

Dedicat Centenarului Gazetei Matematice (1985-1995)

ASUPRA UNEI METODE DE COMPRIMARE/DECOMPRIMARE
A FIȘIERELOR

1. Introducere

Problema de la care se pornește este următoarea: un beneficiar achiziționează un produs soft cu anumite caracteristici (număr de utilizări, tipuri de teste etc.). În timp el și-a epuizat total sau parțial aceste facilități, dar dorește utilizarea în continuare a programului. Este necesar în acest caz un upgrade pentru acest produs soft, iar acest upgrade trebuie să fie cât mai eficient posibil. Desigur că este posibil ca beneficiarului să i se livreze un nou pachet conținând programul cu noii parametri solicitați, dar această operație este costisitoare atât ca timp cât și ca material, și în plus este superfluă, programul propriu-zis este același. Într-o asemenea situație beneficiarul are nevoie doar de modificarea "contoarelor" unde sunt păstrate caracteristicile ce trebuiesc modificate. În cazul a numeroase produse soft aceste caracteristici sunt păstrate în cheia de protecție hard (hardlock), cheie pe care beneficiarul o primește în momentul achiziționării produsului (prima achiziționare). Revenind, pentru un upgrade de acest fel este nevoie doar de modificarea unor chei din hardlock, programul rămânând nemodificat. Utilizatorul trebuie să trimită distribuitorului programului doar conținutul cheii de protecție (sau o parte a sa), distribuitorul face modificările (de obicei incrementări) de rigoare și returnează noua cheie (noul conținut). Această operație trebuie făcută cu cheltuieli minime, atât din punctul de vedere al beneficiarului cât și al distribuitorului. Se

va studia în continuare modul în care se va face acest schimb prin modem sau prin telefon.

2. "Schimbul" de informație

Se va porni în continuare de la premisa că cheia hard are mărimea de 496 octeți. Schimbul prin telefon se va face în modul următor:

- 1) utilizatorul se hotărăște să facă upgrade-ul
- 2) selectează opțiunea respectivă din program (eventual poate folosi un program de sine stătător)
- 3) pe ecran apare conținutul hardlock-ului comprimat, criptat, codificat; se va numi acest conținut **șir de intrare**
- 4) utilizatorul telefonează distribuitorului căruia îi comunică șirul de intrare și noii parametri solicitați
- 5) distribuitorul lansează un program special care cere ca date de intrare șirul de intrare și noii parametri
- 6) distribuitorul introduce datele de intrare
- 7) programul returnează un **șir de ieșire** (tot comprimat, criptat), șir care va fi comunicat beneficiarului
- 8) beneficiarul tastează acest șir ca răspuns la opțiunea selectată
- 9) din acest moment programul va funcționa cu noile contoare.

2.1. Observații

- 1) Este de dorit ca atât șirul de intrare cât și cel de ieșire să fie cât mai scurt posibil.
- 2) Alfabetul folosit pentru scrierea acestor șiruri să fie cât mai inteligibil, având în vedere că șirurile se transmit prin viu grai.

2.2. Schimbul via modem

În cazul transmiterii datelor prin modem nu mai este necesară codificarea cheii comprimate, care implicit înseamnă creșterea dimensiunii șirului de intrare. Ambiguitățile care pot apărea în cazul schimbului prin telefon aici nu mai există. Acest schimb este mai avantajos din multe puncte de vedere, dar presupune că utilizatorul are modem!

3. Comprimarea propriu-zisă

Problema care se pune este comprimarea unui fișier de 496 octeți. Acest fișier de intrare are (sau ar trebui să aibă) anumite caracteristici, de exemplu conține multe secvențe FFFF și 0000. Pornind de la această observație, ca prim pas în comprimare, se va efectua o compresie cu pattern pe 16 biți.

