

## Lucrarea 4

# Procedeul de codare Huffman

### 4.1 Obiectivul lucrării

În această lucrare este studiată codarea surselor de informație discrete pentru canale fără perturbații prin procedeul Huffman. Vor fi analizate codurile Huffman binare și  $n$ -are și se va studia, de asemenea, aplicarea algoritmului Huffman pentru compresia datelor.

### 4.2 Introducere teoretică

Algoritmul de codare Huffman este un algoritm optimal de codare a surselor discrete de informație, în sensul că prin niciun alt algoritm nu se poate obține o lungime a cuvintelor de cod mai mică. Comparativ cu algoritmul Shannon-Fano, util numai pentru codarea surselor care generează simboluri ale căror probabilități sunt puteri negative ale lui 2, procedeul lui Huffman este aplicabil pentru orice distribuție de probabilități ale simbolurilor sursei de informație. El permite obținerea unui cod ale cărui cuvinte au lungimi variabile, fără prefix. Proprietatea codurilor fără prefix este aceea că niciun cuvânt de cod nu este un prefix pentru un alt cuvânt de cod.

Un cod fără prefix este un cod instantaneu, în cazul căruia decodarea unui cuvânt se efectuează pe măsura recepționării lui, fără a fi necesară examinarea cuvântului următor.

În cazul codării surselor pentru canale fără perturbații, sursa de informație primară generează un număr de simboluri sau mesaje care alcătuiesc alfabetul sursei a cărei dimensiune este  $N$ :

$$[\mathbf{S}] = [s_1 \ s_2 \ \dots \ s_N]. \quad (4.1)$$

Aceste mesaje  $s_i$  sunt puse în corespondență biunivocă cu tot atâtea cuvinte de cod  $c_i$  având lungimile  $l_i$ . Lungimea  $l_i$  a unui cuvânt de cod  $c_i$  este dată de numărul de simboluri  $x_i$  din structura sa. Pentru construcția acestor cuvinte de cod se folosesc literele  $x_i$  ale unei surse secundare, care aparțin alfabetului ce conține  $D$  litere:

$$[\mathbf{X}] = [x_1 \ x_2 \ \dots \ x_D] \quad (4.2)$$

cu probabilitățile  $P[\mathbf{X}] = [p(x_1) \ p(x_2) \ \dots \ p(x_D)]$ .

Pentru a evalua performanțele codului

$$[\mathbf{C}] = [c_1 \ c_2 \ \dots \ c_N] \quad (4.3)$$

se definește lungimea medie a cuvintelor de cod:

$$\bar{l} = \sum_{i=1}^N p(s_i) \cdot l_i. \quad (4.4)$$

Valoarea minimă a acestei mărimi, pentru o distribuție dată a probabilităților sursei primare, se obține atunci când probabilitățile simbolurilor sursei secundare sunt egale între ele; ea este dată de relația:

$$\bar{l}_{\min} = \frac{H(\mathbf{S})}{\log_2 D}. \quad (4.5)$$

unde  $H(\mathbf{S})$  reprezintă entropia sursei, calculată cu formula:

$$H(\mathbf{S}) = - \sum_{i=1}^N p(s_i) \log_2 p(s_i). \quad (4.6)$$

Este de dorit ca lungimea medie  $\bar{l}$  să fie cât mai apropiată de  $\bar{l}_{\min}$  și, în consecință, eficiența unui cod este definită ca:

$$\eta = \frac{\bar{l}_{\min}}{\bar{l}}. \quad (4.7)$$

Atunci când eficiența unui cod are valoarea 1, se spune că acesta este absolut optimal, iar dacă eficiența sa are valoarea maximă posibilă, dar este subunitară, se spune că este optimal.

În cazul general, al unei distribuții oarecare de probabilități ale simbolurilor sursei primare, prin procedeul Huffman se obține un cod optimal, iar dacă aceste probabilități sunt puteri negative ale lui 2, se obține un cod absolut optimal.

