

CAPITOLUL 3

CODAREA SURSĂ

3.1. CUVINTE DE COD DE LUNGIME FIXĂ

Unele surse de informație generează semnale analogice ce pot fi descrise drept variabile aleatoare de tip continuu. În această clasă intră semnalul vocal generat de un microfon, unda seismică generată de un geofon și, în general, semnalele generate de diverse traductoare (un traductor fiind un dispozitiv care convertește o mărime neelectrică în una electrică și invers). Alte surse de informație generează șiruri de simboluri dintr-o mulțime discretă și finită. De exemplu, o dactilografă bătând la mașina de scris generează un șir de caractere din setul de caractere pentru care sunt prevăzute taste. Acest exemplu pare a nu mai fi de actualitate, căci mașina de scris a fost trimisă la Muzeul de Istorie a Tehnicii prin larga răspândire a calculatoarelor personale. Dar aceeași dactilografă, adaptându-se vremurilor, generează de la claviatura unui PC un text ce poate include un set mult mai larg de simboluri, de asemenea însă discret și finit. Mai mult decât atât, cineva care scrie un text utilizând calculatorul personal generează de fapt un fișier, adică un șir compus numai din două simboluri, 0 și 1, șir ce poate fi memorat sau transmis la distanță și pe care o imprimantă îl poate oricând transforma în textul tipărit. De la cursul de „Semnale, circuite și sisteme“, știm că, prin eșantionare și cuantizare, un semnal analogic poate fi transformat într-un semnal digital din care poate fi recuperat cu o distorsiune ce poate fi făcută acceptabil de mică. Prin urmare, indiferent dacă sursa de informație generează un semnal analogic sau unul digital, acest semnal poate fi transformat într-unul echivalent utilizând numai două simboluri, fie ele 0 și 1. Interesul nostru special pentru această posibilitate,

de a exprima orice informație cu numai două simboluri, se explică prin faptul că tranzistorul, care este baza electronicii, funcționează în regim de comutație între două stări, de conducție și de blocare, cărora le facem să corespundă cele două simboluri. Remarcabil însă este că această idee a avut-o filosoful englez Francis Bacon (1561–1626) cu circa trei secole și jumătate înainte de inventarea tranzistorului. El scria: „Un om poate exprima și comunica intențiile minții sale, la orice distanță prin obiecte... capabile numai de o diferență binară“.

Recapitulând, orice sursă de informație poate fi transformată într-o sursă de simboluri binare. O problemă importantă în memorarea și în transmiterea digitală a informației este reprezentarea eficientă a datelor generate de sursă. Procesul prin care se realizează această reprezentare se numește *codare sursă*. Dispozitivul care efectuează această reprezentare se numește *codor sursă*. Dacă variabila aleatoare X este ieșirea unei surse discrete, entropia $H(X)$ a sursei reprezintă cantitatea medie de informație emisă de sursă. În acest capitol, vom considera problema codării semnalului de ieșire a sursei, adică, reprezentarea semnalului de ieșire printr-un șir de simboluri binare, numite biți. O măsură a eficienței metodei de codare sursă se poate obține comparând numărul mediu de biți pe simbol de ieșire dintr-o sursă de entropie $H(X)$.

Un codor sursă eficient trebuie să satisfacă două cerințe funcționale:

1. Cuvintele de cod produse de codor trebuie să fie în formă *binară*.
2. Codul sursă să fie *univoc decodabil*, astfel încât șirul sursă original să poată fi reconstruit perfect din șirul binar codat.

Presupunem că o sursă discretă fără memorie (SDFM) produce un simbol, sau o literă, de ieșire la fiecare T_s secunde. Fiecare simbol este selectat dintr-un alfabet finit de simboluri $x_i, i = 1, 2, \dots, K$, ce apar cu probabilități $P(x_i), i = 1, 2, \dots, K$. Entropia acestei SDFM în biți pe simbol sursă este

$$H(X) = -\sum_{i=1}^K P(x_i) \log_2 P(x_i) \leq \log_2 K \quad (3.1)$$

EXEMPLUL 3.1: La cursul de „Semnale, circuite și sisteme“, se studiază eșantionarea, un proces prin care un semnal analogic este convertit într-un șir corespunzător de eșantioane, care sunt valori ale semnalului analogic prelevate în momente discrete $t_k = kT_e$, unde T_e este perioada de eșantionare, iar k este un număr natural. În cazul foarte important al semnalului vocal, semnalul generat de un microfon este mai întâi filtrat

astfel încât să nu conțină componente de frecvență în afara intervalului [300, 3400 Hz] și este apoi eșantionat cu o frecvență $f_e = 1/T_e = 8$ kHz, rezultând astfel un șir de valori indexate după k . Fie D gama dinamică a lui $s(t)$, definită ca diferență între valoarea maximă și valoarea minimă pe care le poate lua $s(t)$: $D = s_{\max}(t) - s_{\min}(t)$. În cazul semnalului vocal, gama dinamică D se împarte în $2^8 = 256$ de nivele. Cuantizarea este procedeul prin care mulțimea infinită a valorilor pe care un eșantion dat de mărime $s(kT_e)$ le poate avea în interiorul gamei dinamice D este redusă la mulțimea discretă și finită a celor 256 de nivele. În acest caz, alfabetul sursă este alcătuit din cele $K = 256$ de nivele. Fiecărui nivel de eșantionare îi corespunde un cuvânt binar de $\log_2 256 = 8$ biți.

