

TEORIA DEPENDENȚEI CONCEPTUALE – INSTRUMENT PENTRU SUMARIZARE

Autor: Victoria LAZU

Universitatea Tehnică a Moldovei

Email: lazu_vica@mail.utm.md

Abstract: Sumarizarea reprezintă un aspect foarte important în ceea ce privește coordonarea numărului vast de texte citite de oameni. Sumarizarea are rolul de a reduce cantitatea de informație care trebuie citită și oferă posibilitatea de a decide gradul de relevanță al informației necesare. Odată cu utilizarea calculatoarelor în cadrul procesării textelor scrise, una dintre primele sarcini apărute este cea de a realiza sumarul unui text care va prezenta cât mai succint conținutul, prin reducerea lungimii documentului, menținând totodată sensul și păstrând ideile de bază.

Cuvinte cheie: sumarizare, abstracție, reprezentarea cunoștințelor, dependența conceptuală, entități, acțiuni primitive, attribute, reguli semantice.

1. Introducere

Sumarizarea automată a textelor este un domeniu destul de vast și variat și include sarcini ce nu au fost încă exploatate sau nu au adus rezultate destul de relevante, deoarece limbajul natural este destul de ambiguu și întâmpină dificultăți în procesul de procesare a lui. În era calculatoarelor este firesc să apară dorința de a face computerele să înțeleagă limba vorbită și scrisă, să caute în texte întrebările și răspunsurile, să deducă concluzii logice asupra conținutului lor, să le poată sumariza și generaliza, adică să realizeze ceea ce face, o persoană, în mod normal cu textul. Sumarizarea are rolul de a reduce cantitatea de informație care trebuie citită și oferă posibilitatea de a decide gradul de relevanță al informației necesare. Dependența conceptuală este o teorie despre modul de reprezentare a diferitelor tipuri de cunoștințe, despre evenimentele conținute în mod uzual în propoziții exprimate în limbaj natural. Ea exprimă relațiile dintre concepte într-o manieră independentă de limba în care au fost exprimate la origine propozițiile.

2. Sumarizarea

Sumarizarea automată este crearea unei versiuni mai scurte a unui text de către un program. Rezultatul acestei operații conține totuși majoritatea punctelor importante din textul original. Tehnologiile care construiesc un sumar coerent, dintr-un text de orice natură, trebuie să ia în considerare diverse variabile precum lungimea, stilul de scriere și sintaxa pentru a realiza un sumar util. Un sumarizator este un sistem de prelucrare automată a unui sau mai multor texte, cu scopul obținerii unui rezumat (sumar) util unui utilizator uman. După forma de rezumare se disting două posibilități de sumarizare: *extracție* și *abstracție*.

Tehnicile de extracție copiază informațiile considerate cele mai importante din text în sumar adică sumarul este complet compus din secvențe de cuvinte (propoziții, fraze sau cuvinte cheie), copiate din documentul original în timp ce *tehnicele de abstracție* implică parafrizarea unor secțiuni din documentul sursă, ele conțin cuvinte, secvențe ce nu sunt prezente în originale. În general, abstracția poate condensa un text mai bine decât extracția, dar astfel de programe sunt mai greu de dezvoltat și necesită utilizarea tehnologiei de generare de limbaj natural, domeniu care este în dezvoltare.

Modelul de abstracție conține atât analiza morfologică a propozițiilor cât și cea sintactică și semantică.

Așa cum arată figura 1, abstracția are două abordări de bază. Prima (partea de sus a figurii) folosește o metodă tradițională lingvisticii care analizează propoziții sintactice. Această metodă utilizează informații semantice pentru a adnota arborii de analiză. Procedurile de compactare funcționează direct pe baza arborilor pentru a elimina și regrupa o parte din propoziții de exemplu, prin tăiere arborilor în conformitate cu criteriile structurale.

A doua abordare își are rădăcinile în inteligența artificială și se concentrează asupra înțelegerii limbajului natural. Inteligența artificială în special necesită cunoaștere.

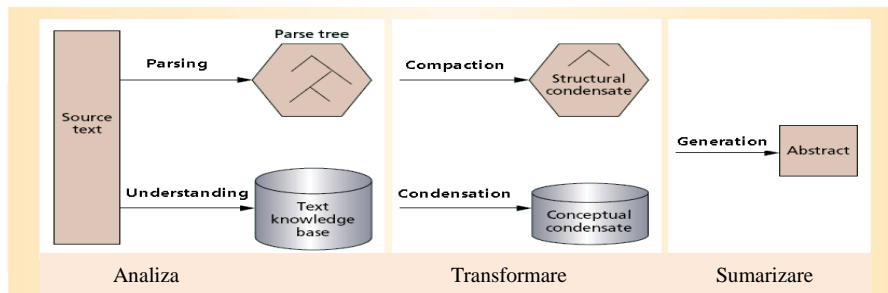


Fig. 1. Arhitectura de abstracție

Reprezentarea cunoștințelor poate fi dificilă: cunoașterea e voluminoasă, e greu de caracterizat cu precizie și e în permanentă schimbare.