Redundanța unui cod este dată de relația:

$$\rho = 1 - \eta. \quad (4.8)$$

Pentru ca un cod fără prefix să existe, trebuie să fie îndeplinită inegalitatea Kraft-McMillan:

$$K = \sum_{i=1}^N D^{-l_i} \leq 1. \quad (4.9)$$

Această relație exprimă condiția suficientă pentru un cod fără prefix și condiția necesară pentru un cod unic decodabil.

Dacă inegalitatea nu este îndeplinită, atunci codul nu este unic decodabil.

Un cod optimal fără prefix are următoarele proprietăți:

1. Simbolurile cu probabilitatea cea mai mare vor corespunde cuvintelor de cod cele mai scurte. Pentru  $p(s_i) > p(s_j)$ , rezultă  $l_i \leq l_j$ .
2. Ultimele  $D$  simboluri cu cea mai mică probabilitate vor corespunde unor cuvinte de cod de lungime egală.

Codul Huffman se creează după următorul algoritm:

1. Se ordonează simbolurile sursei primare în ordine descrescătoare a probabilităților:  $p(s_1) \geq p(s_2) \geq \dots$

2. Ultimele  $D$  simboluri (unde  $D$  este dimensiunea alfabetului codului) cu cele mai mici probabilități se grupează într-un nou simbol a cărui probabilitate este egală cu suma probabilităților simbolurilor grupate:  $p(r_1) = \sum_{i=1}^D p(s_i)$ . Se obține o sursă restrânsă cu  $N - D + 1$  simboluri.
3. Se repetă pașii 1 și 2 pentru sursele restrânse astfel obținute, până când se obține o sursă restrânsă cu  $D$  simboluri.
4. Cuvântul de cod pentru fiecare simbol este obținut prin parcurgerea schemei de codare, în sens opus restrângerii (echivalent cu parcurgerea unui arbore de la nodul rădăcină către frunze). La fiecare ramificație se va adăuga o literă din cele  $D$  litere ce compun alfabetul codului.

**Exemplu:** Se consideră o sursă discretă de informație cu următoarea distribuție de probabilități:

$$P(\mathbf{S}) = \begin{bmatrix} 0.30 & 0.10 & 0.05 & 0.25 & 0.20 & 0.10 \end{bmatrix}.$$

Sursa va fi codată cu ajutorul algoritmului Huffman binar ( $D = 2$ ) și Huffman ternar ( $D = 3$ ).

a. **Codarea binară** ( $D = 2$ )

Schema de codare și codul obținut sunt prezentate în Fig. 4.1 și în Tabelul 4.1.

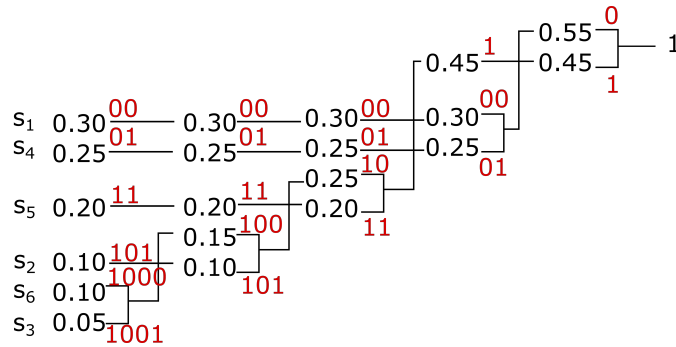


Figura 4.1: Procedeul de codare Huffman folosind un alfabet binar

Simbol	Probabilitate	Cuvânt de cod	Lungimea cuvântului de cod
$s_1$	0.30	00	2
$s_2$	0.10	101	3
$s_3$	0.05	1001	4
$s_4$	0.25	01	2
$s_5$	0.20	11	2
$s_6$	0.10	1000	4

Tabelul 4.1: Codul Huffman binar pentru sursa  $S$

Entropia sursei este calculată cu ajutorul relației (4.6):

$$H(\mathbf{S}) = 2.37 \text{ [biti/simbol]}.$$

Lungimea medie a cuvintelor de cod, calculată cu formula (4.4) este:

$$\bar{l} = 2.4,$$

iar lungimea medie minimă, calculată cu relația (4.5) este:

$$\bar{l}_{min} = 2.37.$$

În aceste condiții, eficiența codului este, conform relației (4.7):

$$\eta = 0.9875.$$

În acest caz, relația Kraft-McMillan (4.9) este îndeplinită cu egalitate:

$$K = 1.$$

**b. Codarea ternară ( $D = 3$ )**

Pentru a realiza codarea ternară și, în general, codarea  $n$ -ară, trebuie ca ultima sursă restrânsă să genereze exact  $D$  simboluri. După  $n$  pași de restrângere, o sursă va genera  $N - nD + n$  simboluri.