EXEMPLUL 3.2: Cu ajutorul calculatorului personal, putem introduce de la claviatură mult mai multe semne decât permite o mașină de scris obișnuită. Codul ASCII (de la American National Standard Code for Information Interchange) are 8 biți pe caracter. Șapte biți sunt utilizați pentru a specifica un caracter, iar al optulea este un bit de paritate, care poate fi *pară* sau *impară*, la alegere. Deci, în acest cod, alfabetul sursă are $2^7 = 128$ caractere.

Numărul mediu de biți pe simbol emiși de sursă este $H(X)$ dată de (3.1), iar viteza cu care debitează sursa în biți/s este prin definiție $H(X)/T_s$.

Un *cod sursă* C pentru o variabilă aleatoare X este o aplicație de la domeniul de existență al lui X la mulțimea șirurilor de lungime finită de simboluri dintr-un alfabet cu B litere. Cel mai adesea, vom considera alfabetul binar $B = \{0, 1\}$.

Notă. Reamintim de la cursul de Analiză matematică ce înțelegem printr-o *aplicație*; aceasta este o *funcție* pentru care domeniul de definiție și cel de existență nu sunt neapărat mulțimi numerice, ci și mulțimi abstracte (de exemplu, mulțimi de evenimente).

Fie C_i cuvântul de cod corespunzător simbolului sursă x_i și să notăm cu $l(x_i)$ lungimea lui C_i . Prin definiție, *lungimea medie* L a unui cod sursă C pentru o variabilă aleatoare X având funcția masă de probabilitate $P(x)$ este dată de

$$L = \sum_{i=1}^K P(x_i)l(x_i) \quad (3.2)$$

EXEMPLUL 3.3 (*Codul Morse*): Pictorul american Samuel Finley Breese Morse (1791 – 1872) a inventat un cod relativ eficient pentru alfabetul latin utilizat în limba engleză cu numai patru simboluri: un punct, o linie, un spațiu între litere și un spațiu între cuvinte. În codul Morse, șiruri scurte reprezintă litere frecvente (de exemplu, un singur punct reprezintă E), iar șiruri lungi reprezintă litere ce apar rareori (de exemplu, Q se reprezintă prin „linie, linie, punct, linie“).

Să considerăm o schemă de codare bloc în care fiecare simbol din alfabetul sursă se reprezintă printr-un vector binar unic cu L componente. Fiindcă alfabetul sursă are K simboluri, iar K poate fi sau nu o putere a lui 2, distingem două cazuri:

1. $K = 2^k$

Avem deci

$$L = \log_2 K = \log_2 2^k = k \quad (3.3)$$

2. K nu este o putere a lui 2. Avem în acest caz $2^k < K < 2^{k+1}$ și deci

$$L = k + 1 \quad (3.4)$$

unde k este cel mai mare întreg mai mic decât $\log_2 K$.

Definim eficiența codării pentru o sursă discretă fără memorie (SDFM) drept raportul

$$\eta = \frac{H(X)}{L}. \quad (3.5)$$

Se vede că, dacă mărimea alfabetului K este o putere a lui 2 iar literele debitate de sursă sunt egal probabile, $L = H(X)$ și eficiența este maximă: $\eta = 1$. Dacă, însă, K nu este o putere a lui 2, dar simbolurile sursă sunt egal probabile, L diferă de $H(X)$ cu cel mult 1 bit pe simbol. Pentru $\log_2 K \gg 1$, eficiența acestei scheme de codare este ridicată. Pe de altă parte, dacă alfabetul are un număr K mic de simboluri, se poate mări eficiența unui cod de lungime fixă codând împreună un șir de J simboluri. Codul sursă constă din vectori binari cu N componente, ce permit codarea a 2^N șiruri de J simboluri dintr-un alfabet cu K litere. Condiția pentru N se scrie:

$$N \geq J \log_2 K \quad (3.6)$$

Prin urmare, valoarea întreagă minimă a lui N este

$$N = J \log_2 2^k + 1 = Jk + 1 \quad (3.7)$$

Numărul mediu de biți pe simbol sursă este acum $N/J = k + 1/J$ și astfel ineficiența s-a redus cu aproximativ un factor de $1/J$ în raport cu schema de codare simbol cu simbol descrisă mai înainte. Făcând J suficient de mare, eficiența procedului de codare, măsurată în raportul $JH(X)/N$, poate fi făcută oricât de apropiată de 1 se dorește.

Metodele de codare descrise mai sus nu introduc distorsiune întrucât codarea simbolurilor sursă sau a blocurilor de simboluri de cuvinte de cod este univocă. Acest tip de codare se numește *fără zgomot*.

