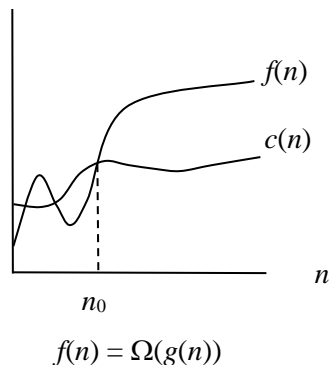
Exemplul 3.2.2. Fie funcțiile $f(n), g(n) \in \mathfrak{R}$ definite prin relațiile

$f(n) = 10n$ și $g(n) = n^2$. Dacă calculăm $f(n)$ și $g(n)$ pentru $1 \leq n \leq 9$ se obține: $f(1) = 10, g(1) = 1; f(2) = 20, g(2) = 4; f(3) = 30, g(3) = 9; f(4) = 40, g(4) = 16, \dots, f(9) = 90, g(9) = 81$. Pentru $n \geq 10$ se observă că $n^2 \geq 10n$, adică $0 \leq f(n) \leq g(n)$.

3.3. Ω - notația

Definiția 3.3.1 $f(n) = \Omega(g(n)) (n \rightarrow \infty)$ dacă $g(n) = O(f(n)) (n \rightarrow \infty)$.

$$\Omega(g(n)) = \{f(n) \in \mathfrak{R} | \exists c > 0, \exists n_0 \in \mathbb{N} \text{ astfel încât } f(n) \geq cg(n) \geq 0, \forall n \geq n_0\}.$$



În studiul algoritmilor, simbolul Ω este util pentru a pune în evidență un timp minim posibil de execuție. Așadar simbolul Ω este util în precizarea unei margini asimptotice inferioare nu neapărat tare.

Exercițiu 3.3.1:

Arătați că $\theta(f(n)) = \{g(n) \in \mathfrak{R} | g \in O(f(n)) \cap \Omega(f(n))\}$.

În practică se demonstrează existența marginilor asimptotice tari pornind de la margini asimptotice superioare și inferioare și nu invers.

3.4. o - notația

Definiția 3.4.1 $f(n) = o(g(n)) (n \rightarrow \infty)$ dacă $\forall c > 0, \exists n_0 \in \mathbb{N}$ astfel încât $\forall n \geq n_0$, are loc relația $0 \leq f(n) < cg(n)$.

Prin urmare, $o(g(n)) = \{f(n) \in \mathfrak{R} | \forall c > 0, \exists n_0 \in \mathbb{N} \text{ astfel ca } 0 \leq f(n) < cg(n), \forall n \geq n_0\}$ Prin notația „ o ” vom semnala o margine asimptotică superioară care nu este o margine asimptotică superioară tare.

Exemplul 3.4.1.

$10n = o(n^2)$, iar $10n^2 \neq o(n^2)$.

Conform definiției simbolului o putem spune că $f(n) = o(g(n)) (n \rightarrow \infty)$ dacă $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)}$ există și este egală cu 0.

O astfel de notație exprimă faptul că funcția $f(n)$ crește strict mai lent decât funcția $g(n)$ pentru valori ale lui n suficient de mari.

Să observăm, de exemplu, că $\sin(n) = o(n)$ implică $\sin(n) = O(n)$, dar $\sin(n) = O(1)$ nu implică $\sin(n) = o(1)$. Ce puteți afirma în acest sens ?

Propoziția 3.4.2. O condiție necesară și suficientă ca $\lim_{n \rightarrow \infty} f(n) = 0$ este ca $f(n) = o(1)$

3.5. ω - notația

Definiția 3.5.1. $f(n) = \omega(g(n)) (n \rightarrow \infty)$ dacă $\forall c > 0, \exists n_0 \in \mathbb{N}$ astfel încât $\forall n \geq n_0$, are loc relația $f(n) > cg(n) \geq 0$.

Simbolul ω se raportează la Ω într-un mod asemănător lui o față de O .

Așadar, $\omega(g(n)) = \{f(n) \in \mathfrak{R} | \forall c > 0, \exists n_0 \in \mathbb{N} \text{ astfel ca } f(n) > cg(n) \geq 0, \forall n \geq n_0\}$ Se observă că $f(n) = \omega(g(n))$ dacă și numai dacă $g(n) = o(f(n))$.

Așadar ω desemnează o margine asimptotică inferioară care nu este o margine asimptotică inferioară tare.

Exemplul 3.5.1.

$n^2 = \omega(n)$, iar $n^2 \neq \omega(n^2)$.

Se observă că $f(n) = \omega(g(n))$ implică $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \infty$ (dacă limita există).

3.6. ~ - notația

Definiția 3.6.1.

$f(n) \sim g(n) (n \rightarrow \infty)$ dacă există $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 1$.

Propoziția 3.6.1. O condiție necesară și suficientă ca $f(n) \sim g(n) (n \rightarrow \infty)$ este ca $f(n) = g(n) (1 + o(1))$.

Demonstrație:

Dacă $f(n) \sim g(n)$ atunci $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 1$ implică

$$\lim_{n \rightarrow \infty} \frac{f(n) - g(n)}{g(n)} = \lim_{n \rightarrow \infty} \left(\frac{f(n)}{g(n)} - 1 \right) = \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} - 1 = 0$$

Așadar $(f(n) - g(n))/g(n) = o(1)$, adică $f(n) = g(n) (1 + o(1))$.

Reciproc se demonstrează asemănător.

Propoziția 3.6.2. Dacă $f(n) \sim g(n) (n \rightarrow \infty)$, atunci $f(n) = O(g(n)) (n \rightarrow \infty)$.

Demonstrație: exercițiu!

4. Compararea asimptotică a funcțiilor

Multe dintre proprietățile relațiilor dintre numerele reale se aplică și la compararea asimptotică a funcțiilor. Pentru cele ce urmează vom presupune că $f(n)$ și $g(n)$ sunt asimptotic pozitive.

Tranzitivitate :

$f(n) = \theta(g(n))$ și $g(n) = \theta(h(n))$ implică $f(n) = \theta(h(n))$,

$f(n) = O(g(n))$ și $g(n) = O(h(n))$ implică $f(n) = O(h(n))$,

$f(n) = \Omega(g(n))$ și $g(n) = \Omega(h(n))$ implică $f(n) = \Omega(h(n))$,

$f(n) = o(g(n))$ și $g(n) = o(h(n))$ implică $f(n) = o(h(n))$,

$f(n) = \omega(g(n))$ și $g(n) = \omega(h(n))$ implică $f(n) = \omega(h(n))$.

Reflexivitate:

$f(n) = \theta(f(n))$,

$f(n) = O(f(n))$,

$f(n) = \Omega(f(n))$,

Simetrie:

$f(n) = \theta(f(n))$, dacă și numai dacă $g(n) = \theta(f(n))$.

Antisimetrie:

$f(n) = O(f(n))$, dacă și numai dacă $g(n) = \Omega(f(n))$.

$f(n) = o(f(n))$, dacă și numai dacă $g(n) = \omega(f(n))$.

Deoarece aceste proprietăți sunt valide pentru notații asimptotice, se poate trasa o analogie între compararea asimptotică a două funcții f și g și compararea asimptotică a două numere reale a și b :

$f(n) = O(f(n)) \approx a \leq b$,

$f(n) = \Omega(f(n)) \approx a \geq b$,

$f(n) = \theta(f(n)) \approx a = b$,

$f(n) = o(f(n)) \approx a < b$,

$f(n) = \omega(f(n)) \approx a > b$.

O proprietate a numerelor reale, totuși, nu se transpune la notații asimptotice:

Trihotomia: pentru orice două numere reale a și b exact una dintre următoarele relații este adevărată: $a < b$, $a = b$, $a > b$.

Deși orice două numere reale pot fi comparate, nu toate funcțiile sunt asimptotic comparabile. Cu alte cuvinte, pentru două funcții $f(n)$ și $g(n)$, se poate întâmpla să nu aibă loc nici $f(n) = O(g(n))$, nici $f(n) = \Omega(g(n))$. De exemplu, funcțiile n și $n^{1+\sin n}$ nu pot fi comparate utilizând notații asimptotice, deoarece valoarea exponentului în $n^{1+\sin n}$ oscilează între 0 și 2, luând toate valorile intermediare.

EXERCITII:

Exemplul 3.1.1.

$(n+1)^3 = \Theta(n^3)$

(1)

$(5 + (2n)^{0.5})^{0.5} = \Theta(n^{0.25})$

(2)

$$(1+7/n)^n = \Theta(1) \quad (3)$$

Exemplul 3.4.2.

$$n^4 = o(n^7) \quad (1)$$

$$5 \log n = o(n^{0.03}) \quad (2)$$

$$\sin n = o(n) \quad (3)$$

$$5.2 n^{0.5} = o(n/8 + 9 \sin n) \quad (4)$$

Exemplul 1.3.9.

$$n^4 + 3n^3 - 2 \sim n^4 \quad (1)$$

$$(5n + 1)^2 \sim 25n^2 \quad (2)$$

$$\frac{(6n^5 + 7n^2 + 1)}{(3n^3 + 5)} \sim 2n^2 \quad (3)$$

$$4^n + 8 \log n + 78 \sim 4^n \quad (4)$$

$$\sin 1/n \sim 1/n \quad (5)$$

$$n + \sin n \sim n \quad (6)$$

5. RELAȚII RECURENTE

Cel mai important câștig al exprimării recursive este faptul că ea este naturală și compactă, fără să ascundă esența algoritmului prin detaliile de implementare. Pe de altă parte, apelurile recursive trebuie folosite cu discernământ, deoarece solicită și ele resursele calculatorului (timp și memorie). Analiza unui algoritm recursiv implică rezolvarea unui sistem de recurențe. Vom vedea în continuare cum pot fi rezolvate astfel de recurențe.

Când un algoritm conține o apelare recursivă la el însuși, timpul său de execuție poate fi descris adesea printr-o recurență. O *recurență* este o ecuație sau o inegalitate care descrie întregul timp de execuție al unei probleme de dimensiune n cu ajutorul timpilor de execuție pentru date de intrare de dimensiuni mici. Putem, apoi, folosi instrumente matematice pentru a rezolva problema de recurență și pentru a obține margini ale performanței algoritmului.

O recurență pentru timpul de execuție al unui algoritm de tipul divide și stăpânește se bazează pe cele trei etape definite în descrierea metodei. La fel ca până acum, vom nota cu $T(n)$ timpul de execuție al unei probleme de dimensiune n . Dacă dimensiunea problemei este suficient de mică, de exemplu $n \leq c$ pentru o anumită constantă c , soluția directă ia un timp constant de execuție, pe care îl vom nota cu $\Theta(1)$. Să presupunem că divizăm problema în a subprobleme, fiecare dintre acestea având dimensiunea de $1/b$ din dimensiunea problemei originale. Dacă $D(n)$ este timpul necesar pentru a divide problema în subprobleme, iar $C(n)$ este timpul necesar pentru a combina soluțiile subproblemelor în soluția problemei originale, obținem recurența

$$T(n) = \begin{cases} \Theta(1), & n \leq c \\ aT\left(\frac{n}{b}\right) + D(n) + C(n), & n > c \end{cases}$$

5.1 Metoda master

Metoda master furnizează o „rețetă” pentru rezolvarea recurențelor de forma

$$T(n) = aT(n/b) + f(n) \quad (5.1)$$

unde $a \geq 1$ și $b > 1$ sunt constante, iar $f(n)$ este o funcție asimptotic pozitivă. Metoda master pretinde memorarea a trei cazuri, dar apoi soluția multor recurențe se poate determina destul de ușor, de multe ori fără creion și hârtie.

Recurența (1.4.1) descrie timpul de execuție al unui algoritm care împarte o problemă de dimensiune n în a subprobleme, fiecare de dimensiune n/b , unde a și b sunt constante pozitive. Cele a subprobleme sunt rezolvate recursiv, fiecare în timp $T(n/b)$. Costul divizării problemei și al combinării rezultatelor subproblemelor este descris de funcția $f(n)$ (Adică, utilizând notația $f(n) = D(n) + C(n)$). Din punctul de vedere al corectitudinii tehnice, recurența nu este, de fapt, bine definită, deoarece n/b ar putea să nu fie întreg. Înlocuirea fiecăruia dintre cei a termeni $T(n/b)$ fie cu $T(\lfloor n/b \rfloor)$ fie cu $T(\lceil n/b \rceil)$ nu afectează, totuși, comportamentul asimptotic al recurenței. În mod normal, ne va conveni, de aceea, să omitem funcțiile parte întreagă inferioară și superioară când scriem recurențe divide și stăpânește de această formă.

Teorema master

Metoda master depinde de următoarea teoremă.

Teorema 5.1.1. (teorema master) Fie $a \geq 1$ și $b > 1$ constante, fie $f(n)$ o funcție și fie $T(n)$ definită pe întregii nenegativi prin recurență

$$T(n) = aT(n/b) + f(n)$$

unde interpretăm n/b fie ca $\lfloor n/b \rfloor$ fie ca $\lceil n/b \rceil$. Atunci $T(n)$ poate fi delimitată asimptotic după cum urmează.

1. Dacă $f(n) = O(n^{\log_b^a - \varepsilon})$ pentru o anumită constantă $\varepsilon > 0$, atunci $T(n) = \Theta(n^{\log_b^a})$.
2. Dacă $f(n) = \Theta(n^{\log_b^a})$, atunci $T(n) = \Theta(n^{\log_b^a} \lg n)$.
3. Dacă $f(n) = \Omega(n^{\log_b^a + \varepsilon})$, pentru o anumită constantă $\varepsilon > 0$ și dacă $af(n/b) \leq cf(n)$ pentru o

anumită constantă $c < 1$ și toți n suficient de mari, atunci $T(n) = \Theta(f(n))$.

Înainte de a aplica teorema master la câteva exemple, să ne oprim puțin ca să înțelegem ce spune. În fiecare dintre cele trei cazuri, comparăm funcția $f(n)$ cu funcția $n^{\log_b^a}$. Intuitiv, soluția recurenței este determinată de cea mai mare dintre cele două funcții. Dacă funcția $n^{\log_b^a}$ este mai mare, ca în cazul 1, atunci soluția este $T(n) = \Theta(n^{\log_b^a})$. Dacă funcția $f(n)$ este mai mare, ca în cazul 3, atunci soluția este $T(n) = \Theta(f(n))$. Dacă cele două funcții sunt de același ordin de mărime, ca în cazul 2, înmulțim cu un factor logaritmic, iar soluția este

$$T(n) = \Theta(n^{\log_b^a} \lg n) = \Theta(f(n) \lg n).$$

Dincolo de această intuiție, există niște detalii tehnice care trebuie înțelese.

În primul caz, $f(n)$ trebuie nu numai să fie mai mică decât $n^{\log_b^a}$, trebuie să fie polinomial mai mică. Adică $f(n)$ trebuie să fie asimptotic mai mică decât $n^{\log_b^a}$ cu un factor n^ε pentru o anumită constantă $\varepsilon > 0$. În al treilea caz, $f(n)$ trebuie nu numai să fie mai mare decât $n^{\log_b^a}$, trebuie să fie polinomial mai mare și, în plus, să verifice condiția de „regularitate” $af(n/b) \leq cf(n)$. Această condiție este satisfăcută de majoritatea funcțiilor polinomiale mărginite pe care le vom întâlni.

Este important de realizat că cele trei cazuri nu acoperă toate posibilitățile pentru $f(n)$.

Există un gol între cazurile 1 și 2 când $f(n)$ este mai mic decât $n^{\log_b^a}$ dar nu polinomial mai mic. Analog există un gol între cazurile 2 și 3 când $f(n)$ este mai mare decât $n^{\log_b^a}$ dar nu polinomial mai mare. Dacă funcția $f(n)$ cade într-unul dintre aceste goluri, sau când condiția de regularitate din cazul 3 nu este verificată, metoda master nu poate fi utilizată pentru a rezolva recurența.

