

Laborator PS

Algoritmi de compresie a datelor

Lucrarea 1: Introducere în Compresia Datelor

Lucrarea 2: Algoritmul Shannon-Fano

Lucrarea 3: Algoritmul Huffman static

Lucrarea 4: Algoritmul Huffman dinamic

Prof. Dan Ștefănoiu <dan.stefanoiu@acse.pub.ro>

Ș.I. Alexandru Dumitrașcu <alexandru.dumitrascu@acse.pub.ro>

A. Obiective

Introducerea principalelor concepte din domeniul Compresiei Datelor. Însușirea unui mic vocabular de termeni specifici, care va fi folosit și în lucrările următoare.

B. Semnale utilizate

Pentru a testa algoritmi din acest capitol, sunt disponibile trei categorii de fișiere de date:

- **Documente.** Din această categorie fac parte fișierele:

- **Documentatie_UMAPID.doc;**
- **Documentatie_UMAPID.pdf;**
- **Prefata_Undine.txt;**
- **show_audio.m** (sursă MATLAB);
- **Y04.m** (serie de date în format ASCII).

Acestea pot fi extrase din arhiva **Documents.rar**.

- **Semnale audio reale.** Arhiva **Audio.rar** include 3 semnale înregistrate în formatul **wave**:

- **instr_01.wav** – sunet emis de un instrument muzical;
- **sound_01.wav** – sunet fabricat artificial;
- **speech_01.wav** – semnal *vocal* (emis de vocea umană).

- **Executabile.** Fișierele din această categorie sunt incluse în arhiva **Executables.rar** și conțin cod executabil sau date împachetate în formatul **mat** din MATLAB. Ele sunt următoarele:

- **KARMA_DATA482#1_5_V7.mat;**
- **quartz.dll;**
- **WinRar.exe;**
- **WINZIP32.EXE.**

C. Suport teoretic

Compresia Datelor (CoDa), este un domeniu de cercetare care face parte din câmpul Teoriei Informației – ramură a Matematicii născută către sfârșitul anilor '40, odată cu lucrările publicate de către cercetătorul american C.E. Shannon (de la Bell Laboratories) și matematicianul ungar D. Gabor (unul dintre pionierii Teoriei Comunicației, laureat al premiului Nobel). Aceste lucrări prezintă o serie de soluții (mai mult de ordin teoretic) ale unor probleme legate de emisia, culegerea, prelucrarea, transmisia, recepția și interpretarea informației cu ajutorul unui mijloc automat de calcul. În această lucrare, sunt prezentate principalele concepte care definesc domeniul CoDa, fără a intra în detalii de natură teoretică.



Concepte de bază

Obiectul de studiu al CoDa îl constituie metodele și tehnicile destinate reducerii volumului unui set de date reprezentate sub formă binară, cu posibilitatea reconstituirii acestuia fie exact, fie aproximativ (conform unor criterii definite a priori). Posibilitatea de a reduce volumul unui set de date se bazează pe câteva fenomene importante legate de reprezentarea binară a informației.

Astfel, setul de date este privit ca un șir finit de *simboli*. Acesta este un termen consacrat (nu foarte diferit de termenul colocvial de *simboluri*). Fiecare dintre simbolii are o reprezentare numerică unică numită *cod (numeric) original*. Numărul simbolilor ce compun un set de date, adică **volumul** acestuia, este, în general, mult mai mare decât numărul de simbolii distincți ce participă la construcția sa. **De regulă, un simbol este reprezentat în baza 2, pe 8 biți (1 octet), astfel că setul de date se construiește folosind maximum $2^8 = 256$ simbolii distincți.** Mulțimea simbolilor distincți dintr-un set de date se numește **alfabet**, denumire sugerată de conceptul uzual de *alfabet literar*, din vorbirea colocvială. Alfabetul din CoDa include literele curente (A, B, C etc.), dar nu se reduce doar la acestea; el poate conține și alți simbolii (cifre numerice, semne de punctuație, semne speciale care nu au o interpretare directă etc.). Reprezentarea simbolilor pe 8 biți este realizată

conform unui standard internațional numit *ASCII*. Tabela ASCII pune în corespondență biunivocă mulțimea celor 256 de simboluri cu mulțimea celor 256 de coduri numerice (de la 0 la 255). De exemplu, codul (numeric) original asociat literei 'A' este 65, în bază zecimală. Acest cod poate fi exprimat prin **0x41** în baza 16 (prin convenție, numerele precedate de **0x** sunt reprezentate în baza 16 sau în *cod hexazecimal*) sau prin **0100 0001** în baza 2 sau în cod binar.