Modalitățile tradiționale de reprezentare a cunoștințelor sunt: - Rețelele semantice (apropiate de limbajul natural, permițând ierarhizarea conceptelor); - Dependența conceptuală (independentă de limbajul natural); - Cadrele (pentru situații stereotipe); - Scenariile (pentru secvențe de evenimente stereotipe).

3. Dependența conceptuală

Construirea bazelor de cunoștințe mari și a sistemelor care exploatează aceste cunoștințe trebuie să țină cont de eficiența reprezentărilor interne asociate unui model, de capacitatea de acces și modificare a cunoștințelor. În plus, s-a constatat ca și la nivel simbolic, achiziția cunoștințelor, înțelegerea, utilizarea și întreținerea bazelor de cunoștințe este dificilă în lipsa unei organizări conceptuale adecvate. Aceste motive au dus la apariția unei noi clase de modele de reprezentare a cunoștințelor: cunoștințele structurate.

Cele mai importante reprezentări structurate ale cunoștințelor sînt: rețelele semantice și grafurile conceptuale, unitățile, numite în diverse abordări și cadre sau prototipuri, dependențele conceptuale și scenariile. Toate aceste modele prezintă o caracteristică comună: oferă posibilitatea de a organiza și structura cunoștințele în funcție de relațiile existente între obiectele universului problemei.

Paradigmele de reprezentare se numesc modele de cunoștințe structurate tocmai datorită accentului fundamental pe care îl pun pe structura cunoștințelor în reprezentare, prin indexarea cunoștințelor în funcție de obiectele importante în sistem. Dacă sistemul bazat pe cunoștințe are nevoie de una din informațiile asociate unui astfel de obiect, structura de cunoștințe asociată obiectului este regăsită și toate faptele relevante despre acel obiect sunt identificate deodată. Modelele cunoștințelor structurate au asociate metode de inferență specifice, incluse direct în reprezentare.

Dependențele conceptuale și scenariile pot fi caracterizate drept "structuri tari" de reprezentare a cunoștințelor, deoarece ele includ noțiuni specifice despre tipul obiectelor și relațiilor existente între acestea. Structurile tari de reprezentare a cunoștințelor pot fi puse în corespondență, la nivelul controlului, cu strategiile de căutare informată, în care fiecare algoritm de căutare include informație euristică specifică domeniului problemei.

Dependențele conceptuale sunt un model structurat care permite reprezentarea cunoștințelor conținute în propozițiile limbajului natural. Scopul dependențelor conceptuale este acela de a reprezenta cunoștințele astfel încît:

- reprezentarea să fie independentă de limbajul în care au fost formulate propozițiile;
- propoziții diferite, dar avînd aceeași semnificație, să aibă aceeași reprezentare sub forma dependențelor conceptuale;
- reprezentarea semnificației propozițiilor să fie neambiguă;
- să faciliteze execuția inferențelor determinate de informațiile conținute în propoziții.

Din aceste motive, dependențele conceptuale folosesc o reprezentare a propozițiilor limbajului natural care nu se bazează pe cuvintele propoziției ci pe o mulțime de entități primitive conceptuale care pot fi combinate pentru formarea semnificației cuvintelor și a propoziției. Acest model a fost propus de Schank în 1972 și a fost folosit în diverse programe de înțelegere a limbajului natural.

În teoria dependențelor conceptuale se disting cinci tipuri de elemente componente (primitive ontologice) care constituie componentele constructive ale reprezentării. Aceste elemente de bază sunt: entități, acțiuni, cazuri conceptuale, timpuri conceptuale și dependențe conceptuale, fiecare element avînd o serie de subtipuri.

Entități

Obiectele sau producătorii de scenarii, numite pe scurt PP, reprezintă persoanele (actorii) sau obiectele fizice (inclusiv memoria umană) care acționează în universul discursului.

Atributele sau ucenicii scenariilor, numite pe scurt PA, reprezintă proprietățile producătorilor de scenarii.

