

OPTIMALITATEA CODULUI HUFFMAN

Stan Bianca-Mihaela

Sa ne amintim mai intai ce era un cod Huffman...

- Codul Huffman este un mod de a comprima un set de date astfel incat sa avem 0 pierderi.

Cum face asta codul Huffman?

- Folosind o codificare de tip “prefix-free”=niciun cod nu este prefix pentru un alt cod obtinut.
- exemplu:
 - {00, 01, 111} este o codificare in care niciun element nu este prefix pentru un alt element din sir
 - {00, 01, 011} nu respecta proprietatea de mai sus (01 este prefix pentru 011)

Cum se realizeaza aceasta proprietate utilizand codurile Huffman?

Sa luam un exemplu:

- Vrem sa codificam mesajul: BCCABBDDECCBBAEDDCC.

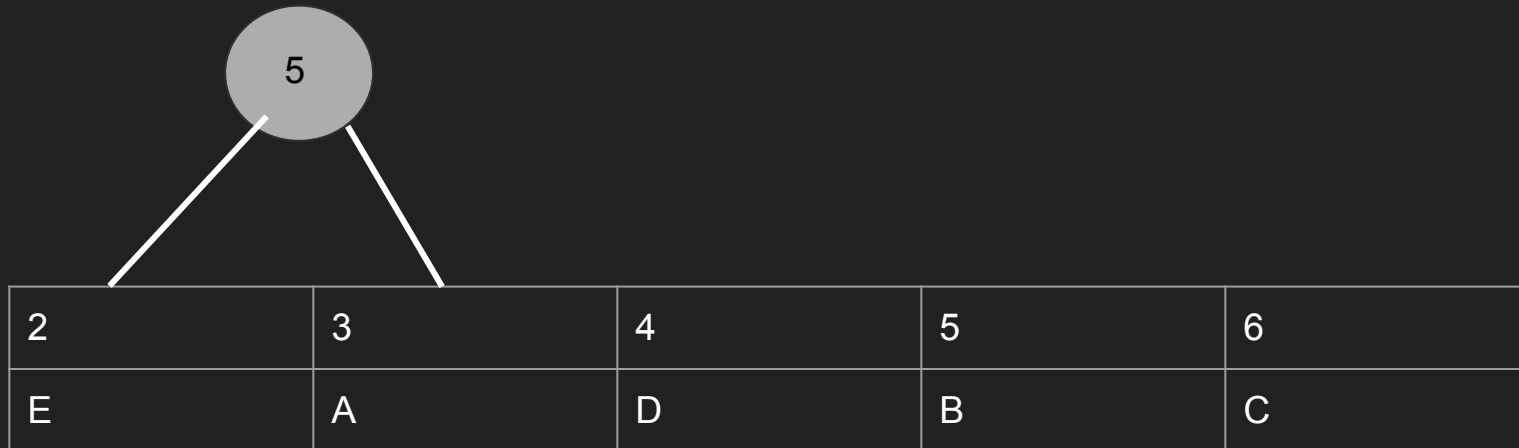
Fiecare litera este retinuta prin codul sau ASCII, care este un cod de 8 biti.

Avem 20 de litere => dimensiunea mesajului este $8 \times 20 = 160$ biti.

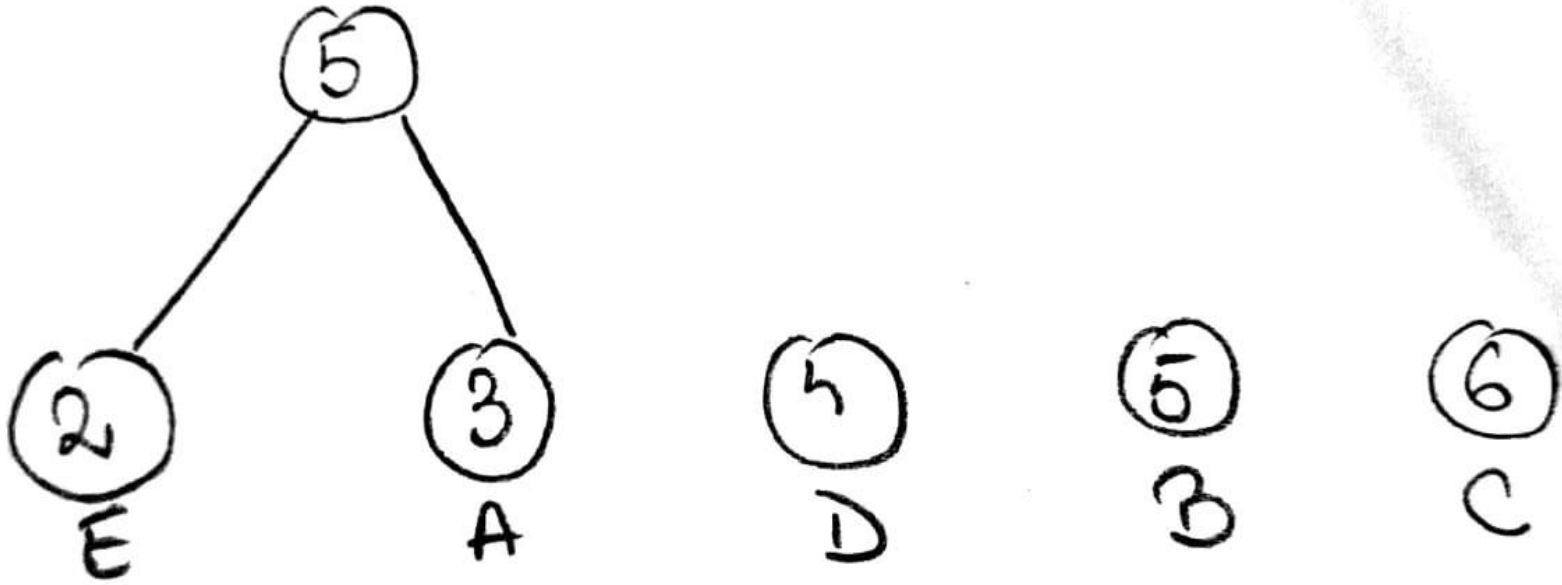
- Observam ca vectorul de frecventa arata:

A	B	C	D	E
3	5	6	4	2

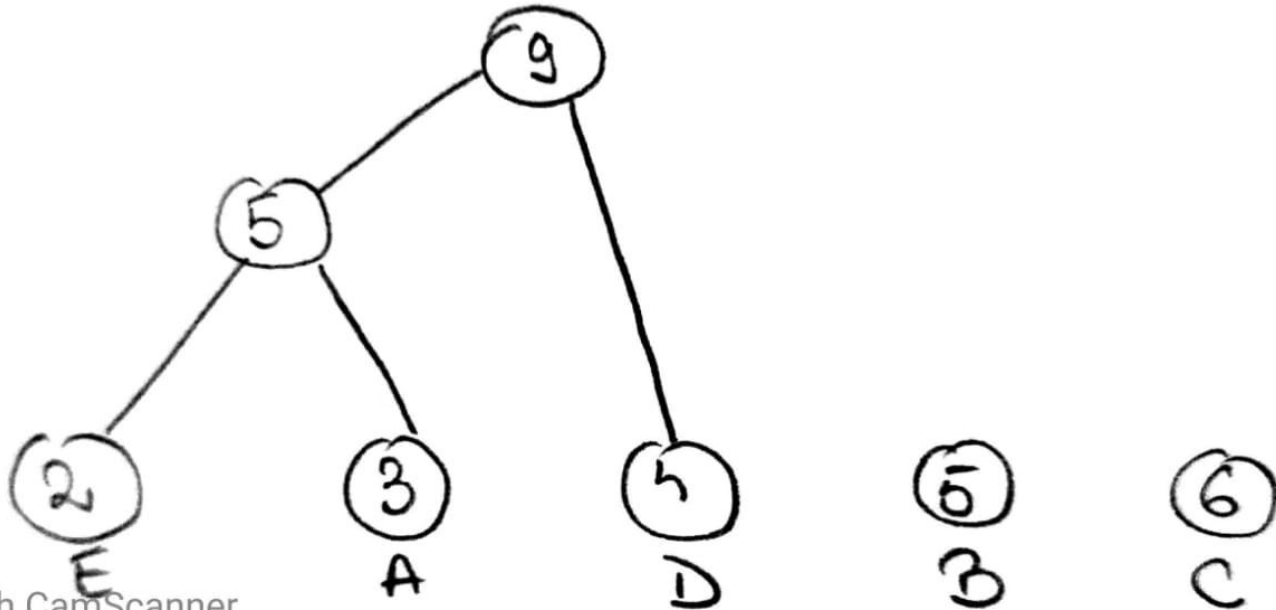
- Sortez literele in ordinea crescatoare dupa numaru de aparitii in mesaj.
- Urmaresc ceea ce se numeste “optimal merge parttern” care spune:
 - pun fiecare litera (reprezentata de numarul sau de aparitii) ca frunza intr-un arbore binar
 - la fiecare pas adaug in arbore un nou nod, format prin insumarea celor mai mici 2 noduri din arbore (nu luam in considerare nodurile care sunt deja fii)



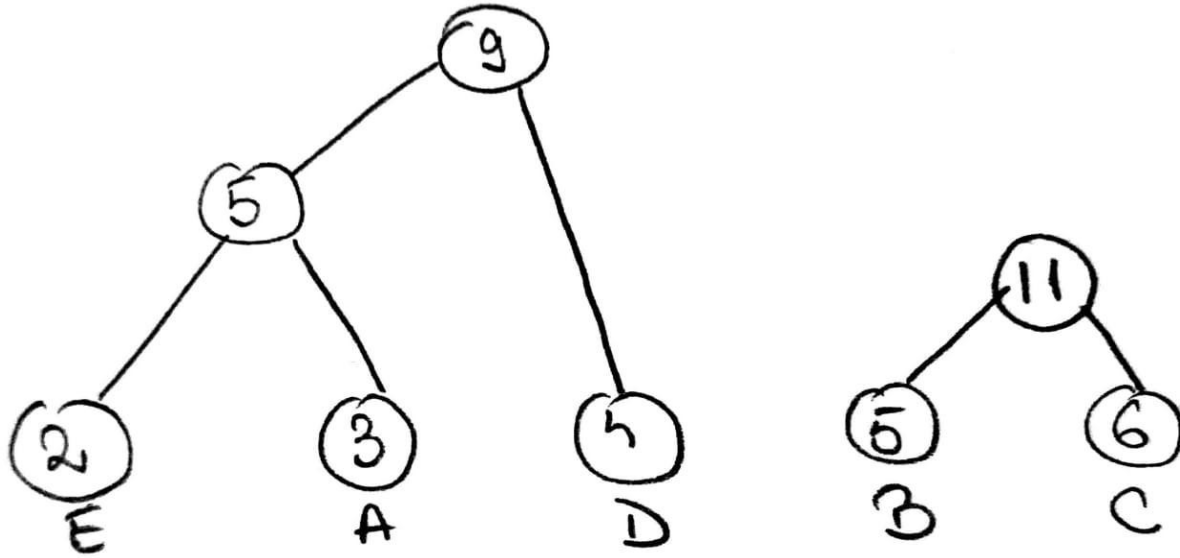
Pasul 1: cele mai mici frunze sunt 2 si 3 =>5



