

# Preprocesarea textelor în sistemele de tip „întrebare-răspuns”.

Victoria BOBICEV

Technical University of Moldova

victoria\_bobicev@rol.md

**Rezumat** — Sisteme de întrebare-răspuns (Question-Answering systems) au nevoie de procesare intensivă a textului, etapele de preprocesare fiind analiza morfologică și sintactică. În lucrarea dată este descris procesul de analiză a textelor în sistemul de întrebare-răspuns ce este creat în cadrul proiectului „Cercetarea în domeniul de Regăsire a Informației în scopul creării sistemului electronic de informare publică”. Prima etapă descrisă este marcarea morfologică în baza dicționarului, a doua este dezambiguizarea morfologică utilizând metoda statistică. A treia etapă descrisă este analiza sintactică parțială, numită „chunking” ce definește grupurile sintactice în propoziții.

**Cuvinte cheie** — analiza morfologică automată, analiza sintactică parțială, dezambiguizarea morfologică, procesarea textului, sisteme de întrebare-răspuns.

## I. INTRODUCERE

Sistemele de întrebare - răspuns (în engleză "question answering systems", sau sisteme QA) sunt caracterizate prin faptul că primesc întrebări formulate în limbaj natural și, în baza unei colecții de documente, extrag răspunsul sau un set de răspunsuri găsite în documentele date. Astfel de sisteme sunt considerate ca fiind următorul pas în evoluția motoarelor de căutare a informației în surse textuale [1].

Ca regulă astfe de sistem constă din următoarele module:

1. Modulul de analiză a întrebării – transformă întrebările formulate în limbaj natural uman în interogări pentru motorul de achiziție de documente;
2. Modulul de achiziție de articole – caută în colecția de articole articolele relevante pentru întrebarea formulată de utilizator, în baza datelor primite de la modulul de analiză a întrebării;
3. Modulul de extragere a răspunsului – din colecția de articole returnate de modulul de achiziție de articole, extrage un răspuns succint care constituie răspunsul în limbaj natural uman la întrebarea utilizatorului.

Astfel de sistem are nevoie de pre-procesări elementare asupra textelor din colecția: segmentare lexicală, etichetare morfo-sintactică (POS tagging), recunoașterea grupurilor frazale (chunking), determinarea dependențelor între cuvinte sau grupuri frazale (dependency linking) [7].

Sistemul ce se crează în cadrul proiectului „Cercetarea în domeniul de Regăsire a Informației în scopul creării sistemului electronic de informare publică” la prima etapă va funcționa în baza documentelor prestate de IDSI (Institutul de Dezvoltare a Societății Informaționale), și anume:

1. Codul cu privire la știință și inovare al Republicii Moldova
2. H O T Ă R Î R E cu privire la aprobarea Acordului de parteneriat între Guvern și Academia de Științe a Moldovei pentru anii 2009-2012
3. Cu privire la parcurile științifico-tehnologice și incubatoarele de inovare
4. Știință și inovare.

Documentele acestea au fost preprocesate cu scopul simplificării căutării și extragerii răspunsului.

## II. MARCAREA MORFOLOGICĂ

Problema marcării automate morfologice a fost studiată intens pe parcursul anilor. Primele încercări de a rezolva problema aceasta sunt bazate pe cunoștințe [8]. Mai târziu au fost efectuate încercările de a aplica metode statistice pentru dezambiguizarea morfologică a cuvintelor textului. Au fost executate încercările de a aplica diferite sisteme statistice în baza lanțului Markov.

În [6] este descrisă o aplicație a metodei statistice pentru marcarea automată a textului român cu codurile morfologice. Metoda aceasta este bazată pe modelul Markov și este practic independentă de limba textului la care este aplicată. Modelul statistic este antrenat în baza unui text marcat cu codurile morfologice, obținând statistica secvențelor din coduri morfologice. Apoi, având statistica obținută metoda aplică algoritmul lui Viterbi pentru dezambiguizarea marcării cuvintelor textului cu codurile respective. Inițial metoda a fost aplicată pentru textul englez, în [6] aceasta este aplicată pentru marcarea textelor române.