Acțiuni

Acțiunile primitive, numite pe scurt ACT. În reprezentarea dependențelor conceptuale acțiunile sunt formate dintr-o mulțime de acte primitive. Un set tipic de astfel de acțiuni primitive este următorul:

- ATRANS - Transferul făcut de o relație abstractă, de exemplu "a da"
- PTRANS - Transferul locației fizice a unui obiect, de exemplu "a merge"
- PROPEL - Aplicarea unei forțe fizice unui obiect, de exemplu "a împinge"
- MOVE - Mișcarea părții unui corp de către acel corp, de exemplu "a lovi"
- GRASP - Apucarea unui obiect de către un actor, de exemplu "a smulge"
- INGEST - Ingerarea unui obiect de către o ființă, de exemplu "a mânca"
- EXPEL - Expulzarea unui obiect din corpul unei ființe, de exemplu "a plînge"
- MTRANS - Transferul unei informații mentale, de exemplu "a destăinui"
- MBUILD - Construirea informațiilor noi din cele vechi, de exemplu "a decide"
- SPEAK - Producerea sunetelor, de exemplu "a vorbi"
- ATTEND - Concentrarea atenției unui organ de simț asupra unui stimul, de exemplu "a asculta"

Ucenicii acțiunilor, numiți pe scurt AA, reprezintă atribute sau proprietăți ale acțiunilor primitive.

Cazuri conceptuale

- Caz obiect (o)
- Caz direct (D)
- Caz instrument (I)
- Caz recipient (R)

Timperi conceptuale

- | | |
|------------------------|---------------------------|
| Condițional (c) | Interogativ (?) |
| Continuu (k) | Tranziție (t) |
| Prezent (nil) | Începutul tranziției (ts) |
| Trecut (p) | Sfârșitul tranziției (tf) |
| Viitor (f) Negativ (/) | Atemporar (delta) |

Timperile conceptuale oferă o modalitate de modificare a descrierii evenimentelor prin indicarea timpului, modului sau aspectului unui verb, deci a informațiilor care apar de obicei în limbajul natural.

Dependențe conceptuale

Aceste dependențe reprezintă reguli semantice de formare a structurilor de dependențe pe baza entităților și a acțiunilor, cum ar fi relația între un actor și un eveniment sau relația între o acțiune primitivă și un instrument. Există o serie de astfel de dependențe conceptuale stabilite, o parte dintre cele mai importante fiind prezentate în Figura 2. În această figură prima coloană conține forma dependenței conceptuale, a doua coloană conține un exemplu de utilizare a dependenței conceptuale și a treia coloană conține exprimarea în limbaj natural a acestei dependențe.

Elementele componente prezentate sunt folosite pentru crearea *structurilor conceptuale*, numite și structuri de dependențe sau conceptualizări, care pot descrie adecvat înțelesul unei propoziții în limbaj natural. O astfel de structură conceptuală se reprezintă printr-o formă particulară de graf. În reprezentarea grafică a dependențelor conceptuale nodurile reprezintă entități sau acțiuni, săgețile indică direcția dependențelor, săgețile duble indică legături bidirecționale între actori (entități) și acțiuni. În plus, toate legăturile (arcele) din graf au asociate etichete corespunzătoare relațiilor de caz, deci cazurilor conceptuale, sau relații temporare, deci timpurilor conceptuale.

Regula 1 descrie relația între un actor și evenimentul cauzat de acesta. Relația de dependență este bidirecțională deoarece nici actorul nici evenimentul nu pot fi considerate mai importante în relație. Litera p cu care este etichetată legătura indică timpul trecut.

Regula 2 descrie relația între un producător de scenarii și un ucenic de scenarii. Multe descrieri de stări, cum ar fi înălțime, sunt reprezentate în dependențele conceptuale prin valori cuantificate.

Regula 3 descrie relația între o acțiune primitivă și instrumentul cu care se execută acea acțiune. Instrumentul trebuie să fie întotdeauna o conceptualizare, i.e. trebuie să conțină o acțiune și nu un simplu obiect fizic.

Regula 4 descrie relația între o conceptualizare și o altă conceptualizare care a cauzat-o pe prima. Se observă că săgețile indică dependența unei conceptualizări de o alta și de aceea au o direcție opusă direcției implicației existente între concepte. Cele două forme ale regulii descriu cauza unei acțiuni (a) și cauza unei schimbări de stare (b).

Regula 5 descrie relația între o conceptualizare și o altă conceptualizare ce reprezintă timpul în care s-a produs prima.