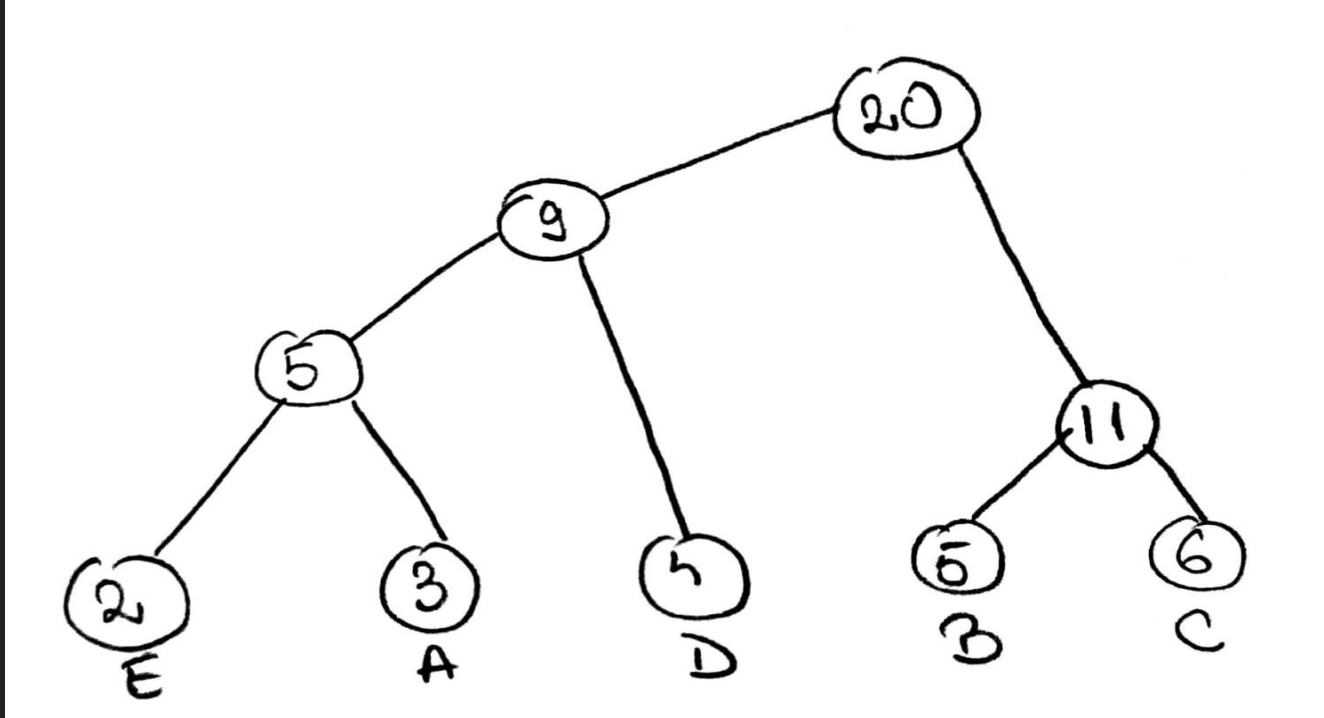
Pasul 2: cele mai mici frunze sunt 5 si 4 \Rightarrow 9



