

## Metrici pentru compresia fișierelor

Prof. dr. Ion IVAN, ec. Daniel VERNIȘ, ec. Petrișor OPREA  
 Catedra de Informatică Economică, A.S.E., București

*Pentru a se putea compara algoritmii de compresie, trebuie definit un sistem de indici pentru măsurarea performanțelor acestora. În această lucrare sunt prezentate principali indicatori folosiți în compresia de date.*

**Cuvinte cheie :** metricile fișierelor, compresie de date, entropia.

### 1. Introducere

Compresia de date este de multe ori privită ca un exercițiu de reducere a redundanței în fișier. De fapt, reprezintă cu mult mai mult. Compresia de date rezolvă una din problemele majore ale teoriei informației : cum se poate codifica în mod eficient un mesaj, din punct de vedere al numărului de biți pe simbol. Când un mesaj a fost exprimat în cât mai puțini biți pe simbol putem spune că a avut loc o compresie optimală. Nu sunt biți folosiți inefficient. Acesta este scopul compresiei de date.

În literatura de specialitate, sunt folosiți deseori termenii de compresie și compactare. Între acești doi termeni există o diferență.

Prin definiție, **compactarea datelor** este o tehnică de reducere a lungimii reprezentărilor fizice a datelor, cu păstrarea numai a unui subset din setul de date inițial. Acest subset poartă denumirea de **informație relevantă**.

Prin **compresie de date** se înțelege un procedeu de compactare care este complet reversibil. Adică din fișierul compresat, se poate obține prin aplicarea decompresiei fișierul inițial. În continuare, F este fișierul inițial iar F' este fișierul compresat.

Deci, dacă la compresie, fișierul se poate recupera integral după compresie, la compactare acest lucru nu mai este valabil. De exemplu, prin compactarea unui text se poate înțelege eliminarea spațiilor nesemnificative. Pentru

măsurarea lungimii fișierelor, în compresia de date, se folosesc două tipuri de lungimi:

- lungimea fișierului ca număr de simboluri,  $L(F)$  ;
- lungimea fișierului ca număr de biți,  $I(F) = \sum I(x_i)$ , unde  $x_i$  reprezintă simbolul din fișierul F de pe poziția i.

### 2. Indicatori de performanță a algoritmilor

Performanța unui algoritm de compresie este privită sub trei aspecte semnificative :

- modul în care se face compresia fiecarui fișier în parte ;
- comportamentul statistic global al produsului care implementează algoritmul ;
- resursele necesare pentru realizarea compresiei ;

#### Gradul de compresie

O expresie cantitativă a puterii algoritmului de a reduce lungimea  $I(F')$  față de  $I(F)$ , dată prin :

$$g = (I(F) - I(F')) / I(F)$$

#### Rata de compresie

Reprezintă ponderea pe care o are fișierul compresat în fișierul inițial.

$$r = I(F') / I(F)$$

#### Lungimea minimală

Reprezintă lungimea sub care fișierul compresat devine mai lung decât fișierul inițial. Se identifică deci

situația limită cînd algoritmul de compresie își pierde eficiența.

Pentru coduri de compresie cu lungime fixă, se consideră :

- $n$  numărul de simboluri distincte din fișier;
- $k$  lungimea codului ;
- $m$  lungimea fișierului ca număr de simboluri ;

Presupunem ca această lungime se memorează pe patru biți .

$h(a)$  - lungimea codului inițial, avînd valoarea 8 ;

Din inegalitatea:

$$n \cdot h(a) + n \cdot k + 4 + I(F') < I(F)$$

rezultă:

$$n \cdot (8 + k) + 4 + m \cdot k < m \cdot 8$$

și se obține lungimea minimală:

$$l_{\min} = \frac{n \cdot (8 + k) + 4}{8 - k} < m.$$

Dacă se consideră coduri de lungime variabilă asociate simbolurilor, se calculează o lungime medie a codurilor  $\bar{k}$  și se determină intervalul  $[\bar{k} - 3\sigma, \bar{k} + 3\sigma]$ . Vor rezulta valori limită precum și un nivel mediu de încadrare a lungimii peste care compresia devine eficientă.

