

COMPRESIA ENTROPICĂ

2.1 CONCEPTE FUNDAMENTALE

2.1.1 Entropia unei surse staționare

Se consideră un *alfabet* finit **X** constituit din *M* simboluri și se definește un *mesaj* ca o secvență de simboluri. O sursă discretă de informație este numită stohastică dacă selecția simbolurilor care formează mesajele este făcută din cadrul alfabetului **X** în conformitate cu o distribuție de probabilitate $p_i \triangleq P(x_i)$. Din punct de vedere probabilistic, mulțimea mesajelor poate fi privită ca un proces aleatoriu discret adică, precum o secvență $(\xi_n)_{n=0}^{\infty}$ de variabile aleatoare, fiecare dintre ele luând valori în mulțimea **X** cu probabilitatea de distribuție $\{p_i\}$.

Se presupune în plus că sursa este *staționară*, adică

$$P\{\xi_{i1} = x_1, \dots, \xi_{ik} = x_k\} = P\{\xi_{i1+h} = x_1, \dots, \xi_{ik+h} = x_k\} \quad (2.1)$$

pentru toate numerele naturale i_1, \dots, i_k, h și toate $x_1, \dots, x_k \in \mathbf{X}$.

O caracteristică globală a unei astfel de surse din punctul de vedere al informației este *entropia sursei*.

$$H(x) \triangleq \sum_{i=1}^M p_i \log_2 \frac{1}{p_i} \quad (2.2)$$

care este exprimată în biți/simbol și indică cantitatea medie de informație a alfabetului **X**.

Teorema 2.1. - Entropia $H(x)$ a unei surse cu un alfabet conținând **M** simboluri, satisface inegalitatea

$$H(x) \leq \log_2 M \quad (2.3)$$

egalitatea fiind satisfăcută pentru simboluri echiprobabile.

În scopul transmiterii la distanță sau pentru stocarea informației sursei este eficient să se treacă de la reprezentarea oferită de alfabetul sursei **X** la o reprezentare printr-un alt alfabet **Y**. Procesul este numit *codarea datelor*. Un cod este o corespondență între mulțimea mesajelor sursei și mulțimea cuvintelor de cod.

Se consideră alfabetul sursei **X** constituit din simbolurile x_1, \dots, x_M . Fiecărui simbol *i* se asociază o secvență finită de simboluri ale alfabetului codului **Y**. În practică, alfabetul **Y** este constituit din două simboluri **0** sau **1**. În acest fel fiecărui simbol *i* se asociază o secvență finită de biți numită *cuvânt de cod*.

O codare eficientă (compresie) urmărește minimizarea *lungimii medii* a cuvintelor de cod.

$$\bar{n} \triangleq E\{n\} = \sum_{i=1}^M p_i n_i \quad (2.4)$$

unde n_i este lungimea cuvântului de cod (număr de biți) care reprezintă simbolul x_i , iar *n* este o variabilă aleatoare care reprezintă lungimea cuvântului de cod (luând valoarea n_i cu probabilitatea p_i , $i = 1, 2, \dots, M$).

Codurile realizează o corespondență între mesajele sursei (secvențe de simboluri - cuvinte ale alfabetului sursei) și cuvintele de cod (secvențe de biți - cuvinte ale alfabetului **Y**). După modul de stabilire a acestei corespondențe codurile pot fi clasificate în:

- i) bloc-bloc;
- ii) bloc-variabile;
- iii) variabile-bloc;
- iv) variabile-variabile;

Codurile bloc-bloc sunt coduri care stabilesc o corespondență între mesaje ale sursei de lungime fixată și cuvintele de cod de lungime fixată.

Codurile variabile-variabile stabilesc o corespondență între mesaje ale sursei de lungime variabilă și cuvinte de cod variabile.

EXEMPLUL 2.1

Se consideră

X = { a, b, c, d, e, f, g, h, spațiu }

Y = { 0, 1 }

Simboluri	Cuvinte de cod
a	000
b	001
c	010
d	011
e	100
f	101
g	110
Spațiu	111

Tab. 2.1 Tabel de codare bloc-bloc

Codul din tabelul 2.1 este un cod bloc-bloc și în aceeași categorie se încadrează codurile ASCII care pun în corespondență un alfabet de 64 (sau 256) simboluri cu cuvinte de cod de 6 (sau 8) biți.