Pasul 3: Cele mai mici noduri sunt acum 5 si 6 \Rightarrow 11

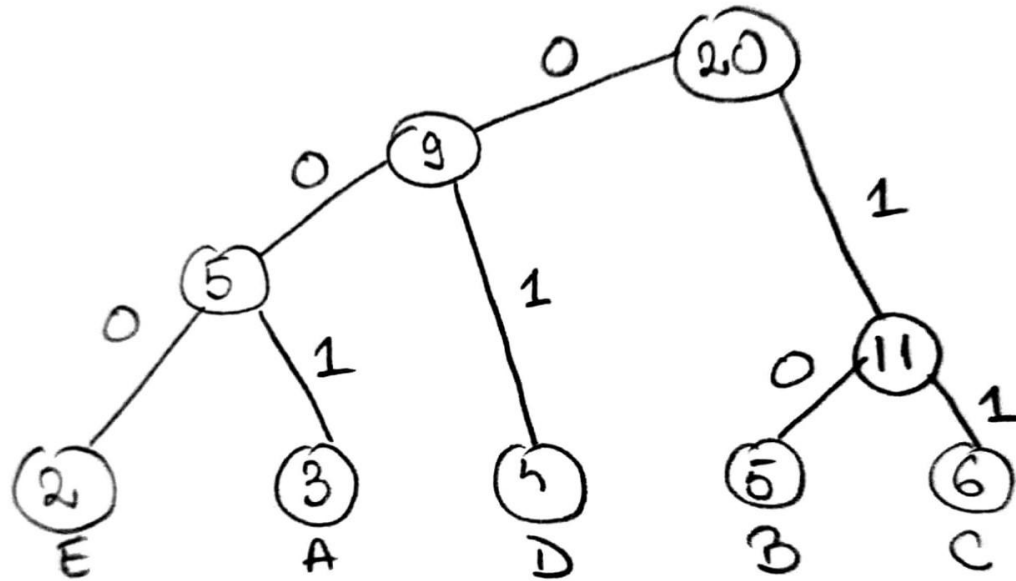


Pasul 4: cele mai mici noduri(care nu sunt parinti)
sunt 9 si 11 =>20

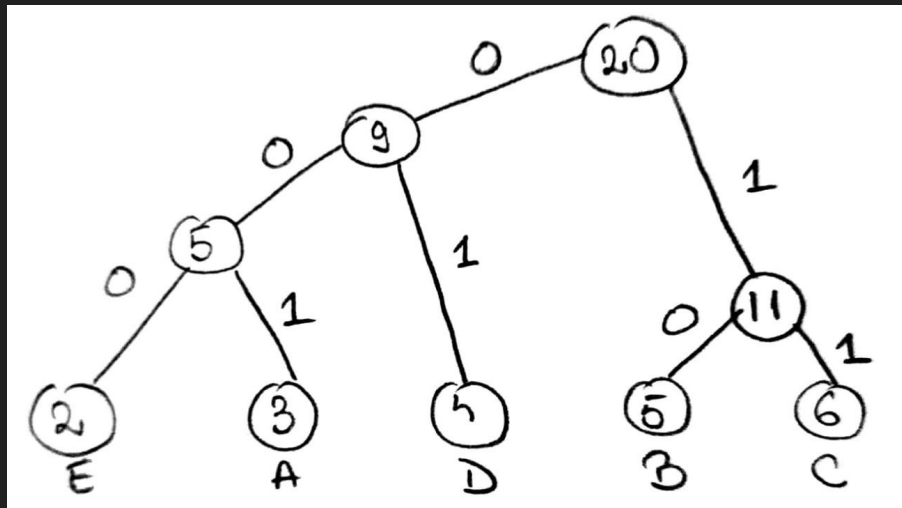


Pasul 5:

- legatura dintre tata si fiu drept o marcheaz cu 1
- legatura dintre tata si fiu stang o marcheaz cu 0



Pasul 6: codul pentru fiecare litera va fi format din literele intalnite in drumul din varful arborelui spre frunza corespunzatoare literei



	A	B	C	D	E
Frecventa	3	5	6	4	2
cod	001	10	11	01	000

Care este dimensiunea actuala a mesajului?

	A	B	C	D	E
frecventa	3	5	6	4	2
cod	001	10	11	01	000
dimensiune	$3b \cdot 3$ "9 ↑ $1b = 1 \text{ bit}$	$2b \cdot 5$ "10	$2b \cdot 6$ "12	$2b \cdot 4$ "8	$3b \cdot 2$ "6

=> dimensiunea totală a mesajului
= $9 + 10 + 12 + 8 + 6 = 45 \text{ biti}$

Imi trebuie totusi si un tabel pentru a putea decodifica mesajul. Ce dimensiune va avea tabelul?

A	B	C	D	E	= 5 cifre \cdot 8 biti	
001	10	11	01	000	40 biti	
3 biti	2 biti	2 biti	2 biti	3 biti		
12 biti						

$40 + 12 = 52$ de biti va avea tabelul

=> in final vom avea $45 + 52 = 97$ biti

Se observa ca aceasta dimensiune este mult mai mica decat dimensiunea initiala a mesajului, care era de 160 de biti.

Este insa aceasta dimensiune minima? Cu alte cuvinte, este codul Huffman o modalitate de a obtine o codificare optima?

Vom demonstra mai intai optimalitatea pe coduri binare, urmand sa generalizam foarte simplu din asta. Fie:

$B = 2 \leftarrow$ alfabetul codului nostru este binar

$A = \{0,1\} \leftarrow \text{alphabet}$

$$X = \{x_1, x_2, \dots, x_m\} \leftarrow \text{alphabet initial, } m \in \mathbb{N}$$
$$p = (p_1, \dots, p_m) \text{ pmf on } X$$

a e um pmf?

probability mass function: pentru x_k, p_k este probabilitatea ca alegând orice element din X , valoarea elementului să fie egală cu x_k .

ex: A A B A C $P_A = \frac{3}{5}$ ← $\frac{\text{nr de A-uri}}{\text{nr total de litere}}$

$e = (e_1, \dots, e_m) \leftarrow$ lungimea codurilor

Ce înseamnă că codul Huffman e optim?

Codifică mesajul corect și ea o dimensiune minimă.

Am văzut în exemplul de mai sus că dimensiunea finală a mesajului este:

$$\underbrace{\sum_{k=1}^{m_0} e_k \cdot p_k}_{\text{mesajul}} + \underbrace{\sum e_k + m_0 \cdot m_{\text{biti-un-caracter}}}_{\text{tabelele}}$$

Lema 1 (a ordonarii inverse)