Exemplu 5.1.1.

$$T(n) = 9T(n/3) + n.$$

Pentru această recurență, avem $a = 9$, $b = 3$, $f(n) = n$ și astfel $n^{\log_b^a} = T(n) = n^{\log_3^9} = \Theta(n^2)$.

Deoarece $f(n) = O(n^{\log_3^9 - \varepsilon})$, unde $\varepsilon = 1$, putem să aplicăm cazul 1 al teoremei master și să considerăm că soluția este $T(n) = \Theta(n^2)$.

Exemplu 5.1.2.

$$T(n) = T(2n/3) + 1.$$

Pentru această recurență, avem $a = 1$, $b = 3/2$, $f(n) = 1$ și $n^{\log_b^a} = n^{\log_{3/2}^1} = n^0 = 1$. Cazul 2 este cel care se aplică deoarece $f(n) = \Theta(n^{\log_b^a}) = \Theta(1)$ și astfel soluția recurenței este

$$T(n) = \Theta(\lg n).$$

Exemplu 5.1.3.

$$T(n) = 3T(n/4) + n \lg n,$$

Pentru această recurență, avem $a = 3$, $b = 4$, $f(n) = n \lg n$ și $n^{\log_b^a} = n^{\log_4^3} = O(n^{0.793})$. Deoarece $f(n) = \Omega(n^{\log_4^3 + \varepsilon})$, unde $\varepsilon \approx 0.2$, cazul 3 se aplică dacă putem arăta că pentru $f(n)$ este verificată condiția de regularitate. Pentru n suficient de mare,

$$af(n/b) = 3(n/4) \lg(n/4) \leq (3/4)n \lg n = cf(n) \text{ pentru } c = 3/4.$$

În consecință, din cazul 3, soluția recurenței este $T(n) = \Theta(n \lg n)$.

Exemplu 5.1.4.

Metoda master nu se aplică recurenței

$$T(n) = 2T(n/2) + n \lg n,$$

chiar dacă are forma potrivită: $a = 2$, $b = 2$, $f(n) = n \lg n$ și $n^{\log_b^a} = n$. Se pare că ar trebui să se aplice cazul 3, deoarece $f(n) = n \lg n$ este asimptotic mai mare decât $n^{\log_b^a} = n$, dar nu polinomial mai mare. Raportul $f(n)/n^{\log_b^a} = (n \lg n)/n$ este asimptotic mai mic decât n^ε pentru orice constantă pozitivă ε . În consecință, recurența cade în golul dintre cazurile 2 și 3.

5.2. Metoda ecuațiilor caracteristice

5.2.1. Recurențe liniare omogene

Există tehnici care pot fi folosite aproape automat pentru a rezolva anumite clase de recurențe. Vom considera recurențe liniare omogene, de forma

$$a_0 t_n + a_1 t_{n-1} + \dots + a_k t_{n-k} = 0 \quad (5.2.1)$$

unde t_i sunt valorile pe care le căutăm, iar coeficienții a_i sunt constante.

Vom căuta soluții de forma

$$t_n = x^n$$

unde x este o constantă (deocamdată necunoscută). Încercăm aceasta soluție în (5.2.1) și obținem

$$a_0 x^n + a_1 x^{n-1} + \dots + a_k x^{n-k} = 0.$$

Soluțiile acestei ecuații sunt fie soluția trivială $x = 0$, care nu ne interesează, fie soluțiile ecuației

$$a_0 x^k + a_1 x^{k-1} + \dots + a_k = 0$$

care este ecuația caracteristică a recurenței (1.4.2).

Presupunând deocamdată ca cele k rădăcini r_1, r_2, \dots, r_k ale acestei ecuații caracteristice sunt distincte, orice combinație liniară

$$t_n = \sum_{i=1}^k c_i r_i^n$$

este o soluție a recurenței (5.2.1), unde constantele c_1, c_2, \dots, c_k sunt determinate de condițiile inițiale. Este remarcabil că (1.4.2) are numai soluții de aceasta formă.

Exemplu 5.2.1.

Recurența care definește șirul lui Fibonacci:

$$t_n = t_{n-1} + t_{n-2} \quad n \geq 2$$

iar $t_0 = 0, t_1 = 1$. Putem să rescriem această recurență sub forma $t_n - t_{n-1} - t_{n-2} = 0$ care are ecuația caracteristică $x^2 - x - 1 = 0$ cu rădăcinile $r_{1,2} = (1 \pm \sqrt{5})/2$. Soluția generală are forma

$$t_n = c_1 r_1^n + c_2 r_2^n$$

Impunând condițiile inițiale, obținem

$$\begin{aligned} c_1 + c_2 &= 0, & n &= 0 \\ r_1 c_1 + r_2 c_2 &= 1, & n &= 1 \end{aligned}$$

de unde determinăm

$$c_{1,2} = \pm 1/\sqrt{5}$$

Deci, $t_n = 1/\sqrt{5} (r_1^n - r_2^n)$. Observăm că $r_1 = \phi = (1 + \sqrt{5})/2, r_2 = -\phi^{-1}$ și obținem $t_n = 1/\sqrt{5} (\phi^n - (-\phi)^{-n})$

Putem să arătăm acum că timpul pentru algoritmul care determină șirul Fibonacci este în $\Theta(\phi^n)$.

Se poate arăta că, dacă r este o rădăcină de multiplicitate m a ecuației caracteristice, atunci

$$t_n = r^n, t_n = nr^n, t_n = n^2 r^n, \dots, t_n = n^{m-1} r^n$$

sunt soluții pentru (1.4.2). Soluția generală pentru o astfel de recurență este atunci o combinație liniară a acestor termeni și a termenilor proveniți de la celelalte rădăcini ale ecuației caracteristice. Din nou, sunt de determinat exact k constante din condițiile inițiale.

Exemplu 5.2.2.

Fie recurența $t_n = 5t_{n-1} - 8t_{n-2} + 4t_{n-3}$ iar $t_0 = 0, t_1 = 1, t_2 = 2$. Ecuația caracteristică

$$x^3 - 5x^2 + 8x - 4 = 0$$

are rădăcinile 1 (de multiplicitate 1) și 2 (de multiplicitate 2). Soluția generală este:

$$t_n = c_1 1^n + c_2 2^n + c_3 n 2^n$$

Din condițiile inițiale, obținem $c_1 = -2, c_2 = 2, c_3 = -1/2$.

5.2.2. Recurențe liniare neomogene

Considerăm acum recurențe de următoarea formă mai generală

$$a_0 t_n + a_1 t_{n-1} + \dots + a_k t_{n-k} = b^n p(n) \quad (5.2.2)$$

unde b este o constantă, iar $p(n)$ este un polinom în n de grad d . Ideea generală este că, prin manipulări convenabile, să reducem un astfel de caz la o formă omogenă.

Exemplu 5.2.3.

De exemplu, o astfel de recurență poate fi:

$$t_n - 2t_{n-1} = 3^n$$

În acest caz, $b = 3$ și $p(n) = 1$, un polinom de grad 0. O simplă manipulare ne permite să reducem acest exemplu la forma (1.4.2). Înmulțim recurența cu 3, obținând $3t_n - 6t_{n-1} = 3^{n+1}$. Înlocuind pe n cu $n+1$ în recurența inițială, avem $t_{n+1} - 2t_n = 3^{n+1}$. În fine, scădem aceste două ecuații și obținem $t_{n+1} - 5t_n + 6t_{n-1} = 0$.

Am obținut o recurență omogenă pe care o putem rezolva ca în secțiunea precedentă. Ecuația caracteristică este:

$$x^2 - 5x + 6 = 0$$

adică $(x-2)(x-3) = 0$.

Intuitiv, observăm că factorul $(x-2)$ corespunde părții stângi a recurenței inițiale, în timp ce factorul $(x-3)$ a apărut ca rezultat al manipularilor efectuate, pentru a scăpa de partea dreaptă.

Generalizând acest procedeu, se poate arăta că, pentru a rezolva (1.4.3), este suficient să luăm următoarea ecuație caracteristică:

$$(a_0x^k + a_1x^{k-1} + \dots + a_k)(x-b)^{d+1} = 0$$

Odată ce s-a obținut această ecuație, se procedează ca în cazul omogen.

5.2.3. Schimbarea variabilei

Uneori, printr-o schimbare de variabilă, putem rezolva recurențe mult mai complicate. În exemplele care urmează, vom nota cu $T(n)$ termenul general al recurenței și cu t_k termenul noii recurențe obținute printr-o schimbare de variabilă. Presupunem pentru început ca n este o putere a lui 2.

Exemplu 5.2.4.

Fie recurența $T(n) = 4T(n/2) + n$, $n > 1$ în care înlocuim pe n cu 2^k , notăm $t_k = T(2^k) = T(n)$ și obținem $t_k = 4t_{k-1} + 2^k$. Ecuația caracteristică a acestei recurențe liniare este

$(x-4)(x-2) = 0$ cu răși deci,

$$t_k = c_14^k + c_22^k.$$

Înlocuim la loc pe k cu $\lg n$ și obținem

$$T(n) = c_1n^2 + c_2n$$

Rezultă că

$T(n) \in O(n^2 \mid n \text{ este o putere a lui } 2)$.

EXERCITII:

II. PROIECTAREA ALGORITMILOR PRELIMINARII

1. TEHNICA DIVIDE ET IMPERA.

În această secțiune vom examina o abordare numită *divide și stăpânește*. Vom utiliza această abordare pentru a construi un algoritm de sortare. Unul din avantajele algoritmilor de tipul *divide și stăpânește* este acela că timpul lor de execuție este adesea ușor de determinat folosind tehnici date mai sus.

Mulți algoritmi au o structură recursivă: pentru a rezolva o problemă dată, aceștia sunt apelați o dată sau de mai multe ori pentru a rezolva subproblemele apropiate. Acești algoritmi folosesc de obicei o abordare de tipul *divide și stăpânește*: ei rup problema de rezolvat în mai multe probleme similare problemei inițiale, dar de dimensiune mai mică, le rezolvă în mod recursiv și apoi le combină pentru a crea o soluție a problemei inițiale.

Paradigma *divide și stăpânește* implică trei pași la fiecare nivel de recursivitate:

- **Divide** problema într-un număr de subprobleme.
- **Stăpânește** subproblemele prin rezolvarea acestora în mod recursiv. Dacă dimensiunile acestora sunt suficient de mici, rezolvă subproblemele în mod uzual, nerecursiv.
- **Combină** soluțiile tuturor subproblemelor în soluția finală pentru problema inițială.

1.1. Sortarea prin interclasare.

Algoritmul de sortare prin interclasare urmează îndeaproape paradigma *divide și stăpânește*. Intuitiv acesta operează astfel:

Divide: Împarte șirul de n elemente care urmează a fi sortat în două subșiruri de câte $n/2$ elemente.

Stăpânește: Sortează recursiv cele două subșiruri utilizând sortarea prin interclasare.

Combină: Interclasează cele două subșiruri sortate pentru a produce rezultatul final.

Să observăm că recursivitatea se oprește când șirul de sortat are lungimea 1, caz în care nu mai avem nimic de făcut, deoarece orice șir de lungime 1 este deja sortat.

Operația principală a algoritmului de sortare prin interclasare este interclasarea a două șiruri sortate, în pasul denumit mai sus *combină*. Pentru aceasta vom utiliza o procedură auxiliară **MERGE**(A, p, q, r), unde A este vector, și p, q, r sunt indici ai vectorului, astfel încât $p \leq q < r$. Procedura presupune că subvectorii $A[p..q]$ și $A[q+1..r]$ sunt sortați. Ea interclasează pentru a forma un subvector sortat care înlocuiește subvectorul curent $A[p..r]$.

Deși vom lăsa pseudocodul pentru această procedură ca exercițiu, este ușor de imaginat o procedură de tip **MERGE** al cărei timp de execuție este de ordinul $\Theta(n)$, în care $n = r - p + 1$ este numărul elementelor interclasate. Revenind la exemplul nostru cu cărțile de joc, să presupunem că avem două pachete de cărți de joc așezate pe masă cu fața în sus. Fiecare dintre aceste două pachete de cărți este sortat, cartea cu valoarea cea mai mică fiind deasupra. Dorim să amestecăm cele două pachete într-un singur pachet sortat, care să rămână așezat pe masă cu fața în jos. Pasul principal este acela de a selecta cartea cu valoarea cea mai mică dintre cele două aflate deasupra pachetelor (fapt care va face ca o nouă carte să fie deasupra pachetului respectiv) și de a o pune cu fața în jos pe locul în care se va forma pachetul sortat final. Repetăm acest procedeu pînă cînd unul din pachete nu este ieput. În această fază, este suficient să luăm pachetul rămas și să-l punem pe pachetul deja sortat, întorcînd toate cărțile cu fața în jos. Din punctul de vedere al timpului de execuție, fiecare pas de bază durează un timp constant, deoarece comparăm de fiecare dată doar două cărți. Deoarece avem de făcut cel mult n astfel de operații elementare, timpul de execuție pentru procedura **MERGE** este $\Theta(n)$.

MERGE-SORT (A, p, r)

1. **if** $p < r$
2. **then** $q \leftarrow [(p + r)/2]$
3. **MERGE-SORT** (A, p, q)
4. **MERGE-SORT** ($A, q+1, r$)
5. **MERGE** (A, p, q, r)

Acum putem utiliza procedura **MERGE** ca subrutină pentru algoritmul de sortare prin interclasare.

Procedura **MERGE-SORT** (A, p, r) sortează elementele din subvectorul $A[p..r]$. Dacă $p \geq r$, subvectorul are cel mult un element și este, prin urmare, deja sortat. Astfel pasul de divizare este prezent aici prin simplul calcul al unui indice q care împarte $A[p..r]$ în doi subvectori $A[p..q]$, conținînd $\lfloor n/2 \rfloor$ elemente și $A[q+1..r]$ conținînd $\lfloor n/2 \rfloor$ elemente. Pentru a sorta întregul șir $A = \langle A[1], A[2], \dots, A[n] \rangle$, vom apela procedura **MERGE-SORT** ($A, 1, \text{lungime}[A]$), unde din nou, $\text{lungime}[A] = n$. Dacă analizăm modul de operare al procedurii, de jos în sus, cînd n este o putere a lui 2, algoritmul constă din interclasarea perechilor de șiruri de lungime 1, pentru a forma șiruri sortate de

lungime 2, interclasarea acestora în șiruri sortate de lungime 4, și așa mai departe, pînă cînd două șiruri sortate de lungime $n/2$ sunt interclasate pentru a forma șirul sortat de dimensiune n . Figura 2.1 ilustrează acest proces.

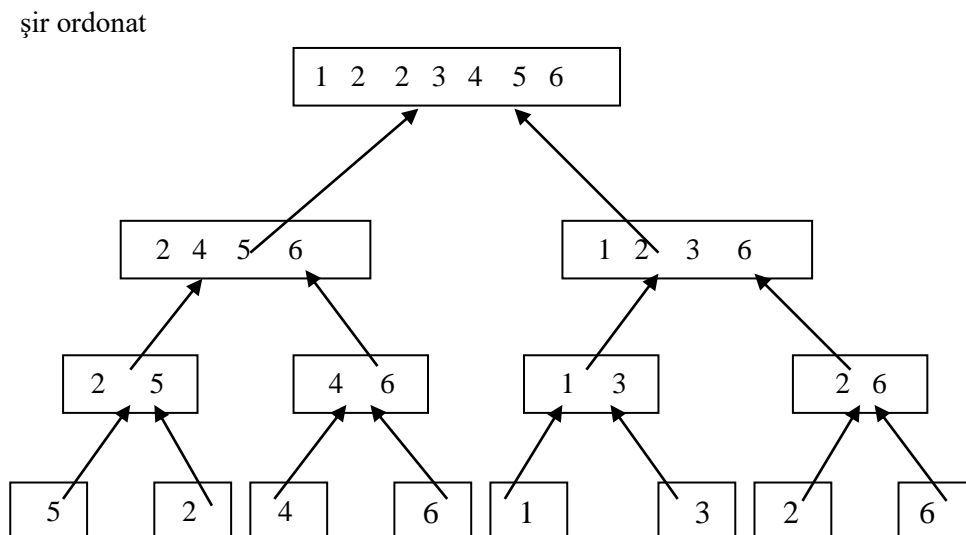


Fig. 1.1. Modul de operare al sortării prin interclasare asupra vectorului $A = \langle 5, 2, 4, 6, 1, 3, 2, 6 \rangle$. Lungimile șirurilor sortate, în curs de interclasare, cresc pe măsură ce algoritmul avansează de jos în sus.