Să presupunem că încercăm să reducem lungimea cuvântului de cod L slăbind condiția ca procesul de codare să fie univoc. În acest scop, selectăm cele mai probabile $2^N - 1$ blocuri de câte J simboluri pe care le codăm univoc, urmând ca celelalte $K^J - (2^N - 1)$ blocuri de câte J simboluri să se reprezinte prin singurul cuvânt de cod rămas nefolosit. La recepție, acest procedeu are drept rezultat o decodare greșită, ce intervine cu o anumită probabilitate de eroare, de fiecare dată când un bloc de probabilitate scăzută este pus în corespondență cu acel unic cuvânt de cod rămas nefolosit la codarea blocurilor mai probabile. Să notăm cu P_e această probabilitate de eroare. Pe baza acestui procedeu de codare bloc, Shannon a demonstrat în anul 1948 următoarea teoremă de codare sursă.

Teorema codării sursă I

Fie X variabila aleatoare reprezentând simbolurile emise de o SDFM având entropie finită $H(X)$. Blocuri de câte J simboluri debitate de sursă se codează în cuvinte de cod de lungime N dintr-un alfabet binar. Pentru orice $\varepsilon > 0$, probabilitatea P_e a unei decodări eronate a unui bloc poate fi făcută arbitrar de mică dacă

$$L \equiv \frac{N}{J} \geq H(X) + \varepsilon \quad (3.8)$$

iar J este suficient de mare. Reciproc, dacă

$$L \leq H(X) - \varepsilon, \quad (3.9)$$

P_e devine arbitrar de apropiată de 1 dacă J este suficient de mare.

Din această teoremă, rezultă că numărul mediu de biți pe simbol necesari pentru a coda ieșirea unei SDFM cu o probabilitate de decodare

eronată arbitrar de mică este mărginit inferior de entropia sursei. Pe de altă parte, dacă $L < X(X)$, frecvența unei decodări eronate se apropie de 100% atunci când J crește arbitrar.

3.2 CUVINTE DE COD DE LUNGIME VARIABILĂ

Dacă simbolurile emise de sursă nu sunt egal probabile, o metodă mai eficientă de codare este de a utiliza cuvinte de cod de lungime variabilă. Codul Morse menționat în **Exemplul 3.3** este cea mai veche ilustrare a acestei idei. Acest tip de codare se numește *codare entropică*.

Pentru a vedea ce probleme ridică o codare sursă în care cuvintele de cod au lungime variabilă în funcție de probabilitatea simbolurilor sursei, să considerăm exemplul unei SDFM cu litere de ieșire a_1, a_2, a_3, a_4 și probabilități corespunzătoare $P(a_1) = \frac{1}{2}$, $P(a_2) = \frac{1}{4}$ și $P(a_3) = P(a_4) = \frac{1}{8}$. În Tabelul 3.1 se arată trei variante de codare.

Tabelul 3.1

Coduri de lungime variabilă

Literă	$P(a_k)$	Cod I	Cod II	Cod III
a_1	$\frac{1}{2}$	1	0	0
a_2	$\frac{1}{4}$	00	10	01
a_3	$\frac{1}{8}$	01	110	011
a_4	$\frac{1}{8}$	10	111	111

Codul I are un „cusur“ fundamental. Să presupunem că recepționăm șirul 001001... Primii doi biți 00 corespund în mod clar simbolului a_2 , însă următorii patru biți prezintă ambiguitate, căci nu sunt univoc decodabili: ei pot fi decodați fie ca a_4a_3 , fie ca $a_1a_2a_1$. Este posibil ca această ambiguitate să fie rezolvată așteptând să ajungă la receptor biți

suplimentari, dar o astfel de întârziere de decodare este nedorită. De aceea, vom considera numai coduri care sunt decodabile *instantaneu*, adică, fără nici o întârziere de decodare.

Codul II din Tabelul 3.1 este *univoc decodabil* și *decodabil instantaneu*. Este convenabil să reprezentăm cuvintele de cod grafic drept frunze (noduri terminale) ale unui arbore, după cum se arată în figura 3.1.

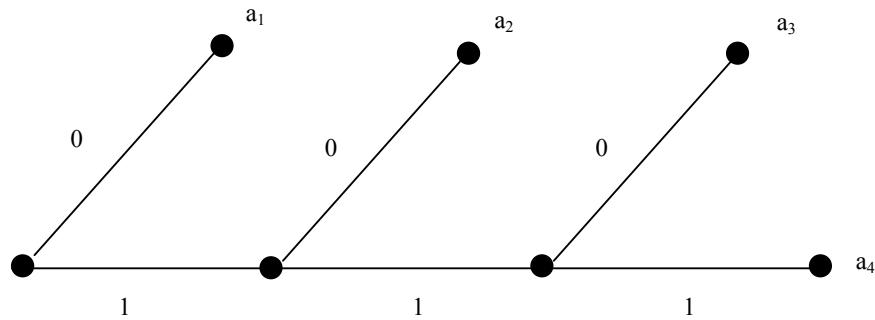


Fig. 3.1. Arbore de cod pentru codul II din Tabelul 3.1.

