

O METODĂ DE REZOLVARE A PROBLEMEI DE SUMARIZĂRE A INFORMAȚIEI TEXTUALE

¹Victoria LAZU, ¹Stanislav SANDUȚA, ¹Valeria UNGUREANU,
²Ghenadie SAFONOV

¹Universitatea Tehnică a Moldovei, ²Academia Militară a Forțelor Armate „Alexandru cel Bun”

Abstract: În lucrarea de față sunt prezentate rezultatele dezvoltării unei metode de rezolvare a problemei de sumarizare a informației textuale bazată pe selectarea cuvintelor cheie care sunt întâlnite cel mai frecvent în text și selectarea propozițiilor care includ aceste cuvinte pentru modelarea informației de sumarizare a textului.

Cuvinte cheie: Logica Fuzzy, teoria mulțimilor, sumarizarea textului, procesarea limbajului natural, calculul natural, calculul membranelor, P-systems.

Introducere

Procesul de sumarizare este o operație efectuată de un sistem hardware-software cu scopul extragerii unui rezumat, care include punctele importante, cu un conținut ce descrie esența documentului original [1,2]. Complexitatea algoritmică ale sistemelor dedicate pentru crearea unui rezumat includ operații ce depind de lungimea, stilul și sintaxa textului procesat. Sinteza unui rezumat face parte din domeniul inteligenței artificiale și include operații din *machine learning* și *data mining*.

Există două abordări generale ale metodelor de sumarizare automată: sumarizare prin extragere și sumarizare prin abstractizare. Metodele bazate pe extragere funcționează prin selectarea unui subset de cuvinte, fraze sau propoziții existente în textul original pentru a forma rezumatul. Pe când, metodele bazate pe abstractizare construiesc o reprezentare semantică internă și apoi folosesc tehnici de generare a limbajului natural pentru a crea un rezumat mai aproape de gândirea umană. Un astfel de rezumat ar putea include inovații verbale obținute în rezultatul auto-învățării [2,3].

Una din metodele cele mai eficiente pentru sporirea performanțelor ale sistemelor dedicate pentru procesarea limbajelor naturale este utilizarea procesoarelor specializate cu procesare paralelă a datelor. Totodată, dezavantajul acestor sisteme este: domeniul limitat de utilizare, lipsa unor modele matematice bazate pe paralelism, cost mare pentru dezvoltarea aplicațiilor hardware și software [4].

Un avantaj deosebit poate fi oferit de utilizarea metodelor netradiționale de procesare a datelor care după conținut și structură impun un paralelism în procesul de calcul.

În lucrare se propune dezvoltarea unui sistem specializat destinat pentru rezolvarea problemei de sumarizare a informației textuale bazat pe paralelismul oferit de calculul natural în special calculul celular (P-systems) [5,6].

1. Dezvoltarea modelului matematic pentru rezolvarea problemei de sumarizare a informației textuale

Fie este dat textul $T^I = \bigcup_{i=1}^I (c_i)$, unde c_i este mulțimea de cuvinte și I numărul total de cuvinte în

text. Textul T^I este structurat în aliniate $T^I = \bigcup_{j=1}^J (A_j)$ și respectiv, fiecare aliniat A_j este structurat în

propoziții $A_j = \bigcup_{l=1}^{L_j} (P_{j,l})$. Pentru fiecare propoziție $P_{j,l}$ este definită mulțimea de cuvinte $P_{j,l} = \bigcup_{i=1}^{I_{j,l}} (c_{j,l,i})$,

unde $\forall c_{j,l,i} \in T^I$.

În scopul obținerii textului sumarizat T^S din textul T^I sunt definite următoarele operații de procesare a datelor:

$T^I \xrightarrow{Q} T^K$ - extragerea cuvintelor cheie, unde: Q - este algoritmul de extragere a cuvintelor cheie;

$T^K = \bigcup_{m=1}^M (c_m^K)$ și $\forall c_m^K \in T^I$;

$T^I \xrightarrow{G(T^K)} T^S$ - extragerea rezumatului, unde $G(T^K)$ este algoritmul de extragere a textului pentru sumarizare în baza cuvintelor cheie T^K .