Lungimea minimală se poate obține în acest caz prin tehnici de simulare.

### 3. Entropia fișierelor

Pentru fișiere, entropia reprezintă numărul de biți pe simbol ce s-ar obține printr-o compresie optimală a mesajului.

Este foarte ușor să se determine, din acest moment, și gradul de redundanță în fișier. Dacă un fișier de date se reprezintă în cod ASCII, adică 8 biți pe simbol, se poate obține gradul de redundanță în fișier astfel:

$$R_F = 8 - H$$

măsurată în biți pe simbol, unde  $H$  se determină prin relația lui Shannon [1] astfel:

$$H = -\sum f_i \log(f_i)$$

Entropia poate fi considerată o metrică pentru fișierele de date, datorită proprietăților :

$$H(X) \geq 0 ; H(X) \leq -\sum(1/n)\log_2(1/n).$$

unde  $1/n$  reprezintă probabilitatea de apariție a simbolului  $x_i$  în condiții de echiprobabilitate.

Considerăm două alfabete  $A = \{x_1, x_2, \dots, x_n\}$  și  $B = \{y_1, y_2, \dots, y_p\}$  și  $A \cap B = \emptyset$ . Se construiesc pe cele două alfabete multimi de fișiere  $F_A, F_B$ .

Fie  $F_1 \in F_A$  și  $F_2 \in F_B$  cărora li se asociază cîmpurile de probabilitate  $X_1$  respectiv  $X_2$ . În aceste condiții

$$H(X_1, X_2) = H(X_1) + H(X_2).$$

Dacă alfabetele  $A$  și  $B$  sunt oarecare (nu mai au proprietatea de independentă )

$$H(X_1, X_2) = H(X_1) + H(X_2/X_1).$$

unde  $H(X_2/X_1)$  este entropia alfabetului sursă  $B$  condiționat de apariția simbolurilor din alfabetul  $A$ . Deci  $H(X_1, X_2) \leq H(X_1) + H(X_2)$

$$\begin{aligned} H(X_1/X_2) &= H(X_2/X_1) + H(X_1) - H(X_2) \\ 0 &\leq H(X_1) - H(X_1/X_2) \leq H(X_1) \end{aligned}$$

Aceste proprietăți permit efectuarea de cuantificări pe fișiere și o ierarhizare a acestora în raport cu gradul de nedeterminare.

De exemplu, se consideră fișierele :

$$F_1 = \{a, a, a, a, b, b, b, b, c, c, c, c, d, d, d, d\}$$

$$F_2 = \{a, a, b, b, b, b, c, c, c, c, c, c, c, c, d, d\}$$

$$F_3 = \{a, b, b, b, b, c, c, c, c, d, d, d, d, d, d, d\}$$

În care  $Ig(F_1) = Ig(F_2) = Ig(F_3)$ , unde  $Ig()$  este funcția de lungime a unui fișier care indică numărul de simboluri conținute.

Fișierele  $F_1, F_2, F_3$  sunt definite pe alfabetul  $A = \{a, b, c, d\}$ . Frecvențele de apariție a simbolurilor în cele 3 fișiere sunt date în tabelul 1. :

Tabelul 1.

Simbol	F1	F2	F3
a	4	2	1
b	4	4	4
c	4	8	4
d	4	2	7

Probabilitățile, ca frecvențe relative de apariție sunt date în tabelul 2.

Tabelul 2.

Simbolul	F1	F2	F3
a	1/4	1/8	1/16
b	1/4	1/4	1/4
c	1/4	1/2	1/4
d	1/4	1/8	7/16
Entropia	1.999	1.749	1.771

Prin concatenarea fișierelor  $F_1$ ,  $F_2$ ,  $F_3$  rezultă fișierul  $F_C$  caracterizat prin valorile date în tabelul 3.

Tabelul 3.

Simbolul	Frecvență	$f_i$
a	7	7/48
b	12	12/48
c	16	16/48
d	13	13/48
Entropia	-	1.943

Omogenitatea fișierelor în raport cu distribuția de probabilitate a simbolurilor poate fi considerată un criteriu de alegere a unui algoritm de compresie de date.

