

METRICA REPREZENTATIVITĂȚII ÎN ENTITĂȚI TEXT

Ion Ivan

ionivan@ase.ro

Daniel Milodin

daniel.milodin@ase.ro

Mihai Georgescu

mihaita.georgescu@gmail.com

Academia de Studii Economice București

Abstract: Se definesc conceptele de entități text structurate și componentele acestora. Se detaliază conceptul de ortogonalitate a entităților structurate și se aplică acest concept în cadrul textelor componente ale aceleiași entități. Sunt definite criteriile ce stau la baza studiului ortogonalității. Se detaliază operațiunea de normalizare a textelor, în vederea eficientizării analizelor efectuate. Se prezintă conceptul de repetitivitate a subșirurilor text.

Cuvinte cheie: entități structurate, ortogonalitate, repetitivitate.

Abstract: Are defined the concepts of structured text entities and their components. Is presented in detail the orthogonality concept of structured entities and it is applied in the text components of the same entity. Are defined the criteria underlying the study of orthogonality. Is detailed the operation of normalizing the texts in order to make the performed analysis more efficient. Is presented the concept of text substrings repetition.

Keywords: structured entities, orthogonality, repetition.

1. Entități text structurate

Entitățile structurate creează și implementează forme de stocare și prelucrare care sunt folosite pentru gestiunea informației.

Organizarea informației se realizează în funcție de caracteristicile care o definesc și care o individualizează.

Structurarea entităților reflectă asocierile existente între date.

Entitățile structurate sintetizează modalitatea de grupare a caracteristicilor, precum și a legăturilor dintre acestea.

Între elementele componente ale unei entități structurate există o relație de dependență.

Fiecare componentă a entității are un corespondent în plan real, materializat prin niveluri sau subniveluri de organizare.

Structurarea entităților este realizată pe baza dispunerii elementelor ce compun sistemul analizat, structurile construite lucrând cu informații care sunt stocate în cadrul arborescenței implementate. Condițiile ce trebuie respectate în cadrul structurării implică dispunerea entităților pe nivele și prelucrarea informațiilor stocate pe fiecare nivel.

Organizarea este realizată pe diferite domenii de activitate, definirea entităților text făcându-se în funcție de criteriul de structurare și de informațiile cu care operează.

Entitățile text structurate sunt construite pe baza unor concepte de bază, precum [5]:

- simbolul constituie o modalitate de reprezentare; prin folosirea simbolurilor se pun în corespondență obiecte, noțiuni, imagini cu reprezentarea acestora; simbolul sugerează elemente reale prin înlocuirea acestora cu anumite forme considerate reprezentative; simbolul este definit printr-un singur semn, α reprezentând litera grecească alfa, sau printr-un grup de semne, H_2O reprezentând simbolul chimic al apei; cu ajutorul simbolului sunt reprezentate cantități, fenomene, imagini, idei, operații, formule; de-a lungul timpului modalitatea de utilizare a simbolului a evoluat, în funcție de gradul de dezvoltare a societății umane; inițial, simbolul a definit obiecte, oameni, animale, rolul lui fiind acela de a constitui elemente ale unui alfabet în procesul de comunicare; treptat, o dată cu apariția literelor și alfabetelor, simbolul și-a pierdut din importanță, acesta trecând în plan secundar;
- alfabetul este o mulțime finită, alcătuită din literele folosite în scrierea unei limbi; caracteristica de bază a literelor ce redau cuvintele unei limbi este de a respecta o ordine

convențională; funcție de simbolurile utilizate alfabetele se clasifică în alfabete fonetice și alfabete simbolice precum alfabetul Morse, alfabetul Braille; plecând de la rolul de bază al alfabetului, acela de a reda sunetele unei limbi, implicațiile alfabetului s-au extins asupra unor arii noi de interes; în informatică alfabetul definește totalitatea simbolurilor care stau la baza unui limbaj de programare;