EXEMPLUL 2.2

Se consideră

X = { a, b, c, d, e, f, g, h, spațiu }

Y = { 0, 1 }

Simboluri	Cuvinte de cod
aa	0
bbb	1
cccc	10
ddddd	11
eeeeee	100
ffffff	101
ggggggg	110
spațiu	111

Tab. 2.2. Tabel de codare variabile-variabile

Se observă că pentru un mesaj :

aa_bbb_cccc_ddddd_eeeeee_ffffffggggggg (40 simboluri in total)
se obține o lungime cod de 30 față de o lungime de cod de 120 dacă codarea s-ar fi făcut conform
tabelului 2.1 .

Observație : S-a folosit simbolul " _ " (underscore) pentru a nota caracterul spațiu

Un cod este *distinct* dacă fiecare cuvânt de cod este discernabil în raport cu celelalte
(corespondența între mesajele sursei și cuvintele de cod este biunivocă).

Un cod distinct este *unic decodabil* dacă fiecare cuvânt de cod este identificabil
când se află imersat într-o secvență de cuvinte de cod.

Codurile din Exemplul 2.1 și Exemplul 2.2 sunt ambele distincte, dar codul din
Exemplul 2.2 nu este unic decodabil. Spre exemplu codul 11 poate fi decodat fie ca dddd
fie ca bbbbbb.

Un cod unic decodabil este un *cod prefix* dacă are proprietatea de prefix care
pretinde ca nici un cuvânt de cod să nu fie prefixul altui cuvânt de cod. Codul care are
drept cuvinte de cod {1, 100000, 00} este un exemplu de cod care este unic decodabil dar
care nu are proprietatea de prefix.

Codurile prefix sunt *decodabile instantaneu* adică se bucură de proprietatea că
mesajul codat poate fi despărțit în cuvinte de cod fără a fi necesară examinarea în avans
(lookahead) a următorilor biți. Deci decodarea este posibilă imediat ce s-a primit ultimul bit
al cuvântului de cod fără a mai fi necesară recepționarea altor biți. Spre exemplu dacă se
folosește codul care are cuvintele de cod {1, 100000, 00}, decodarea mesajului 1000000001
și identificarea lui 1 ca fiind primul cuvânt de cod necesită examinarea celui de-al zecelea
simbol (era posibil ca primul cuvânt de cod să fie 100000).

Un *cod prefix minimal* este un cod prefix cu proprietatea ca dacă x este un prefix al
unui cuvânt de cod atunci și x sigma este fie un cuvânt de cod fie un prefix al unui cuvânt
de cod oricare ar fi litera sigma aparținând alfabetului Y. Spre exemplu considerând codul
{00, 01, 10} se observă că acesta este un cod prefix care nu este minimal. Astfel 1 este
prefix al cuvântului de cod 10 și ar impune ca și 11 să fie un cuvânt de cod fie un prefix al
unui cuvânt de cod valid. Dacă în exemplul de mai sus se înlocuiește cuvântul de cod 10 cu
cuvântul de cod 1 se obține un cod prefix minimal cu cuvinte de cod fie mai scurte.

2.1.2 Codarea alfabetului sursei

O reprezentare grafică utilă a codurilor prefix constă în asocierea fiecărui cuvânt de
cod cu un nod terminal (frunză) dintr-un arbore binar. Pornind de la rădăcina arborelui
primele două ramuri corespund selecției între 0 și 1 ca primă cifră a cuvântului de cod. Cele
două ramuri pornind din nodurile de ordinul întâi corespund celei de a doua cifre a
cuvântului de cod și procesul continuă mai departe (a se vedea fig.2.1).

Cuvintele de cod fiind asignate numai nodurilor terminale, nici un cuvânt de cod nu
poate fi prefix al unui alt cuvânt de cod.