1.1.1 Analiza algoritmilor de tipul divide și stăpînește

Cînd un algoritm conține un apel recursiv la el însuși, timpul său de execuție poate fi, adesea, descris printr-o **relație de recurență** sau mai simplu, **recurență**, care descrie întregul timp de execuție al unei probleme de dimensiune n cu ajutorul timpilor de execuție pentru date de intrare de dimensiuni mai mici. Putem, apoi, folosi instrumente matematice pentru a rezolva problema de recurență și pentru a obține margini ale performanței algoritmului.

O recurență pentru timpul de execuție al unui algoritm de tipul divide și stăpînește se bazează pe cele trei etape definite în descrierea metodei. La fel ca pînă acum, vom nota cu $T(n)$ timpul de execuție al unei probleme de dimensiune n . Dacă dimensiunea problemei este suficient de mică, de exemplu $n \leq c$ pentru o anumită constantă c , soluția directă ia un timp constant de execuție, pe care îl vom nota cu $\Theta(1)$. Să presupunem că divizăm problema în a subprobleme, fiecare dintre acestea avînd dimensiunea de b ori mai mică decât dimensiunea problemei originale. Dacă $D(n)$ este timpul necesar pentru a divide problema în subprobleme, iar $C(n)$ este timpul necesar pentru a combina soluțiile subproblemelor în soluția problemei originale, obținem recurența

$$T(n) = \begin{cases} \Theta(1), & n \leq c \\ aT(n/b) + D(n) + C(n), & n > c \end{cases}$$

1.1.2. Analiza sortării prin interclasare

Deși algoritmul **MERGE-SORT** funcționează corect și cînd numărul de elemente nu este par, analiza bazată pe recurență se simplifică dacă presupunem că dimensiunea problemei originale este o putere a lui 2. Fiecare pas de împărțire generează deci două subșiruri avînd dimensiunea exact $n/2$. În continuare vom vedea că această presupunere nu afectează ordinul de creștere a recurenței.

Pentru a determina recurența pentru $T(n)$, timpul de execuție al sortării prin interclasare a n numere, în cazul cel mai defavorabil, vom raționa în felul următor. Sortarea prin interclasare a unui singur element are nevoie de un timp constant. Cînd avem $n > 1$ elemente, vom descompune timpul de execuție după cum urmează:

Divide: La acest pas, se calculează doar mijlocul subvectorului, calcul care are nevoie de un timp constant de execuție. Astfel, $D(n) = \Theta(1)$.

Stăpînește: Rezolvăm recursiv două probleme fiecare de dimensiune $n/2$, care contribuie cu $2T(n/2)$ la timpul de execuție.

Combină: Am observat, deja că procedura **MERGE** pentru un subvector cu n elemente consumă $\Theta(n)$ timp de execuție, deci $C(n) = \Theta(n)$.

Cînd adunăm funcțiile $D(n)$ și $C(n)$ pentru analiza sortării prin interclasare, adunăm o funcție cu timpul de execuție $\Theta(1)$. Această sumă este o funcție lineară în raport cu n , adică are timpul de execuție $\Theta(n)$. Adăugînd

aceasta la termenul $2T(n/2)$ de la pasul *stăpînește*, obținem timpul de execuție $T(n)$ în cazul cel mai defavorabil pentru sortarea prin interclasare:

$$T(n) = \begin{cases} \Theta(1), & n = 1 \\ 2T(n/2) + \Theta(n), & n > 1 \end{cases}$$

Se poate arăta că $T(n)$ este $\Theta(n \lg n)$ unde $\lg n$ reprezintă $\log_2 n$. Pentru numere suficient de mari, sortarea prin interclasare, are timpul de execuție $\Theta(n \lg n)$.

2. TEHNICA GREEDY

Algoritmii *greedy* (greedy = lacom) sunt în general simpli și sunt folosiți pentru rezolvarea problemelor de optimizare, cum ar fi: să se găsească cea mai bună ordine de executare a unor lucrări pe calculator, să se găsească cel mai scurt drum într-un graf etc. Algoritmii aplicați problemelor de optimizare sunt compuși dintr-o secvență de pași, la fiecare pas existând mai multe alegeri posibile. Un *algoritm greedy* va alege la fiecare moment de timp soluția ce pare a fi cea mai bună la momentul respectiv. Deci este vorba despre o alegere optimă, făcută local, cu speranța că ea va conduce la un optim global. Acest capitol tratează probleme de optimizare ce pot fi rezolvate cu ajutorul algoritmilor greedy.

Algoritmii greedy conduc în multe cazuri la soluții optime, dar nu întotdeauna.

2.1. O problemă de selectare a activităților

Primul exemplu pe care îl vom considera este o problemă de repartizare a unei resurse mai multor activități care concurează pentru a obține resursa respectivă. Vom vedea ca un algoritm de tip greedy reprezintă o metodă simplă și elegantă pentru selectarea unei mulțimi maximale de activități mutual compatibile.

Să presupunem că dispunem de o mulțime $S = 1, 2, \dots, n$ de n activități care doresc să folosească o aceeași resursă (de exemplu o sală de lectură). Această resursă poate fi folosită de o singură activitate la un anumit moment de timp. Fiecare activitate i are un timp de pornire s_i și un timp de oprire f_i , unde $s_i \leq f_i$. Dacă este selecționată activitatea i , ea se desfășoară pe durata intervalului $[s_i, f_i)$. Spunem că activitățile i și j sunt *compatibile* dacă intervalele $[s_i, f_i)$ și $[s_j, f_j)$ nu se intersectează (adică i și j sunt compatibile dacă $s_i \geq f_j$ sau $s_j \geq f_i$). *Problema selectării activităților* constă din selectarea unei mulțimi maximale de activități mutual compatibile.

Un algoritm greedy pentru această problemă este descris de următoarea procedură, prezentată în pseudocod. Vom presupune că activitățile (adică datele de intrare) sunt ordonate crescător după timpul de terminare:

$$f_1 \leq f_2 \leq \dots \leq f_n \quad (2.1)$$

În caz contrar această ordonare poate fi făcută în timp $O(n \lg n)$. Algoritmii de mai jos presupune că datele de intrare s și f sunt reprezentate ca vectori.

SELECTOR-ACTIVITĂȚI-GREEDY (s, f)

```

1:  $n \leftarrow \text{lungime}[s]$ 
2:  $A \leftarrow \{1\}$ 
3:  $j \leftarrow 1$ 
4: for  $i \leftarrow 2$  to  $n$  do
5:   if  $s_i \geq f_j$  then
6:      $A \leftarrow A \cup \{i\}$ 
7:      $j \leftarrow i$ 
8: returnează  $A$ 

```

Operațiile realizate de algoritm pot fi vizualizate în figura 2.1. în mulțimea A se introduc activitățile selectate. Variabila j identifică ultima activitate introdusă în A . Deoarece activitățile sunt considerate în ordinea nedescrescătoare a timpilor lor de terminare, f_j , va reprezenta întotdeauna timpul maxim de terminare a oricărei activități din A . Aceasta înseamnă că

$$f_j = \max\{f_k : k \in A\} \quad (2.2)$$

În liniile 2-3 din algoritm se selectează activitatea 1, se inițializează A astfel încât să nu conțină decât această activitate, iar variabila j ia ca valoare această activitate. În continuare liniile 4-7 consideră pe rând fiecare activitate i și o adaugă mulțimii A dacă este compatibilă cu toate celelalte activități deja selectate. Pentru a vedea dacă activitatea i este compatibilă cu toate celelalte activități existente la momentul curent în A , este suficient, conform formulei (2.2), să fie îndeplinită condiția din linia 5 adică momentul de pornire s_i , să nu fie mai devreme decât momentul de oprire f_j , al activității cel mai recent adăugate mulțimii A . Dacă activitatea i este compatibilă, atunci în liniile 6-7 ea este adăugată mulțimii A , iar variabila j este actualizată. Procedura **SELECTOR-ACTIVITĂȚI-GREEDY** este foarte eficientă. Ea poate planifica o mulțime S de n activități în $\Theta(n)$, presupunând că activitățile au fost deja ordonate după timpul lor de terminare. Activitatea aleasă de procedura **SELECTOR-ACTIVITĂȚI-**

GREEDY este întotdeauna cea cu primul timp de terminare care poate fi planificată legal. Activitatea astfel selectată este o alegere "greedy" (lacomă) în sensul că, intuitiv, ea lasă posibilitatea celorlalte activități rămase pentru a fi planificate. Cu alte cuvinte, alegerea greedy maximizează cantitatea de timp neplanificată rămasă.

i	s_i	f_i
1	1	4
2	3	5
3	0	6
4	5	7
5	3	8
6	5	9
7	6	10
8	8	11
9	8	12
10	2	13
11	12	14

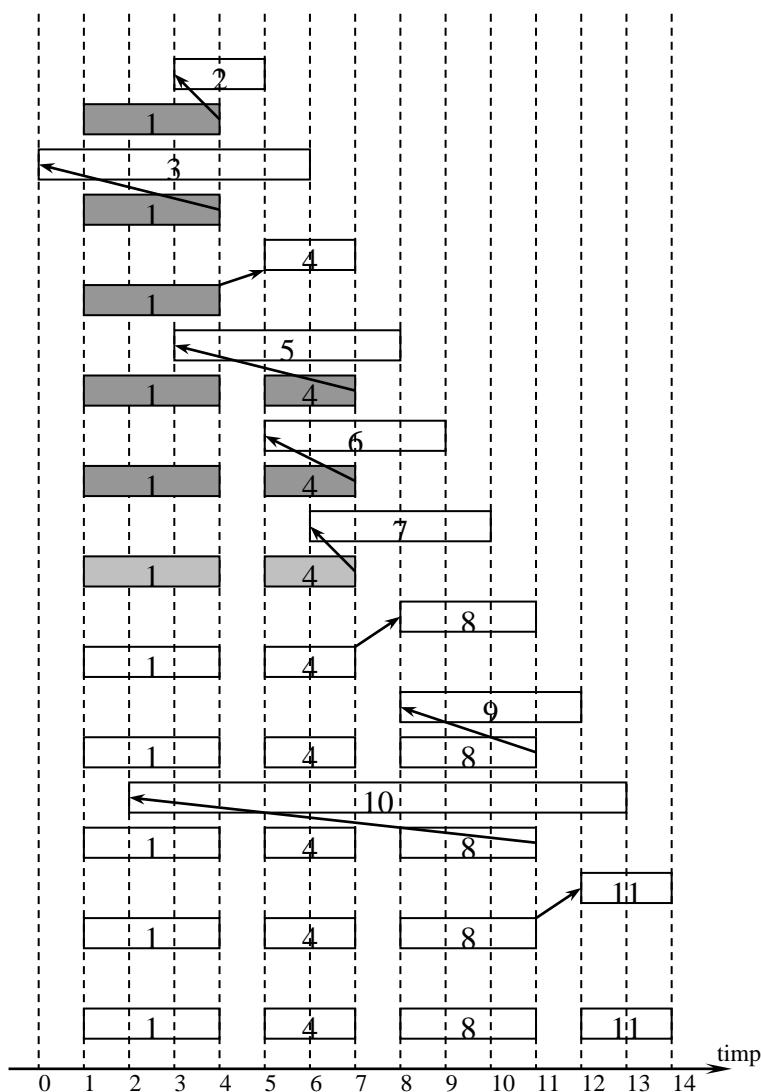


Figura 2.1. Operațiile algoritmului **SELECTOR-ACTIVITĂȚI-GREEDY** asupra celor 11 activități date în stânga. Fiecare linie a figurii corespunde unei iterații din ciclul pentru din liniile 4-7. Activitățile care au fost selectate pentru a fi incluse în mulțimea A sunt hașurate, iar activitatea curentă i este nehașurată. Dacă timpul de pornire s_i , al activității i este mai mic decât timpul de terminare al activității j (săgeata dintre ele este spre stânga), activitatea este ignorată, în caz contrar (săgeata este îndreptată spre dreapta), activitatea este acceptată și este adăugată mulțimii A .

2.1.1. Demonstrarea corectitudinii algoritmului greedy

Algoritmii de tip greedy nu furnizează întotdeauna soluțiile optime. Cu toate acestea, algoritmul **SELECTOR-ACTIVITĂȚI-GREEDY** determină întotdeauna o soluție optimă pentru o instanță a problemei selectării activităților.

Teorema 2.1. Algoritmul **SELECTOR-ACTIVITĂȚI-GREEDY** furnizează soluții de dimensiune maximă pentru problema selectării activităților.

Demonstrație. Fie $S = \{1, 2, \dots, n\}$ mulțimea activităților care trebuie planificate. Deoarece presupunem că activitățile sunt ordonate după timpul de terminare, activitatea 1 se va termina cel mai devreme. Vrem să arătăm că există o soluție optimă care începe cu activitatea 1, conform unei alegeri greedy.

Să presupunem că $A \subseteq S$ este o soluție optimă pentru o instanță dată a problemei selectării activităților. Vom ordona activitățile din A după timpul de terminare. Mai presupunem că prima activitate din A este activitatea k . Dacă $k = 1$, planificarea mulțimii A începe cu o alegere greedy. Dacă $k \neq 1$ vrem să arătăm că există o altă soluție optimă B a lui S care începe conform alegerii greedy, cu activitatea 1. Fie $B = A - \{k\} \cup \{1\}$. Deoarece $f_1 \leq f_k$,

activitățile din B sunt distincte, și cum B are același număr de activități ca și A , B este de asemenea optimă. Deci B este o soluție optimă pentru S , care conține alegerea greedy a activității 1. Am arătat astfel că există întotdeauna o planificare optimă care începe cu o alegere greedy.

Mai mult, o dată ce este făcută alegerea greedy a activității 1, problema se reduce la determinarea soluției optime pentru problema selectării acelor activități din S care sunt compatibile cu activitatea 1. Aceasta înseamnă că dacă A este o soluție optimă pentru problema inițială, atunci $A' = A - \{1\}$ este o soluție optimă pentru problema selectării activităților

$S' = \{i \in S, s_i \geq f_1\}$. De ce? Dacă am putea găsi o soluție B' pentru S' cu mai multe activități decât A' , adăugarea activității 1 la B' va conduce la o soluție B a lui S cu mai multe activități decât A , contrazicându-se astfel optimalitatea mulțimii A . În acest mod, după ce este făcută fiecare alegere greedy, ceea ce rămâne este o problemă de optimizare de aceeași formă ca problema inițială. Prin inducție după numărul de alegeri făcute, se arată că realizând o alegere greedy la fiecare pas, se obține o soluție optimă.

2.2. Coduri Huffman

Codurile Huffman reprezintă o tehnică foarte utilizată și eficientă pentru compactarea datelor; în funcție de caracteristicile fișierului care trebuie comprimat, spațiul economisit este între 20% și 90%. Algoritmul greedy pentru realizarea acestei codificări utilizează un tabel cu frecvențele de apariție ale fiecărui caracter. Ideea este de a utiliza o modalitate optimă pentru reprezentarea fiecărui caracter sub forma unui șir binar.