3.1. Comprimarea cu pattern pe 16 biți

Fie un șir de pornire de dublu octeți. Se va construi alfabetul asociat acestui șir:

s_1, s_2, \dots, s_k - șirul de dublu octeți (nu neapărat distincți)

a_1, a_2, \dots, a_t - alfabetul asociat acestui șir

p_1, p_2, \dots, p_t - frecvențele elementelor alfabetului.

Atunci: a) $1 \leq k$

$$b) \sum_{i=1}^t p_i = k$$

c) lungimea șirului de intrare are $16 \times k$ biți.

Se va comprima șirul s (de intrare) la șirul z (de ieșire) astfel:

a) z este un șir de octeți

b) în z_1 și z_2 se va păstra lungimea în biți a șirului de intrare s ($16 \times k$)

c) în z_1 și z_2 se va păstra dublu octetul care apare cel mai frecvent (acel a_i pentru care $p_i = \max \{ p_j, j=1, \dots, t \}$)

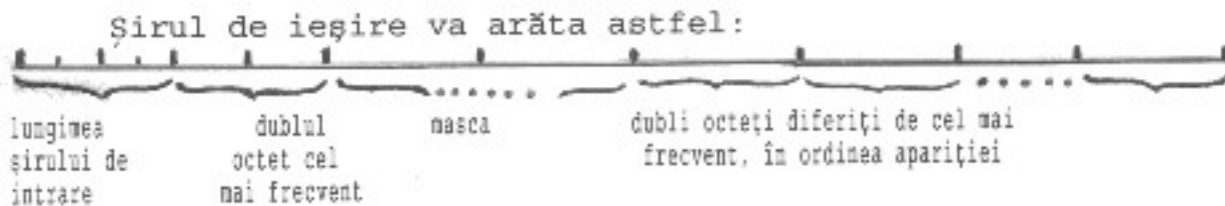
d) se creează o mască de biți astfel:

- bit pe 1 dacă dublul octet (caracterul) respectiv este a_i
- bit pe 0 dacă dublul octet (caracterul) respectiv nu este a_i
- această mască se scrie în z_1, z_2, \dots, z_t
- lungimea acestui subșir este $t-5+1$ octeți și

$$t-5+1 = \begin{cases} \frac{k}{8} & \text{dacă } 8 \text{ divide } k \\ \left[\frac{k}{8} \right] + 1 & \text{dacă } 8 \text{ nu divide } k \end{cases}$$

- dacă 8 nu divide pe k se completează ultimii biți (neseemnificativi) din z , cu 0.

e) se completează z cu acei dubli octeți din s care sunt diferiți de a_1 în ordinea apariției lor în z cu repetiție (dacă se repetă).



La decompresare se procedează astfel:

- se citește lungimea șirului de intrare
- se citește dublul octet cel mai frecvent
- se parcurge masca bit cu bit (mai puțin eventualii biți neseemnificativi din ultimul octet din mască; aceștia se cunosc pentru că se cunoaște lungimea șirului) și pentru fiecare bit 1 se scrie în șirul de ieșire dublul octet cel mai frecvent, iar pentru fiecare bit 0 se scrie dublul octet următor din secvența de după mască.

Problema care se pune este următoarea: este șirul de ieșire z mai scurt decât cel de intrare s ? Evident, pentru a se realiza o compresie propriu-zisă răspunsul trebuie să fie afirmativ.

Fie $l(s)$ și $l(z)$ lungimile celor două șiruri.

Atunci $l(z) = 2 + 2 + t - 4 + 2 \cdot (p_1 + \dots + p_{i-1} + p_{i+1} + \dots + p_l) = k + 2 \cdot (k - p_1)$.

Pentru cazul considerat $k = 496/2 = 248$. Fie $x = \frac{p_1}{8}$ procentul în care cel mai frecvent dublu octet apare în șirul de intrare.

Lungimea șirului de ieșire devine $l(z) = 35 + 496 \cdot (1 - x)$. Lungimea șirului de intrare este $l(s) = 496$ (octeți).