Se observă că bitul 0 indică sfârșitul unui cuvânt de cod pentru primele trei cuvânt de cod și că nici un cuvânt de cod nu este prefixul vreunui alt cuvânt de cod. În general, *condiția de prefix* cere ca, pentru un cuvânt de cod C_i de lungime i având componentele (b_1, b_2, \dots, b_i) , să nu existe nici un alt cuvânt de cod de lungime $j < i$ cu componente (b_1, b_2, \dots, b_j) pentru $1 \leq j \leq i-1$. Cu alte cuvinte, să nu existe nici un cuvânt de cod de lungime $j < i$ care să fie identic cu primii j biți ai altui cuvânt de cod de lungime $i > j$. Această proprietate face ca un cod să fie decodabil instantaneu.

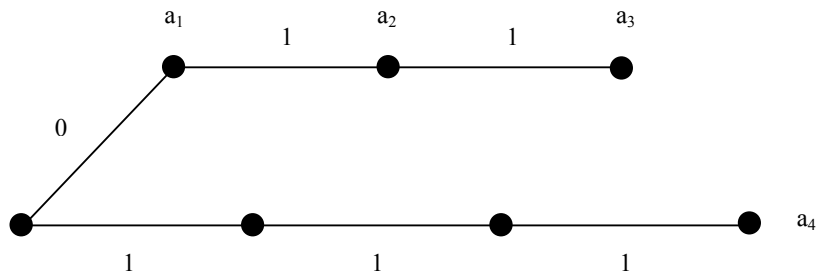


Fig. 3.2. Arbore de cod pentru codul III din Tabelul 3.1.

Codul III dat în Tabelul 3.1 are graful arbore arătat în figura 3.2. Se vede că este decodabil univoc dar nu și instantaneu. Este clar că acest cod nu satisface condiția de prefix.

Principalul nostru obiectiv este de a imagina o procedură sistematică pentru construcția codurilor de lungime variabilă decodabile univoc care să fie eficiente în sensul că lungimea medie (3.2) în biți pe literă debitată de sursă să fie minimizată. Condițiile pentru existența unui cod ce satisface condiția de prefix sunt date de teorema următoare.

Inegalitatea lui Kraft-McMillan

O condiție necesară și suficientă pentru existența unui cod binar instantaneu (adică, ce îndeplinește condiția de prefix) este ca lungimile cuvintelor de cod $l_1 \leq l_2 \leq \dots \leq l_K$ să satisfacă inegalitatea

$$\sum_{i=1}^K 2^{-l_i} \leq 1 \quad (3.10)$$

DEMONSTRAȚIE

Fie l mai mare decât oricare dintre lungimile l_i , în biți. Construim un arbore binar complet de ordin l care are 2^l noduri terminale (frunze) și două noduri de ordin i provenind din fiecare nod de ordin $i-1$, pentru orice i , $1 \leq i \leq l$. Ramurile arborelui reprezintă simbolurile binare ale cuvântului de cod. Spre exemplu, cele două ramuri ce ies din nodul rădăcină reprezintă cele două valori posibile ale primului bit. Fiecare cuvânt de cod este reprezentat de o frunză a arborelui. Drumul de la rădăcină la frunză ne dă biții ce compun respectivul cuvânt de cod. Ilustrăm aceasta în figura 3.3 pentru un arbore având 16 noduri terminale și o sursă care generează cinci litere cu $l_1 = 1, l_2 = 2, l_3 = 3$ și $l_4 = l_5 = 4$.

Condiția de prefix impusă cuvintelor de cod face ca nici un cuvânt de cod să nu fie strămoșul vreunui alt cuvânt de cod din arbore. De aceea, fiecare cuvânt de cod își elimină descendenții ca posibile cuvinte de cod.

Să considerăm toate cele 2^l noduri ale arborelui de la nivelul l . Unele dintre ele sunt cuvinte de cod, altele sunt descendenți ai unor cuvinte de cod, iar celelalte nu sunt nici una, nici alta. Dacă una din primele două ramuri ce ies din rădăcină (adică, $i = 1$) formează ea singură un cuvânt de cod ($l_1 = 1$), jumătate din ramurile arborelui, adică 2^{l-1} , sunt eliminate.