- cuvântul reprezintă asocierea dintre un sens și un complex sonor; cuvântul este unitatea de bază a vocabularului; în funcție de modalitatea de formare cuvintele se clasifică în: cuvinte de bază – cuvinte care servesc ca elemente de bază în formarea altor cuvinte, cuvinte compuse – cuvinte formate prin compunerea mai multor cuvinte și cuvinte derivate – cuvinte formate plecând de la un alt cuvânt; unui cuvânt îi corespund unul sau mai multe sensuri;
- vocabularul este format din mulțimea cuvintelor specifice unei limbi; frecvența de utilizare a cuvintelor împarte vocabularul în:
 - vocabular activ alcătuit din totalitatea cuvintelor specifice unei limbi folosite în mod efectiv în exprimare;
 - vocabular pasiv ce conține cuvintele specifice unei limbi care nu sunt utilizate în mod frecvent;
 - vocabular fundamental alcătuit din fondul principal de cuvinte folosite;
 - vocabular secundar ce formează masa vocabularului; de asemenea, termenul de vocabular definește cuvintele specifice categoriilor sociale, domeniilor de activitate, de cercetare;
- separatorul este un simbol care nu aparține alfabetului, cu rolul de a delimita cuvintele din alfabet când sunt folosite pentru transmiterea de informații; textele sunt construite prin utilizarea de separatori; separatorii : ; ? ! @ # \$ % ^ & * () ` ~ , . < > / \ - _ + = | " ' [] " } { sunt alocați mulțimii separatorilor;
- frecvența de apariție este un indicator care furnizează informații despre nivelul de utilizare a termenilor din vocabularul de bază în cadrul entităților structurate; în funcție de numărul de apariții, termenii se împart în:
 - termeni cu frecvență scăzută sunt termenii care la un anumit număr de cuvinte au un număr scăzut de apariții;
 - termeni cu frecvență ridicată de apariție sunt termenii care la un număr dat de cuvinte au un număr ridicat de apariții;
- textul este alcătuit din mai multe cuvinte separate care transmit o anumită informație; textul este caracterizat de lungime, măsurată în număr de cuvinte sau în număr de simboluri folosite, de frecvența de apariție a cuvintelor sau simbolurilor, precum și de gradul de apartenență la un vocabular prin care se stabilește în ce măsură textul respectiv folosește termeni dintr-un domeniu; textul *Ortogonalitatea textelor oferă informații despre nivelul de originalitate al textelor și despre nivelul de asemănare al textelor* conține un număr de 17 cuvinte, din care 11 cuvinte distincte și 112 caractere;
- tezaurul este mulțimea formată din totalitatea cuvintelor considerate definitorii pentru un domeniu;
- sistemul referit definește ansamblul de elemente structurate, studiat și implementat folosind entitățile structurate.

Gradul de dificultate a conceptelor folosite în cadrul entităților oferă indicii cu privire la nivelul de reprezentare folosit. Cuvântul, simbolul și separatorii sunt folosiți în cuprinsul textelor și în construirea de entități text. Vocabularul și tezaurul sunt create ținând seama de apartenența cuvintelor folosite la o arie de interes definită cu ajutorul tezaurului și a vocabularului [1].

Modelele de stocare a textelor construite în mod flexibil și cu deschidere pentru elementele de noutate permit prelucrarea de diverse tipuri de texte, orientate spre diverse domenii de activitate.

Entitățile text structurate sunt în general folosite pentru lucrul cu proiectele.

Un text este structurat când este alcătuit din mai multe subtexte ce au o existență de sine stătătoare.

Șirurile de cuvinte sunt caracterizate prin poziții ale cuvintelor, prin legăturile existente între acestea, prin definirea de contexte și prin existența unei corespondențe între aceste cuvinte și lumea reală.

Entitățile text au o arie largă de aplicabilitate, pornind de la articole simple, de ziar, până la produse program scrise în diverse limbaje.

Dimensiunea unui vocabular este dată de numărul cuvintelor componente. Vocabularul limbii române este compus dintr-un număr relativ mic de cuvinte, aproximativ 200.000 de cuvinte, comparativ cu limba engleză care a depășit recent un număr de 1.000.000 cuvinte.

DEX - Dicționarul explicativ al limbii române 2009 conține un număr de 65000 de cuvinte (cele mai utilizate în limba română contemporană și în operele literare clasice), precum și un mare număr de regionalisme, arhaisme și neologisme, sensuri și unități frazeologice noi.