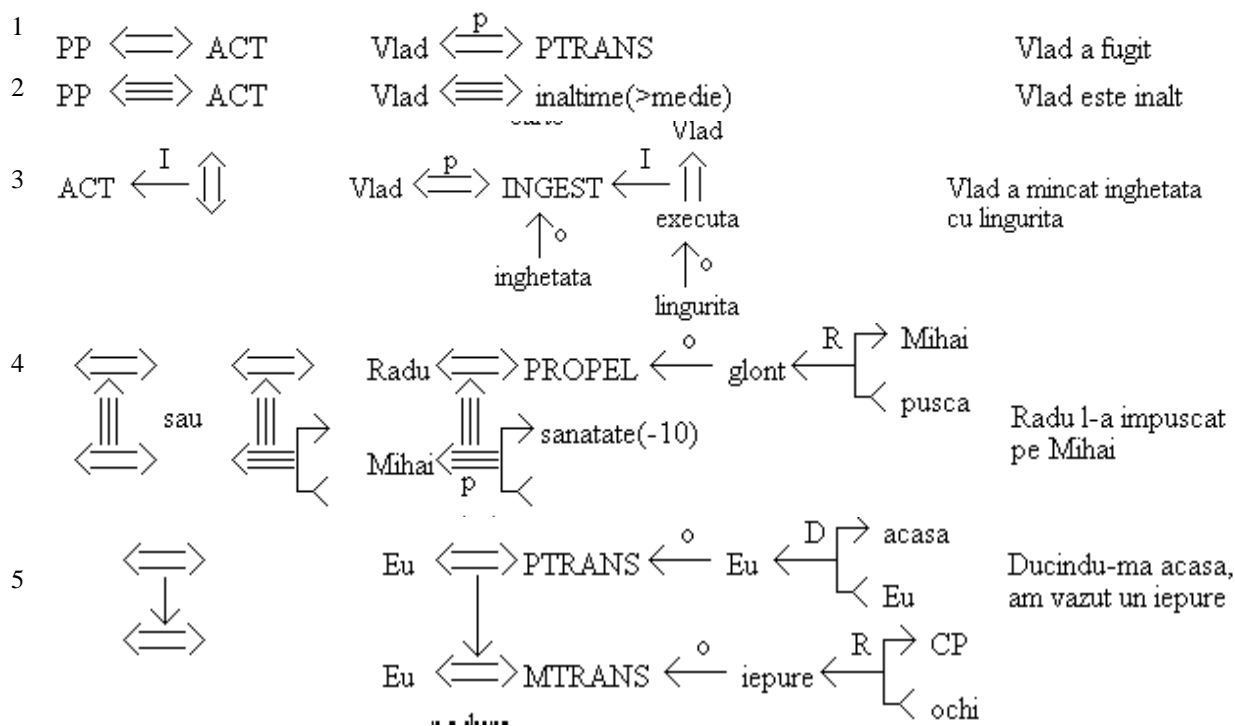


Fig. 2. Dependente conceptuale

Exemplul acestei reguli pune în evidență modul în care dependențele conceptuale imită o prelucrare informațională specifică modelului mental uman: acțiunea "a vedea" este reprezentată ca un transfer de informație între ochi și un element de prelucrare conștient.

Utilizând aceste structuri conceptuale de bază, se pot construi structuri complexe corespunzătoare semnificației propozițiilor în limbaj natural.

Din punct de vedere al reprezentării cunoștințelor și a facilităților de raționament, dependențele conceptuale oferă următoarele avantaje:

- (1) Numărul de reguli de inferență necesare rezolvării problemei scade.
- (2) Multe reguli de inferență sunt deja incluse în reprezentare.
- (3) Structura inițial construită pentru reprezentarea informației dintr-o propoziție conține șabloane care trebuie completate. Aceste șabloane pot fi folosite ulterior ca un mecanism de focalizare a atenției programului pentru înțelegerea propozițiilor următoare.

4. Concluzie

Sumarizarea automată a textului ia unul sau mai multe texte și extrage "cea mai importantă" informație sau informația legată de aspectul ales de către utilizator. Există sisteme care ar putea fi foarte utile pentru cercetători care au nevoie să stabilească rapid conținutul unui articol. Din păcate, cu tehnologia actuală, este dificil să se producă automat un rezumat, care să înlocuiască întregul document.

Reprezentarea cunoștințelor prin dependență conceptuală presupune o abordare decompozițională, care analizează propoziția în limbaj natural pentru a reține structura sa de bază, elementele semantice

fundamentale și facilitează execuția inferențelor determinate de informațiile conținute în propoziții. Ceea ce va duce la ușurarea sumarizării informației.

Bibliografie

1. Mani, Inderjeet (2001). *Automatic Summarization*. ISBN 1-58811-060-5
2. Marcu, Daniel (2000). *The Theory and Practice of Discourse Parsing and Summarization*. ISBN 0-262-13372-5
3. Schank, R., "A Conceptual Dependency Representation for a Computer-Oriented Semantics", Ph.D. Thesis University of Texas, Austin 1969
4. Schank, R.C. (1975). *Conceptual Information Processing*. New York: Elsevier.
5. Roger C. Schank, *The Cognitive Computer, On Language, Learning and Artificial Intelligence*. Addison, Wesley Publishing Company, Inc., Reading. 1984