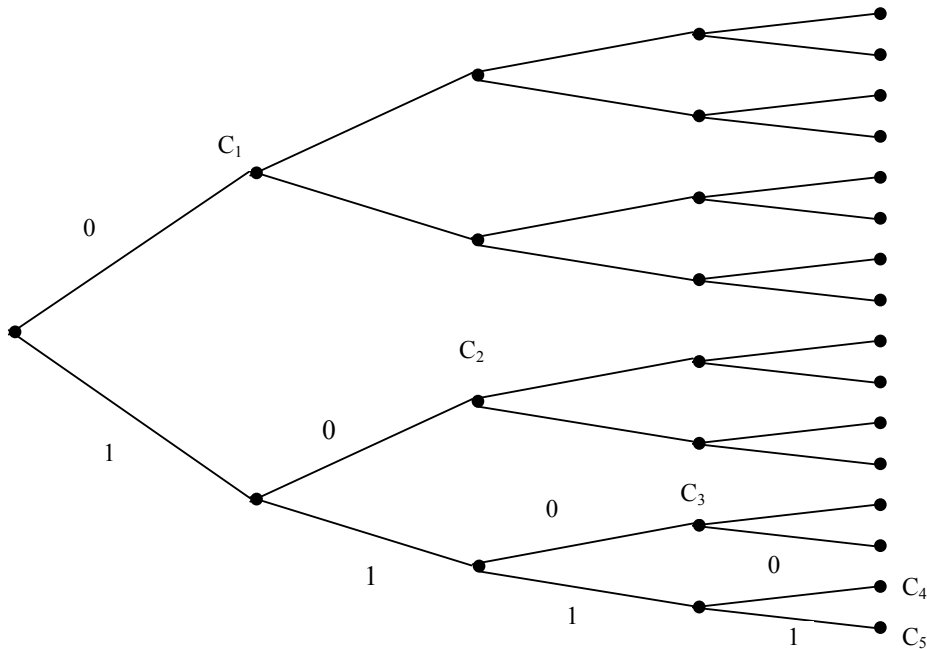


Fig. 3.3. Construcția unui cod arbore binar cuprins într-un arbore complet.

Conform aceluiași raționament, un cuvânt de cod de la nivelul l_i are 2^{l-l_i} descendenți la nivelul maxim l . Aceste mulțimi de descendenți trebuie să fie disjuncte. De asemenea, numărul total al nodurilor din aceste mulțimi trebuie să fie mai mic sau egal cu 2^l . Prin urmare, sumând peste toate cuvintele de cod, obținem

$$\sum_{i=1}^K 2^{l-l_i} \leq 2^l \quad (3.11)$$

din care, prin împărțire la 2^l , rezultă inegalitatea lui Kraft-McMillan (3.10).

Reciproc, dacă se dă orice mulțime de lungimi ale cuvintelor de cod l_1, l_2, \dots, l_K ce satisfac inegalitatea lui Kraft-McMillan, putem întotdeauna construi un arbore asemenea celui din figura 3.3. Etichetăm primul nod de adâncime l_1 drept cuvântul de cod C_1 și eliminăm descendenții săi din arbore. Apoi etichetăm primul nod rămas de adâncime l_2 drept cuvântul de

cod C_2 , etc. Procedând în acest mod, construim un cod prefix cu lungimile specificate l_1, l_2, \dots, l_K .

Teorema codării sursă II

Fie X variabila aleatoare reprezentând simbolurile $x_i, 1 \leq i \leq K$, emise de o SDFM cu entropie finită $H(X)$ cu probabilitățile corespunzătoare de apariție $p_i = P(x_i), 1 \leq i \leq K$. Este posibil să construim un cod ce satisface condiția de prefix și are o lungime medie L ce satisface inegalitățile

$$H(X) \leq L < H(X) + 1 \quad (3.12)$$

DEMONSTRAȚIE

Vom aproxima distribuția $p_i, 1 \leq i \leq K$ cu distribuția

$$r_i = \frac{2^{-l_i}}{c}, \quad 1 \leq i \leq K \quad (3.13)$$

unde

$$c = \sum_{i=1}^K 2^{-l_i} \quad (3.14)$$

Datorită inegalității lui Kraft-McMillan, $c \leq 1$. Scriem diferența dintre lungime medie și entropie astfel:

$$\begin{aligned} L - H(X) &= \sum_{i=1}^K p_i l_i - \sum_{i=1}^K p_i \log_2 \frac{1}{p_i} \\ &= \sum_{i=1}^K p_i \log_2 p_i - \sum_{i=1}^K p_i \log_2 2^{-l_i} \\ &= \sum_{i=1}^K p_i \log_2 \frac{p_i}{r_i} - \log_2 c \\ &= D(p \| r) + \log_2 \frac{1}{c} \\ &\geq 0 \end{aligned} \quad (3.15)$$

din cauză că entropia relativă este nenegativă iar $c \leq 1$. Rezultă că $L \geq H(X)$ cu egalitate dacă și numai dacă $p_i = 2^{-l_i}$, adică, dacă și numai dacă $-\log_2 p_i$ este un întreg pentru toți i .

Alegerea lungimilor $l_i = \log_2 \frac{1}{p_i}$ ne dă deci $L = H(X)$. Dar fiindcă $\log_2 \frac{1}{p_i}$ s-ar putea să nu fie un întreg, îl rotunjim prin majorare pentru a obține lungimi întregi ale cuvintelor de cod:

$$l_i = \left\lceil \log_2 \left(\frac{1}{p_i} \right) \right\rceil, \quad (3.16)$$

unde $\lceil x \rceil$ este cel mai mic întreg $\geq x$. Aceste lungimi satisfac inegalitatea lui Kraft-McMillan deoarece

$$\sum_{i=1}^K 2^{-\lceil \log_2 \frac{1}{p_i} \rceil} \leq \sum_{i=1}^K 2^{-\log_2 \frac{1}{p_i}} = \sum_{i=1}^K p_i = 1. \quad (3.17)$$

Această alegere a lungimilor cuvintelor de cod satisface

$$\log_2 \frac{1}{p_i} \leq l_i < \log_2 \frac{1}{p_i} + 1. \quad (3.18)$$

Înmulțind cu p_i și sumând peste i , obținem (3.12).