Fiecare cuvânt al unei limbi reprezintă o legătură între un obiect, o stare, o conjunctură și un limbaj universal acceptat de vorbitorii acelei limbi. Prin intermediul cuvântului sunt create premisele unei interfețe care asigură comunicarea între membri societății.

Astfel, cu cât vocabularul unei limbi este mai variat, cu atât comunicarea se realizează mai eficient, mai exact.

2. Ortogonalitatea internă a entităților structurate

Ortogonalitatea studiază gradul de asemănare între două sau mai multe entități. Prin intermediul acestei caracteristici de calitate este determinată măsura în care entitățile diferă una de cealaltă [2].

Conceptul de ortogonalitate provine din domeniul matematicii, unde are în vedere aspectele următoare:

- două plane sunt ortogonale dacă unghiul format la intersecția lor are cosinusul egal cu valoarea zero; o mulțime finită de plane este ortogonală dacă planele sunt perpendiculare două câte două;
- două drepte sunt ortogonale dacă formează unghiuri adiacente congruente; o mulțime finită de drepte este ortogonală dacă dreptele sunt perpendiculare două câte două;
- doi vectori sunt ortogonali dacă produsul scalar al acestora este nul; o mulțime finită de vectori este ortogonală dacă produsul scalar al oricăror doi vectori diferiți este nul.

Ortogonalitatea este studiată pe baza criteriilor de ortogonalitate. Cu ajutorul acestor criterii sunt evidențiate caracteristicile care au aceeași valoare pentru entitățile studiate și sunt determinate nivelurile de asemănare.

Ortogonalitatea entităților text este studiată ținând cont de:

- conținutul lor informațional;
- semnele utilizate pentru a reda informația, respectiv de semiotica entităților.

Pentru a fi studiată ortogonalitatea a două sau mai multe entități, ca o primă condiție ce trebuie respectată, acestea trebuie să aibă aceeași structură, adică să fie definite prin intermediul aceluiași caracteristici.

Plecând de la caracteristicile entităților, sunt definiți indicatori de calitate ai acestora. Pe baza acestor indicatori este construit un indicator ce ține cont de valorile indicatorilor componenți.

Un domeniu în care ortogonalitatea este foarte importantă este programarea. Limbajele de programare sunt astfel proiectate încât ortogonalitatea acestora să fie maximă, în sensul că noțiunile implementate, termenii și cuvintele cheie folosite sunt unice pentru a nu crea confuzie, atât în rândul utilizatorilor, cât mai ales pentru a nu aduce compilatoarele în situația de a nu cunoaște ce presupune o linie cod, din cauza mai multor opțiuni posibile [3].

Compararea a două entități se reduce la a raporta o entitate la cealaltă entitate, respectiv la a identifica părțile comune și părțile care diferă. Sunt astfel, comparate caracteristicile

corespondente ale celor două entități.

Pentru entitățile text un rol important îl are frecvența de apariție a cuvintelor în cadrul entităților. Prin intermediul frecvenței este determinată importanța cuvintelor în cadrul entităților și este determinat gradul de utilizare a cuvintelor, precum și felul cum acestea influențează construirea entităților.

Ortogonalitatea textelor se determină ca ortogonalitate internă și ortogonalitate externă. Ortogonalitatea internă stabilește măsura în care, în cadrul textului, cuvintele folosite se repetă, identificând legăturile existente între cuvintele din cadrul textului, precum și modalitatea de dispunere a acestora.

Ortogonalitatea externă arată care sunt diferențele sau asemănările ce există între texte.

Ortogonalitatea entităților text este studiată pe baza următoarelor criterii [4]:

- dimensiunea fișierelor care stochează entitățile text, exprimată ca număr de biți;
- frecvențele de apariție a cuvintelor în cadrul textelor, calculate prin identificarea numărului de apariții ale cuvintelor în cadrul textelor;
- frecvențele de apariție a cuvintelor de legătură, calculate prin identificarea numărului de apariții ale cuvintelor de legătură în cadrul textelor;
- dimensiunea textelor privită ca număr de cuvinte sau ca număr de litere care alcătuiesc textele;
- dispunerea cuvintelor în cadrul textelor, respectiv stabilirea poziției fiecărui cuvânt și identificarea similitudinilor;
- distanțele dintre cuvinte definite ca număr de cuvinte care se interpun între două cuvinte considerate;
- numărul de paragrafe ale textelor.