Sa presupunem că avem un fișier ce conține 100.000 de caractere, pe care dorim să îl memorăm într-o formă compactată. Frecvențele de apariție ale caracterelor în text sunt date de figura 2.2: există doar șase caractere diferite și fiecare dintre ele apare de 45.000 de ori.

	a	b	c	d	e	f
Frecvența (în mii)	45	13	12	16	9	5
Cuvânt codificat; lungime fixă	000	001	010	011	100	101
Cuvânt codificat; lungime variabilă	0	101	100	111	1101	1100

Figura 2.2. O problemă de codificare a caracterelor. Un fișier de 100.000 de caractere conține doar caracterele a–f, cu frecvențele indicate. Dacă fiecărui caracter îi este asociat un cuvânt de cod pe 3 biți, fișierul codificat va avea 300.000 de biți. Folosind codificarea cu lungime variabilă, fișierul codificat va ocupa doar 224.000 de biți.

Există mai multe modalități de reprezentare a unui astfel de fișier. Vom considera problema proiectării unui cod binar al caracterelor (pe scurt cod) astfel încât fiecare caracter este reprezentat printr-un șir binar unic. Dacă utilizăm un cod de lungime fixă, avem nevoie de 3 biți pentru a reprezenta șase caractere: a=000, b=001, ..., f=101. Această metodă necesită 300.000 de biți pentru a codifica tot fișierul. Se pune problema dacă se poate face o compactare și mai bună.

O codificare cu lungime variabilă poate îmbunătăți semnificativ performanțele, atribuind caracterelor cu frecvențe mai mari cuvinte de cod mai scurte iar celor cu frecvențe mai reduse cuvinte de cod mai lungi. Figura 2.3(b) prezintă o astfel de codificare; șirul 0 având lungimea de 1 bit reprezintă caracterul a în timp ce șirul 1100 de lungime 4 reprezintă caracterul f. Această codificare necesită $(45 \cdot 1 + 13 \cdot 3 + 12 \cdot 3 + 16 \cdot 3 + 9 \cdot 4 + 5 \cdot 4) \cdot 1.000 = 224.000$ biți pentru a reprezenta un fișier și economisește aproximativ 25% din spațiu. Vom vedea că aceasta este de fapt o codificare-caracter optimă pentru acest fișier.

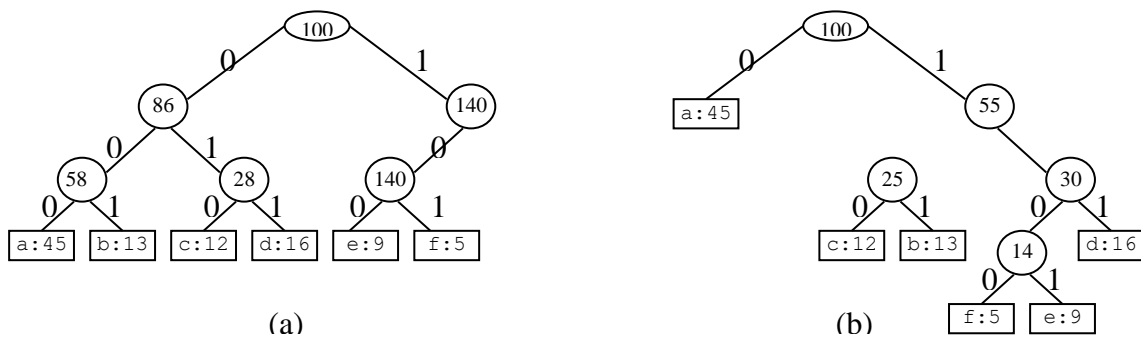


Figura 2.3. Arborii corespunzătorii schemelor de codificare din figura 2.2. Fiecare frunză este etichetată cu un caracter și cu frecvența de apariție a acestuia. Fiecare nod intern este etichetat cu suma ponderilor frunzelor din

subarborele aferent. (a) Arborele corespunzător codificării de lungime fixă $a=000, \dots, f=101$. (b) Arborele asociat codificării prefix optime $a=0, b=101, \dots, f=1100$.

2.2.1. Coduri prefix

Vom considera în continuare doar codificările în care nici un cuvânt de cod nu este prefixul altui cuvânt. Astfel de codificări se numesc *codificări prefix*. Se poate arăta că o compresie optimă a datelor, realizată prin codificarea caracterelor, poate fi realizată și prin codificarea prefix, deci considerarea codificării prefix nu scade din generalitate.

Codificarea prefix este utilă deoarece simplifică atât codificarea (deci compactarea) cât și decodificarea. Codificarea este întotdeauna simplă pentru orice codificare binară a caracterelor; se concatenează cuvintele de cod reprezentând fiecare caracter al fișierului. De exemplu, prin codificarea prefix având lungime variabilă din figura 17.3, putem codifica un fișier de 3 caractere abc astfel: $0 \cdot 101 \cdot 100 = 0101100$, unde semnul " \cdot " reprezintă operația de concatenare.

Decodificarea este de asemenea relativ simplă pentru codurile prefix. Cum nici un cuvânt de cod nu este prefixul altuia. Începutul oricărui fișier codificat, nu este ambiguu. Putem deci să identificăm cuvântul inițial de cod, să îl "traducem" în caracterul original, să-l îndepărtăm din fișierul codificat și să repetăm procesul pentru fișierul codificat rămas. În exemplul nostru, șirul 001011101 se translatează automat în $0 \cdot 0 \cdot 101 \cdot 1101$, secvență care se decodifică în $aabe$.

Procesul de decodificare necesită o reprezentare convenabilă a codificării prefix astfel încât cuvântul inițial de cod să poată fi ușor identificat. O astfel de reprezentare poate fi dată de un arbore binar ale cărui frunze sunt caracterele date. Interpretăm un cuvânt de cod binar pentru un caracter ca fiind drumul de la rădăcina până la caracterul respectiv, unde 0 reprezintă "mergi la fiul stâng" iar 1 "mergi la fiul drept". În figura 3 sunt prezentați arborii pentru cele două codificări ale exemplului nostru. Observați că aceștia nu sunt arbori binari de căutare, deoarece frunzele nu trebuie să fie ordonate, iar nodurile interne nu conțin chei pentru caractere.

O codificare optimă pentru un fișier este întotdeauna reprezentată printr-un arbore binar complet, în care fiecare vârf care nu este frunză are doi fii.

Codificarea cu lungime fixă din exemplul de mai sus nu este optimă deoarece arborele asociat, prezentat în figura 3(a), nu este un arbore binar complet: există două cuvinte de cod care încep cu $10\dots$, dar nici unul care să înceapă cu $11\dots$. Conform celor de mai sus ne putem restrânge atenția numai asupra arborilor binari compleți, deci dacă C este alfabetul din care fac parte caracterele, atunci arborele pentru o codificare prefix optimă are exact $|C|$ frunze, una pentru fiecare literă din alfabet, și exact $|C| - 1$ noduri interne.

Dându-se un arbore T , corespunzător unei codificări prefix, este foarte simplu să calculăm numărul de biți necesari pentru a codifica un fișier. Pentru fiecare caracter c din alfabet, fie $f(c)$ frecvența lui c în fișier și să notăm cu $d_T(c)$ adâncimea frunzei în arbore (adică nivelul pe care se află). Să observăm că $d_T(c)$ reprezintă de asemenea cuvântul de cod pentru caracterul c . Numărul de biți necesari pentru a codifica un fișier este

$$B(T) = \sum_{c \in C} f(c) d_T(c)$$

Vom numi acest număr *costul* arborelui T .

2.2.2. Construcția unui cod Huffman

Huffman a inventat un algoritm greedy care construiește o codificare prefix optimă numită *codul Huffman*. Algoritmul construiește arborele corespunzător codificării optime, într-o manieră "bottom-up". Se începe cu o mulțime de $|C|$ frunze și se realizează o secvență de $|C| - 1$ operații de interclasare pentru a crea arborele final.

În algoritmul care urmează, vom presupune că C este o mulțime de n caractere și fiecare caracter $c \in C$ este un obiect având o frecvență dată $f[c]$. Va fi folosită o coadă de priorități pentru a identifica cele două obiecte cu frecvența cea mai redusă care vor fuziona. Rezultatul interclasării celor două obiecte este un nou obiect a cărui frecvență este suma frecvențelor celor două obiecte care au fost interclasate.

HUFFMAN (C)

1: $n \leftarrow |C|$

2: $Q \leftarrow C$

3: **pentru** $i \leftarrow 1, n - 1$ **execută**

4: $z \leftarrow \text{ALOCĂ-NOD}()$

5: $x \leftarrow \text{stânga}[z] \leftarrow \text{EXTRAGE-MIN}(Q)$

6: $y \leftarrow \text{dreapta}[z] \leftarrow \text{EXTRAGE-MIN}(Q)$

7: $f[z] \leftarrow f[x] + f[y]$

8: **INSEREAZĂ** (Q, z)

9: **returnează** $\text{EXTRAGE-MIN}(Q)$

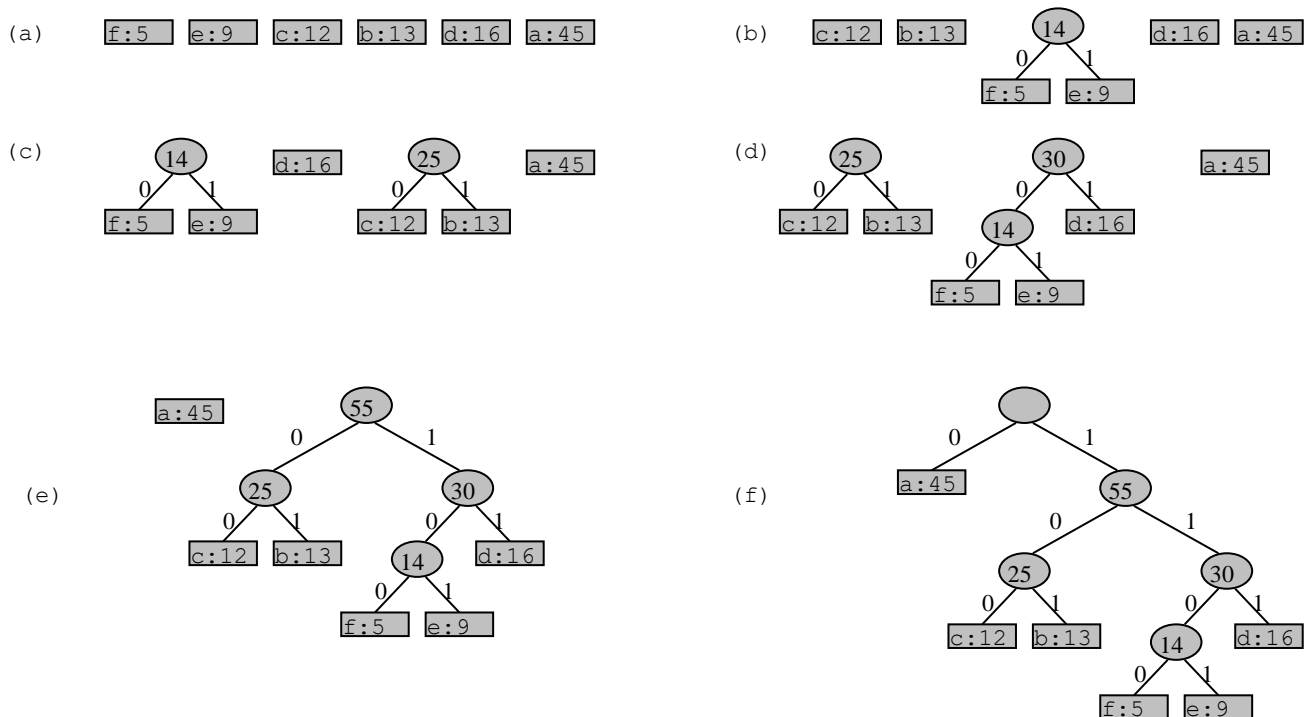


Figura 2.4 Pași algoritmului Huffman pentru frecvențele date în figura 2.2. Fiecare parte ilustrează conținutul cozii ordonate crescător după frecvență. La fiecare pas, cei doi arbori cu cele mai scăzute frecvențe fuzionează. Frunzele figurează ca dreptunghiuri ce conțin un caracter și frecvența sa. Nodurile interne figurează ca cercuri ce conțin suma frecvențelor fiilor lor. O muchie care leagă un nod intern cu fiii ei este etichetată cu 0 dacă este o muchie către un fiu stâng, respectiv cu 1 dacă este către fiul drept. Cuvântul de cod pentru o literă este secvența etichetelor de pe muchiile care leagă rădăcina de frunza asociată caracterului respectiv. (a) Mulțimea inițială de $n=6$ noduri, unul pentru fiecare literă. (b)-(e) Etape intermediare. (f) Arborele final.

Pentru exemplul nostru, algoritmul lui Huffman lucrează ca în figura 2.3. Deoarece există 6 litere în alfabet, dimensiunea inițială a cozii este $n = 6$, iar pentru a construi arborele sunt necesari 5 pași de interclasare. Arborele final reprezintă codificarea prefix optimă. Codificarea unei litere este secvența etichetelor nodurilor ce formează drumul de la rădăcină la literă.

În linia 2 coada de priorități Q este inițializată cu caracterele din C . Ciclul *for* din liniile 3-8 extrage în mod repetat două noduri x și y cu cea mai mică frecvență din coadă, și le înlocuiește în coadă cu un nou nod z , reprezentând fuziunea lor. Frecvența lui z este calculată în linia 7 ca fiind suma frecvențelor lui x și y . Nodul z are pe x ca fiu stâng și pe y ca fiu drept. (Această ordine este arbitrară; interschimbarea fiilor stâng și drept ai oricărui nod conduce la o codificare diferită, dar de același cost.) După $n - 1$ fuzionări, unicul nod rămas în coadă, rădăcina arborelui de codificare, este returnat în linia 9.

2.2.3. Corectitudinea algoritmului lui Huffman

Pentru a demonstra că algoritmul de tip greedy al lui Huffman este corect, vom arăta că problema determinării unei codificări prefix optime implică alegeri greedy și are o substructură optimă. Următoarea lemă se referă la proprietatea alegerii greedy.

Lema 2.1. Fie C un alfabet în care fiecare caracter $c \in C$ are frecvența $f[c]$. Fie x și y două caractere din C având cele mai mici frecvențe. Atunci există o codificare prefix optimă pentru C în care cuvintele de cod pentru x și y au aceeași lungime și diferă doar pe ultimul bit.

Demonstrație. Ideea demonstrației este de a lua arborele T reprezentând o codificare prefix optimă și a-l modifica pentru a realiza un arbore reprezentând o altă codificare prefix optimă. În noul arbore, caracterele x și y vor apărea ca frunze cu același tată și se vor afla pe nivelul maxim în arbore. Dacă putem realiza acest lucru, atunci cuvintele lor de cod vor avea aceeași lungime și vor diferi doar pe ultimul bit.