3.3. CODURI HUFFMAN

În anul 1952, Huffman a conceput un algoritm de codare cu lungime variabilă, optim în sensul că numărul mediu de biți necesari pentru a reprezenta simbolurile sursă este minim astfel încât cuvintele de cod să satisfacă acea condiție de prefix descrisă mai sus care permite ca șirul recepționat să fie decodabil univoc și instantaneu. Codarea Huffman este cel mai ușor de înțeles prin câteva exemple. Fie o sursă discretă fără memorie (SDFM) cu probabilitățile literelor $p_i = P(x_i)$, $i = 1, 2, \dots, K$.

EXEMPLUL 3.4: Considerăm o SDFM ce debitează șapte simboluri posibile x_1, x_2, \dots, x_7 cu probabilitățile ilustrate în figura 3.4.

Se ordonează simbolurile sursă în ordine descrescătoare a probabilităților: $P(x_1) > P(x_2) > \dots > P(x_7)$. Se începe procesul de codare cu cele două

simboluri care sunt cel mai puțin probabile x_6 și x_7 . Cele două simboluri sunt frunzele unui arbore; ramurile terminale respective se întâlnesc într-un nod în care simbolurile x_6 și x_7 se combină într-un simbol, să spunem, x'_6 . Următorul pas este de a uni cele două simboluri mai puțin probabile din mulțimea $x_1, x_2, x_3, x_4, x_5, x'_6$. Acestea sunt x_5 și x'_6 , care au împreună probabilitatea de 0,05. Atribuim ramurii de la x_5 un 0, iar ramurii de la x'_6 un 1. Această procedură continuă până când epuizăm mulțimea literelor sursă. Rezultatul este un arbore de cod ale cărui ramuri conțin cuvintele de cod dorite. Cuvintele de cod se obțin începând de la nodul cel mai din dreapta și mergând spre stânga. Cuvintele de cod sunt listate în Tabelul 3.2.

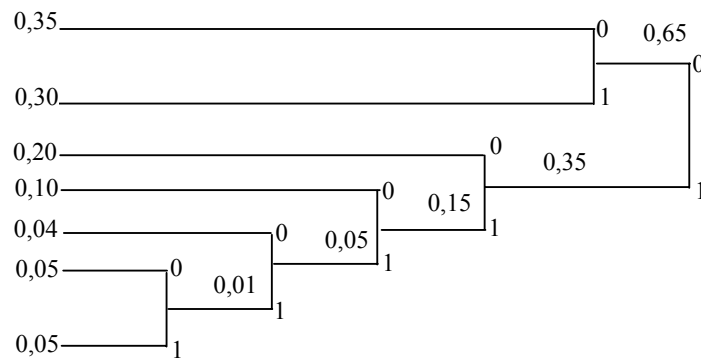


Fig. 3.4 . Exemplu de codare sursă pentru o SDFM cu lungime variabilă a cuvintelor.

Tabelul 3.2

Codul Huffman corespunzător figurii 3.4

Literă	Probabilitate	Autoinformație	Cod
x_1	0,35	1,5146	00
x_2	0,30	1,7370	01
x_3	0,20	2,3219	10
x_4	0,10	3,3219	110
x_5	0,04	4,6439	1110
x_6	0,005	7,6439	11110
x_7	0,005	7,6439	11111

Pentru acest cod, $H(X) = 2,11$ $L = 2,21$.

Acest cod nu este în mod necesar unic. Spre exemplu, la penultimul pas din procedura de codare, avem egalitate între x_1 și x_3 , căci aceste simboluri sunt egal probabile. În acest punct, facem alegerea de a împerechea x_1 cu x_2 . O alternativă este de a împerechea x_2 cu x_3 . Codul ce rezultă dacă facem această alegere este ilustrat în figura 3.5. Numărul mediu de biți pe simbol sursă pentru acest cod este tot 2,21. Codurile sunt, deci, la fel de eficiente. În al doilea rând, atribuirea unui 0 ramurii superioare și a unui 1 ramurii inferioare (mai puțin probabile) este arbitrară. N-avem decât să inversăm atribuirea lui 0 și a lui 1 pentru a obține un cod eficient care îndeplinește condiția de prefix.

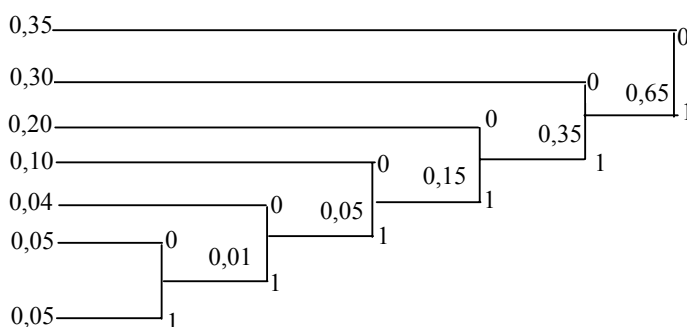


Fig. 3.5. Un alt cod, echivalent, pentru Exemplul 3.4.