Dacă după  $n$  restrângeri, sursa trebuie să genereze  $D$  simboluri, atunci trebuie să fie îndeplinită următoarea relație:

$$D = (N - nD) + n,$$

de unde rezultă următoarea relație:

$$n = \frac{N - D}{D - 1}. \quad (4.10)$$

În mod evident, numărul de restrângeri este întotdeauna natural ( $n \in \mathbb{N}$ ).

Pentru a îndeplini această condiție, în cazul codării  $n$ -are poate fi necesară adăugarea de simboluri artificiale sursei primare, astfel încât să fie îndeplinită această condiție.

Pentru sursa propusă, se cunoaște  $N = 6, D = 3$ . Atunci:

$$n = \frac{6 - 3}{3 - 1} = \frac{3}{2} \notin \mathbb{N}.$$

Se va adăuga un simbol artificial cu probabilitatea 0, astfel încât  $N' = 7$ , și atunci:

$$n = \frac{7 - 3}{3 - 1} = 2 \in \mathbb{N}.$$

Schema de codare și codul obținut sunt prezentate în Fig. 4.2 și în Tabelul 4.2.

Pentru a analiza performanțele codului, se vor calcula aceleași mărimi ca în cazul codării binare.

Entropia sursei primare rămâne neschimbată, ea depinzând doar de probabilitățile simbolurilor sursei:



```

import heapq

class Node:
    def __init__(self, p, s):
        self.left = None
        self.right = None
        self.prob = p
        self.symbol = s

    def __lt__(self, other):
        return self.prob < other.prob

    def __repr__(self):
        return "Node({}, {}, {})".format(repr([self.prob, self.symbol]), ...,
                                           repr(self.left), repr(self.right))

def HuffmanTree(SP):
    pq = []
    for symbol, prob in SP.items():
        pq.append(Node(prob, symbol))
    heapq.heapify(pq)

    while len(pq) > 1:
        n1 = heapq.heappop(pq)
        n2 = heapq.heappop(pq)
        top = Node(n1.prob+n2.prob, n1.symbol+n2.symbol)
        top.left = n1
        top.right = n2
        heapq.heappush(pq, top)
        # print(n1)
        # print(n2)
        print(pq, '\n')
    return pq

def encode(dic_code, root, code):
    if root.left is None and root.right is None:
        dic_code[root.symbol] = code
    else:
        encode(dic_code, root.left, code+'0')
        encode(dic_code, root.right, code+'1')

S = ['s1', 's2', 's3', 's4', 's5', 's6']
P = [0.30, 0.10, 0.05, 0.25, 0.20, 0.10]
SP = dict((S[i], P[i]) for i in range(len(S)))
dic_code = dict((S[i], '') for i in range(len(S)))

PQ = HuffmanTree(SP)
encode(dic_code, root=PQ[0], code='')
print(dic_code)

```

Rezolvarea problemei propuse folosește *heapq*, un modul Python ce permite implementarea structurilor de date *heap*, ce permit efectuarea unor operații de bază asupra cozilor cu prioritate. Structura *heap* este o structură de date în care fiecare cheie este mai mare decât cheile de pe două poziții determinate, adică:

$$\begin{aligned} \text{heap}[k] &\leq \text{heap}[2*k+1] \\ \text{heap}[k] &\leq \text{heap}[2*k+2] \end{aligned}$$

pentru orice  $k$  pornind de la 0. O proprietate interesantă a unei structuri *heap* este faptul că cel mai mic element este întotdeauna *heap*[0].