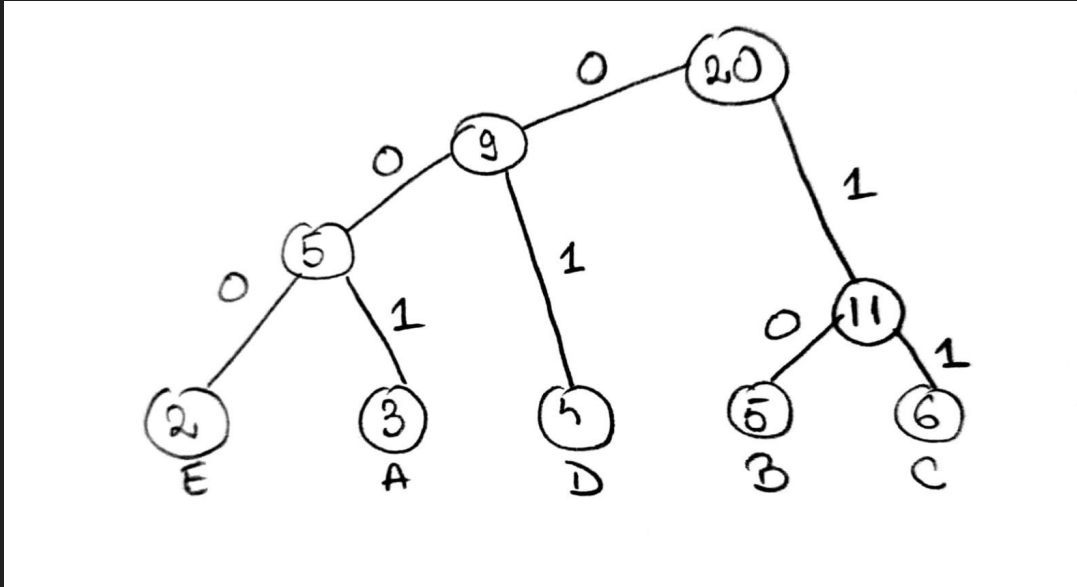
Intuitiv: Pentru orice cod optim de prefixe, lungimile codurilor sunt in ordinea inversa a probabilitatii fiecaredui cod. Putem observa asta din tabel:

- cum C apare de 6 ori, codul sau are 2 cifre
- E apare doar de 2 ori, deci are 3 cifre

	A	B	C	D	E
frecventa	3	5	6	4	2
cod	001	10	11	01	000

Lema 1 (a ordonarii inverse)

Acest lucru devine intuitiv cand ne gandim la arbore. Vrem ca o litera utilizata des sa fie usor de accesat, adica aproape de varful arborelui.



Lema 1 (a ordonarii inverse)

Formal:

Pentru orice "prefix code" optim :

$$\forall j, k, i \neq j, p_j > p_k, \text{ atunci } e_j \leq e_k$$

Demonstratie:

→ Fie C un "prefix code" optimal, cu lungimile
codurilor $w = (w_1, \dots, w_m) \in A^*$.

→ Stim ca $p_j > p_k$.

$$\text{ex: } w_2 = 01101$$

→ Fie C' astfel $\cdot w_i' = w_i \quad \forall i$ astfel $i \neq j'$
 $i \neq k$

$$\begin{cases} \cdot w_j' = w_k \\ \cdot w_{k'}' = w_j \end{cases}$$

(folosim C , doar inversam w_k cu w_j)

Curr c e optimal $\Rightarrow L_c \leq L_{c'}$
 \uparrow \nwarrow
 suma $e_i c_i$ suma $e_i c'_i$
 cu $e_i c_i$ cu $e_i c'_i$

$$\Rightarrow 0 \leq L_{c'} - L_c$$

$$\Rightarrow 0 \leq \sum_{i=1}^m e'_i p_i - \sum_{i=1}^m e_i p_i$$

evident $\left\{ \begin{array}{l} e'_i \geq e_i \quad \forall i \neq j, k \\ e'_j \geq e_k \\ e_k \geq e'_j \end{array} \right.$ (deci am inversat e_k cu e'_j)

$$0 \leq \sum_{\substack{i=1 \\ i \neq j, k}}^m e_i p_i + e'_j p_j + e'_k p_k - \sum_{i=1}^m e_i p_i$$

$$\Rightarrow 0 \leq \underbrace{e'_j p_j}_{= e_k} + \underbrace{e'_k p_k}_{= e'_j} - e_j p_j - e_k p_k$$

$$(2) \quad 0 \leq e_k p_j + e_j p_k - e_j p_j - e_k p_k$$

$$(2) \quad 0 \leq e_k (p_j - p_k) + e_j (p_k - p_j)$$

$$(2) \quad 0 \leq (e_k - e_j) (p_j - p_k)$$

$$\text{ntirm ca } p_j > p_k \Rightarrow p_j - p_k > 0$$

$$\Rightarrow 0 \leq e_k - e_j$$

$$(2) \quad e_j \leq e_k$$

Lema 2 (lema slaba a fratilor)

Formal:

Pentru orice $p \neq \emptyset$ $p = (p_1, p_2, \dots, p_m)$
există un „prefix code” optimal BINAR
cu 2 frati.

Ce sunt 2 frati în acest context?

Două coduri cu lungime maximă a1 :

a) au același lungime

b) se diferentiază doar prin ultimul bit

exemplu: 000
 001

în ordonat
după lungime

$\begin{bmatrix} 10100 \\ 11101 \\ 10110 \\ 10111 \end{bmatrix}$ asta 2 sunt frati

Demonstratie:

Fie C un "prefix code" optim al Σ e' e minimal.

Demonstrăm a): Dacă dintre cele mai lungi coduri au același lungime.

Presupunem că a) nu e adevărat.

\Leftrightarrow Există un cod care e mai lung decât toate celelalte.

Să ricem că lungimile sunt ordonate:

$$e_1 \leq \dots \leq e_m$$

↑
ce mai lung e strict mai mare ca celelalte

Lema 2 (lema slaba a fratilor)

Dacă metam codurile din C astfel:

$$W = (w_1, \dots, w_{m-1}, w_m)$$

putem să construim C' astfel:

$$W' = (w_1, \dots, w_{m-1}, w_m')$$

$$\text{unde } \underbrace{w_m' \alpha}_{\uparrow} = w_m, \alpha \in A$$

w_m' concatenat cu o literă din alfabetul A

$$e_m' = e_{m-1}$$

Lema 2 (lema slaba a fratilor)

Este doar acest c' mai construit, un "prefix-code"? Verifică proprietatea ca orice cod din c' să nu fie prefix pentru un alt cod din c' ?