Modelul matematic al algoritmului de extragere a cuvintelor cheie.

Algoritmul de extragere a cuvintelor cheie T^K din textul T^I s-a elaborat în baza modelului matematic:

$$T^K = \bigcup_{m=1}^M \left(c_i \mid \max^{N \geq N^*} \left\{ \sum (c_i^R), \forall i = \overline{1, I} \right\} \right),$$

unde: $\sum (c_i^R)$ - numărul de rădăcini ale cuvântului c_i în textul T^I ; $\max^{N \geq N^*}$ - selectarea cuvintelor c_i care se întâlnesc mai frecvent de N^* ori în textul T^I .

Modelul matematic al algoritmului de sumarizare în baza cuvintelor cheie.

Algoritmul de sumarizare $G(T^K)$ în baza cuvintelor cheie T^K s-a elaborat în baza modelului matematic:

$$T^S = \bigcup_{l=1}^{L_s} \left(\max^{N \geq N^*} \left\{ \bigcup_{h=1}^H (P_h) \mid c_m^K \in P_h \right\} \right), \text{ unde: } \bigcup_{h=1}^H (P_h) - \text{reuniunea propozițiilor care includ cuvântul}$$

cheie c_m^K ; $\max^{N \geq N^*}$ - selectarea propozițiilor care includ cuvântul cheie c_m^K de mai multe ori ca N^* .

Mențiuni

Cercetările efectuate fac parte din tematica tezelor de doctorat planificate în cadrul Departamentului Informatică și Ingineria Sistemelor, FCIM, UTM. Modelările și testările experimentale au fost efectuate în baza dispozitivelor oferite de CSCT „Hard & Soft” și ORANGE Cafee.

Bibliografie

1. Marcu, Daniel. *The Theory and Practice of Discourse Parsing and Summarization*. MIT Press Cambridge, MA, USA, 2000, 248p., ISBN 0-262-13372-5.
2. Mani, Inderjeet. *Automatic Summarization*. John Benjamins Publishing, 2001, 285p., ISBN 1-58811-060-5.
3. Louis, A. & Nenkova, A., *Performance Confidence Estimation for Automatic Summarization*. Proceedings of the 12th Conference of the European Chapter of the ACL, pp. 541–548, Athens, Greece, 30 March – 3 April, 2009.
4. Н.И. Червяков, П.А. Сахнюк, А.В. Шапошников, С.А. Ряднов. *Модулярные параллельные вычислительные структуры нейтропроцессорных систем*. - М.: ФИЗМАТЛИТ, 2003. - 288 с., ISBN: 5-9221-0327-X.
5. Gh. Păun. *A quick introduction to membrane computing*. The Journal of Logic and Algebraic Programming, 79, 2010, pp. 291–294.
6. G. Zhang, J. Cheng, T. Wang, X. Wang, J. Zhu. *Membrane Computing: Theory and Applications*, Science Press, Beijing, China, 2015.
7. ABABII, V.; SUDACEVSCHI, V.; LAZU, V.; MUNTEANU, S.; UNGUREANU, V. Distributed computing system based on mobile program code. *In The XIV International Conference “Measurement and Control in Complex Systems” – MCCS-2018, October 15-17, 2018, Vinnitsia, Ukraine*. { <http://ir.lib.vntu.edu.ua/handle/123456789/22769> }, file: 145.pdf.
8. АБАБИЙ, В.; СУДАЧЕВСКИ, В.; ЛАЗУ, В.; МУНТЯНУ, С.; УНГУРЯНУ, В. Организация распределенных вычислительных процессов для решения задачи многокритериального поиска информации в неструктурированных текстовых документах. *Proceedings of the Seventh International Conference on Informatics and Computer Technics Problems PICT-2018, 11-14 October, 2018, Chernivtsi, Ukraine*. pp. 71-74.
9. ABABII, V.; SUDACEVSCHI, V.; LAZU, V.; MUNTEANU, S.; UNGUREANU, V. Distributed data processing method for extracting knowledge. *Book of Abstracts of the 26th Conference on Applied and Industrial Mathematics - CAIM-2018, 20-23 September, 2018, Technical University of Moldova, Chisinau, Moldova*. pp. 114.