Fie b și c două caractere reprezentând noduri terminale (frunze) situate pe nivelul maxim al arborelui T . Fără a restrânge generalitatea, vom presupune că $f[b] \leq f[c]$ și $f[x] \leq f[y]$. Cum $f[x]$ și $f[y]$ sunt frunzele cu frecvențele cele mai scăzute, în această ordine, iar $f[b]$ și $f[c]$ sunt deci frecvențe arbitrare, în ordine, avem $f[x] \leq f[b]$ și $f[y] \leq f[c]$. După cum se vede în figura 2.5, vom schimba în T pozițiile lui b și x pentru a produce arborele T' , iar apoi vom schimba în T' pozițiile lui c și y pentru a obține arborele T'' . Conform ecuației (17.3), diferența costurilor lui T și T' este

$$\begin{aligned} B(T) - B(T') &= \sum_{c \in C} f(c)d_T(c) - \sum_{c \in C} f(c)d_{T'}(c) \\ &= f[x]d_T(x) + f[b]d_T(b) - f[x]d_{T'}(x) - f[b]d_{T'}(b) \\ &= f[x]d_T(x) + f[b]d_T(b) - f[x]d_{T'}(b) - f[b]d_{T'}(x) \\ &= (f[b] - f[x])(d_T(b) - d_T(x)) \geq 0, \end{aligned}$$

deoarece atât $f[b] - f[x]$ și $d_T[b] - d_T[x]$ sunt nenegative. Mai precis, $f[b] - f[x]$ este nenegativ deoarece x este frunză având frecvența minimă iar $d_T[b] - d_T[x]$ este nenegativ pentru că b este o frunză aflată pe nivelul maxim în T . În mod analog, deoarece interschimbarea lui y cu c nu mărește costul, diferența $B(T) \leq B(T'')$ este nenegativă. Astfel $B(T) \leq B(T'')$, ceea ce implică $B(T'') = B(T)$. Așadar T'' este un arbore optim în care x și y apar ca noduri terminale frați, și se află pe nivelul maxim, ceea ce trebuia demonstrat.

Din lema 2.1 se deduce faptul ca procesul de construcție a unui arbore optim prin fuzionări poate, fără a restrânge generalitatea, să înceapă cu o alegere greedy a interclasării acelor două caractere cu cea mai redusă frecvență. De ce este aceasta o alegere greedy? Putem interpreta costul unei singure interclasări ca fiind suma frecvențelor celor două obiecte care fuzionează. La fiecare pas, dintre toate interclasările posibile, algoritmul HUFFMAN o alege pe aceea care determină costul minim.

Lema 2.2. Fie T un arbore binar complet reprezentând o codificare prefix optimă peste un alfabet C , unde frecvența $f[c]$ este definită pentru fiecare caracter $c \in C$. Considerăm două caractere x și y oarecare care apar ca noduri terminale frați în T , și fie z tatăl lor. Atunci, considerând z ca un caracter având frecvența $f[z] = f[x] + f[y]$, arborele $T' = T - \{x, y\}$ reprezintă o codificare prefix optimă pentru alfabetul $C' = C - \{x, y\} \cup \{z\}$.

Demonstrație. Vom arăta mai întâi că $B(T)$, costul arborelui T , poate fi exprimat în funcție de costul $B(T')$ al arborelui T' considerând costurile componente din ecuația (17.3). Pentru fiecare $c \in C - \{x, y\}$, avem $d_T[c] = d_{T'}[c]$, și $f[c]d_T[c] = f[c]d_{T'}[c]$. Cum $d_T[x] = d_T[y] = d_{T'}[z] + 1$, avem

$$f[x]d_T(x) + f[y]d_T(y) = f[x] + f[y]d_{T'}(z) + 1 = f[z]d_{T'}(z) + (f[x] + f[y]),$$

de unde deducem că

$$B(T) = B(T') + f[x] + f[y].$$

Dacă T' reprezintă o codificare prefix care nu este optimă pentru alfabetul C' , atunci există un arbore T'' ale cărui frunze sunt caractere în C' astfel încât $B(T'') < B(T')$. Cum z este tratat ca un caracter în C' , el va apare ca frunză în T'' . Dacă adăugăm x și y ca fiind fiii lui z în T'' , atunci vom obține o codificare prefix pentru C având costul $B(T'') + f[x] + f[y] < B(T)$, ceea ce intră în contradicție cu optimalitatea lui T . Deci T' trebuie să fie optim pentru alfabetul C' .

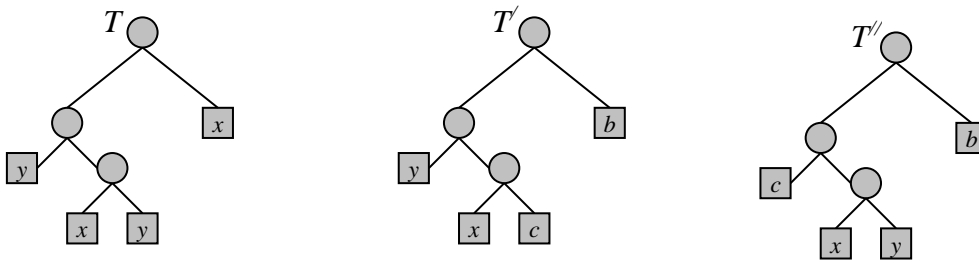


Figura 2.5 O ilustrare a pasului cheie din demonstrația lemei 2.1. În arborele optim T , b și c sunt două dintre frunzele aflate pe cel mai de jos nivel în arbore; aceste noduri se consideră frați. x și y sunt două frunze pe care algoritmul Huffman le interclasează primele; ele apar în poziții arbitrare în T . Frunzele b și x sunt interschimbate pentru a obține arborele T' . Apoi, frunzele c și y sunt interschimbate pentru a obține arborele T'' . Cum cele două interschimbări nu măresc costul, arborele rezultat T'' este de asemenea un arbore optim.

2.3. Arbori parțiali de cost minim

Fie $G = \langle V, M \rangle$ un graf neorientat conex, unde V este mulțimea vârfurilor și M este mulțimea muchiilor. Fiecare muchie are un *cost* nenegativ (sau o *lungime* nenegativă). Problema este să găsim o submulțime $A \subseteq M$, astfel încât toate vârfurile din V să rămână conectate atunci când sunt folosite doar muchii din A , iar suma

lungimilor muchiilor din A sa fie cat mai mica. Căutam deci o submulțime A de cost total minim. Aceasta problema se mai numește și *problema conectării orașelor cu cost minim*, având numeroase aplicații.

Graful parțial $\langle V, A \rangle$ este un arbore și este numit *arborele parțial de cost minim* al grafului G (*minimal spanning tree*). Un graf poate avea mai mulți arbori parțiali de cost minim și acest lucru se poate verifica pe un exemplu.

Vom prezenta doi algoritmi greedy care determina arborele parțial de cost minim al unui graf. În terminologia metodei greedy, vom spune că o mulțime de muchii este o *soluție*, dacă constituie un arbore parțial al grafului G , și este *fezabilă*, dacă nu conține cicluri. O mulțime fezabilă de muchii este *promițătoare*, dacă poate fi completată pentru a forma soluția optimă. O muchie *atinge* o mulțime dată de vârfuri, dacă exact un capăt al muchiei este în mulțime. Următoarea proprietate va fi folosită pentru a demonstra corectitudinea celor doi algoritmi.

Proprietatea 2.1 Fie $G = \langle V, M \rangle$ un graf neorientat conex în care fiecare muchie are un cost nenegativ. Fie $W \subset V$ o submulțime strictă a vârfurilor lui G și fie $A \subseteq M$ o mulțime promițătoare de muchii, astfel încât nici o muchie din A nu atinge W . Fie m muchia de cost minim care atinge W . Atunci, $A \cup \{m\}$ este promițătoare.

Demonstrație: Fie B un arbore parțial de cost minim al lui G , astfel încât $A \subseteq B$ (adică, muchiile din A sunt conținute în arborele B). Un astfel de B trebuie să existe, deoarece A este promițătoare. Dacă $m \in B$, nu mai rămâne nimic de demonstrat. Presupunem că $m \notin B$. Adăugându-l pe m la B , obținem exact un ciclu (Exercițiul 3.2). În acest ciclu, deoarece m atinge W , trebuie să mai existe cel puțin o muchie m' care atinge și ea pe W (altfel, ciclul nu se închide). Eliminându-l pe m' , ciclul dispăre și obținem un nou arbore parțial B' al lui G . Costul lui m este mai mic sau egal cu costul lui m' , deci costul total al lui B' este mai mic sau egal cu costul total al lui B . De aceea, B' este și el un arbore parțial de cost minim al lui G , care include pe m . Observăm că $A \subseteq B'$ deoarece muchia m' , care atinge W , nu poate fi în A . Deci, $A \cup \{m\}$ este promițătoare.

Mulțimea inițială a candidaților este M . Cei doi algoritmi greedy aleg muchiile una câte una într-o anumită ordine, această ordine fiind specifică fiecărui algoritm.

2.3.1. Algoritmul lui Kruskal

Arborele parțial de cost minim poate fi construit muchie, cu muchie, după următoarea metoda a lui Kruskal (1956): se alege întâi muchia de cost minim, iar apoi se adaugă repetat muchia de cost minim nealeasă anterior și care nu formează cu precedentele un ciclu. Alegem astfel $\#V - 1$ muchii. Este ușor de dedus că obținem în final un arbore. Este însă acesta chiar arborele parțial de cost minim căutat?

Înainte de a răspunde la întrebare, să considerăm, de exemplu, graful din Figura 2.6(a). Ordonăm crescător (în funcție de cost) muchiile grafului: $\{1, 2\}$, $\{2, 3\}$, $\{4, 5\}$, $\{6, 7\}$, $\{1, 4\}$, $\{2, 5\}$, $\{4, 7\}$, $\{3, 5\}$, $\{2, 4\}$, $\{3, 6\}$, $\{5, 7\}$, $\{5, 6\}$ și apoi aplicăm algoritmul. Structura componentelor conexe este ilustrată, pentru fiecare pas, în Figura 2.7.

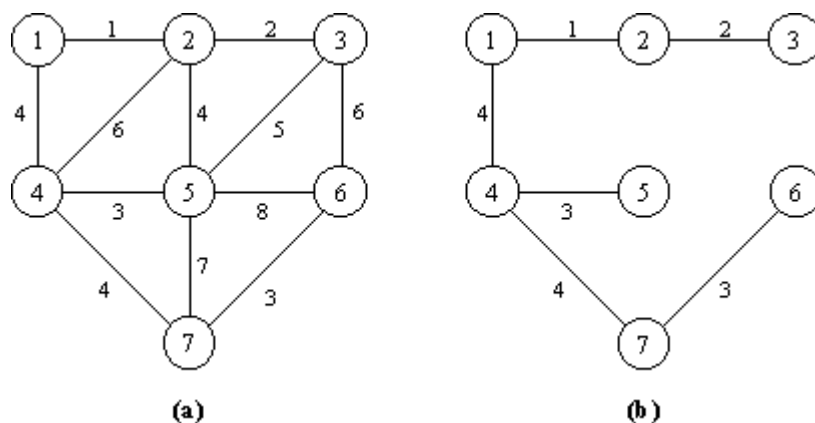


Figura 2.6 Un graf și arborele său parțial de cost minim.

Pasul	Muchia considerata	Componentele subgrafului $\langle V, A \rangle$	conexe	ale
inițializare	—	{1}, {2}, {3}, {4}, {5}, {6}, {7}		
1	{1, 2}	{1, 2}, {3}, {4}, {5}, {6}, {7}		
2	{2, 3}	{1, 2, 3}, {4}, {5}, {6}, {7}		
3	{4, 5}	{1, 2, 3}, {4, 5}, {6}, {7}		
4	{6, 7}	{1, 2, 3}, {4, 5}, {6, 7}		
5	{1, 4}	{1, 2, 3, 4, 5}, {6, 7}		
6	{2, 5}	respinsa (formează ciclul)		
7	{4, 7}	{1, 2, 3, 4, 5, 6, 7}		

Figura 2.7 Algoritmul lui Kruskal aplicat grafului din Figura 2.6(a)

Mulțimea A este inițial vidă și se completează pe parcurs cu muchii acceptate (care nu formează un ciclu cu muchiile deja existente în A). În final, mulțimea A va conține muchiile $\{1, 2\}, \{2, 3\}, \{4, 5\}, \{6, 7\}, \{1, 4\}, \{4, 7\}$. La fiecare pas, graful parțial $\langle V, A \rangle$ formează o pădure de componente conexe, obținută din pădurea precedentă unind două componente. Fiecare componentă conexă este la rândul ei un arbore parțial de cost minim pentru vârfurile pe care le conectează. Inițial, fiecare vârf formează o componentă conexă. La sfârșit, vom avea o singură componentă conexă, care este arborele parțial de cost minim căutat (Figura 2.6(b)).

Ceea ce am observat în acest caz particular este valabil și pentru cazul general, din Proprietatea 2.1 rezultând:

Proprietatea 2.2 În algoritmul lui Kruskal, la fiecare pas, graful parțial $\langle V, A \rangle$ formează o pădure de componente conexe, în care fiecare componentă conexă este la rândul ei un arbore parțial de cost minim pentru vârfurile pe care le conectează. În final, se obține arborele parțial de cost minim al grafului G .

Pentru a implementa algoritmul, trebuie să putem manipula submulțimile formate din vârfurile componentelor conexe. Folosim pentru aceasta o structură de mulțimi disjuncte și procedurile de tip *find* și *merge* (Vezi disciplina „Structuri de date și algoritmi”). În acest caz, este preferabil să reprezentăm graful ca o listă de muchii cu costul asociat lor, astfel încât să putem ordona această listă în funcție de cost. Iată algoritmul:

KRUSKAL($G = \langle V, M \rangle$)

1. {inițializare}
2. sortează M crescător în funcție de cost
3. $n \leftarrow \#V$
4. $A \leftarrow \emptyset$ {va conține muchiile arborelui parțial de cost minim}
5. {inițializează n mulțimi disjuncte conținând fiecare câte un element din V }

{bucla greedy}

repeat

{ u, v } \leftarrow muchia de cost minim care
încă nu a fost considerată

$u_{comp} \leftarrow find(u)$

$v_{comp} \leftarrow find(v)$

if $u_{comp} \neq v_{comp}$ **then** $merge(u_{comp}, v_{comp})$

$A \leftarrow A \cup \{\{u, v\}\}$

until $\#A = n-1$

return A

Pentru un graf cu n vârfuri și m muchii, presupunând că se folosesc procedurile *find3* și *merge3*, numărul de operații pentru cazul cel mai nefavorabil este în:

- $O(m \log m)$ pentru a sorta muchiile. Deoarece $m \leq n(n-1)/2$, rezulta $O(m \log m) \subseteq O(m \log n)$. Mai mult, graful fiind conex, din $n-1 \leq m$ rezulta și $O(m \log n) \subseteq O(m \log m)$, deci $O(m \log m) = O(m \log n)$.
- $O(n)$ pentru a inițializa cele n mulțimi disjuncte.
- Cele cel mult $2m$ operații *find3* și $n-1$ operații *merge3* necesită un timp în $O((2m+n-1) \lg^* n)$, după cum am specificat în Capitolul 3. Deoarece $O(\lg^* n) \subseteq O(\log n)$ și $n-1 \leq m$, acest timp este și în $O(m \log n)$.
- $O(m)$ pentru restul operațiilor.

Deci, pentru cazul cel mai nefavorabil, algoritmul lui Kruskal necesită un timp în $O(m \log n)$.

O altă variantă este să păstrăm muchiile într-un min-heap. Obținem astfel un nou algoritm, în care inițializarea se face într-un timp în $O(m)$, iar fiecare din cele $n-1$ extrageri ale unei muchii minime se face într-un timp în $O(\log m) = O(\log n)$. Pentru cazul cel mai nefavorabil, ordinul timpului rămâne același cu cel al vechiului algoritm.

Avantajul folosirii min-heap-ului apare atunci când arborele parțial de cost minim este găsit destul de repede și un număr considerabil de muchii rămân netestate. În astfel de situații, algoritmul vechi pierde timp, sortând în mod inutil și aceste muchii.