Tabelul 3.3

Codul Huffman corespunzător figurii 3.5.

Literă	Cod
x_1	0
x_2	10
x_3	110
x_4	1110
x_5	11110
x_6	111110
x_7	111111

Lungimea medie a acestei variante de cod este tot $L = 2,21$.

EXEMPLUL 3.5: O sursă discretă fără memorie (SDFM) are un alfabet de ieșire cu opt litere ale căror probabilități sunt date în figura 3.6. Entropia acestei surse este $H(X) = 2,63$ biți pe simbol. Codul Huffman ilustrat în figura 3.6 are o lungime medie de $L = 2,70$ biți pe simbol, ceea ce înseamnă o eficiență de 0,97.

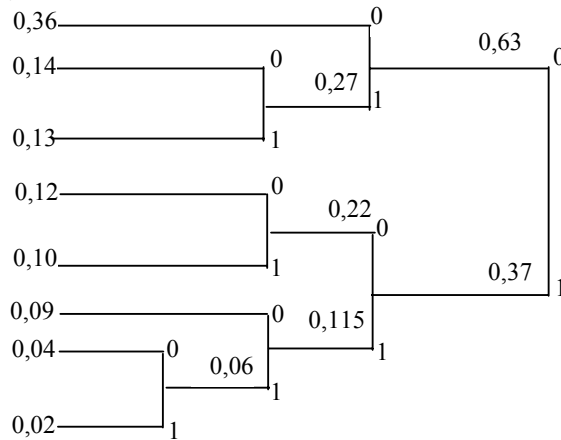


Fig. 3.6. Codul Huffman pentru Exemplul 3.5.

Rezultă codul Huffman dat în Tabelul 3.4.

Tabelul 3.4

Codul Huffman corespunzător figurii 3.6

Literă	Cod
x_1	00
x_2	010
x_3	011
x_4	100
x_5	101
x_6	110
x_7	1110
x_8	1111

Algoritmul Huffman descris în Exemplele 3.4 și 3.5, codând simbol cu simbol, generează un cod prefix având o lungime medie L ce satisface inegalitatea (3.12). O procedură mai eficientă este de a coda blocuri de J simboluri considerate ca „supersimboluri“, ceea ce face ca lungimea medie L_J să satisfacă inegalitatea

$$JH(X) \leq L_J < JX(X) + 1. \tag{3.19}$$

Împărțind la J , obținem

$$H(X) \leq \frac{L_J}{J} < H(X) + \frac{1}{J} \tag{3.20}$$

Numărul mediu de biți pe simbol sursă este

$$L = \frac{L_J}{J} \tag{3.21}$$

Luând J suficient de mare, se vede că putem face ca lungimea medie L să fie oricât de apropiată de entropia sursei $H(X)$.

EXEMPLUL 3.6: O SDFM are alfabetul de ieșire compus din trei litere x_1, x_2 și x_3 pe care le generează cu probabilitățile date în Tabelul 3.5. Entropia sursei este $H(X) = 1,518$. Codul Huffman dat în acest tabel necesită $L = 1,55$ biți pe simbol, ceea ce reprezintă o eficiență de 97,9%.

Tabelul 3.5

Cod Huffman pentru Exemplul 3.6

Literă	Probabilitate	Autoinformație	Cod
x_1	0,45	1,156	1
x_2	0,35	1,520	00
x_3	0,20	2,330	01

Să codăm acum perechi de litere ($J = 2$). Entropia variabilei aleatoare care ia ca valori perechi de litere este $2H(X) = 3,036$ biți pe perechi de simboluri. Codul Huffman necesită 3,0675 biți pe pereche de simboluri. Eficiența codării crește deci la $2H(X)/L_2 = 0,990$, adică, la 99%.

Tabelul 3.6

Cod Huffman pentru codarea perechilor de litere

Literă	Probabilitate	Autoinformație	Cod
x_1x_1	0,2025	2,312	10
x_1x_2	0,1575	2,675	001
x_2x_1	0,1575	2,675	010
x_2x_2	0,1225	3,039	011
x_1x_3	0,09	3,486	111
x_3x_1	0,09	3,486	0000
x_2x_3	0,07	3,850	0001
x_3x_2	0,07	3,850	1100
x_3x_3	0,04	4,660	1101

3.4. ALGORITMUL DE CODARE LEMPEL-ZIV

Deși algoritmul de codare Huffman este optim în sensul că este satisfăcută condiția de prefix iar lungimea medie a blocurilor este minimă, el are neajunsul că necesită cunoașterea distribuției de probabilitate a alfabetului sursei. În practică, statisticile surselor nu se cunosc întotdeauna *a priori*. Acesta este și cazul când transmitem un mesaj prin Internet, generând un lung șir de litere din alfabetul latin și de alte semne alfanumerice.