Un al doilea fenomen important al reprezentării binare a informațiilor este cel al frecvenței diferite de apariție a simbolurilor alfabetului de-a lungul setului de date. Nu toți simbolii participă în mod uniform la construcția setului de date: unii sunt mai frecvenți decât alții. Acest fapt a sugerat ideea că **alfabetul ar putea fi recodificat cu ajutorul unor coduri de lungime variabilă**, în așa fel încât simbolurilor ce apar frecvent să li se aloce coduri mai scurte decât simbolurilor ce apar mai rar în setul de date.



Rată de compresie, factor de compresie

Un prim parametru care privește performanțele metodei de compresie utilizate este ***rata de compresie***. Aceasta se notează cu ρ și are următoarea exprimare:

$$\rho = \left[1 - \frac{\#dc}{\#do} \right] \times 100 [\%], \quad (1.1)$$

unde: $\#dc$ este lungimea setului de date comprimat, iar $\#do$ este lungimea setului de date original.

Această definiție arată că:

- dacă $\rho = 100 \%$, atunci metoda de compresie aferentă este *ideală*, câștigul obținut prin compresie fiind maxim; în practică, însă, $\rho < 100 \%$;

- dacă $\rho = 0 \%$, atunci metoda nu realizează nici o compresie efectivă, câștigul fiind nul; în practică, se spune, în acest caz, că setul de date a fost *povestit (stored)*, adică redat ca atare, fără a suferi o compresie efectivă;
- dacă $\rho < 0 \%$, atunci metoda aferentă nu realizează o compresie efectivă, ci o *expandare* a datelor (adică o reprezentare a acestora cu mai mulți biți decât originalul), fiind ineficientă.

Termenul de *rată de compresie* este adesea înlocuit de cel de ***factor de compresie***, notat cu φ . Acesta se definește ca raportul procentual dintre lungimea setului de date comprimat și cea a setului de date original:

$$\varphi = 1 - \rho = \frac{\#dc}{\#do} \times 100 [\%]. \quad (1.2)$$

Interpretarea sa este considerată mai directă decât cea a ratei de compresie. *O metodă de compresie este cu atât mai eficientă cu cât factorul de compresie este mai mic.*



Entropie, cantitate de informație, redundanță

Un concept care leagă CoDa de Teoria Informației este cel de ***redundanță***. Înainte de a-l defini, trebuie subliniat faptul că obiectivul principal al unei metode de compresie a datelor este detectarea, localizarea și eliminarea *redundanțelor* (adică a repetabilității informațiilor) din setul de date.

O măsură a *redundanței* unui set de date nu este dată de rata de compresie (așa cum am fi tentați să credem), deoarece rate diferite de compresie se pot obține pentru același set de date, dacă se apelează la metode de compresie diferite. *Redundanța* este o constantă a setului de date și nu depinde de metoda de compresie utilizată.

Volumul unui set de date este numărul simbolilor (nu neapărat distincți) care îl compun. **Lungimea** unui set de date este numărul total de biți pe care este el reprezentat. Cele două noțiuni sunt diferite.

Vom nota prin $v(s)$ **frecvența de apariție a simbolului** s . Ea constituie o estimatie a *probabilității de apariție* a simbolului s , notată cu $\mathcal{P}(s)$. Plecând de la $v(s)$, se poate calcula numărul minim de biți suficient pentru a recodifica simbolul s , fără a confuzie cu alți simbolii. Acest număr se numește **entropie** și poate fi fracționar (nu neapărat întreg). Practic, însă, el va fi adesea rotunjit la întregul superior, din motive lesne de înțeles. Entropia va fi cel mult egală cu numărul de biți ai codului original (de exemplu, 8). Diferența dintre ele constituie o măsură a *redundanței* intrinseci transportate de acel simbol în cadrul setului de date (adică numărul de biți inutili prezenți în codul original).