Evident, primele $m-1$ coduri respectă proprietatea pentru că sunt identici cu codurile din c .

Se pune problema pt $w_{m'}$.

Să spunem că $w_{m'}$ e prefix pt alt cod din w' , să îi spunem w_k .

⇒ are aceiași lungime cu w_k (ca
(w_m era strict mai lung ca toate)

- $w_{m'}$ prefix pt w_k
- w_m are aceiași lungime cu w_k

⇒ $w_{m'} = w_k$

dar asta înseamnă că w_k era prefix pt w_m (contradicție cu proprietatea lui c)

Lema 2 (lema slaba a fratilor)

$\Rightarrow C'$ este un "prefix code"

Cum optimalitatea este data de $\sum e_i p_i$ și am micșorat suma totală a lungimilor $\Rightarrow \sum e_i' p_i \leq \sum e_i p_i$

$$e_1 p_1 + \dots + e_m' p_m \leq e_1 p_1 + \dots + e_m p_m$$

↑
e posibilă în egalitatea
pentru că p_m poate fi 0,
caz în care nu am modificat
deoc optimalitatea

Dar presupunem în definiția lui C că $\sum e_i$ e minimă (motiv
pentru care $\sum e_i p_i$ era minimă).

Dar $\sum e_i' = \sum e_i - 1 \Rightarrow$ contradicție asupra optimalității lui C
 \Rightarrow a) e adevărat

Lema 2 (lema slabă a fratilor)

Demonstrăm b):

Presupunem că b) e falsă.

(\Rightarrow) Fie K lungimea maximă a codurilor.

\Rightarrow oricărui 2 coduri de lungime K diferă undeva în primele $K-1$ cifre.

Ce vom face? Pentru:

000
001
:
1000

Luăm unul dintre codurile de lungime maximă și îi tăiem ultimul bit!

1000
1010
11101
10010
10111

\Downarrow
• noul cod va fi optim
• vom verifica că are propr. de prefix

Formal: Avem c-ue moștenite cu $W = (w_1, \dots, w_{m-1}, w_m)$ ordonate cu $\ell_1 \leq \dots \leq \ell_m$

Construiesc C' cu $W' = (w_1, \dots, w_{m-1}, w_{m'})$

unde:

- $w_{m'} = w_m$, $\alpha \in A$
- $\ell_{m'} = \ell_m - 1$

Cum am văzut ea demonstrează de ea a), C' este optim.
Este încă prefix-cod? 4

Lema 2 (lema slaba a fratilor)

La fel ca la a), stim ca w_1, \dots, w_n respecta proprietatea de prefix.

Trebuie doar sa verificam doar ca:

- w_m nu e prefix pentru niciun alt cod
- niciun alt cod nu e prefix pentru w_m

• este w_m prefix pentru un alt cod? // lungime K

→ daca e prefix pt unul din codurile mai lungi decat ee

\Rightarrow cele $K-1$ cifre ale lui w_m sunt identice cu primele $K-1$ cifre din celalalt cod

\Downarrow

contradictie cu proprietatea ca
orice 2 coduri de lungime maxima diferă
în primele $K-1$ cifre.

→ daca e prefix pt. unul din codurile de lungime $K-1$,
notat cu w_g

$|e_m| = K-1 \Rightarrow$ este identic cu celalalt cod

$|e_g| = K-1$

w_m prefix al lui w_g

$\Rightarrow w_m$ identic cu w_g

$\Rightarrow w_g$ e prefix pt w_m

\Downarrow

contradictie
cu proprietatea de prefix a
lui c

Lema 2 (a fratilor)

- micum alt cod nu e prefix pt. wm' ?

Fie $wg \in W'$ at wg prefix pt wm' .

$\Rightarrow wg$ prefix pt $wm' \cdot \alpha$, $\alpha \in A$

$\Leftrightarrow wg$ prefix pt wm

\Downarrow

contradicție ea
proprietatea de prefix a
lui c

\Rightarrow c' este optim și are proprietatea de prefix

\Rightarrow b) e adevărat

\Rightarrow Lema 2 este adevărată \cup

Lema 3 (lema puternica a fratilor)

Lema 3:

Fie p unu $p_m \neq a$ $p_1 \geq \dots \geq p_m$. Există un prefix code binar w $w = (w_1, \dots, w_m)$ a :

a) $e_1 \leq \dots \leq e_{m-1} = e_m$

b) w_{m-1} și w_m diferă doar la ultimul bit

Lema 3 (lema puternica a fratilor)

Demonstratie:

Ma' folosesc de urma neaba' a fratilor.

Fie c un prefix code optimal cu 2 frati, iar pmf-ul lui c este ordonat $\hat{p}_1 \geq \dots \geq \hat{p}_m$.

Demonstrat a):

Fie $p_k \geq p_{k+1}$.

\Rightarrow fie : 1. $p_k > p_{k+1}$

2. $p_k = p_{k+1}$

Cazul 1: $p_k > p_{k+1}$

Lema 1
 $\Rightarrow e_k \leq e_{k+1}$

Cazul 2: $p_k = p_{k+1}$

\Rightarrow pot ordona p_k m', p_{k+1} \hat{a}

$$\min(e_k, e_{k+1}) \leq \max(e_k, e_{k+1})$$

\Rightarrow pt $p_1 \geq \dots \geq p_m$, $e_1 \leq e_2 \leq \dots \leq e_{m-1} \leq e_m$
din urma neaba' a fratilor: $e_m = e_{m-1}$

$\Rightarrow e_1 \leq e_2 \leq \dots \leq e_{m+1} = e_m$ c.c.t.d.