Condiția ca compresia să fie eficientă este ca $l(s) > l(z)$, adică $x > 35/496 \Rightarrow x > 0.07$. În concluzie, dacă procentul în care cel mai frecvent dublu octet este de cel puțin 7%, atunci această metodă este eficientă.

Pentru a se comprima un șir se va folosi această metodă

recursiv. Evident există un număr critic de aplicare a metodei începând de la care dimensiunea șirului nu mai scade. Se pot determina condiții matematice riguroase pentru aplicarea eficientă recursivă a metodei de un anumit număr de ori. De exemplu, pentru ca aplicarea recursivă de două ori a acestei metode de comprimare asupra unui șir de 512 octeți, condiția este $p_2 + \frac{p_1}{16} > \frac{1}{16} + \frac{1}{256}$, unde

p_1 și p_2 sunt procentajele celor mai frecvenți doi dublu octeți din șir. Bineînțeles, nu este nevoie să se verifice asemenea condiții înainte de aplicarea recursivă a metodei de comprimare, este suficient să se aplice metoda atâta timp cât dimensiunea șirului scade.

3.2. Comprimarea cu pattern pe 8 biți

Este similară cu cea prezentată anterior (pe 16 biți), diferența fiind aceea că elementele șirului de intrare sunt octeți. Ca și consecință a acestui fapt, dimensiunea măștii se dublează, condiția de eficiență fiind mai slabă ($x > 13\%$).

3.3. Comprimare combinată

Pentru unele fișiere (șiruri) se obțin rezultate eficiente dacă aceste metode de comprimare recursivă se aplică alternativ.

4. Codificarea

Odată șirul de intrare comprimat, se pune problema codificării acestuia, în sensul scrierii lui cu ajutorul unui alfabet accesibil (să nu uităm că șirul va fi transmis prin viu grai). Este dificil pentru o persoană să transmită prin telefon un șir format din caractere de genul: # \$ & etc. Se va folosi pentru codificare un alfabet format din 64 de caractere {a, ..., z, A, ..., Z, +, -}. Fiecare grup de 6 biți se codifică obținându-se printr-un caracter din acest alfabet. Există în acest caz un dezavantaj și anume că dimensiunea fișierului codificat crește de 4/3 ori.

Pentru a realiza compresia, decompresia (pe 8 biți și pe 16 biți), codificarea și decodificarea s-au scris patru programe în limbajul C, iar în continuare sunt prezentate câteva dintre rezultatele obținute.

Comprimare:

Fișierul	Tipul compresiei	Lungime inițială	1	2	3	4	5	6	7
R.HAS	16 biți	496	121	107	96	94	96		
R.HAS	8 biți	496	146	113	105	90	82	83	
TEST.HAS	16 biți	496	137	124	106	97	102		
TEST.HAS	8 biți	496	168	145	115	105	99	90	94

Codificare:

Fișierul	Lungime necodificat (octeți)	Lungime codificat (octeți)
R.HAS	94	127
R.HAS	82	111
TEST.HAS	97	131
TEST.HAS	90	121

BIBLIOGRAFIE

1. Ionescu T.C., Zsako I., Structuri arborescente, Ed.Tehnică, București, 1990
2. Knuth D.E., Tratat de programare a calculatoarelor, Algoritmi fundamentali, Ed.Tehnică, București, 1976
3. Aho A.V., Hopcroft J.E., Data Structures and Algorithms, Addison Wesley Reading, Massachussets, 1983

ON A COMPRESSION/DECOMPRESSION METHOD

Abstract. In this paper are presented methods for compression, codification, decompression, decodification generally for certain relatively short files that have to be transmitted "by voice" on the phone (usually the contain of a hardlock) in order to upgrade a certain soft product.

Universitatea din Baia Mare
Str. Victoriei nr. 76, 4800 Baia Mare
ROMÂNIA