5.2.2 Algoritmul lui Prim

Cel de-al doilea algoritm greedy pentru determinarea arborelui parțial de cost minim al unui graf se datorează lui Prim (1957). În acest algoritm, la fiecare pas, mulțimea A de muchii alese împreună cu mulțimea U a vârfurilor pe care le conectează formează un arbore parțial de cost minim pentru subgraful $\langle U, A \rangle$ al lui G . Inițial, mulțimea U a vârfurilor acestui arbore conține un singur vârf oarecare din V , care va fi rădăcina, iar mulțimea A a muchiiilor este vidă. La fiecare pas, se alege o muchie de cost minim, care se adaugă la arborele precedent, dând naștere unui nou arbore parțial de cost minim (deci, exact una dintre extremitățile acestei muchii este un vârf în arborele precedent). Arborele parțial de cost minim crește “natural”, cu câte o ramură, până când va atinge toate vârfurile din V , adică până când $U = V$. Funcționarea algoritmului, pentru exemplul din Figura 6(a), este ilustrată în Tabelul 6.2. La sfârșit, A va conține aceleași muchii ca și în cazul algoritmului lui Kruskal. Faptul că algoritmul funcționează întotdeauna corect este exprimat de următoarea proprietate, pe care o puteți demonstra folosind Proprietatea 6.2.

Pasul	Muchia considerata	U
inițializare	—	{1}
1	{2, 1}	{1, 2}
2	{3, 2}	{1, 2, 3}
3	{4, 1}	{1, 2, 3, 4}
4	{5, 4}	{1, 2, 3, 4, 5}
5	{7, 4}	{1, 2, 3, 4, 5, 6}
6	{6, 7}	{1, 2, 3, 4, 5, 6, 7}

Tabelul 5.2 Algoritmul lui Prim aplicat grafului din Figura 6.4a.

Proprietatea 5.4 În algoritmul lui Prim, la fiecare pas, $\langle U, A \rangle$ formează un arbore parțial de cost minim pentru subgraful $\langle U, A \rangle$ al lui G . În final, se obține arborele parțial de cost minim al grafului G .

Descrierea formală a algoritmului este dată în continuare.

```

function Prim-formal( $G = \langle V, M \rangle$ )
  {inițializare}
   $A \leftarrow \emptyset$  {va conține muchiile arborelui parțial de cost minim}
   $U \leftarrow$  {un vârf oarecare din  $V$ }
  {bucla greedy}
  while  $U \neq V$  do
    găsește  $\{u, v\}$  de cost minim astfel ca  $u \in V \setminus U$  și  $v \in U$ 
     $A \leftarrow A \cup \{\{u, v\}\}$ 
     $U \leftarrow U \cup \{u\}$ 
  return  $A$ 

```

Pentru a obține o implementare simplă, presupunem că: vârfurile din V sunt numerotate de la 1 la n , $V = \{1, 2, \dots, n\}$; matricea simetrică C da costul fiecărei muchii, cu $C[i, j] = +\infty$, dacă muchia $\{i, j\}$ nu există. Folosim două tablouri paralele. Pentru fiecare $i \in V \setminus U$, $vecin[i]$ conține vârful din U , care este conectat de i printr-o muchie de cost minim; $mincost[i]$ da acest cost. Pentru $i \in U$, punem $mincost[i] = -1$. Mulțimea U , în mod arbitrar inițializată cu $\{1\}$, nu este reprezentată explicit. Elementele $vecin[1]$ și $mincost[1]$ nu se folosesc.

```

function Prim( $C[1 .. n, 1 .. n]$ )
  {inițializare; numai vârful 1 este în  $U$ }
   $A \leftarrow \emptyset$ 
  for  $i \leftarrow 2$  to  $n$  do  $vecin[i] \leftarrow 1$ 
     $mincost[i] \leftarrow C[i, 1]$ 
  {bucla greedy}
  repeat  $n-1$  times
     $min \leftarrow +\infty$ 
    for  $j \leftarrow 2$  to  $n$  do
      if  $0 < mincost[j] < min$  then  $min \leftarrow mincost[j]$ 

```

```

        k ← j
    A ← A ∪ { {k, vecin[k]} }
    mincost[k] ← -1 {adaugă vârful k la U}
    for j ← 2 to n do
        if C[k, j] < mincost[j] then    mincost[j] ← C[k, j]
                                       vecin[j] ← k
    return A

```

Bucula principală se execută de $n-1$ ori și, la fiecare iterație, buclele **for** din interior necesită un timp în $O(n)$. Algoritmul *Prim* necesită, deci, un timp în $O(n^2)$. Am văzut că timpul pentru algoritmul lui Kruskal este în $O(m \log n)$, unde $m = \#M$. Pentru un graf *dens* (adică, cu foarte multe muchii), se deduce că m se apropie de $n(n-1)/2$. În acest caz, algoritmul *Kruskal* necesită un timp în $O(n^2 \log n)$ și algoritmul *Prim* este probabil mai bun. Pentru un graf *rar* (adică, cu un număr foarte mic de muchii), m se apropie de n și algoritmul *Kruskal* necesită un timp în $O(n \log n)$, fiind probabil mai eficient decât algoritmul *Prim*.

6. PROGRAMAREA DINAMICĂ

Programarea dinamică, asemenea metodelor *divide și stăpânește*, rezolvă problemele combinând soluțiile unor subprobleme. (În acest context, termenul de "programare" se referă la o metodă tabelară și nu la scrierea codului pentru un calculator.) După cum am arătat mai sus, algoritmi *divide și stăpânește* partiționează problema în subprobleme independente, rezolvă recursiv subproblemele și apoi combina soluțiile lor pentru a rezolva problema inițială. Spre deosebire de această abordare, programarea dinamică este aplicabilă atunci când subproblemele nu sunt independente, adică subproblemele au în comun sub-subprobleme. În acest context, un algoritm de tipul *divide și stăpânește* ar presupune mai multe calcule decât ar fi necesar dacă s-ar rezolva în mod repetat sub-subproblemele comune. Un algoritm bazat pe programare dinamică rezolvă fiecare sub-subproblemă o singură dată și memorează soluția acesteia într-un tablou, prin aceasta evitând recalcularea soluției ori de câte ori respectiva sub-subproblemă apare din nou.

În general, metoda programării dinamice se aplică *problemelor de optimizare*. Asemenea probleme pot avea mai multe soluții posibile. Fiecare soluție are o valoare iar ceea ce se dorește este determinarea soluției a cărei valoare este optimă (minimă sau maximă). O asemenea soluție se numește o *soluție* optimă a problemei prin contrast cu *soluția* optimă, deoarece pot fi mai multe soluții care realizează valoarea optimă.

Dezvoltarea unui algoritm bazat pe programarea dinamică poate fi împărțită într-o secvență de patru pași.

1. Caracterizarea structurii unei soluții optime.
2. Definirea recursivă a valorii unei soluții optime.
3. Calculul valorii unei soluții optime într-o manieră de tip "bottom-up".
4. Construirea unei soluții optime din informația calculată.

Pașii 1-3 sunt baza unei abordări de tip programare dinamică. Pasul 4 poate fi omis dacă se dorește doar calculul unei singure soluții optime. În vederea realizării pasului 4, deseori se păstrează informație suplimentară de la execuția pasului 3, pentru a ușura construcția unei soluții optime.

6.1. Înmulțirea unui șir de matrice

Primul nostru exemplu de programare dinamică este un algoritm care rezolvă problema înmulțirii unui șir de matrice. Se dă un șir (o secvență) (A_1, A_2, \dots, A_n) de n matrice care trebuie înmulțite, și se dorește calcularea produsului

$$A_1 A_2 \dots A_n \quad (6.1)$$

Expresia (3.1) se poate evalua folosind algoritmul standard de înmulțire a unei perechi de matrice ca subrutină, o dată ce expresia a fost parantezată (parantezarea este necesară pentru a exclude toate ambiguitățile privind înmulțirile de matrice). Un produs de matrice este *complet parantezat* fie dacă este format dintr-o unică matrice, fie dacă este produsul a două produse de matrice care sunt la rândul lor complet parantezate. Cum înmulțirea matricelor este asociativă, toate parantezările conduc la același produs. De exemplu, dacă șirul de matrice este (A_1, A_2, A_3, A_4) , produsul $A_1 A_2 A_3 A_4$ poate fi complet parantezat în cinci moduri distincte, astfel avem:

$$(A_1(A_2(A_3A_4))), (A_1((A_2A_3)A_4)), ((A_1A_2)(A_3A_4)), ((A_1(A_2A_3))A_4), (((A_1A_2)A_3)A_4).$$

Să considerăm mai întâi costul înmulțirii a două matrice. Algoritmul standard este dat de următoarea procedură, descrisă în pseudocod. Prin *linii* și *coloane* sunt referite numărul de linii, respectiv de coloane ale matricei.

ÎNMULȚIRE-MATRICE(A, B)

```

1. if coloane[A] ≠ linii [B]
2.   then mesaj de eroare: „dimensiuni incompatibile”
3. else
4.   for i ← 1 to linii[A]
5.     do for j ← 1 to coloane[B]
6.       do C[i, j] ← 0
7.         for k ← 1 to coloane[A]
8.           do C[i, j] ← C[i, j] + A[i, k] * B[k, j]
9.   return C

```

Două matrice A și B se pot înmulți numai dacă numărul de coloane din A este egal cu numărul de linii din B . Astfel, dacă A este o matrice având dimensiunea $p \times q$ și B este o matrice având dimensiunea $q \times r$, matricea produsă C este o matrice având dimensiunea $p \times r$. Timpul necesar calculului matricei C este determinat de numărul de înmulțiri scalare (a se vedea linia 8 din algoritm) care este pqr . În cele ce urmează, vom exprima timpul de execuție în funcție de numărul de înmulțiri scalare.

Pentru a ilustra modul în care apar costuri diferite la parantezări diferite ale produsului de matrice considerăm problema șirului (A_1, A_2, A_3) de trei matrice. Presupunem că dimensiunile matricelor sunt 10×100 , 100×5 și, respectiv, 5×50 . Dacă efectuăm înmulțire conform parantezării $((A_1A_2)A_3)$, atunci vom avea $10 \cdot 100 \cdot 5 = 5000$ înmulțiri scalare pentru a calcula matricea A_1A_2 de dimensiune 10×5 plus alte $10 \cdot 5 \cdot 50 = 2500$ înmulțiri scalare pentru a înmulți această matrice cu matricea A_3 . Astfel rezultă un total de 7500 înmulțiri scalare. Dacă însă vom efectua înmulțirile conform parantezării $(A_1(A_2A_3))$, vom avea $100 \cdot 5 \cdot 50 = 25.000$ înmulțiri scalare pentru a calcula matricea A_2A_3 de dimensiuni 100×50 plus alte $10 \cdot 100 \cdot 50 = 50.000$ înmulțiri scalare pentru a înmulți A_1 cu această matrice. Astfel rezultă un total de 75.000 înmulțiri scalare. Deci, calculând produsul conform primei parantezări, calculul este de 10 ori mai rapid.

Problema înmulțirii șirului de matrice poate fi enunțată în următorul mod: dându-se un șir (A_1, A_2, \dots, A_n) de n matrice, unde pentru $i = 1, 2, \dots, n$, matricea A_i are dimensiunile $p_{i-1} \times p_i$ să se parantezeze complet produsul $A_1A_2 \dots A_n$, astfel încât să se minimizeze numărul de înmulțiri scalare.

6.1.1. Evaluarea numărului de parantezări

Înainte de a rezolva problema înmulțirii șirului de matrice prin programare dinamică, trebuie să ne convingem că verificarea exhaustivă a tuturor parantezărilor nu conduce la un algoritm eficient. Fie $P(n)$ numărul de parantezări distincte ale unei secvențe de n matrice. Deoarece secvența de matrice o putem diviza între matricele k și $(k+1)$ pentru orice $k = 1, 2, \dots, n-1$ și apoi putem descompune în paranteze, în mod independent, fiecare dintre cele două subsecvențe, astfel obținem recurența:

$$P(n) = \begin{cases} 1 & \text{dacă } n = 1, \\ \sum_{k=1}^{n-1} P(k)P(n-k) & \text{dacă } n \geq 2. \end{cases}$$

Soluția acestei recurențe este secvența de **numere Catalan**:

$$P(n) = C(n-1),$$

unde

$$C(n) = \frac{1}{n+1} \binom{2n}{n} = \Omega(4^n/n^{3/2})$$

Numărul de soluții este exponențial în n și metoda forței brute a căutării exhaustive este, prin urmare, o strategie slabă de determinare a modalității de parantezare a șirului de matrice.

6.1.2. Structura unei parantezări optime

Primul pas din schema generală a metodei programării dinamice constă în caracterizarea structurii unei soluții optime. Pentru problema înmulțirii șirului de matrice, aceasta este descrisă în continuare. Pentru simplitate, vom adopta convenția de notare $A_{i..j}$ pentru matricea care rezultă în urma evaluării produsului $A_iA_{i+1} \dots A_j$. O parantezare optimă a produsului $A_1A_2 \dots A_n$ împarte produsul între A_k și A_{k+1} pentru un anumit întreg k din intervalul $1 \leq k < n$. Aceasta înseamnă că, pentru o valoare a lui k , mai întâi calculăm matricele $A_{1..k}$ și $A_{k+1..n}$ și apoi, le înmulțim pentru a produce rezultatul final $A_{1..n}$. Costul acestei parantezări optime este costul calculului matricei $A_{1..k}$ plus costul calculului matricei $A_{k+1..n}$, plus costul înmulțirii celor două matrice.

Observația cheie este că parantezarea subșirului "prefix" $A_1A_2\cdots A_k$, în cadrul parantezării optime a produsului $A_1A_2\cdots A_n$, trebuie să fie o parantezare *optimă* pentru $A_1A_2\cdots A_k$ deoarece dacă ar fi existat o modalitate mai puțin costisitoare de parantezare a lui $A_1A_2\cdots A_k$, înlocuirea respectivei parantezări în parantezarea lui $A_1A_2\cdots A_n$ ar produce o altă parantezare pentru $A_1A_2\cdots A_n$ al cărei cost ar fi mai mic decât costul optimal, ceea ce este o contradicție. O observație asemănătoare este valabilă și pentru parantezarea subșirului $A_{k+1}A_{k+2}\cdots A_n$ în cadrul parantezării optime a lui $A_1A_2\cdots A_n$: aceasta trebuie să fie o parantezare optimă pentru $A_{k+1}A_{k+2}\cdots A_n$.

Prin urmare, o soluție optimă a unei instanțe a problemei înmulțirii șirului de matrice conține soluții optime pentru instanțe ale subproblemelor. Existența substructurilor optime în cadrul unei soluții optime este una dintre caracteristicile cadrului de aplicare a metodei programării dinamice.

6.1.3. O soluție recursivă

Al doilea pas în aplicarea metodei programării dinamice este definirea valorii unei soluții optime în mod recursiv, în funcție de soluțiile optime ale subproblemelor. În cazul problemei înmulțim șirului de matrice, o subproblemă constă în determinarea costului minim al unei parantezări a șirului $A_iA_{i+1}\cdots A_j$, pentru $1 \leq i \leq j \leq n$. Fie $m[i, j]$ numărul minim de înmulțiri scalare necesare pentru a calcula matricea $A_{i..j}$; costul modalității optime de calcul a lui $A_{1..n}$ va fi, atunci, $m[1, n]$.