Marcarea morfologică a textelor în rîndul său are nevoie de unele etape preliminare. Una din etapele importante preliminare este împărțirea textului în cuvinte. În cazul textelor procesate etapa dată nu prezintă dificultăți. În multe cazuri etapa împărțirii textului se unește cu marcarea cuvintelor cu codurile morfosintactice cu ajutorul dicționarului. Dicționarul utilizat a fost creat la catedra “Informatica Aplicată” și conține aproximativ 90 000 de forme de cuvinte. Fiecare formă de cuvînt este însoțită de forma inițială a lui (lema) și un cod morfosintactic. Mai jos este prezentat fragmentul dicționarului dat:

alcoolica	alcoolic	Afsry
abrevierilor	abreviere	Ncfpoy
alegînd	alege	Vmg
ales	alege	Vmp--sm

Fiecare intrare în dicționarul dat ocupă un rînd și conține o formă a cuvîntului, forma lui inițială (care este utilizată în dicționare) și codul morfoogic ce descrie forma dată a cuvîntului. De exemplu, codul Afsry descrie Adjectiv,

feminin, singular, cazul nominativ-acuzativ, articulat, fiecare literă a codului descrie o caracteristică morfologică. Codificarea dată a fost propusă în [4].

Folosind astfel de dicţionar obţinem textul marcat cu codurile morfosintactice. În textul marcat sunt trei tipuri de elemente: cuvintele marcate cu un cod morfologic, cuvintele marcate cu câteva coduri morfologice – cuvinte ambigue şi cuvintele care nu au fost marcate cu nici un cod morfologic - cuvintele care nu sunt găsite în dicţionar se marchează ca 'unknown'. Un fragment de text marcat este prezentat mai jos:

Condiţiile	condiţie/Ncfpry
de	=/S
eligibilitate	unknown
a	avea/Va--3s al/Tsfs =/Q
rezidenţilor	rezident/Ncmppy
parcului	parc/Ncmsoy
ştiinţifico-tehnologic	=/Amsrn =/Amson

Cuvîntul „a avea/Va--3s|al/Tsfs|=/Q” în exemplul dat este ambiguu, el poate fi verb, articol sau particula ce este indicat de toate codurile posibile delimitate cu linia verticală. În tabelul 1 este prezentată statistica obţinută pentru textele analizate. Din tabel se observă că cuvintele problematice (ambigue şi necunoscute) formează în jur de 20% din text ce înseamnă că din fiecare 10 cuvinte aproximativ unul este ambiguu şi unul necunoscut.

TABELUL I. STATISTICA CUVINTELOR AMBIGUE ŞI NECUNOSCUTE ÎN TEXTELE ANALIZATE DUPĂ PRIMA ETAPĂ DE MARCARE

Document	N de cuvinte ambigue	N de cuvinte necunoscute	total cuvinte
Codul cu privire la ştiinţă şi inovare	5340 (15%)	2856 (8%)	35552
Hotărîre cu privire la aprobarea Acordului de parteneriat	4179 (10,5%)	3928 (9,9%)	39723
parcurile ştiinţifico-tehnologice şi incubatoarele de inovare	651 (16%)	283 (7%)	4034
Ştiinţă şi inovare	5127 (15%)	2670 (7,9%)	33622
Total	15297 (13,5%)	9737 (8,6%)	112931

Pasul următor este dezambiguizarea morfologică.

### III. DEZAMBIGUIZAREA MORFOLOGICĂ

Modele statistice pentru rezolvarea problemelor prelucrării automate a limbajului natural sunt create în așa mod ca să reprezinte regularitățile statistice prezentate în limbă. Cercetătorii au observat că modelele create pentru comprimarea textului sunt foarte asemănătoare cu modelele create pentru prelucrarea lui. Pentru comprimare se crează un model în baza aceluiași lanț lui Markov. Astfel este logic de presupus că modelele create pentru comprimare

pot fi folosite pentru rezolvarea problemelor prelucrării limbajului natural. Metoda PPM (prediction by partial matching – precizare prin corespundere parțială) este o metodă de comprimare în baza contextului limitat (finite-context modeling), care estimează probabilitatea simbolurilor în baza contextului - simbolurilor precedente în text. În [5] a fost descrisă implementarea practică algoritmului optimizat PPMC, care este considerat cel mai bun algoritm de comprimare a textelor. PPM prezintă o variantă de amestecare când probabilitățile obținute în baza contextelor cu lungimea diferită sunt unite într-o probabilitate comună. În cazul general probabilitatea amestecată a simbolului curent  $p(s)$  poate fi calculată ca:

$$p(s) = \prod_{i=1}^0 p'(c_i | c_{i-0} c_{i-0+1} c_{i-0+2} \dots c_{i-1}) \quad (1)$$

unde

$p'(c_i | c_{i-0} c_{i-0+1} c_{i-0+2} \dots c_{i-1})$  – probabilitățile simbolului curent  $c_i$ , determinate de model pentru toate contextele, începînd cu contextul maximal  $\hat{t}$ .

Modelul acesta poate fi creat în baza literelor textului, în baza cuvintelor sau codurilor morfologice. În cazul nostru modelul este creat în baza secvențelor de coduri în text marcat cu codurile morfologice.

Modelul PPM descris este aplicat pentru rezolvarea problemei ambiguității cuvintelor, selectarea codului morfologic potrivit pentru fiecare cuvînt. Cuvintele ambigue ca regulă au două sau trei coduri atașate din care sistemul trebuie să selecteze unul corect. Cuvintele necunoscute prezintă cea mai mare problemă pentru sistem fiindcă pentru cuvintele acestea el trebuie să selecteze din toate codurile posibile. În experimentele efectuate noi intenționat am restrîns numărul de coduri posibile pentru cuvintele necunoscute cu scopul micșorării timpului lucrului programului. În marea majoritate cuvintele necunoscute aparțin părților de vorbire 'deschise' - substantivelor, adjectivelor și verbelor. Acestea sunt afișate de numeroase, că practic nu pot fi prezentate în dicționar în întregime.

În metoda descrisă pentru găsirea codului potrivit este folosită o variantă a algoritmului Viterbi [3]. În urmare este prezentat algoritmul folosit în experimentele executate.

Inițializarea:

Se introduc patru semne speciale care formează contextul primului cuvînt. Pentru primul cuvînt al propoziției se calculează probabilitățile tuturor codurilor posibile pentru cuvîntul dat:

$$P(t_{i-1}^k) = P_{PPM}(t_{i-1}^k | B B B B) \text{ pentru toate } t_{i-1}^k \in t_{i-1}^1 \dots t_{i-1}^n;$$

Iterații:

Pentru toate cuvintele propoziției  $w_i$  de la al doilea pînă la ultimul  $w_2 \dots w_p$ :

Pentru toate codurile posibile cuvîntului dat  $t_{i-1}^k \in t_{i-1}^1 \dots t_{i-1}^n$ :

Pentru toate codurile posibile cuvîntului precedent  $t_{i-2}^k \in t_{i-2}^1 \dots t_{i-2}^n$ :  
Calculăm recursiv probabilitatea codului curent

$$P(t_{i-1}^k) = P_{PPM}(t_{i-2}^k | P_{PPM}(t_{i-1}^k));$$

Dacă probabilitatea aceasta este maximală, memorizăm valoarea codului precedent în vectorul codurilor finale:

$$P(t_{i-1}^k) = t_{i-1}^k \text{ pentru care } \max P(t_i);$$

Sfîrșit

Sfîrșit

Sfîrșit

Pentru toate codurile posibile ultimului cuvînt al propoziției  $t_p^k \in t_p^1 \dots t_p^n$ :

Aflăm valoarea maximală a probabilității  $P(t_p^k)$

Şi memorizăm codul corespunzător în vectorul rezultatelor finale.  
Sfârşit

Folosind metoda PPM cu contextul maximal patru coduri precedente am efectuat antrenarea sistemului în baza corpusului de texte marcate morfologic, verificate şi corectată manual. În baza probabilităţilor calculate cu ajutorul algoritmului descris mai sus am executat dezambiguizarea automată a fişierelor deja analizate. Rezultatele obţinute verificate manual sunt prezentate în tabelul 2.