O condiție necesară și suficientă ca un cod să fie cod prefix este dată de teorema
următoare:

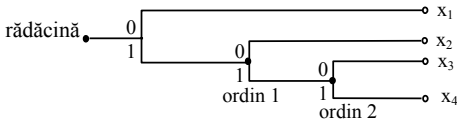


Fig.2.1 Arbore binar de codare

Simboluri	Cuvinte de cod
x ₁	0
x ₂	10
x ₃	110
x ₄	111

Tab.2.3 Tabelul de codare asociată arborelui

TEOREMA 2.2 Inegalitatea lui Kraft

Un cod binar care este un cod prefix având lungimile cuvintelor de cod n₁, n₂, ..., n_M
există dacă și numai dacă

$$\sum_{i=1}^M 2^{-n_i} \leq 1 \quad (2.5)$$

Observație:

Dacă p_i=2^{-n_i}, i=1, 2, ..., M condiția (2.5) devine o egalitate.

În acest caz se obține pe baza relațiilor (2.2) și (2.4) că $\bar{n} = H(x)$.

În general însă $\bar{n} \geq H(x)$ și se urmărește ca acest n să fie cât mai aproape de H(x).

Un cod prefix poate fi determinat astfel încât n să satisfacă teorema următoare:

TEOREMA 2.3

Se poate determina un cod binar prefix pentru orice alfabet al sursei de entropie
H(x); codul având o lungime medie care să satisfacă inegalitatea:

$$H(x) \leq \bar{n} < H(x)+1 \quad (2.6)$$

2.1.3 Clasificarea metodelor

Există mai multe criterii de clasificare a metodelor de compresie entropică.Toate
pleacă însă de la ideea de a reduce redundanța naturală a mesajului inițial, și acest lucru se
poate face doar dacă anumite simboluri sau grupe (subșiruri) de simboluri, pe care le vom
numi în general *forme* se repetă. În acest sens, putem clasifica metodele de compresie
entropică în două categorii:

- metode bazate pe frecvența de apariție a unei forme; aceste metode se realizează în principiu în două trecuri asupra textului, prima trecere permițând tocmai stabilirea acestor frecvențe

$$\varepsilon = \frac{\Delta [vH(x)]}{\frac{\bar{n}_v}{v}} \quad (2.9)$$

și redundanța $\eta \triangleq 1 - \varepsilon$

EXEMPLUL 2.5

Se consideră un alfabet al sursei $X = \{x_1, x_2, x_3\}$ cu $p_1=0,5$; $p_2=0,3$ și $p_3=0,2$. Să construim un nou alfabet $Y = X^2 = \{x_1x_1, x_1x_2, \dots, x_3x_3\}$. Să se construiască tabelul simbolurilor y_i , $i=1, \dots, 9$ cu probabilitățile ordonate descrescător și apoi pentru acest alfabet să se aplice algoritmul Huffman.

Soluție:

x_i	p_i	c_i
x_1	0,5	0
x_2	0,3	10
x_3	0,2	11

Tab. 2.5 Tabelul de codare pentru sursa X

$$\bar{n} = 0,5 \cdot 2 + 0,5 \cdot 1 = 1,5 \text{ biți /simbol}$$

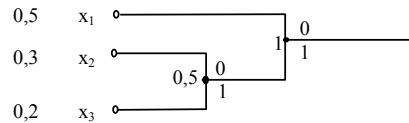


Fig. 2.3 Arbore de codare pentru alfabetul X

y_i	p_i	c_i
$Y_1 = x_1 x_1$	0,25	00
$Y_2 = x_1 x_2$	0,15	010
$Y_3 = x_1 x_3$	0,1	100
$Y_4 = x_2 x_1$	0,15	011
$Y_5 = x_2 x_2$	0,09	1100
$Y_6 = x_2 x_3$	0,06	1101
$Y_7 = x_3 x_1$	0,1	101
$Y_8 = x_3 x_2$	0,06	1110
$Y_9 = x_3 x_3$	0,04	1111