În cazul algoritmului de codare Huffman binar prezentat mai sus, fiecare pas al implementării presupune (1) extragerea din *heap* a două noduri ce au probabilitățile cele mai mici (folosind metoda *heappop(heap)*) (2) gruparea lor într-un nod restrâns având probabilitatea egală cu suma celor două probabilități corespunzătoare nodurilor extrase, și (3) inserarea nodului restrâns înapoi în coada cu priorități (folosind metoda *heappush(heap, item)*) pentru a participa la pașii următori de codare. Restrângerea se finalizează în momentul în care în coadă rămâne un singur simbol restrâns. În Python, scrierea *n1.symbol + n2.symbol* reprezintă concatenarea a două elemente de tip *string*. Pentru a înțelege modul de creare al arborelui, se poate urmări crearea sa folosind instrucțiunea *print*. În ceea ce privește nodul, acesta conține în structura sa atât simbolul, cât și probabilitate asociată. Metoda *\_\_lt\_\_(self, other)* este suprascrisă pentru a permite compararea directă a nodurilor prin valorile probabilităților aferente. Metoda *\_\_repr\_\_(self)* întoarce reprezentarea obiectului într-o anumită formă, de obicei un *string* ce conține o reprezentare a obiectului ce poate fi afișată în consolă.

Funcția *encode* începe codarea de la rădăcina arborelui ce are asociat drept cod un *string* gol *''*. Pornind de la rădăcină către frunzele arborelui, prin frunze înțelegând acele noduri ce nu au descendenți, se va alocă câte un bit de 0 și 1 pe fiecare ramură. În final, codul asociat sursei va fi reținut în *dic\_code*.

## 4.4 Exerciții

Pornind de la algoritmul prezentat în platformă, rezolvați următoarele exerciții.

1. Implementați funcția:

$$p1(S: List[str], P: List[float]) \rightarrow List[float]$$

Funcția primește ca parametri:

- S - o listă de simboluri
- P - o listă de probabilități pentru fiecare simbol.

Se garantează:

- toate simbolurile sunt diferite între ele;
- cele 2 liste au aceeași lungime;
- suma tuturor probabilităților este 1.

Funcția codează simbolurile folosind codarea Huffman binară și returnează, într-o listă, următoarele 3 valori în această ordine:

- Lungimea medie a cuvântului de cod (calculată după formula (4.4));
- Eficiența codului (calculată după formula (4.7));
- Valoarea lui K pentru inegalitatea Kraft-McMillan (calculată după formula (4.9)).

2. Implementați funcția:

$$p2(S: List[str], P: List[float]) \rightarrow float$$

Funcția primește ca parametri:

- S - o listă de simboluri.
- P - o listă de probabilități pentru fiecare simbol.

Se garantează:

- toate simbolurile sunt diferite între ele;
- cele 2 liste au aceeași lungime;
- suma tuturor probabilităților este 1.

Funcția codează simbolurile folosind codarea Huffman ternară și returnează lungimea medie a cuvântului de cod (calculată după formula (4.4)).

Sfaturi:

- Modificați clasa Node astfel încât să permită 3 copii (descendenți) în loc de 2.
- Modificați funcția HuffmanTree astfel încât să țină cont de relația (4.10) și să adauge un nou simbol ( ' ' ) cu probabilitatea 0.0 atunci când relația nu este îndeplinită.
- Modificați funcția HuffmanTree astfel încât să realizeze suma ultimelor 3 noduri cu cele mai mici probabilități în loc de suma ultimelor 2.
- Modificați funcția encode astfel încât să țină cont că un nod poate avea 3 copii (descendenți) în loc de 2.
- Modificați funcția encode astfel încât să țină cont că poate exista și simbolul '2' pentru codare.

### 3. Implementați funcția:

$$p_4(\text{text}: \text{str}) \rightarrow \text{List}[\text{int}]$$

Funcția primește ca parametru un text.

Funcția returnează o listă de 2 întregi: numărul de biți necesari codării întregului text folosind codarea Huffman binară precum și numărul de biți necesari codării întregului text dacă fiecare caracter se codează pe  $\lceil \log_2(n) \rceil$  biți.