Putem defini $m[i, j]$ în mod recursiv, după cum urmează. Dacă $i = j$ șirul constă dintr-o singură matrice $A_{i..i} = A_i$ și, pentru calculul produsului nu este necesară nici o înmulțire scalară. Atunci, $m[i, i] = 0$ pentru $i = 1, 2, \dots, n$. Pentru a calcula $m[i, j]$ când $i < j$, vom folosi structura unei soluții optimale găsite la pasul 1. Să presupunem că descompunerea optimă în paranteze împarte produsul $A_iA_{i+1}\cdots A_j$ între A_k și A_{k+1} , cu $i \leq k < j$. Atunci $m[i, j]$ este egal cu costul minim pentru calculul subprodusului $A_{i..k}$ și $A_{k+1..j}$, plus costul înmulțirii acestor două matrice rezultat. Deoarece evaluarea produsului $A_{i..k}A_{k+1..j}$ necesită $p_{i-1}p_kp_j$ înmulțiri scalare, vom obține

$$m[i, j] = m[i, k] + m[k + 1, j] + p_{i-1}p_kp_j.$$

Aceasta ecuație recursivă presupune cunoașterea valorii lui k , lucru care nu este posibil. Există doar $j - i$ valori posibile pentru k , și anume $k = i, i + 1, \dots, j - 1$. Deoarece parantezarea optimă trebuie să folosească una dintre aceste valori pentru k , va trebui să le verificăm pe toate pentru a o putea găsi pe cea mai bună. Atunci, definiția recursivă a costului minim a parantezării produsului $A_iA_{i+1}\cdots A_j$, devine

$$m(i, j) = \begin{cases} 0 & i = j \\ \min\{m[i, k] + m[k + 1, j] + p_{i-1}p_kp_j\} & i < j \end{cases} \quad (6.2)$$

Valorile $m[i, j]$ exprimă costul soluțiilor optime ale subproblemelor. Pentru a putea urmări modul de construcție a soluției optime, să definim $s[i, j]$ care va conține valoarea k pentru care împărțirea produsului $A_iA_{i+1}\cdots A_j$ produce o parantezare optimă. Aceasta înseamnă că $s[i, j]$ este egal cu valoarea k pentru care $m[i, j] = m[i, k] + m[k + 1, j] + p_{i-1}p_kp_j$.

6.1.4. Calculul costului optimal

În acest moment, este ușor să scriem un algoritm recursiv bazat pe recurența (6.2) care calculează costul minim $m[1, n]$ al produsului $A_1A_2\cdots A_n$. După cum vom vedea mai târziu, acest algoritm necesită un timp exponențial –nu mai bun decât metoda forței brute folosită pentru verificarea fiecărui mod de descompunere în paranteze a produsului de matrice.

Observația importantă care se impune în acest moment este că avem relativ puține subprobleme: o problemă pentru fiecare alegere a lui i și j ce satisfac $1 \leq i \leq j \leq n$, adică un total de $\binom{n}{2} + n = \Theta(n^2)$. Un algoritm recursiv poate întâlni fiecare subproblemă de mai multe ori pe ramuri diferite ale arborelui său de recurență. Această proprietate de suprapunere a subproblemelor este a doua caracteristică a programării dinamice.

În loc să calculăm recursiv soluția recurenței (4.2), vom aplica cel de-al treilea pas al schemei generale a metodei programării dinamice și vom calcula costul optimal cu o abordare "bottom-up". Algoritmul următor presupune că matricea A , are dimensiunile $p_{i-1} \times p_i$, pentru $i = 1, 2, \dots, n$. Intrarea este secvența (p_0, p_1, \dots, p_n) , unde $lungime[p] = n + 1$. Procedura folosește un tablou auxiliar $m[1..n, 1..n]$ pentru costurile $m[i, j]$ și un tablou auxiliar $s[1..n, 1..n]$ care înregistrează acea valoare a lui k pentru care s-a obținut costul optim în calculul lui $m[i, j]$.

ORDINE –ȘIR –MATRICE(p)

```

1: n ← lungime[p] - 1
2: for i ← 1 to n do
3:   m[i, i] ← 0
4: for l ← 2 to n do
5:   for i ← 1 to n - l + 1 do
6:     j ← i + l - 1

```

```

7:    $m[i,j] \leftarrow \infty$ 
8:   for  $k \leftarrow i$  to  $j-1$  do
9:      $q \leftarrow m[i, k] + m[k+1, j] + p_i p_k p_j$ 
10:    if  $q < m[i, j]$  then
11:       $m[i, j] \leftarrow q$ 
12:       $s[i, j] \leftarrow k$ 
13: return  $m, s$ 

```

Algoritmul completează tabloul m într-un mod ce corespunde rezolvării problemei parantezării unor șiruri de matrice de lungime din ce în ce mai mare. Ecuația (6.2) arată că $m[i,j]$, costul de calcul al produsului șirului de $j-i+1$ matrice, depinde doar de costurile calculării produselor șirurilor de mai puțin de $j-i+1$ matrice. Aceasta înseamnă că, pentru $k = i, i+1, \dots, j-1$, matricea $A_{i..k}$ este un produs de $k-i+1 < j-i+1$ matrice, iar matricea $A_{k+1..j}$ este un produs de $j-k < j-i+1$ matrice.

Algoritmul inițializează, mai întâi, $m[i, i] \leftarrow 0$ pentru $i = 1, 2, \dots, n$ (costul minim al șirurilor de lungime 1) în liniile 2-3. Se folosește, apoi, recurența (6.2) pentru a calcula $m[i, i+1]$ pentru $i = 1, 2, \dots, n-1$ (costul minim al șirurilor de lungime 2) la prima execuție a ciclului din liniile 4-12. La a doua trecere prin ciclu, se calculează $m[i, i+2]$ pentru $i = 1, 2, \dots, n-2$ (costul minim al șirurilor de lungime 3) etc. La fiecare pas, costul $m[i, j]$ calculat pe liniile 9-12 depinde doar de intrările $m[i, k]$ și $m[k+1, j]$ ale tabloului, deja calculate.

Figura 6.1 ilustrează această procedură pe un șir de $n = 6$ matrice. Deoarece am definit $m[i, j]$ doar pentru $i \leq j$, din tabloul m folosim doar porțiunea situată strict deasupra diagonalei principale. Figura arată tablourile rotite astfel încât diagonala principală devine orizontală, iar șirul de matrice este înscris de-a lungul liniei de jos. Folosind acest model, costul minim $m[i, j]$ de înmulțire a subșirului $A_i A_{i+1} \dots A_j$, de matrice poate fi găsit la intersecția liniilor care pornesc

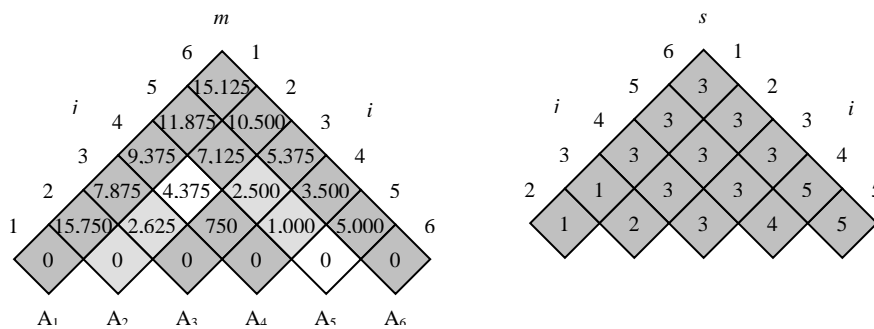


Figura 6.1 Tabelele m și s calculate de ORDINE-ȘIR-MATRICE pentru $n = 6$ și matricele cu dimensiunile date mai jos:

matrice	dimensiuni
A_1	30x35
A_2	35x15
A_3	15x5
A_4	5x10
A_5	10x20
A_6	20x25

Tabelele sunt rotite astfel încât diagonala principală devine orizontală. În tabloul m sunt utilizate doar elementele de deasupra diagonalei principale, iar în tabloul s sunt utilizate doar elementele de sub diagonala principală. Numărul minim de înmulțiri scalare necesare pentru înmulțirea a șase matrice este $m[1,6] = 15.125$. Dintre pătratele hașurate cu gri deschis, perechile care sunt colorate cu aceeași nuanță sunt considerate împreună în linia 9, când calculăm pe direcția nord-est din A_i și pe direcția nord-vest din A_j . Fiecare rând orizontal din tablou conține valorile pentru șiruri de matrice de aceeași lungime. Procedura ORDINE-ȘIR-MATRICE calculează rândurile de jos în sus și de la stânga la dreapta în cadrul fiecărui rând. O valoare $m[i, j]$ este calculată pe baza folosirii produselor $p_i \dots p_k p_j$ pentru $k=i, i+1, \dots, j-1$ și a tuturor valorilor aflate pe direcția sud-est și sud-vest față de $m[i, j]$.

$$m[2,5] = \min \left\{ \begin{array}{l} m[2,2] + m[3,5] + p_1 p_2 p_5 = 0 + 2500 + 35 \cdot 15 \cdot 20 = 13000 \\ m[2,3] + m[4,5] + p_1 p_3 p_5 = 2625 + 1000 + 35 \cdot 5 \cdot 20 = 7125 \\ m[2,4] + m[5,5] + p_1 p_4 p_5 = 4375 + 0 + 35 \cdot 10 \cdot 20 = 11375 \end{array} \right\} = 7125$$

O simplă examinare a structurii de ciclu imbricat din ORDINE-ȘIR-MATRICE conduce la constatarea că timpul de execuție a algoritmului este $O(n^3)$. În fiecare dintre cele trei cicluri imbricate, indicii acestora (l, i și k) iau cel mult n valori.

6.1.5. Construirea unei soluții optime

Deși algoritmul ORDINE-ȘIR-MATRICE determină numărul optim de înmulțiri se necesare pentru calculul produsului șirului de matrice, acesta nu prezintă în mod direct, trebuie făcută înmulțirea. Pasul 4 al schemei generale a metodei programării dinamice urmează construirea unei soluții optime din informația disponibilă.

În acest caz particular, vom folosi tabloul $s[l..n, l..n]$ pentru a determina modul optim de înmulțire a matricelor. Fiecare element $s[i, j]$ conține valoarea lui k pentru care parantezarea optimă a produsului $A_i A_{i+1} \dots A_j$ împarte produsul între A_k și A_{k+1} . Atunci știm că, în produsul final de calcul al matricei $A_{l..n}$, optimul este $A_{l..s[l, n]} A_{s[l, n]+1..n}$. Înmulțirile anterioare determinate recursiv, deoarece $s[l, s[l, n]]$ determină ultima înmulțire matriceală din calculul lui $A_{l..s[l, n]}$ și $s[s[l, n]+1, n]$ determină ultima înmulțire matriceală din calculul lui $A_{s[l, n]+1..n}$. Următoarea procedură recursivă calculează produsul șirului de matrice $A_{i..j}$, fiind date matricea $A = (A_1, A_2, \dots, A_n)$, tabloul s calculat de ORDINE-ȘIR-MATRICE și indicii i și j . Apelul este ÎNMULȚIRE-ȘIR-MATRICE (A, s, i, n).

ÎNMULȚIRE-ȘIR-MATRICE (A, s, i, n)

```

1: if  $j > i$  then
2:    $X \leftarrow$  ÎNMULȚIRE-ȘIR-MATRICE ( $A, s, i, s[i, j]$ )
3:    $Y \leftarrow$  ÎNMULȚIRE-ȘIR-MATRICE ( $A, s, s[i, j]+1, j$ )
4:   return ÎNMULȚIRE-ȘIR-MATRICE ( $X, Y$ )
5: else
6:   return  $A_i$ 

```

În exemplul din figura 6.1, apelul ÎNMULȚIRE-ȘIR-MATRICE ($A, s, 1, 6$) calculează produsul șirului de matrice conform descompunerii în paranteze $((A_1(A_2A_3))(A_4A_5)A_6)$.

6.2. Elemente de programare dinamică

Cu toate că am prezentat un exemplu de aplicare a metodei programării dinamice, este posibil să existe, în continuare, întrebări asupra situațiilor în care ea este aplicabilă, sau, altfel spus, când trebuie căutată o soluție de programare dinamică pentru o anumită problemă? În această secțiune se vor examina două componente de bază pe care trebuie să le prezinte problema de optimizare pentru ca programarea dinamică să fie aplicabilă: substructura optimă și subproblemele suprapuse. Se va urmări, însă, și o variantă a metodei, numită memoizare, în care se utilizează proprietatea de suprapunere a subproblemelor.

Substructura optimă

Primul pas al rezolvării unei probleme de optimizare prin programare dinamică constă în caracterizarea structurii unei soluții optime. Spunem că problema prezintă o **substructură optimă** dacă o soluție optimă a problemei include soluții optime ale subproblemelor. Ori de câte ori o problemă prezintă o structură optimă, acesta este un bun indiciu pentru posibilitatea aplicării programării dinamice.

În secțiunea 6.1 am văzut că problema înmulțirii șirului de matrice prezintă o substructură optimă. Am observat că o parantezare optimă pentru $A_1 A_2 \dots A_n$, care separă produsul între A_k și A_{k+1} , include soluții optime pentru problemele parantezării matricelor $A_1 A_2 \dots A_k$ și $A_{k+1} A_{k+2} \dots A_n$. Tehnica utilizată pentru a arăta că subproblemele au soluții optime este clasică. Presupunem că există o soluție mai bună pentru o subproblemă și arătăm că această presupunere contrazice optimalitatea soluției problemei inițiale.

Substructura optimă a unei probleme sugerează, deseori, un spațiu potrivit al subproblemelor pentru care se poate aplica programarea dinamică. De obicei, există mai multe clase de subprobleme care pot fi considerate "naturale" pentru o problemă. De exemplu, spațiul subproblemelor pe care le-am considerat la înmulțirea șirului de matrice conținea toate subsecvențele șirului de intrare. La fel de bine, am fi putut alege ca spațiu al subproblemelor secvențe arbitrare de matrice din șirul de intrare, dar acest spațiu al subproblemelor ar fi fost prea mare. Un algoritm de programare dinamică bazat pe acest spațiu al subproblemelor ar rezolva mult mai multe probleme decât ar fi necesar.

Investigarea substructurii optime a unei probleme prin iterarea instanțelor subproblemelor acestea este o metodă bună pentru a găsi un spațiu potrivit de subprobleme pentru programarea dinamică. De exemplu, după examinarea structurii unei soluții optime a problemei înmulțirii șirului de matrice, putem itera și examina structura soluțiilor optime ale subproblemelor, a subsubproblemelor și așa mai departe. Descoperim, astfel, că toate subproblemele sunt compuse din subsecvențe ale șirului ($A_1 A_2 \dots A_n$). Deci, mulțimea șirurilor de forma $(A_i A_{i+1} \dots A_j)$ cu $1 \leq i \leq j \leq n$ un spațiu de subprobleme natural și rezonabil de utilizat.

6.2.1. Suprapunerea subproblemelor

O a doua componentă pe care trebuie să o aibă o problemă de optimizare pentru ca programarea dinamică să fie aplicabilă este ca spațiul subproblemelor să fie "restrâns", în sensul că un algoritm recursiv rezolvă mereu aceleași subprobleme, în loc să genereze subprobleme noi.

De obicei, numărul total de subprobleme distincte este dependent polinomial de dimensiunea datelor de intrare. Când un algoritm recursiv abordează mereu o aceeași problemă, se spune că problema de optimizare are *subprobleme suprapuse*. Spre deosebire de aceasta, o problemă care se rezolvă cu un algoritm de tip divide și stăpânește generează, la fiecare nouă etapă, probleme noi. În general, algoritmii de programare dinamică folosesc suprapunerea subproblemelor prin rezolvarea o singură dată a fiecărei subprobleme, urmată de stocarea soluției într-un tablou unde poate fi regăsită la nevoie, necesitând pentru regăsire un timp constant.

Pentru a ilustra proprietatea de subprobleme suprapuse, vom reexamina problema înmulțirii șirului de matrice. Revenind la figura 4.1. observăm că algoritmul ORDINE-ȘIR-MATRICE preia, în mod repetat, soluțiile subproblemelor din liniile de jos atunci când rezolvă subproblemele din liniile superioare. De exemplu, valoarea $m[3,4]$ este referită de patru ori: în timpul calculului pentru $m[2,4]$, $m[1,4]$, $m[3,5]$ și $m[3,6]$. Dacă $m[3,4]$ ar fi recalculată de fiecare dată în loc să fie preluată din tablou, creșterea timpului de execuție ar fi dramatică.