Tab. 2.6 Tabelul de codare pentru sursa Y

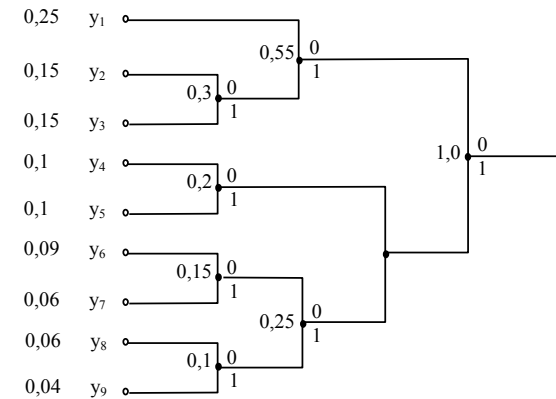


Fig. 2.4 Arbore de codare pentru alfabetul $Y = X^2$

Se obțin astfel lungimile medii ca fiind

$$\bar{n}_2 = 0,25 \cdot 2 + 0,5 \cdot 3 + 0,16 \cdot 4 = 2,24 \text{ biți/simbol}$$

$$\bar{n}_2/2 = 1,12 \text{ biți /simbol}$$

Exemplul 2.6

Metoda Huffman valorifică ideea că într-un cod optimal la $p_i > p_j$ corespunde $n_i < n_j$, și adaugă cerința ca cele mai puțin probabile două mesaje să aibă aceeași lungime. Tehnica codării constă în rescrierea tabelului de probabilități intercalând în ordine descrescătoare suma ultimelor două mesaje (cele mai puțin probabile), iterația oprindu-se când rămân în tabel două mesaje. Combinația de cod se citește urmând traseul săgeților de la dreapta la stânga (figura 5.4):

Mesaj	p_i	$\log(1/p_i)$	n_i	s_i^* Huffman
s_1	0,4	1,32	2	1
s_2	0,18	2,47	3	001
s_3	0,10	3,32	4	011
s_4	0,10	3,32	4	0000
s_5	0,07	3,83	4	0100
s_6	0,06	4,06	5	0101
s_7	0,05	4,32	5	00010
s_8	0,04	4,64	5	00011
			\bar{n}/η	2,61 97,8%

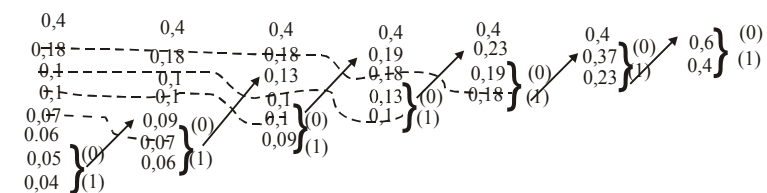


Fig 5.4

2.2.2 Algoritmul Shannon-Fano

Se presupune că cele M simboluri ale alfabetului sursei sunt ordonate în ordine descrescătoare a probabilității: $p_1 \geq p_2 \geq \dots \geq p_M$

- Se notează cu $F_i = \sum_{k=i}^M p_k$, unde $i = 1, 2, \dots, M$ și $F_1 = 1$

- Se determină un întreg n_i astfel încât

$$\log_2(1/p_i) \leq n_i < 1 + \log_2(1/p_i) \quad (2.10)$$

pentru fiecare $i = 1, 2, \dots, M$

- Cuvântul de cod asociat simbolului x_i este reprezentarea binară a fracției F_i pe n_i biți

Observație:

Prin reprezentarea binară a unei fracții se înțelege

$b_1 b_2 \dots b_k = b_1 \cdot 2^{-1} + b_2 \cdot 2^{-2} + \dots + b_k \cdot 2^{-k}$, unde $b_i \in \{0, 1\}$, $i = 1, \dots, k$

Se remarcă următoarele proprietăți ale algoritmului:

- (i) Mesaje cu probabilitate mare de apariție sunt reprezentate prin cuvinte cod scurte
- (ii) Cuvântul de cod pentru x_i va diferi de toate celelalte cuvinte de cod prin unul sau mai mulți biți. Prin aceasta se asigură decodarea unică. Pentru a demonstra această afirmație, se rescrie relația (2.10) sub forma

$$1/2^{n_i} \leq p_i < 1/2^{n_i-1}$$

Reprezentarea binară a lui F_i va fi diferită de toate celelalte prin unul sau mai mulți biți. Spre exemplu F_i va fi diferită față de F_{i+1} prin cel puțin bitul n_i deoarece $p_i \geq 1/2^{n_i}$. Prin aceasta cuvântul de cod pentru x_{i+1} va diferi de cel pentru x_i prin cel puțin un bit.

