

INTRODUCERE

Compresia Datelor – metode si tehnici de reducere a volumului unui set de date reprezentate sub forma binara, cu posibilitatea reconstituirii lor fie exact, fie aproximativ.

Set de date = sir finit de simbolii, cu reprezentare numerica unica (*cod numeric original*), in care nr. simbolilor ce-l compun (*volumul setului de date*) > nr. simbolilor distincti ce construiesc setul de date.

1 simbol \rightarrow 8 biti (standard ASCII) \Rightarrow setul de date se construiesc din max. $2^8 = 256$ simbolii distincti
Alfabet = multimea simbolilor distincti dintr-un set de date

Frecventa diferita de aparitie a simbolilor de-a lungul setului de date \Rightarrow recodificarea alfabetului folosind coduri de lungime variabila (simbolilor ce apar frecvent le sunt asociate coduri mai scurte decat simbolilor ce apar mai rar in setul de date)

Lungimea setului de date = nr. total de biti pe care este reprezentat setul de date

Volumul unui set de date = numarul de simbolii (nu neaparat distincti) care compun setul de date

Compresie $\begin{cases} \rightarrow \text{conservativa} - \text{studiaza metodele de reconstructie exacta a setului de date din versiunea sa comprimata (text, progr. sursa)} \\ \rightarrow \text{neconservativa} - \text{studiaza metode de reconstructie aproximative (inexacte) a setului de date original (semnal audio, imagine)} \end{cases}$

Rata de compresie = parametru prin care se masoara performantele metodei de compresie utilizate

$$\gamma = \left[1 - \frac{\#c}{\#o} \right] \times 100\%, \#c \text{ lungimea setului de date comprimat, } \#o \text{ lungimea setului de date original.}$$

$\gamma=100\%$, metoda de compresie ideala, castig maxim

$\gamma=0\%$, metoda nu realizeaza compresie, castig nul (storing)

$\gamma<0\%$, metoda nu realizeaza compresie, expandarea datelor

γ masoara castigul relativ de lungime obtinut prin utilizarea unei metode de compresie $\Rightarrow \gamma \in (0,100)\%$ pt o metoda de compresie eficienta.

Viteza de compresie/decompresie $\sim 1/\gamma$ (depinde de complexitatea metodei; compromis intre cei doi parametri).

Compresie conservativa - rata de compresie mare, vit. de compresie scazuta

Compresie neconservativa - rata de compresie mica, vit. de compresie mare

Expandare = cresterea lungii setului de date procesat fata de cel original (fenomen produs de metode de compresie ineficiente)

Decompresie = actiunea de reconstituire (exacta sau aproximativa) a setului de date original, plecand de la un set de date original

$$\text{Factor de compresie } \gamma' = \frac{\#c}{\#o} \times 100\% = 100\% - \gamma$$

Metoda de compresie eficienta \Leftrightarrow factor de compresie mic

Obiectivul unei metode de compresie a datelor = detectarea si eliminarea *redundantelor* din setul de date (adica a repetabilitatii datelor)

Redundanta – constanta a setului de date, care nu depinde de metoda de compresie utilizata

- se masoara cu ajutorul unui *model statistic* asociat setului de date, placand de la frecventele de aparitie ale simbolilor alfabetului setului de date

Frecventa de aparitie = raportul intre numarul aparitiilor simbolului si numarul total de simbolii ai setului de date, notata cu $\nu(s)$ (frecventa de aparitie a simbolului s).

- constituie o estimatie a probabilitatii de aparitie a simbolului s, $P(s)$.

Entropie = numarul minim de biti suficient pentru a recodifica simbolul s, fara a se realiza confuzia cu alti simbolii.

- numarul poate fi fractionar, nu neaparat intreg, dar de obicei se rotunjeste la intregul superior
- va fi cel mult egala cu nr. de biti ai codului original
- diferenta dintre ele este o masura a redundantei intrinseci (nr. de biti inutili prezenti in codul original)
- $H(s) = -\log_2 \nu(s) \Rightarrow \nu(s) = 2^{-H(s)} \in (0,1]$
- $0 < \nu(s) \leq 1 \Rightarrow H(s) > 0$

Ex: nr. aparitii(a) = 38 intr-un text cu Vol = 100 caractere $\Rightarrow \nu(a)=0,38$ si $H(a) = -\log_2(0,38) \approx 1,4$ biti, deci 6,6 biti (din 8) din codajul initial sunt inutili (redundanti) \Rightarrow sunt suficienti 2 biti pentru alocarea literei „a” la o prima comprimare a textului.

Cantitatea de informatie $I(s) = \frac{C}{\nu(s)}$, C – constanta de proportionalitate

Un simbol care apare frecventa in setul de date transporta o cantitate de informatie mica si va primi un numar mic de biti de recodificare: $H(s) = \log_2 I(s)$.

$R(D) = \sum_{s \in D} [8 - H(s)]$, D – setul de date asupra caruia se opereaza; alocare initiala = 8 biti/simbol

Redundanta depinde direct de entropiile simbolilor si este cu atat mai mare cu cat entropiile sunt mai mici

$R(D) = \sum_{s \in D} [8 + \log_2 \nu(s)] \Rightarrow$ frecvente mari de aparitie ale simbolilor conduc la o redundanta mare a acestuia.

Ex: sirul „aaaaa” e reprezentat pe $5 \cdot 8 = 40$ biti. Daca el face parte dintr-un set mai mare de date, in care $\nu(a) = 1/16 \Rightarrow H(a) = 4$ si $R(aaaaa) = 40 - 5 \cdot 4 = 20$ biti \Rightarrow doar 20 de biti de reprezentare sunt utili (50%). Codificarea ideala a cuvintului este pe 7 biti (in loc de 20), din care 4 biti pentru reprezentarea literei „a” si 3 biti pentru numarul de aparitii consecutive ($5 = 101$).