

Codare Shannon Fano sau Huffman Static (2019-2020)

A. Tema

Sa se implementeze coderul si decodor ul pentru unul din algoritmi considerați: Shannon Fano sau Huffman static.

B. Aspecte teoretice

Ambii algoritmi considerați intră în categoria celor **statistici**. Implementarea cerută aici este cea consacrată adică cea cu **model semistatic** (notația Huffman static este incorectă dar consacrată, de aceea o vom folosi totuși în continuare).

Codorul parcurge următoarele **etape**:

1. Determinare a statisticii
2. Construire model (coduri)
3. Codare pe baza modelului

Decodorul parcurge următoarele etape:

1. Preluare model
2. Decodare pe baza modelului

Pentru a simplifica transmiterea modelului spre decodor propun o soluție în care de fapt se transmite spre decodor explicit statistica, **etapele** parcurse de **decodor** fiind în acest caz următoarele:

1. Preluare statistică (din antet)
2. Construire model (coduri)
3. Decodare pe baza modelului

B1. Codificare Shannon-Fano - exemplu

Să considerăm fluxul: ABCDECDECD

1. Determinare a statisticii

Obținem frecvența de apariție a fiecărui simbol:

A	B	C	D	E
1	1	3	3	2

2. Construire model

Primul pas: sortarea (de obicei crescătoare) după frecvența de apariție.

A	B	E	C	D
1	1	2	3	3

Se împarte apoi mulțimea de simboluri în două submulțimi având suma frecvențelor cât mai apropiată și atribuirea codului 0 respectiv 1 celor două "jumătăți" (se încearcă astfel atribuirea biților 0 și 1 cu o probabilitate cât mai apropiată de 0.5, cât este optimul teoretic)

A	B	E	C	D	Simbol
1	1	2	3	3	Frecvență
-8	-6	-2	4		Diferența în cazul împărțirii în fiecare din poziții

Diferența minimă (în valoare absolută) este dată la împărțirea între E și C. Se împarte astfel și se atribuie codul 0 primei "jumătăți" și codul 1 celei de a doua "jumătăți". Se continuă apoi (de obicei recursiv) asemănător pe fiecare din "jumătățile" obținute până se ajunge la simboluri individuale.

În figura următoare prezentăm întregul proces.

A	B	E	C	D	Simbol
1	1	2	3	3	Frecvență
0			1		
0		1	0	1	
0	1				

Rezultă următoarele coduri:

A = 000
B = 001
C = 10
D = 11
E = 01

Observăm că avem de a face cu o **abordare top-down** de construcție a codurilor.

3. Codare pe baza modelului

Se parcurge fluxul (din nou) și se înlocuiește fiecare simbol cu codul aferent. Obținem astfel:

ABCDECDECD

șir inițial

0000011011011011011011

flux de date comprimate

Din păcate delimitarea între simboluri se pierde (codurile având lungime variabilă), ceea ce va complica puțin decodarea.

B2. Codificare Huffman Static - exemplu

Să considerăm fluxul: ABCDECDECD

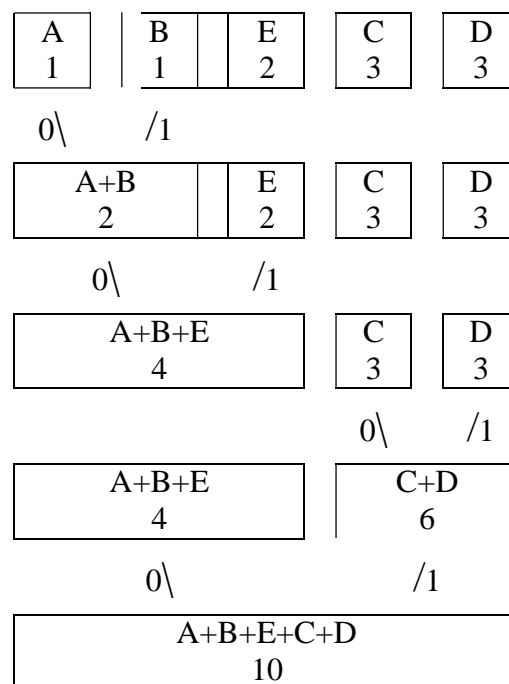
1. Determinare a statisticii

Obținem (la fel):

A	B	C	D	E
1	1	3	3	2

2. Construire model

În acest caz se grupează la fiecare pas cele două simboluri cu frecvența cea mai mică și se înlocuiesc cu un simbol având frecvența egală cu suma acestora. Celor două simboluri care s-au grupat li se atribuie biții 0 respectiv 1. Se repetă acest procedeu până se grupează toate simbolurile existente.