Se definește **indicatorul normat al entropiei** fișierului  $F$ , fișier definit pe un alfabet cu  $n$  simboluri.

$$I_H = \frac{H(F)}{\lg_2 n}$$

În continuare se prezintă programul **entropia.cpp**, care calculează entropia fișierelor și gradul de redundanță a acestora.

#### Calculul entropiei fișierelor - entropia.cpp

```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#define LOG(x) 3.32*log10(x)
#define ENTROPIA(x) -(x*LOG(x))
FILE *input;
unsigned int TabelaContor[256];
void CitireFisier(void);
double Analiza(void);
void Utilizare (void);
main(int argc,char *argv[])
{double Rezultat;
if (argc==1) Utilizare();
input=fopen(argv[1],"rb");
if (!input) printf("\nFisier de intrare inexistent.");
if (!input) exit(-1);
printf("\n Citire fisier ...");
CitireFisier();
printf("\n Analiza date ...");
Rezultat=Analiza();
printf("\nFisierul \"%s\" are o entropie de %3.3f
biti pe bait.",argv[1],Rezultat);
printf("\nEste posibil sa se obtina o rata medie
de compresie");
```

```
printf("\nfolosind algoritmul Huffman standard
de : %2.2f%%\n\n",100 - (Rezultat*100)/8);
fclose(input);
return(1);
}
void CitireFisier()
{int ch;
while((ch=getc(input))!=EOF)
TabelaContor[ch]++;
}
double Analiza()
{double accum=0.0;
double freq;
long fsize=0L;
unsigned int z;
fsize=fseek(input);
for(z=0;z<256;z++)
if (TabelaContor[z])
{
freq=(double)
TabelaContor[z]/fsize;
accum+=(double)
ENTROPIA(freq);
}
return accum;
}
void Utilizare()
{printf("\n\n\n");
printf(" 'Entropia fisierelor' v1.0 by Daniel
Vernis\n");
printf(" Sintaxa : entropia [fisier]\n\n");
printf(" Acest program calculeaza entropia
fisierelor\n");
printf(" folosind ecuația entropiei introdusa de
Shannon.\n"); exit(1);
}
```

#### 4. Erori ale compresiei de date

Transmisia datelor folosind canalele discrete fără zgomot este din nefericire, un model nu prea realistic în sistemele de comunicație. Actualele sisteme de transmisii de date, sunt inclinate spre două tipuri de erori :

- **erori de fază**, în care codul unui simbol este pierdut sau adăugat ;
- **erori de amplitudine**, în care codul unui simbol este corupt .

Gradul în care erorile de canal degradează transmisia este un parametru important în alegerea metodei de compresie a datelor.

Susceptibilitatea la erori a algoritmilor de compresie, depinde în mare măsură de tipul algoritmului : static sau dinamic.

### Algoritmi statici

Este în general cunoscut că algoritmii de tip Huffman, au proprietatea de autocorecție, prin care erorile de transmisie tind să nu se propage mai departe. Într-un algoritm static, prin sincronizare se înțelege identificarea în același mod, la emisie și la recepție, a cuvîntului transmîs.

În continuare este prezentat un exemplu pentru ilustrarea abilității algoritmului Huffman standard de recuperare a erorii de fază. Se folosește mesajul **bcdaeb**, codificat după algoritmul Huffman astfel :

<b>a</b>	<b>1</b>
<b>b</b>	<b>011</b>
<b>c</b>	<b>010</b>
<b>d</b>	<b>001</b>
<b>e</b>	<b>000</b>

Șirul va fi codificat **0110100011000011**. Figura 1 demonstrează impactul pierderii primului bit, celui de-al doilea sau celui de-al patrulea. Punctele delimităză codurile simbolurilor interpretate.