TABELUL 2. STATISTICA CORECTITUDINII  
DEZAMBIGUIZĂRII MORFOLOGICE ÎN TEXTELE  
ANALIZATE

Document	N de cuvinte prelucrate	N de cuvinte dezambiguizate corect	N de cuvinte dezambiguizate greşit
Codul cu privire la ştiinţă şi inovare	22,5%	18,9%	3,6%
Hotărâre cu privire la aprobarea Acordului de parteneriat	19,5%	13,3%	6,2%
parcurile ştiinţifico-tehnologice şi incubatoarele de inovare	23,8%	20%	3,8%
Ştiinţă şi inovare	18,8%	13,8%	5%
Total	21,15%	16,5%	4,65%

Din tabel se observă că cuvintele prelucrate sunt cele ambigue şi cele necunoscute ( 8% + 13% = 21% ). Numărul de cuvinte prelucrate este cu o fracţie de procent mai mică decât suma celor ambigue şi necunoscute din cauză că unele necunoscute nu sunt de fapt cuvinte. În text se întîlnesc unele elemente care nu au fost recunoscute de sistem şi se considerau cuvinte necunoscute (de exemplu, cifre romane: III, V, X). Aceste erori nu au fost luate în consideraţie pe parcursul evaluării calităţii dezambiguizării din cauză că erorile au fost comise la prima etapă de marcare.

A treia coloniţă a tabelului conţine procentul cuvintelor cu coduri atribuite greşit în procesul de dezambiguizare. Se observă că procentul acesta diferă considerabil de la un document la altul. Media cuvintelor cu codul atribuit greşit este 4,65% ce de fapt presupune aproape un cuvînt greşit din fiecare 20 cuvinte din text. Însă analiza mai detaliată a cuvintelor cu coduri greşit detectate arată că în unele cazuri partea de vorbire a fost determinată corect, eroarea a fost comisă în detectarea unelor caracteristici morfologice, de exemplu, cazului sau numărului substantivului. Astfel de erori afectează analiza etapelor următoare, şi anume, analiza sintactică. Dacă calculăm procentajul erorilor grave, cînd este greşit detectată partea de vorbire, obţinem 3,15% din numărul total al cuvintelor textului. Rezultatul acesta

este destul de bun şi practic la acelaşi nivel ca şi alte sisteme de dezambiguizare morfologică [6].

#### IV. CHUNKING

Etapă următoare după analiza morfologică finalizată este analiza sintactică, însă această sarcină este extrem de complicată şi practic nu poate fi realizată în baza textelor documentelor date. Astfel, în locul analizei sintactice totale este propusă analiza sintactică parţială care a obţinut denumirea “chunking” provenind de la cuvîntul “chunk” – bucată, fragment. Textul este împărţit în fragmente care reprezintă grupuri sintactice, şi anume, grup nominal grup verbal, grup prepoziţional şi altele.

În [2] a fost descrisă metoda de creare a gramaticii pentru analiza sintactică parţială în mod semiautomat. Gramatica dată prezintă un set de reguli sintactice ce sunt aplicate la un text marcat morfologic. Regula sintactică a gramaticii date este formată din coduri morfologice şi reprezintă o secvenţă a codurilor ce se unesc într-un grup. De exemplu, pentru grup nominal se formează astfel de reguli:

Nc Ts Ncmsoy Ams  
Nc Ts Ncmsoy Amp  
Nc Ts Ncmsoy Ams  
Nc Ts Ncmsoy Amp

în care sunt descrise codurile cuvintelor consecutive: Nc (Noun, common) – substantiv comun, Ts – articol posesiv, Ncmsoy – substantiv comun, masculin, singular, cazul dativ-genitiv, articulat, Ams – adjectiv, masculin, singular. Regulele date specifică caracteristicile morfologice ale cuvintelor care indică acordul între cuvinte în grupul dat. De exemplu, în grupul nominal este indicat acordul în gen şi număr între substantiv şi adjectiv.