În final codul este dat de secvența de biți de pe calea spre simbolul inițial (de la final spre simbol).

Rezultă următoarele coduri:

A	=	000
B	=	001
C	=	10
D	=	11
E	=	01

Observăm că avem de a face cu o **abordare down-top** de construcție a codurilor.

Pentru exemple mici (ca cel de față) codurile dau la fel cu cazul Shannon-Fano. Pentru exemple mai mari ele însă vor diferi (teoretic, Huffman este codul optim).

3. Codare pe baza modelului

Se face la fel ca în exemplul de la Shannon-Fano.

C. Indicații de implementare

Sunt doar indicații, nu sunt strict obligatorii, dar cred că simplifică semnificativ implementarea.

C1. Structura fișierului comprimat

După cum am amintit anterior propun salvarea în fișierul comprimat a informațiilor legate de statistica apariției simbolurilor în fluxul de intrare. Schema generală:

Antet (model)	Flux de date comprimate (șirul de biți rezultat)
------------------	--

devine în acest caz:

Antet (frecvențele de apariție)	Flux de date comprimate (șirul de biți rezultat)
------------------------------------	--

C2. Decodarea

1. Preluare statistică – se realizează foarte simplu, prin citirea frecvențelor de apariție din antet
2. Construire model (coduri) – simplu, se folosește aceeași funcție ca la codor (pe baza aceluiași frecvențe se va genera același model).
3. Decodare pe baza modelului. Aici sunt două abordări posibile:

Prima: se citește câte un bit din fluxul de date comprimate până se obține un cod valid. Această soluție este destul de inefficientă (mai ales când sunt multe simboluri și codurile sunt lungi).

A doua: se citește câte un bit și se parcurge arborele (modelul). Se pleacă din rădăcină și se parcurge stânga sau dreapta conform cu bitul următor din flux până se atinge o frunză, obținând astfel simbolul decodificat. Se reia apoi din rădăcină pentru următorul simbol, etc.

C3. Problema sfârșitului fluxului

Deoarece în fluxul de biți comprimați exista în final unii biți adăugați în ultimul octet (maxim 7) care nu trebuie decodificați oprirea decodificării este o problemă. Deși există mai multe soluții posibile precum adăugarea unor informații suplimentare în antet (ex: lungimea fișierului, numărul de biți adăugați în ultimul antet) sau folosirea în model a unui simbol special de tip end-of-file EOF propun soluția simplă bazată pe cunoașterea lungimii fișierului, lungime care poate fi determinată ușor ca sumă a frecvențelor de apariție existente în antet.

C4. Realizarea unei singure aplicații codor-decodor

Deși cele două părți codor respectiv decodor funcționează complet independent propun realizarea lor într-o singură aplicație. În acest fel aplicația se poate realiza și depana mai ușor (de exemplu construirea modelului fiind una și aceeași). Atenție însă să nu avem influențe între codor și decodor altele decât prin intermediul fișierului comprimat.

C5. Colectarea statisticii

Cea mai simplă soluție este aceea de a contoriza apariția fiecărui octet din fișierul de intrare într-un contor corespunzător astfel:

```
Int32 contoare[256] //initializat cu 0
```

Apoi, pentru fiecare octet din fișier, să avem:

```
c = citește octet din fișier
contoare[c]++;
```

C6. Variante de antet

În varianta discutată până acum în antet trebuie **transmise toate cele 256 contoare** a câte 4 octeți.

contoare[0]	contoare[1]	...	contoare[255]	Flux de date comprimate (șirul de biți rezultat)
-------------	-------------	-----	---------------	---

În acest caz antetul va ocupa 1024 octeți. Pentru fluxuri mari acest lucru nu este semnificativ, însă pentru fluxuri mici poate fi semnificativ. Soluția prezentată este **satisfăcătoare** pentru implementarea acestei teme.