Două entități text sunt ortogonale dacă nu au cuvinte comune.

Ortogonalitatea a două texte este calculată prin inventarierea cuvintelor comune.

În cazul studierii ortogonalității foarte important este și sensul cu care sunt folosite datele. Dacă sunt folosite entități pentru a defini un domeniu de activitate și în cadrul entităților sunt anumite restricții privitoare la datele ce sunt disponibile, atunci frecvențele de utilizare vor avea valori apropiate, conținutul va fi asemănător, singurul aspect care deosebește entitățile fiind sensul cu care sunt folosite datele.

Este urmărită modalitatea în care două entități diferă sau sunt asemănătoare, atât ca frecvențe de apariție a cuvintelor, conținut, precum și sens al datelor conținute.

Conceptul de ortogonalitate se aplică pentru a identifica dacă anumite lucrări aparțin sau nu unui domeniu de referință. Este studiat gradul de asemănare, frecvențele de utilizare a cuvintelor, iar pe baza indicatorului rezultat se determină dacă textele sunt asemănătoare sau nu. În situația în care indicatorul de ortogonalitate tinde spre 1, adică lucrările sunt total diferite, este evident că nu aparțin aceluiași domeniu. Pentru a identifica în mod corect dacă două sau mai multe lucrări aparțin aceluiași domeniu trebuie identificate cuvintele de specialitate conținute de lucrări și frecvențele de apariție a acestora.

Comparabilitatea textelor se analizează la următoarele niveluri:

- formă de reprezentare;
- conținut informațional.

Reprezentarea percepțiilor omului privind lumea înconjurătoare se realizează prin intermediul construcțiilor sintactice formate din simboluri alfanumerice, grafice, spațiale.

Criteriile de comparare a două entități vizează:

- lungimea entității text – este exprimată în mai multe moduri:
 - număr de simboluri utilizate;

- număr de caractere utilizate cu posibilitatea de a diferenția pe caractere mici și caractere mari;
 - număr de cuvinte folosite;
 - număr de pagini;
 - număr de bytes ocupați, în cazul entităților text care sunt date în format electronic.
- lungimea vocabulelor entităților – exprimată ca număr de cuvinte distincte utilizate în construcția entităților text;
 - lungimea vocabulelor cuvintelor cheie – sunt luate în considerare cuvintele definiției pentru domeniul abordat;
 - numărul de cuvinte cheie din titlurile entităților;
 - diversitatea vocabulelor entităților – raportează cuvintele din vocabularele entităților cu cele din vocabularul limbii în care au fost realizate;
 - frecvențele de apariție a cuvintelor din entități – pentru fiecare cuvânt din vocabularul unei entități se asociază numărul de apariții în cadrul entității;
 - frecvențele de apariție a cuvintelor din vocabularele entităților în raport cu un vocabular definit de utilizator;
 - construirea matricelor de precedente în vederea analizei comparate a textelor în funcție de pozițiile construcțiilor sintactice.

Comparabilitatea textelor din punctul de vedere al conținutului informațional are în vedere următoarele aspecte:

- delimitarea domeniului abordat prin evidențierea semanticii cuvintelor cheie;
- înțelesul cuvintelor specifice domeniului și sensul în care acestea au fost utilizate în construcția entității text;
- acoperirea cadrului conceptual prin utilizarea cuvintelor specifice domeniului;
- aducerea entităților la o formă comună prin eliminarea cuvintelor de legătură și a prefixelor/sufixelor adăugate conform regulilor gramaticale specifice limbii;
- identificarea sinonimelor pentru cuvintele din vocabularul entităților și modificarea textelor prin utilizarea aceleiași forme a cuvântului;
- identificarea și prezentarea elementelor de natură teoretică și practică exprimate în cele două entități și realizarea de asocieri;
- identificarea funcționalităților descrise în texte, inventarierea lor și evidențierea eventualelor asocieri.

Ortogonalitatea entităților text este studiată prin prisma domeniului căruia îi aparțin textele prelucrate cu ajutorul entităților.