Entropia unui simbol s dintr-un set de date (adică *entropia informatică*) se notează prin $\mathcal{H}(s)$ și se definește astfel:

$$\mathcal{H}(s) = -\log_2 v(s). \quad (1.3)$$

În definiția (1.3) se observă imediat că: $v(s) \in (0,1]$, ceea ce arată că $\mathcal{H}(s) > 0$. În consecință:

$$v(s) = 2^{-\mathcal{H}(s)} \in [0,1]. \quad (1.4)$$

Dacă, de exemplu, litera 'a' apare de 38 de ori într-un text cu volumul de 100 de caractere, atunci $v('a') = 0.38$, ceea ce conduce la entropia informatică $\mathcal{H}('a') = -\log_2 0.38 \cong 1.4$ biți. Rezultă că, în codarea inițială a lui 'a', 6.6 biți (din cei 8 alocăți) sunt inutili, adică *redundanți*. Practic, va fi suficient ca literei 'a' să i se aloce doar 2 biți, pentru a obține o primă compresie a textului.

Un alt concept direct legat de frecvența de apariție a unui simbol s este **cantitatea de informație**, notată cu $i(s)$. Ea este invers proporțională cu frecvența $v(s)$:

$$i(s) = \frac{C}{v(s)}, \quad (1.5)$$

unde $C > 0$ este o constantă de proporționalitate.

Un simbol ce apare frecvent în setul de date transportă o cantitate de informație mică și va primi un număr mic de biți de recodificare. Acest fenomen este mai bine ilustrat de următoarea relație evidentă (obținută prin combinarea definițiilor (1.3) și (1.5)):

$$\mathcal{H}(s) = -\log_2 C + \log_2 i(s), \quad (1.6)$$

în care se poate renunța la termenul constant $-\log_2 C$.

Se notează, mai departe, prin \mathcal{D} setul de date asupra căruia se operează (mulțimea tuturor simbolilor (nu neapărat distincți) ce alcătuiesc setul de date). Reamintim că *alfabetul* setului de date este mulțimea tuturor simbolilor *distincți* ce alcătuiesc setul de date și, în general, diferă de mulțimea \mathcal{D} . *Redundanța* setului de date \mathcal{D} se măsoară prin numărul total de biți utilizați, dar inutili, calculat astfel:

$$\mathcal{R}(\mathcal{D}) = \sum_{s \in \mathcal{D}} [8 - \mathcal{H}(s)]. \quad (1.7)$$

(Am considerat că fiecărui simbol s s-a alocat inițial un cod pe 8 biți.) Redundanța depinde direct de entropiile tuturor simbolilor mesajului. Ea este cu atât mai pronunțată cu cât aceste entropii sunt mai mici. Dacă ținem cont de definiția (1.3) a entropiei, atunci redundanța (1.7) se poate exprima și astfel (în funcție de frecvențele de apariție):

$$\mathcal{R}(\mathcal{D}) = \sum_{s \in \mathcal{D}} [8 - \log_2 v(s)]. \quad (1.8)$$

Astfel, frecvențe mari de apariție ale simbolilor setului de date conduc la o redundanță importantă a acestuia.

Relația (1.8) poate fi utilizată și pentru evaluarea redundanței unor subseturi de date din setul original. De exemplu, șirul de 5 caractere 'aaaaa' (numit ad hoc *cuvânt*) este reprezentat inițial pe un număr de $5 \times 8 = 40$ de biți. Dacă el face parte dintr-un set de date mai larg, în care $v('a') = 1/16$, atunci: $\mathcal{H}('a') = -\log_2 1/16 \cong 4$ și $\mathcal{R}('aaaaa') = 40 - 5 \times 4 = 20$ biți. Rezultă că, în cuvântul 'aaaaa', doar 20 de biți de reprezentare sunt utili, adică 50%. Mai mult, acest cuvânt ar putea fi codificat în mod ideal cu numai 7 biți (în loc de 20), reprezentând codul literei 'a' (pe 4 biți) și numărul de apariții consecutive (5, număr care ocupă 3 biți: 101). Cu toate acestea, o reprezentare pe un număr total de biți scăzut este dificil de realizat, căci, în procesul de decompresie, o serie de informații auxiliare trebuie de asemenea să se regăsească în setul de date comprimat. De exemplu, este necesar să se cunoască numărul de biți alocați simbolului 'a', poziția cuvântului 'aaaaa' în setul de date etc. Ori, această informație poate anula câștigul obținut prin calculul entropiei și redundanței.