Algoritmul de codare Lempel-Ziv este independent de statisticile sursei. El utilizează cuvinte de cod de lungime fixă pentru a reprezenta un număr variabil de simboluri sursă. În practică, se utilizează blocuri fixe având lungimea de 12 biți, ceea ce permite alcătuirea unui dicționar, numit și *carte de cod*, cu 4096 de elemente. Ideea de bază este de a împărți datele emise de sursă în segmente care sunt cele mai scurte subșiruri neîntâlnite încă. Se presupune că simbolurile binare 0 și 1 sunt deja memorate în această ordine în cartea de cod. În scop ilustrativ, vom considera blocuri mult mai scurte, de numai 5 biți.

EXEMPLUL 3.7: O sursă generează șirul binar 1010110100100111010100001100111010110110...

Înainte de începerea codării, scriem

Subșiruri codate: 0, 1

Date de segmentat: 101011010010011101010000110011101011000110110...

Procesul de codare începe de la stânga. Cu simbolurile 0 și 1 deja memorate, *cel mai scurt subșir* întâlnit pentru prima oară este 10, astfel încât avem

Subșiruri codate: 0, 1, 10

Date de segmentat:1011010010011101010000110011101011000110110...

Cel de al doilea subșir cel mai scurt care n-a mai apărut este 101; scriem deci

Subșiruri codate: 0, 1, 10, 101

Date de segmentat:10100100111010100001100111010110001101100...

Următorul subșir cel mai scurt care n-a mai apărut este 1010; deci, scriem

Subșiruri codate: 0, 1, 10, 101, 1010

Date de segmentat:010011101010000110011101011000110110...

Continuăm în modul descris mai sus până când șirul de date va fi fost segmentat complet. Rezultă *cartea de cod* dată în Tabelul 3.7.

Prima coloană din Tabelul 3.7 indică pozițiile numerice ale subșirurilor din cartea de cod. Ultimul simbol din fiecare subșir din cartea de cod este un *simbol de inovație*, denumit astfel întrucât anexarea sa la un subșir particular îl deosebește de toate subșirurile memorate anterior în cartea de cod. De aceea, ultimul bit din cuvântul de cod reprezintă simbolul de inovație pentru respectivul subșir. Restul biților reprezintă în binar „pointerul“ către *subșirul rădăcină* ce corespunde respectivului subșir cu excepția simbolului de inovație. Primul subșir din trenul de date, 10, este compus prin concatenarea celui de al doilea element din cartea de cod, numerotat 1 și având valoarea binară 1, cu simbolul de inovație 0; de aceea, el este reprezentat prin numărul 10. Cel de al doilea subșir din trenul de date, 101, constă din elementul numerotat 2 și din simbolul de inovație 1, așa încât este reprezentat prin numărul 21. Celelalte subșiruri se codează similar. De exemplu, subșirul numerotat 9 este compus prin concatenarea subșirului 6 cu simbolul de inovație 0, astfel încât se scrie 60.

Tabelul 3.7

Carte de cod pentru algoritmul Lempel-Ziv

Poziție numerică	Subșir	Reprezentare numerică	Cuvânt de cod
0	0		
1	1		
2	10	10	00010
3	101	21	00101
4	1010	30	00110
5	01	01	00001
6	00	00	00000
7	11	11	00011
8	10101	41	01001
9	000	60	01100
10	011	51	01011
11	001	61	01101
12	110	70	01110
13	1011	31	00111
14	0001	91	10011
15	10110	130	11010

Decodarea este la fel de simplă. Se utilizează pointerul pentru a identifica subșirul rădăcină și se adaugă apoi simbolul de inovație. Fie, spre exemplu, cuvântul de cod 11010 din poziția 15. Ultimul bit, 0, este bitul de inovație. Restul de biți, 1101, indică subșirul rădăcină 1011 de la poziția 13. Prin urmare, cuvântul de cod 11010 este decodat drept 10110, ceea ce este corect.

Se observă că algoritmul a codat 45 biți sursă în 14 cuvinte de cod de câte 5 biți fiecare, cu un total de 70 de biți. S-ar spune că nu s-a realizat nici o compresie a datelor, ba dimpotrivă. Această ineficiență se explică prin faptul că am considerat un șir foarte scurt. Dacă se mărește lungimea șirului de date, eficiența crește. În cazul unui text obișnuit în limba engleză, codarea Huffman realizează o compactare de aproximativ 43%, în vreme ce algoritmul Lempel-Ziv asigură o compactare de aproximativ 55%.

Desigur, codorul și decodorul trebuie să lucreze în sincronism, astfel încât, după completarea elementelor din *cartea de cod*, a cărei capacitate este mare dar finită, algoritmul să o ia de la început pentru un nou șir de date.

Să nu scăpăm însă din vedere că între codor și decodor este o distanță fizică, iar transmisia informației între ele se face pe un canal de comunicație. Modul în care se efectuează această transmisie se prezintă în capitoul următor.