Putem însă lua în considerare următoarele variații:

- Includerea în începutul antetului a unei **hărți de alocare** a contoarelor (de 256 biți) astfel:

101 ...	contoare[0]	contoare[2]	...	Flux de date comprimate (șirul de biți rezultat)
---------	-------------	-------------	-----	---

Fiecare bit din harta de alocare indică dacă contorul respectiv este nenul și merită pus în antet. În acest caz antetul ocupă 32 octeți pentru harta de alocare și un număr variabil de octeți pentru contoare. Dacă cel puțin 8 simboluri nu există în flux soluția duce la câștig. Dacă nu, nu :)

Studentii de la ISM vor implementa antetul 2

- O altă idee este aceea de a avea **harta de alocare pe câte 2 biți** cu următoarea convenție (de exemplu):
 - 00 contorul nu se reprezintă (este nul)
 - 01 contorul se reprezintă pe 1 octet
 - 10 contorul se reprezintă pe 2 octeți
 - 11 contorul se reprezintă pe 4 octeți
- Studentii C/II vor implementa antetul 3

În acest caz harta de alocare ocupă 64 octeți (pierdere) dar se speră ca câștigul care apare din zona contoarelor să compenseze acest lucru.

C7. Construirea arborelui

La **Shannon-Fano**, în procedura (recursivă) de împărțire a intervalului, în momentul împărțirii se adaugă și doi fii nodului curent. În acest fel la terminarea împărțirilor avem gata și arborele corespunzător.

La **Huffman static** se construiește un listă (sau vector) de arbori, inițial fiecare arbore având doar o frunză (ca rădăcină). Apoi, la fiecare pas, se sortează lista, se scot cei 2 arbori cu contorul minim, se construiește un arbore având cei doi arbori scoși anterior ca fii (și contorul suma) și se introduce înapoi în listă arborele tocmai construit. Astfel, la fiecare pas, lista scade cu un arbore. În final avem un singur arbore în listă, tocmai arborele Huffman căutat.

C8. Codare pe baza arborelui

În cazul în care facem codarea direct pe baza arborelui trebuie avut în vedere că, prin parcurgerea căii din frunză spre rădăcină, se obține codul frunzei în ordinea inversă a biților. Deci, acest șir de biți trebuie inversat înainte de a fi scris în fișierul comprimat.

Pentru a ajunge însă ușor în frunza corespunzătoare unui simbol se recomandă întreținerea unui vector de referințe (pointeri) spre fiecare frunză (vector indexat cu simbolul).

C9. Sugestii interfata

Codor

- Buton “Load file” – cand se va apasa se va deschide case standard the OpenFileDialog si se va putea selecta orice tip de fisier de pe HDD
- Buton “Encode file”– la apasarea acestui buton se va encode fisierul selectat anterior(fie cu codificare Huffman fie cu codificare Shannon).Fisierul encoded va fii salvat in acelasi director cu fisierul original. Filename-ul fisierul encoded va fii identic cu numele fisierul original la care se va adauga extensia .hs (huffman static) sau .sf (Shannon fano). E.x. daca fisierul input se numeste “poza mea.jpeg” -> fisierul encoded -> “poza mea.jpeg.hs”
- Text input – multiline text area unde se va putea introduce free text (Ascii or Unicode format)
- Buton “Encode input text”. La apasarea acestui buton se va encode textul introdus in (fie cu codificare Huffman fie cu codificare Shannon). Numele fisierul generat va fii InputArea.hs sau InputArea.sf (depinzand de algoritmul implementat)
- Show codes – checkbox (va exista atata la coder cat si la decoder functionalitatea). Daca checkbox-ul este checked – in momentul cand se va apasa pe butonul de “Encode file” sau “Encode input text” sau “Decode file” – se va afisa in list box-ul afferent checkbox-ului care este codificarea fiecarui symbol. Per exemplu dat anterior A – 001

Decodor

- Buton “Load encoded file” – cand se va apasa se va deschide case standard the OpenFileDialog si se va putea selecta orice fisier cu extensia .HS sau .SF
- Buton “Decode file” – la apasarea acestui buton se va decoda fisierul selectat anterior.Fisierul decoded va fii salvat in acelasi director cu fisierul selectat anterior. Filename-ul fisierul decoded va fii identic cu numele fisierul selectat la care se va adauga un timestamp + extensia originala. E.x. daca fisierul input se numeste “poza mea.jpeg” -> fisierul encoded -> “poza mea.jpeg.hs”
➔ fisierul decoded se va nume “poza mea.jpeg.hs.13-10-2020-14-54.jpeg”