Două entități text sunt ortogonale dacă nu au cuvinte comune.

Versurile următoare nu conțin elemente repetitive:

*La steaua care-a răsărit
E-o cale-atât de lungă,
Că mii de ani i-au trebuit
Luminii să ne-ajungă.*

Spre deosebire de următoarele versuri în care se repetă primul și ultimul vers:

*Vreme trece, vreme vine,
Toate-s vechi și nouă toate;
Ce e rău și ce e bine
Tu te-ntreabă și socoate;
.....
Tu așează-te deoparte,
Regăsindu-te pe tine,
Când cu zgomote deșarte
Vreme trece, vreme vine.*

Ortogonalitatea a două texte este calculată prin inventarierea cuvintelor comune. Procedul devine complex când textele aparțin aceluiași domeniu sau au valoare științifică.

Compararea a două texte științifice pentru stabilirea ortogonalității implică determinarea schemei de dispunere a cuvintelor. Trebuie identificată modalitatea de folosire a cuvintelor și de poziționare a acestora unele față de altele.

Fiecare autor are un stil propriu de îmbinare a cuvintelor ce aparțin domeniului științific.

Ortogonalitatea entităților text implică aspecte precum:

- determinarea poziției cuvintelor;
- determinarea distanțelor dintre cuvinte;
- stabilirea frecvenței cuvintelor care aparțin vocabularului domeniului.

La studierea ortogonalității trebuie să se țină cont de contribuția pe care autorul o aduce în text.

Pentru o analiză eficientă, textele trebuie normalizate.

Normalizarea are în vedere următoarele operații:

- eliminarea secvențelor de separatori și a caracterelor de control – presupune ștergerea secvențelor redundante de simboluri care marchează zone de text;
- eliminarea diferențierii între caracterele mari și caracterele mici pentru uniformizarea comparațiilor;
- eliminarea folosirii diacriticelor – constă în înlocuirea caracterelor specifice unei limbi cu caractere echivalente;
- înlocuirea simbolurilor din alte limbi – presupune stabilirea corespondenței între cuvintele ce nu aparțin limbii și succesiunea de caractere specifice acesteia.

Prin normalizare este asigurată omogenitatea datelor din punct de vedere al comparabilității lor.

Ortogonalitatea entităților text este studiată din punct de vedere semantic și din punct de vedere semiotic.

Ortogonalitatea semantică studiază cuvintele ca și sens al acestora, respectiv stabilește dacă două sau mai multe cuvinte diferă ca înțeles. Pentru aceasta este determinată rădăcina cuvintelor și se stabilește modul de evoluție al fiecărui cuvânt, pornind de la rădăcină.

Ortogonalitatea semiotică determină în ce măsură textele se aseamănă din punct de vedere al simbolurilor utilizate, al poziției în cadrul textului și a distanței dintre cuvinte în cadrul textului.

3. Repetitivitatea subșirurilor identice în cadrul entităților text

Conceptul de repetitivitate definește calitatea entităților de a fi alcătuite din texte care apar în cadrul acestora la poziții diferite, dar având aceeași structură a dispunerii.

Se consideră entitatea text E_1 alcătuită din textele T_1 , T_2 și T_3 , astfel:

$$T_1 = \{aaa, bbbe, ccc\}$$

$$T_2 = \{ppap, ss, mmm, bbbe, ss, nnc, ss, pprp, nn, ddd, ccc, bbbe, aaa, ii\}$$

$$T_3 = \{ppap, mmm, bbbe, ss, nnc, ss, pprp\},$$

entitatea E_1 fiind formată din reuniunea celor trei texte, astfel:

$$E_1 = T_1 + T_2 + T_3 = \{aaa, bbbe, ccc, ii, ppap, ss, mmm, bbbe, ss, nnc, ss, pprp, nn, ddd, ccc, bbbe, aaa, ii, ppap, mmm, bbbe, ss, nnc, ss, pprp\},$$

unde operația de reuniune este reprezentată prin convenție de simbolul +.

Lungimea entității E_1 este dată de numărul total al cuvintelor componente, respectiv 25.