D. Ghid C++

În continuare, sunt prezentate câteva funcții standard ce pot fi folosite în conceperea programelor C++. Pentru programele MATLAB, funcțiile sunt similare și pot fi interogate cu ajutorul facilităților **Help**, **Doc**, etc.

Pentru a deschide un fișier, se poate utiliza apelul:

```
FILE * fopen ( const char * filename, const char * mode );
```

unde: ieșirea funcției este pointerul către o variabilă de tip **FILE**, **filename** este numele complet al fișierului, scris sub forma unui șir de caractere, iar **mode** este modul de deschidere (unul sau mai multe caractere). Un fișier poate fi deschis în mai multe moduri, dar cele mai frecvente sunt:

- **'r'** = modul *citire* (read); nu se poate înscrie nimic în fișierul deschis;
- **'w'** = modul *scriere* (write); nu se poate citi nimic din fișierul deschis; scrierea se face peste datele deja existente în fișier (dacă există), de la prima poziție a acestuia; modul scriere permite deschiderea un nou fișier, care nu a mai fost creat;
- **'a'** = modul *adăugare* (append); dacă fișierul nu există, el va fi creat și deschis în modul scriere; altfel, el va fi deschis astfel încât la datele deja existente să se poată adăuga altele noi;
- **'r+'** = modul *citire-scriere*; pentru acest mod, fișierul trebuie să fi fost creat anterior, altfel el nu poate fi creat; odată deschis, fișierul poate fi citit și/sau supra-înscriș (adică informația nouă o supra-înscrie pe cea veche);
- **'a+'** = modul *citire-adăugare*; dacă fișierul nu a fost creat anterior, el va fi creat acum și deschis; odată deschis, fișierul poate fi citit și/sau extins (adică informația nouă se adaugă la cea veche).

Dacă funcția returnează pointerul **NULL**, utilizatorul este informat că fișierul nu a putut fi deschis, din diferite cauze (cel mai frecvent, pentru că nu a putut fi găsit pe disc în vreuna dintre locațiile de căutare). Pentru utilizarea acestei funcții, precum și a celorlalte prezentate mai jos, este nevoie de introducerea bibliotecii asociate antetului **stdio.h** (*standard input-output functions*).

Utilizatorul poate accesa informația dintr-un fișier deschis în diferite moduri. Există un pointer virtual care indică locația din fișier ce poate fi citită sau înscrisă. După citirea/scrierea fiecărei date, pointerul se incrementează automat. Utilizatorul poate totuși repoziționa acest pointer astfel:

```
int fseek ( FILE * stream, long int offset, int origin );
```

unde:

- **stream** este indicatorul de fișier, așa cum a fost el returnat de către funcția **fopen**;
- **offset** este un întreg care arată că pointerul trebuie deplasat cu valoarea absolută a acestuia către sfârșitul fișierului (dacă e strict pozitiv), către începutul fișierului (dacă e strict negativ) sau deloc (dacă este nul);
- **origin** este o etichetă de tip șir de caractere, care poate avea una dintre următoarele forme: **SEEK_SET**, **SEEK_CUR** sau **SEEK_END**; ea arată de unde trebuie să plece pointerul: **SEEK_SET** – de la începutul fișierului, **SEEK_CUR** – de la poziția curentă în fișier, respectiv **SEEK_END** – de la sfârșitul fișierului.

Funcția **fseek** returnează un indicator de stare, care arată dacă repoziționarea s-a efectuat cu succes (**0**) sau nu (o valoare diferită de 0).

Dacă poziția indicată de pointerul virtual nu (mai) este cunoscută, ea poate fi (re)găsită astfel:

```
long int ftell ( FILE * stream );
```

Argumentul de ieșire este un întreg care măsoară poziția în octeți, începând cu numărul 0 (începutul de fișier). Dacă există o eroare, este returnată valoarea **-1L**, iar variabila globală **errno** este setată cu o valoare pozitivă. Această valoare poate fi interpretată de funcția **perror**.