6.3. Cel mai lung subșir comun

Următoarea problemă pe care o vom lua în considerare este problema celui mai lung subșir comun. Un subșir al unui șir dat este șirul inițial din care lipsesc unele elemente (eventual nici unul). În mod formal, dându-se șirul $X = (x_1, x_2, \dots, x_m)$, un alt șir $Z = (z_1, z_2, \dots, z_k)$ este un *subșir* al lui X dacă există un șir strict crescător de indici din X , (i_1, i_2, \dots, i_k) , astfel încât, pentru toate valorile $j = 1, 2, \dots, k$, avem $x_{i_j} = z_j$. De exemplu, $Z = (B, C, D, B)$ este un subșir al lui $X = (A, B, C, B, D, A, B)$, secvența de indici corespunzătoare fiind $(2, 3, 5, 7)$.

Dându-se două șiruri X și Y , spunem că șirul Z este un *subșir comun* pentru X și Y dacă Z este un subșir atât pentru X cât și pentru Y . De exemplu, dacă $X = (A, B, C, B, D, A, B)$ și $Y = (B, D, C, A, B, A)$, șirul (B, C, A) este un subșir comun pentru X și Y . Șirul (B, C, A) nu este *cel mai lung* subșir comun (CMLSC) pentru X și Y , deoarece are lungimea 3, iar șirul (B, C, B, A) , care este de asemenea subșir comun lui X și Y , are lungimea 4. Șirul (B, C, B, A) este un CMLSC pentru X și Y , la fel și șirul (B, D, A, B) , deoarece nu există un subșir comun de lungime cel puțin 5.

În problema *celui mai lung subșir comun*, se dau două șiruri $X = (x_1, x_2, \dots, x_m)$ și $Y = (y_1, y_2, \dots, y_n)$ și se dorește determinarea unui subșir comun de lungime maximă pentru X și Y . Această secțiune va arăta că problema CMLSC poate fi rezolvată eficient prin metoda programării dinamice.

6.3.1. Caracterizarea celui mai lung subșir comun

O abordare prin metoda forței brute pentru rezolvarea problemei CMLSC constă în enumerarea tuturor subșirurilor lui X și verificarea dacă acestea constituie un subșir și pentru Y , memorând cel mai lung subșir găsit. Fiecare subșir al lui X corespunde unei submulțimi a indicilor lui X $(1, 2, \dots, m)$. Există 2^m subșiruri pentru X și, deci, această abordare necesită un timp exponențial, ceea ce o face inabordabilă pentru șiruri de lungimi mari.

Problema CMLSC are proprietatea de substructură optimă, așa cum vom arăta în următoarea teoremă. Clasa naturală de subprobleme corespunde unor perechi de "prefixe" ale celor două șiruri de intrare. Mai precis, dându-se un șir $X = (x_1, x_2, \dots, x_m)$, definim *prefixul* i al lui X , pentru $i = 0, 1, \dots, m$, ca fiind $X_i = (x_1, x_2, \dots, x_i)$. De exemplu, dacă $X = (A, B, C, B, D, A, B)$, atunci $X_4 = (A, B, C, B)$ iar X_0 este șirul vid.

Teorema 6.1 (Substructura optimă a unui CMLSC) Fie șirurile $X = (x_1, x_2, \dots, x_m)$ și $Y = (y_1, y_2, \dots, y_n)$ și fie $Z = (z_1, z_2, \dots, z_k)$ un CMLSC pentru X și Y .

1. Dacă $x_m \neq y_n$, atunci $z_k = x_m = y_n$ și Z_{k-1} este un CMLSC pentru X_{m-1} și Y_{n-1} .
2. Dacă $x_m \neq y_n$, atunci, din $z_k \neq x_m$ rezultă că Z este un CMLSC pentru X_{m-1} și Y .
3. Dacă $x_m \neq y_n$, atunci, din $z_k \neq y_n$ rezultă că Z este un CMLSC pentru X și Y_{n-1} .

Demonstrație:

(1) Dacă $z_k \neq x_m$ atunci putem adăuga $x_m = y_n$ la Z pentru a obține un subșir comun a lui X și Y de lungime $k + 1$, contrazicând presupunerea că Z este *cel mai lung* subșir comun pentru X și Y . Deci, va trebui să avem $z_k = x_m = y_n$. Prefixul Z_{k-1} este un subșir comun de lungime $k - 1$ pentru X_{m-1} și Y_{n-1} . Dorim să demonstrăm că este un CMLSC. Să presupunem, prin absurd, că există un subșir W comun pentru X_{m-1} și Y_{n-1} , a cărui lungime este fie mai mare decât $k - 1$. Atunci, adăugând la W elementul $x_m = y_n$ se va forma un subșir comun a lui X și Y a cărui lungime este mai mare decât k , ceea ce este o contradicție.

(2) Dacă $z_k \neq x_m$, atunci Z este un subșir comun pentru X_{m-1} și Y . Dacă ar exista un subșir comun W al lui X_{m-1} și Y , cu lungime mai mare decât k , atunci W ar fi și un subșir comun al lui X_m și Y , contrazicând presupunerea că Z este un CMLSC pentru X și Y .

(3) Demonstrația este simetrică celei de la (2).

Caracterizarea din teorema 6.1 arată că un CMLSC al două șiruri conține un CMLSC pentru prefixele celor două șiruri. Atunci, problema CMLSC are proprietatea de substructură optimală. O soluție recursivă are, de asemenea, proprietatea de suprapunere a problemelor, după cum vom arăta în cele ce urmează.

6.3.2. O soluție recursivă a subproblemelor

Teorema 6.1 implică faptul că există fie una, fie două probleme care trebuie examinate pentru găsirea unui CMLSC pentru $X = (x_1, x_2, \dots, x_m)$ și $Y = (y_1, y_2, \dots, y_n)$. Dacă $x_m = y_n$, trebuie găsit un CMLSC pentru X_{m-1} și Y_{n-1} . Adăugarea elementului $x_m = y_n$ la acest CMLSC produce un CMLSC pentru X și Y . Dacă $x_m \neq y_n$, atunci trebuie rezolvate două subprobleme: găsirea unui CMLSC pentru X_{m-1} și Y și găsirea unui CMLSC pentru X și Y_{n-1} . Cel mai lung CMLSC dintre acestea două va fi CMLSC pentru X și Y .

Se poate observa că problema CMLSC se descompune în subprobleme suprapuse. Pentru a găsi un CMLSC al șirurilor X și Y , va trebui, probabil, calculat CMLSC pentru X și Y_{n-1} respectiv, pentru X_{m-1} și Y . Dar fiecare dintre aceste subprobleme conține sub-subproblema găsirii CMLSC pentru X_{m-1} și Y_{n-1} . Multe alte subprobleme au în comun sub-subprobleme.

Ca și problema înmulțirii șirului de matrice, soluția recursivă pentru problema CMLSC implică stabilirea unei recurențe pentru costul unei soluții optime. Să definim $c[i, j]$ ca lungimea unui CMLSC al șirurilor X_i , și Y_j . Dacă $i=0$ sau $j=0$, CMLSC are lungimea 0. Substructura optimală a problemei CMLSC produce formula recursivă

$$c[i, j] = \begin{cases} 0, & i = 0, j = 0, \\ c[i-1, j-1] + 1, & i, j > 0 \text{ and } x_i = y_j, \\ \max(c[i-1, j], c[i, j-1]), & i, j > 0 \text{ and } x_i \neq y_j. \end{cases} \quad (6.2)$$

6.3.3. Calculul lungimii unui CMLSC

Pe baza ecuației (4.5) se poate scrie un algoritm recursiv având timp exponențial pentru calculul lungimii unui CMLSC a două șiruri. Deoarece există numai $\Theta(mn)$ subprobleme distincte, vom folosi metoda programării dinamice pentru a calcula soluțiile în mod "bottom-up".

LUNGIME-CMLSC(A', R)

1: $m \leftarrow \text{lungime}[X]$

2: $n \leftarrow \text{lungime}[Y]$

3: pentru $i \leftarrow 1, m$ execută

4: $c[i, 0] \leftarrow 0$

5: pentru $j \leftarrow 0, n$ execută

6: $c[0, j] \leftarrow 0$

7: pentru $i \leftarrow 1, m$ execută

8: pentru $j \leftarrow 1, n$ execută

9: dacă $x_i = y_j$ atunci

10: $c[i, j] \leftarrow c[i-1, j-1] + 1$

11: $b[i, j] \leftarrow \text{„}\nwarrow\text{”}$

12: altfel

13: dacă $c[i-1, j] \geq c[i, j-1]$ atunci

14: $c[i, j] \leftarrow c[i-1, j]$

15: $b[i, j] \leftarrow \text{„}\uparrow\text{”}$

16: altfel

17: $c[i, j] \leftarrow c[i, j-1]$

18: $b[i, j] \leftarrow \text{„}\leftarrow\text{”}$

19: returnează c, b

Procedura LUNGIME-CMLSC are ca date de intrare două șiruri $X = (x_1, x_2, \dots, x_m)$ și $Y = (y_1, y_2, \dots, y_n)$. Procedura memorează valorile $c[i, j]$ într-un tablou $c[0..m, 0..n]$ ale cărui elemente sunt calculate în ordinea crescătoare a liniilor. (Aceasta înseamnă că se completează mai întâi prima linie de la stânga la dreapta, apoi a doua linie și așa mai departe.) De asemenea, se construiește un tablou $b[1..m, 1..n]$ care simplifică determinarea unei soluții optime. Intuitiv, $b[i, j]$ indică elementul tabloului care corespunde soluției optime a subproblemei aalese la calculul lui $c[i, j]$. Procedura returnează tablourile b și c ; $c[m, n]$ conține lungimea unui CMLSC pentru X și Y .

Figura 16.3 conține tabloul produs de LUNGIME-CMLSC pentru șirurile $X = (A, B, C, B, D, A, B)$ și $Y = (B, D, C, A, B, A)$. Timpul de execuție al procedurii este $O(mn)$, deoarece calculul fiecărui element al tabloului necesită un timp $O(1)$.

Construirea unui CMLSC

Tabelul b returnat de LUNGIME-CMLSC poate fi folosit pentru construcția rapidă a unui CMLSC pentru $X = (x_1, x_2, \dots, x_m)$ și $Y = (y_1, y_2, \dots, y_n)$. Pur și simplu, se începe cu $b[m, n]$ și se parcurge tabloul conform direcțiilor indicate. Ori de câte ori se întâlnește un element " \wedge " pe poziția $b[i, j]$, aceasta înseamnă că $x_i = y_j$ este un element al CMLSC. Prin această metodă, elementele CMLSC sunt regăsite în ordine inversă. Următoarea procedură recursivă tipărește un CMLSC al șirurilor X și Y în ordinea corectă. Apelul inițial este SCRIE-CMLSC($b, X, lungime[X], lungime[Y]$).

SCRIE-CMLSC (6,A,i,J)

- 1: **dacă** $i = 0$ sau $j = 0$ **atunci**
- 2: revenire
- 3: **dacă** $b[i, j] = "\wedge"$ **atunci**
- 4: SCRIE-CMLSC ($b, X, i-1, j-1$)
- 5: tipărește x_i
- 6: **altfel dacă** $b[i, j] = "\uparrow"$ **atunci**
- 7: SCRIE-CMLSC ($b, X, i-1, j$)
- 8: altfel
- 9: SCRIE-CMLSC ($b, X, i, j-1$)

Pentru tabloul b din figura 16.3, această procedură tipărește "BCBA". Procedura necesită un timp de $O(m + n)$, deoarece, cel puțin unul dintre i și j este decrementat în fiecare etapa a recurenței.

Îmbunătățirea codului

O dată ce s-a construit un algoritm, se descoperă că, deseori, timpul de execuție sau spațiul folosit pot fi îmbunătățite. Acest lucru este adevărat, mai ales, pentru programarea dinamică clasică. Anumite modificări pot simplifica implementarea și pot îmbunătăți factorii constanți, dar nu pot aduce o îmbunătățire asimptotică a performanței. Alte modificări pot aduce îmbunătățiri asimptotice substanțiale atât pentru timpul de execuție cât și pentru spațiul ocupat. De exemplu, se poate elimina complet tabloul b . Fiecare element $c[i, j]$ depinde de numai trei alte elemente ale tabloului c : $c[i-1, j-1]$, $c[i-1, j]$ și $c[i, j-1]$. Dându-se valoarea lui $c[i, j]$, se poate determina, în timp $O(1)$, care dintre aceste trei valori a fost folosită pentru calculul lui $c[i, j]$ fără a mai folosi tabloul b .

j	0	1	2	3	4	5	6
y_j		B	D	C	A	B	A
x_i			0	0	0	0	0
A	0	\uparrow 0	\uparrow 0	\uparrow 0	\wedge 1	\leftarrow 1	\wedge 1
B	0	\wedge 1	\leftarrow 1	\leftarrow 1	\uparrow 1	\wedge 2	\leftarrow 2
C	0	\uparrow 1	\uparrow 1	\wedge 2	\leftarrow 2	\uparrow 2	\uparrow 2
B	0	\wedge 1	\uparrow 1	\uparrow 2	\uparrow 2	\wedge 3	\leftarrow 3
D	0	\uparrow 1	\wedge 2	\uparrow 2	\uparrow 2	\uparrow 3	\uparrow 3
A	0	\uparrow 1	\uparrow 2	\uparrow 2	\wedge 3	\uparrow 3	\wedge 4
B	0	\wedge 1	\uparrow 2	\uparrow 2	\uparrow 3	\wedge 4	\uparrow 4

Figura 4.3 Tabelele c și b calculate de LUNGIME-CMLSC pentru șirurile $X = (A, B, C, B, D, A, B)$

și $Y = (B, D, C, A, B, A)$. Pătratul din linia i și coloana j conține valoarea $c[i, j]$ precum și săgeata potrivită pentru valoarea lui $b[i, j]$. Valoarea 4 a lui $c[7, 6]$ - colțul cel mai din dreapta jos în tablou-este lungimea unui CMLSC (B,C,B, A) al lui X și Y . Pentru $i, j > 0$, valoarea $c[i, j]$ depinde doar de valoarea expresiei $x_i = y_j$ și de valorile elementelor $c[i-1, j]$, $c[i, j-1]$ și $c[i-1, j-1]$, care sunt calculate înainte de $c[i, j]$. Pentru a reconstrui elementele unui CMLSC, se urmăresc săgețile corespunzătoare lui $b[i, j]$ începând cu colțul cel mai din dreapta; drumul este hașurat. Fiecare " \wedge " pe acest drum corespunde unui element (care este evidențiat) pentru care $x_i = y_j$ aparține unui CMLSC.

Ca o consecință, putem reconstrui un CMLSC în timp $O(m + n)$ folosind o procedură similară cu SCRIE-CMLSC. Deși prin această metodă se eliberează un spațiu $\Theta(mn)$, spațiul auxiliar necesar calculului unui CMLSC nu descrește asimptotic, deoarece este oricum necesar un spațiu $\Theta(mn)$ pentru tabloul c .

Totuși, spațiul necesar pentru LUNGIME-CMLSC îl putem reduce deoarece, la un moment dat, sunt necesare doar două linii ale tabloului c : linia care este calculată și linia precedentă. (De fapt, putem folosi doar puțin mai mult spațiu decât cel necesar unei linii a lui c , pentru a calcula lungimea unui CMLSC) Această îmbunătățire

funcționează doar dacă este necesară numai lungimea unui CMLSC; dacă este nevoie și de reconstruirea elementelor unui CMLSC, tabloul mai mic nu conține destulă informație pentru a reface etapele în timp $O(m+n)$.