Textele T_1 , T_2 și T_3 sunt construite pe baza cuvintelor:

$$\langle c_i \rangle = \{aaa\}$$

$\langle c_2 \rangle = \{ bbbe \}$

$\langle c_3 \rangle = \{ ccc \}$

$\langle c_4 \rangle = \{ ppap \}$

$\langle c_5 \rangle = \{ mmm \}$

$\langle c_6 \rangle = \{ nnc \}$

$\langle c_7 \rangle = \{ pprp \}$

$\langle c_8 \rangle = \{ nn \}$.

Pornind de la cuvintele identificate în mod unic în cadrul celor trei texte componente ale entității, se construiește vocabularul reprezentativ al entității, VE, alcătuit din cele 8 cuvinte, astfel:

$VE = \{ aaa, bbbe, ccc, ppap, mmm, nnc, pprp, nn \}$.

De asemenea, cele trei texte conțin și cuvinte cu o frecvență ridicată de apariție, respectiv cuvinte de legătură sau cuvinte cu caracter comun. Aceste cuvinte fac parte din vocabularul VL al cuvintelor de legătură:

$VL = \{ ii, ss, ddd, \}$.

În vederea unei analize eficiente se recomandă implementarea unor algoritmi de normalizare cu rolul de a înlătura atât cuvintele de legătură, cât și cuvintele care au o frecvență ridicată de apariție ce afectează analiza.

După normalizare, entitatea E_1 se prezintă astfel:

$E_1 = \{ aaa, bbbe, ccc, ppap, mmm, bbbe, nnc, pprp, nn, ccc, bbbe, aaa, ppap, mmm, bbbe, nnc, pprp \}$,

noua lungime fiind de 17 cuvinte.

În vederea stabilirii subșirurilor de cuvinte care se repetă în cadrul entității E_1 se propune următorul algoritm:

P_1 : se construiește vocabularul atașat entității text, în cazul prezentat vocabularul VE;

P_2 : se construiește vocabularul pentru cuvintele de legătură și pentru cuvintele cu frecvență ridicată de apariție, în exemplul considerat VL;

P_3 : se normalizează entitatea E_1 ;

P_4 : pornind de la vocabularul VE se parcurge textul entității E_1 și se notează poziția fiecărui cuvânt al entității în cadrul vocabularului, rezultând matricea asociată M;

P_5 : pornind de la prima coloană se parcurge matricea, iar unde se identifică valoarea 1 în cadrul coloanei parcurse, semn că un cuvânt din vocabular este prezent în cadrul textului, se trece la identificarea următoarei coloane care stochează valoarea 1; se construiește în acest fel un subșir care este căutat în matricea M coloană după coloană, pornind de la linia unde a fost identificat primul cuvânt al subșirului; la identificarea de poziții consecutive identice se incrementează contorul asociat subșirului pentru a i se stabili lungimea maximă; căutarea subșirurilor se realizează prin incrementare și prin decrementare, cuvintele formând același subșir și în ordine inversă ca în exemplul: $T_1 = \{ annn arrr, ammm \}$ și $T_2 = \{ ammm, arrr, annn \}$

P_6 : se calculează gradul de repetitivitate al subșirurilor;

Conceptul de repetitivitate este o particularizare a conceptului de ortogonalitate în cadrul aceluiasi text, prin intermediul acestuia identificându-se gradul de originalitate internă a unui text și asigurându-se construirea unor înlănțuiri de cuvinte diferite.

Pornind de la entitatea E_1 se construiește matricea M asociată identificării cuvintelor vocabularului VE în cadrul acestei entități:

Tabelul nr. 1 Matricea M asociată pozițiilor cuvintelor vocabularului VE în cadrul entității E₁

Poziția în text Cuvinte vocabular	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
<i>aaa</i>	1											1					
<i>bbbe</i>		1				1					1				1		
<i>ccc</i>			1							1							
<i>ppap</i>				1									1				
<i>mmm</i>					1									1			
<i>nnnc</i>							1									1	
<i>pprp</i>								1									1
<i>nn</i>									1								

Având la bază algoritmul descris anterior se pornește cu prima coloană a matricei, întrucât pe prima coloană se regăsește primul cuvânt al vocabularului VE. Se observă că primul cuvânt, respectiv *aaa* se regăsește în cadrul textului pe pozițiile 1 și 11. Aceste valori se identifică urmărind pozițiile pe care se află stocată cifra 1 în cadrul liniei corespondente a cuvântului analizat. Pornind de la aceste două coloane, se incrementează sau decrementează valorile pozițiilor unde a fost regăsit cuvântul *aaa*, această operațiune efectuându-se atât timp cât în cadrul matricei M valorile identificate pe pozițiile incrementate / decrementate sunt identice.