Citirea de date dintr-un fișier deschis se poate realiza, în general, cu ajutorul funcției **fread**. Declarația uzuală a acesteia este următoarea:

```
size_t fread ( void * ptr, size_t size, size_t count, FILE * stream );
```

unde:

- **ptr** este un pointer către un bloc de memorie cu dimensiunea minimă de (**size*count**) octeți;
- **size** este dimensiunea în octeți a fiecărui element de citit;
- **count** este un întreg care arată câte date de dimensiune **size** trebuie citite;
- **stream** este indicatorul de fișier (returnat de către funcția **fopen**).

Biblioteca asociată antetului **stdio.h** conține o colecție de funcții importante ce pot fi folosite în realizarea unui program C++. Lista completă de funcții este prezentată mai jos. Pentru mai multe informații cu privire la utilizarea lor se poate accesa pagina web :

<http://www.cplusplus.com/reference/clibrary/cstdio/>

- | | | |
|-------------------|------------------|-------------------|
| • clearerr | • fscanf | • rename |
| • fclose | • fseek | • rewind |
| • feof | • fsetpos | • scanf |
| • ferror | • ftell | • setbuf |
| • fflush | • fwrite | • setvbuf |
| • fgetc | • getc | • sprintf |
| • fgetpos | • getchar | • sscanf |
| • fgets | • gets | • tmpfile |
| • fopen | • perror | • tmpnam |
| • fprintf | • printf | • ungetc |
| • fputc | • putc | • vfprintf |
| • fputs | • putchar | • vprintf |
| • fread | • puts | • vsprintf |
| • freopen | • remove | |

~~E. Sarcini de lucru~~

~~Tema 1 (Operarea cu fișiere, evaluarea dimensiunii acestora)~~

~~Scrieți un program prin care să determinați dimensiunea în octeți și în perechi de octeți pentru fiecare fișier din corpusul celor 12 descris în subsecțiunea de semnale utilizat. Comparați dimensiunea în octeți obținută cu cea afișată de exploratorul de fișiere (*Windows Explorer, Norton Commander, etc.*). Completați apoi programul cu un set de instrucțiuni care să măsoare numărul de biți ai fiecărui fișier. Se potrivește raportul dintre acest număr și cifra 8 cu cel de octeți? Argumentați riguros răspunsul.~~

~~Tema 2 (Histograme)~~

~~Scrieți un program prin care să citiți datele din fiecare fișier al corpusului sub formă de simbolii reprezentați întâi pe 8 biți, apoi pe 16 biți. După citire, evaluați și afișați cele două histogramme ale datelor fiecărui fișier. (Histograma se calculează prin detectarea numărului de apariții ale unui simbol. Dacă simbolul este reprezentat pe 8 biți, atunci alfabetul conține cel mult 256 de simbolii distincți. Lungimea setului de date fiind, de regulă, mai mare, este clar că unii dintre ei (dacă nu toți) sunt obligați să se repete. Un fenomen similar are loc și în cazul simbolilor reprezentați pe 16 biți, numai că, aici, frecvența de apariție a simbolilor este mai mică.) Observați cu atenție aceste 24 de histogramme și indicați 3 dintre ele care par apropiate de distribuția gaussiană. Care credeți că este explicația pentru acest tip de distribuție? Există vreo histogramă care sugerează că datele sunt distribuite uniform? Dacă da, pentru care seturi de date? De ce credeți că acestea sunt distribuite cavi uniform? Argumentați răspunsurile în manieră riguroasă.~~

Tema 3 (Entropii, redundanțe)

Utilizând programul de la tema anterioară, scrieți un nou program care evaluează și trasează grafic entropiile informatice ale simbolilor fiecărui set de date. Pe fiecare pereche de grafice (simboli de 8 biți și simbolii de 16 biți), indicați redundanța setului corespunzător de date. Efectuați o comparație între cele 24 de seturi de date privind redundanța. Enumerați seturile de date în ordinea descrescătoare a compresibilității. (Teoretic, cele mai compresibile sunt și cele mai redundante.)