Lema 3 (lema puternica a fratilor)

Demonstrata b)

w_{m-1} și w_m sunt adiacente și ele din Lema slabă a fratilor

⇒ w_{m-1} și w_m se diferentiază doar prin ultimul bit

⇒ Lema 3 este adevărată

Definitii

Fie $p = (p_1, \dots, p_m)$, $p_1 \geq \dots \geq p_m$
 construieste $p' = (p_1, \dots, p_{m-2}, p_{m-1} + p_m)$

Definitii

- Huffman contraction**: (notată cu HC)

Fiind dat un cod core respectă erma puterica a fratilor, C^S , pentru $p_{m-1} + p_m$, atunci "Huffman contraction", HC , pentru p' are lui C^S au codurile:

$$\left| \begin{array}{l} \omega_1^C = \omega_1^S \\ \omega_{m-2}^C = \omega_{m-2}^S \\ \omega_{m-1}^C = \omega_{m-1}^S \\ \omega_{m-1}^C 1 = \omega_m^S \end{array} \right.$$

exemplu: pentru C^S

0

1 0

$\omega_{m-1}^S = 1 1 0$

$\omega_m^S = 1 1 1$

\Rightarrow

C^C
0

1 0

1 1 ω_{m-1}^C

(evident:

$\omega_{m-1}^C 0 = 1 1 0 = \omega_{m-1}^S$

$\omega_{m-1}^C 1 = 1 1 1 = \omega_m^S$.)

Definitii

- **Huffman extremum** (notată cu HE)

Este procesul invers la a am făcut mai sus.

→ Fiind dat un cod optim C^0 pentru p' , atunci HE (pentru p) este C^E cu codurile:

$$C^0 = (w_1^0, \dots, w_m^0) \Rightarrow C^E = (w_1^0, \dots, w_{m-1}^0, w_m^0, w_m^0 1)$$

exemplu:

$$\begin{array}{c} C^0 \\ \hline 0 \\ 10 \\ 11 \end{array}$$

$$\begin{array}{c} C^E \\ \hline 0 \\ 10 \\ 110 \\ 111 \end{array}$$

Proprietate

Proprietate:

Orică cod Huffman pentru p este HE a unui cod Huffman pentru p' .

(formalizează algoritmul de adăugare a nodurilor în arbore)

Lema 4 (lema extensiei)

Lema 4 (Lema extensiei)

Intuitiv: Dacă Huffmanul codului pentru p' e optim, atunci n_i codul pentru p e optim.

Formal: Fie p un pmf al $p_1, z \dots z p_m$ și p' al $p'_1 z (p_1, \dots, p_{m-2}, p_{m-1} + p_m)$. HE (cațru p) a oricăru cod optim pentru p' este optim pentru p .

Demonstratie:

Fie C^0 optim pentru p' .

exemplu:

C^0	C^E
0	0
10	10
11	110
	111

Fie C^E HE alui C^0 cațru p .

$$L^E = e_1^E p_1 + \dots + e_{m-1}^E p_{m-1} + e_m^E p_m$$

$$\text{Stim ca } \begin{cases} e_i^E = e_i^0, & (\forall) i < m-1 \\ e_{m-1}^E = e_m^E = e_{m-1}^0 + 1 \end{cases}$$

$$\Rightarrow L^E = e_1^0 p_1 + \dots + e_{m-2}^0 p_{m-2} + (e_{m-1}^0 + 1) p_{m-1} + (e_{m-1}^0 + 1) p_m$$

$$\Rightarrow L^E = e_1^0 p_1 + \dots + e_{m-2}^0 p_{m-2} + e_{m-1}^0 (p_{m-1} + p_m) + p_{m-1} + p_m$$

$$\text{Stim ca } L^0 = e_1^0 p_1 + \dots + e_{m-2}^0 p_{m-2} + e_{m-1}^0 (p_{m-1} + p_m)$$

$$\Rightarrow L^E = L^0 + p_{m-1} + p_m$$

Lema 4 (lema extensiei)

Fie C^S un cod core respectând o anumită putere mică a frontierei pe p (implică faptul că e optimal). Fie $C^C \in C$ (câtă p') a lui C^S .

$$L^C = e_1^C p_1 + \dots + e_{m-1}^C (p_{m-1} + p_m)$$

$$\text{Stim că } e_i^C = e_i^S, (\forall) i < m-1$$

$$e_{m-1}^C = e_{m-1}^S - 1 = e_m^S - 1$$

$$\Rightarrow L^C = e_1^S p_1 + \dots + e_{m-2}^S p_{m-2} + (e_{m-1}^S - 1) p_{m-1} + (e_m^S - 1) p_m$$

$$\text{Stiu că } L^S = e_1^S p_1 + \dots + e_{m-1}^S p_{m-1} + e_m^S p_m$$

$$\Rightarrow L^C = L^S - p_{m-1} - p_m$$

Optimalitatea codurilor Huffman

În demonstrație ne vom referi la coduri Huffman standard, adică obținute prin combinarea celor mai mici elemente din pmf-ul dat.



Ce e un pmf?

$$p_x = \frac{\text{nr apariții ale lui } x \text{ în mesaj}}{\text{nr de caractere din mesaj}}$$

Vrem să demonstrăm că orice cod Huffman este optim.



- ① Orice cod Huffman pe un pmf cu k elemente este optim.

Am văzut din primul exemplu că optimalitatea este dată de minimizarea $\sum c_i p_i$.