Rezultă, astfel că, pornind de la cuvântul *aaa* se construiește subșirul *aaa*, *bbbe*, *ccc*, respectiv subșirul în formă inversă *ccc*, *bbbe*, *aaa*, precum și subșirul *ppap*, *mmm*, *bbbe*, *nnnc*, *pprp*, regăsit în forma unui subșir identic în conținutul entității E₁.

De asemenea, mai sunt identificate și alte subșiruri, de lungime mai mică decât cele anterioare, acestea fiind incluse în primele două subșiruri.

Modul de lucru al algoritmului prezentat se bazează pe caracterul unic al apariției unui cuvânt al vocabularului pe o anumită poziție din cadrul textelor ce compun entitatea structurată. Astfel, pe aceeași poziție într-o text va fi identificat un singur cuvânt. Implementarea algoritmului asigură identificarea subșirurilor de cuvinte care se repetă în cadrul unui text, pe baza stabilirii succesiunii de apariții ale cuvintelor. Se identifică în acest mod o caracteristică a modului de utilizare a cuvintelor specifice fiecărui text sau se stabilește dacă un text este format prin utilizarea în mod repetat a acelorași subșiruri.

În cadrul entităților structurate de tip articol, repetitivitatea ia în considerare aspecte care asigură definirea și implementarea noțiunilor prezentate în partea introductivă a textului:

- cuvintele cheie se regăsesc în cadrul articolului;
- unele cuvinte cheie se regăsesc în titlurile articolelor din bibliografie;
- cuvintele cheie se regăsesc în subtitlurile articolului.

Aplicabilitatea directă a conceptului de repetitivitate în cazul entităților structurate de tip articol este dată de următoarele reguli:

- în titlul unui articol cuvintele nu se repetă;
- în cadrul secțiunii *abstract* o frază apare o singură dată;
- în capitolele componente ale articolului o definiție este oferită o singură dată.

Implementarea repetitivității asigură creșterea calității entităților structurate prin definirea de criterii de calitate rezultate din studierea modalității de dispunere a cuvintelor în cadrul textelor. Criteriile de calitate reflectă nivelul de repetare al cuvintelor în cadrul textelor, modalitățile de dispunere a cuvintelor în componența paragrafelor, și gradul de apartenență a cuvintelor la un vocabular de specialitate.

4. Indicatori ai repetitivității

În vederea stabilirii nivelului de ortogonalitate, având la bază algoritmul descris mai sus, se definesc indicatori care certifică originalitatea unui text.

Astfel, se definește un indicator normat de ortogonalitate cuprins în intervalul $[0, 1]$, care ia următoarele valori:

- 1, dacă elementele sunt ortogonale, adică nu au nimic în comun;
- 0, elementele sunt identice, adică nu au valori diferite pentru nicio caracteristică.

Dacă valoarea indicatorului tinde către 1 înseamnă că seturile de date conținute de cele două entități tind către ortogonalitate, iar dacă valoarea indicatorului este apropiată de 0 înseamnă că seturile de date au foarte multe elemente identice.

Problema ortogonalității textelor este luată în considerare pentru validarea textelor. În vederea analizării textelor se implementează reguli ce trebuie respectate pentru construirea textelor. Astfel, textele componente ale entității sunt delimitate de separatori și se încheie cu semnul grafic “.”. Entitatea structurată este compusă din mai multe texte, care se îmbină pentru a evidenția o informație.

Se definește indicatorul de repetitivitate a subșirurilor componente ale unei entități structurate:

$$H(\text{subșir}_A, \text{subșir}_B) = \frac{LC}{\max(L_A, L_B)},$$

unde:

- LC – numărul cuvintelor comune celor două subșiruri;
- L_A – lungime subșir $_A$;
- L_B – lungime subșir $_B$.