011.010.001.1.000.011.	codificarea <b>bcdaeb</b>
<del>1</del> 0.001.1.000.011.	bitul 1 pierdut, interpretare <b>aacdaeb</b>
010.1.000.1.1.000.011.	bitul 2 pierdut, interpretare <b>caeeaaeb</b>
011.1.000.1.1.000.011.	bitul 4 pierdut, interpretare <b>baeaaeb</b>

Fig. 1. Recuperarea erorilor de fază

În cazul pierderii primului bit, are loc resincronizarea bupă cel de-al treilea bit. În sir s-a înlocuit simbolul **b** cu grupul **aa**. Cînd al doilea bit este pierdut, primii 8 biți ai sirului sunt greșit interpretați. Resincronizarea se produce la bitul nouă. În cazul pierderii bitului 4, au loc aceleasi efecte ca în cazul pierderii bitului al doilea. Însă, primul simbol al sirului, **b** este citit corect. Efectul erorilor de amplitudine asupra sirului **bcdaeb** este prezentat în figura 2. La erorile de amplitudine, biții 1, 2 și 4 sunt schimbați, nu pier-

duți ca în cazul erorilor de fază. Astfel, bitul 1 este transformat din **0** în **1**, bitul 2 din **1** în **0**, iar bitul 4 din **0** în **1**.

011.010.001.1.000.011.	
	codificarea <b>bcdaeb</b>
1.1.1.010.001.1.000.011.	bitul 1 alterat, interpretare <b>aaacdaeb</b>
001.010.001.1.000.011.	bitul 2 alterat, interpretare <b>dcdaeb</b>
011.1.1.000.1.1.000.011.	bitul 4 alterat, interpretare <b>baeaaeb</b>

Fig. 2. Recuperarea erorilor de amplitudine

Autosincronizarea poate fi discutată mai bine prin următoarele definiții.

Dacă s este un sufix al unei codificări și există secvențele de cod  $\Gamma$  și  $\Delta$ , putem spune:  $s\Gamma=\Delta$ , unde  $\Gamma$  desemnează secvența de sincronizare a lui s. De exemplu, pentru sirul prezentat anterior, **1** este secvență de sincronizare pentru **01**, în vreme ce **000001** și **011** sunt secvențe de sincronizare pentru sufixul **10**. Dacă fiecare sufix are o secvență de sincronizare, atunci codul este **complet autosincronizat**. Dacă nici unul sau doar cîteva din suficele indicate au secvențe de sincronizare atunci codul este **nesincronizat** sau **partial sincronizat**.

De asemenea, orice cod prefix care este complet autosincronizat se va sincroniza cu probabilitate 1, dacă sirul sursă constă într-o succesiune de mesaje independente.

### Algoritmi dinamici

Algoritmii dinamici sunt de departe cei mai afectați de erorile de transmisie. De exemplu, în cazul algoritmului Huffman adaptiv, deși receptorul se poate resincroniza cu emițătorul, în scopul localizării corecte a unui început de codificare, informația pierdută reprezintă cîțiva biți sau cîteva simboluri.

Din această cauză, la recepție vor fi redefinite greșit codificările simbolu-

rilor din punctul întâlnirii erorii, fiind compromis întreg sirul.

Lempel și Ziv recunosc că principalul dezavantaj al algoritmului LZ78 este susceptibilitatea la eroare.

## 5. Concluzii

În literatura de specialitate recentă, entropia este folosită drept criteriu de neomogenitate a fișierelor. Analiza entropică permite stabilirea apartenenței la o anumită clasă de entropie și în acest fel se va asocia direct componența program corespunzătoare algoritmului care realizează compresie performantă. Se fundamentează construirea metricilor pentru compresie care ține seama de neomogenitatea fișierelor. De asemenea, este necesar studiul corelației dintre factorii care influențează atât compresia cât și decompresia datelor.

## Bibliografie

1. C.Shannon, W.Weaver, Entropy, The Mathematical Theory of Communication, Urbana, Illinois, University of Illinois 1949
2. Abrams, S. Tren, A methodology for interactive computer service measurement, Communication of ACM, vol. 20, nr.12, Decembrie 1977.
3. I.Ivan, D.Verniș, Analiza comparată a algoritmilor de compresie date, PC World, nr.12, Decembrie 1995.
4. I.Ivan, M. Popescu, Metrici software, Byte, nr.5, Mai 1996.
5. I.Gh.Roșca, C.Apostol, B. Ghilic-Micu, Prelucrarea fișierelor în Pascal, Editura Tehnică, București 1994.