Vom face o demonstrație prin inducție.

Pasul 1 (verificarea)

permis $k=2 \Rightarrow p' = (p_1 + p_2)$ optim

Optimalitatea codurilor Huffman

Pasul 2 (inductiv)

Presupunem că ① este valabilă pt $k = m-1$.

Demonstrăm că ① este valabilă pt $k = m$.

→ Fie g un pmf arbitrar cu m elemente.

→ Fie p pmf-ul g sortat descrescător, $p = (p_1, p_2, \dots, p_m)$

→ Fie c un cod Huffman pe p .

Dacă $p' = (p_1, \dots, p_{m-1} + p_m)$ n' c' este un cod Huffman pe p' .

Ipoteză
 $\Rightarrow c'$ este optim

Se vede că c este HE lui c' .

Folosim lemma exterioară:

intuitiv: Dacă Huffman cod-ul pe p' este optim, atunci
n' Huffman cod-ul pe p este optim.

$\Rightarrow c$ este optim \Rightarrow am demonstrat ① pt m
elemente

$$C_n(x_n) = 111 \dots 10.$$

For $i=1, \dots, k$: let $M_i = n/p_i$. So if $l_i \geq M_i$ then $l_i p_i \geq n$.

Let $\mathcal{C} = \{C : C \text{ is U.D. with lengths } l_1, \dots, l_n \text{ where } l_i \leq M_i \text{ for } i=1, \dots, k\}$.

$\Rightarrow C_n \in \mathcal{C}$ since $l_i^{C_n} \leq n = m_i p_i \leq M_i \quad \forall i=1, \dots, k$.

Lemma (Existence): For any pmf p , there exists an optimal prefix code.

Lemma:

Lemma (Siblings): For any prefix $p=(p_1, \dots, p_n)$, there exists an optimal binary prefix code s.t. two of the longest codewords:

- (a) have the same length
and (b) differ only the last bit.

e.g.

$$\#L = \sum l_i p_i$$

$$\begin{matrix} 000 \\ 001 \end{matrix}$$

Proof: Let C be opt. prefix s.t. $\sum l_i$ is minimal.

(a) Suppose (a) is not true $\Rightarrow l_1 \leq \dots \leq l_n$

$\frac{C}{w_1}$
 w_1
 \vdots
 w_{n-1}
 w_n

$\begin{matrix} 1000 \\ \left[\begin{matrix} 10100 \\ 10101 \\ 10110 \\ 101110 \end{matrix} \right] \end{matrix}$

Key Lemma

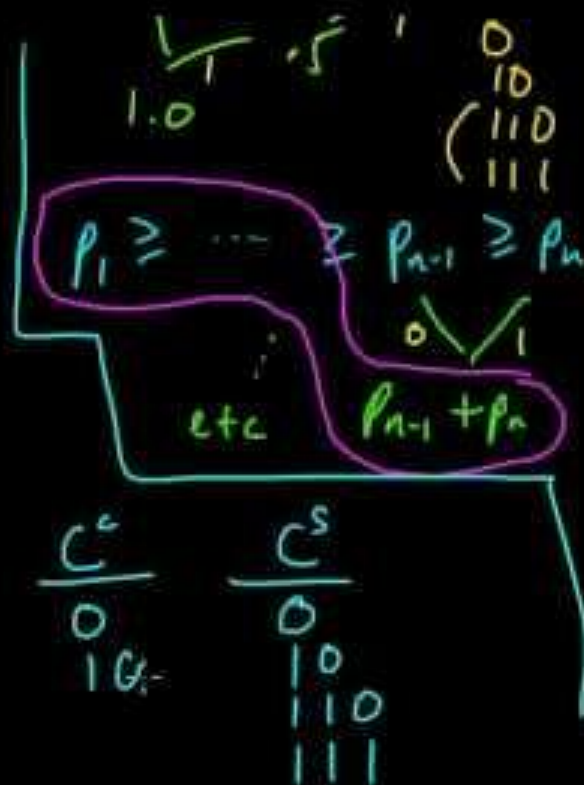
Lemma (Sibling code): Let p be a prob s.t. $p_1 \geq \dots \geq p_n$.

There exists an optimal binary prefix code with
codeword w_1, \dots, w_n ($w_i = C(x_i)$) s.t.

(a) $l_1 \leq \dots \leq l_{n-1} = l_n$

and (b) w_{n-1} and w_n differ only in the last bit.

Defn: Given a sibling code C^s for p , the H-contraction C^c to p' is



Key Lemma

Lemma (Sibling code): Let p be a prob s.t. $p_1 \geq \dots \geq p_n$.

There exists an optimal binary prefix code with
codeword w_1, \dots, w_n ($w_i = (x_i)$) s.t.

(a) $l_1 \leq \dots \leq l_{n-1} = l_n$

and (b) w_{n-1} and w_n differ only in the last bit.

Sibling
code

Let $p = (p_1, \dots, p_n)$, $p_1 \geq \dots \geq p_n$.

0	10	110	111
.5	.3	.1	.1

Lemma (Extension): Let p be a pmf s.t. $p_1 \geq \dots \geq p_n$.
 Let $p' = (p_1, \dots, p_{n-2}, p_{n-1} + p_n)$. The H-extension (to p)
 of any optimal code for p' is optimal for p .

Theorem: Any Huffman code is optimal.

Pf: $(*)_k$ Any Huffman code on a pmf with k elements is optimal.

Base $(*)_2$ is true.



Ind. Suppose $(*)_{n-1}$ is true. Let q be a pmf on n elements