EXEMPLUL 2.6

Fiind dată o sursă care generează simboluri cu probabilitățile p_i specificate în coloana a doua a tabelului prezentat în continuare, să se aplice algoritmul de compresie Shannon-Fano.

Soluție:

Se remarcă ordonarea simbolurilor în sensul descreșterii probabilităților de apariție (coloanele 1 și 2). Se determină numărul de biți ai lungimii de reprezentare, completându-se astfel coloana 3 a tabelului.

Pentru simbolul x_1 numărul de biți se determină pe baza relației:

$$\log_2(128/27) \leq n_1 < 1 + \log_2(128/27)$$

sau

$$2,245 \leq n_1 < 3,245$$

deducându-se că $n_1 = 3$ biți.

Cuvântul de cod c_1 se obține din $F_1 \triangleq 0$. În concluzie $c_1 = 000$. Pentru x_2 se determină n_2

$= 3$ biți și $F_2 = \sum_{i=1}^2 p_i = p_1 = 27/128$. Reprezentarea binară a lui $27/128$ este 0011011. Prin

trunchierea acestei reprezentări la 3 biți se obține cuvântul de cod 001.

x_i	p_i	n_i	F_i	Reprezentare binară	cod c_i
x_1	27/128	3	0	.0000000	000
x_2	27/128	3	27/128	.0011011	001
x_3	9/128	4	54/128	.0110110	0110
x_4	9/128	4	63/128	.0111111	0111
x_5	9/128	4	72/128	.1001000	1001
x_6	9/128	4	81/128	.1001001	1010
x_7	9/128	4	90/128	.1011010	1011
x_8	9/128	4	99/128	.1100011	1100
x_9	3/128	6	108/128	.1101100	110110
x_{10}	3/128	6	111/128	.1101111	110111
x_{11}	3/128	6	114/128	.1110010	111001
x_{12}	3/128	6	117/128	.1110101	111010
x_{13}	3/128	6	120/128	.1111000	111100
x_{14}	3/128	6	123/128	.1111011	111101
x_{15}	2/128	6	126/128	.1111110	111111

Tab. 2.7 Tabelul de codare Shannon-Fano

Se constată că numărul de biți per simbol este $\bar{n}_v = 3,89$ biți / simbol și se obține o eficiență $\epsilon = 62,4\%$, entropia fiind de 2,433 biți / simbol.

Codarea unei secvențe x_8, x_5, x_7, x_2 este făcută, conform tabelului, în șirul de biți 110010011011001.

Pentru decodare, decodorul va examina mai întâi primii 3 biți 110 (lungimea minimă a acestui cod), va concluziona că nu este cuvânt de cod și va încerca un grup de 4 biți 1100 și în acest moment se va decoda simbolul x_8 și procedura va fi reluată începând cu cel de al cincilea bit.

O formulare alternativă a algoritmului Shannon-Fano este:

Pas 1: O listă de simboluri cu probabilități date se sortează în ordine descrescătoare a probabilităților;

Pas 2: Se divide lista în două părți, astfel încât probabilitatea cumulată a porțiunii superioare să fie cât mai apropiată de probabilitatea cumulată a părții inferioare;

Pas 3: Părții superioare i se asignează o cifră binară 0, iar părții inferioare cifra binară 1;

Pas 4: Se aplică recursiv pașii 2 și 3 acestor două jumătăți divizându-se și analizând biții codului până când fiecare simbol a devenit o frunză a arborelui

În exemplul următor este ilustrată aplicarea algoritmului Shannon-Fano în această formulare alternativă.

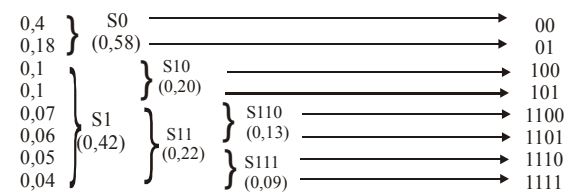


Fig 5.3