Determinarea repetitivității este raportată la dimensiunile textelor care se analizează și nu la dimensiunea totală a textului reunit, fapt care conduce către un rezultat cu un grad de reprezentativitate mai ridicat.

Se consideră cele trei texte, normalizate:

$$T_1 = \{aaa, bbbe, ccc\}$$

$$T_2 = \{ppap, mmm, bbbe, nnc, pprp, nn, ccc, bbbe, aaa\}$$

$$T_3 = \{ppap, mmm, bbbe, nnc, pprp\},$$

precum și cele două subșiruri comune identificate în cadrul celor trei texte:

$$S_1 = \{aaa, bbbe, ccc\}$$

$$S_2 = \{ppap, mmm, bbbe, nnc, pprp\},$$

cele două subșiruri regăsindu-se în cadrul textelor T_1 și T_2 , pentru subșirul S_1 , respectiv T_2 și T_3 , pentru subșirul S_2 .

Dimensiunile celor trei texte normalizate sunt următoarele:

$$L(T_1) = 3;$$

$$L(T_2) = 9;$$

$$L(T_3) = 5,$$

$$L(S_1) = 3;$$

$$L(S_2) = 5, \text{ cele două subșiruri reprezentând cuvintele comune subșirurilor.}$$

$$H(T_1, T_2) = 3 / 9 = 0.33,$$

$$H(T_2, T_3) = 5 / 9 = 0.56.$$

În vederea obținerii unei valori care să cumuleze valorile parțiale ale ortogonalităților, obținute prin compararea textelor luate două câte două, se consideră indicatorul agregat, construit pe baza formulei:

$$HG = 1 - \left(\prod_{i=1}^{NS-1} \prod_{j=i+1}^{NS} HG(\text{șir}_i, \text{șir}_j) \right)^{\frac{1}{NS(NS-1)}}$$

unde:

NS – numărul de șiruri care compun entitatea text analizată;

Șir_i – șirul *i* care este comparat cu un alt șir component al entității.

Pentru exemplul considerat valoarea indicatorului agregat HG este 0.24, ceea ce arată că nivelul ortogonalității interne a entității E₁ este foarte mic, șirurile fiind utilizate în mod excesiv la construirea entității.

Repetitivitatea subșirurilor este urmărită pentru secvențe de texte complete, precum propozițiile și frazele, nu doar pentru șiruri alcătuite din foarte puține cuvinte. Se urmărește astfel identificarea repetitivității care influențează originalitatea textelor.

5. Concluzii

Conceptul de repetitivitate identifică gradul de reutilizare al subșirurilor în cadrul unui text. Sunt identificate atât caracterul original al textelor, cât mai ales tendințele de a refolosi subșirurile în vederea construirii unui text. Algoritmul prezentat asigură determinarea gradului de reutilizare al textelor pe baza vocabularului construit pornind de la cuvintele ce compun textul analizat. Se asigură astfel eficientizarea stabilirii ortogonalității interne a textelor.

În vederea obținerii de rezultate cu o acuratețe ridicată sunt identificate și înlăturate cuvintele de legătură, precum și cuvintele cu grad ridicat de repetitivitate, care influențează indicatorii rezultați.

Studierea conceptului de repetitivitate în cadrul textelor asigură creșterea gradului de diversitate și implicit creșterea gradului de originalitate, prin utilizarea de texte cu caracter unic, orientate spre aceeași problematică.

BIBLIOGRAFIE

1. **IVAN, I., M. POPA:** Text Entities – Development, Evaluation, Analysis, ASE Printing House, Bucharest, 2005.
2. **POPA, M.:** Evaluarea calității entităților text, Teorie și Practică, Editura ASE, București, 2005.
3. **SMEUREANU, I., M. DARDALA:** Programarea în limbajul C/C++, Editura CISON, București, 2004.
4. **IVAN, I., D. MILODIN, M. POPA:** Operation on text entities, Informatica Economica, vol. 11, nr. 1, 2007, pp. 14 – 20.
5. **MILODIN, D.:** – Ortogonalitatea entităților structurate – Teză de doctorat, noiembrie 2009.