În lucrarea menţionată au fost create doar reguli pentru grupuri nominale. Grupuri nominale sunt de fapt elementele principale în text, mai ales în textele formale care sunt procesate în lucrarea dată. Statistica părţilor de vorbire în textele date demonstrează că numărul de substantive în text este aproximativ de cinci ori mai mare decît numărul de verbe. În propoziţie pe lîngă un verb apar aproximativ cîte cinci substantive care formează grupuri nominale complicate cu articole, adjective şi pronume asociate. Specificul documentelor de aşa tip este aşa numită „nominalizarea” verbelor, cînd acţiunea în propoziţie este redată de substantiv format din verb (verb nominalizat). De exemplu: „a accepta condiţiile” se transformă în „acceptarea condiţiilor”, „a asigura integritatea” respectiv în „asigurarea integrităţii”, „a aproba acordarea statutului” – în „aprobarea acordării statutului”. De menţionat că în ultimul exemplu sunt două verbe nominalizate şi astfel de combinaţii pot fi prelungite din care cauză este uneori greu de înţeles mesajul documentelor formale. La fel de grea este procesarea automată a acestui tip de documente. Una din problemele de bază în traducerea automată este anume procesarea grupurilor nominale compuse din mai multe substantive („compound nouns”).

Altă problemă dificilă prezintă grupurile prepoziţionale. Prepoziţiile joacă un rol important în formarea lanţurilor de grupuri nominale, ca de exemplu: „domeniile de activitate prevăzute de proiectele de inovare şi transfer tehnologic”. Însă în procesul unirii noi ne conducem de sens şi deosebit cazul „a fost aprobat proiectul de Consiliu Suprem” şi cazul „a fost aprobat proiectul de transfer

tehnologic”. În primul caz „proiectul” și „de Consilium Suprem” sunt două grupuri atașate ambele la grup verbal: „a fost aprobat proiectul” și „a fost aprobat de Consilium Suprem”, într-a doilea caz este un grup „proiectul de transfer tehnologic”. Problema atașării grupurilor prepoziționale este la fel una din problemele nerezolvate în procesarea limbajului natural. În lucrarea de față problema aceasta la fel nu este rezolvată; rezolvarea ei este lăsată pe viitor.

Astfel, pentru documentele date au fost formate 53 reguli sintactice, 45 pentru grupuri nominale și 8 pentru grupuri verbale. Textele documentelor analizate au fost marcate utilizând regulile date. În urmărire este prezentat un fragment de text marcat.

```
<CHUNK type=NP morph=Ncfpry head= Condițiile len=1>
  Condițiile condiție/Ncfpry
</CHUNK>
  de =/S
<CHUNK type=NP morph=Ncfsmr head=încetare len=3>
  încetare ???/Ncfsmr
  al/Tsfs
  statutului statut/Ncmsoy
</CHUNK>
  de =/S
<CHUNK type=NP morph=Ncmrsm head=rezident len=4>
  rezident =/Ncmrsm
  a al =/Tsms
  parcului parc/Ncmsoy
  științifico-tehnologic =/Amsrn
</CHUNK>
  și =/C
<CHUNK type=NP morph=Ncmsoy head=statutului len=2>
  a al/Tsfs
  statutului statut/Ncmsoy
</CHUNK>
  de =/S
<CHUNK type=NP morph=Ncmrsm head=rezident len=3>
  rezident =/Ncmrsm
  al =/Tsms
  incubatorului incubator/Ncmsoy
</CHUNK>
  de =/S
<CHUNK type=NP morph=Ncfsmr head=inovare len=1>
  inovare =/Ncfsmr
</CHUNK>
```

Din fragmentul de text prezentat se observă că fiecare grup sintactic este marcat cu tagul <CHUNK> </CHUNK>, elementul inițial conține attribute type ce indică tipul grupului, head și morph care indică elementul principal în grupul dat și caracteristicile lui morfologice. Informația aceasta este necesară pentru etapele următoare de analiză sintactică care presupune unirea grupurilor date în grupuri mai mari, repetînd acțiunea de unire pînă ce toată propoziția va fi analizată. În procesul acesta este important de luat în considerație informația morfologică și de respectat acordul morfologic între elementele propoziției. La fel de importantă este detectarea elementului principal în grupurile deja marcate. Din exemplul prezentat se observă că o mare parte a textului este marcată, în afara tagurilor rămînd preponderent prepozițiile. Prepozițiile se unesc cu grupuri nominale care le urmează și formează grupuri prepoziționale. De exemplu: „de încetare al statutului” constă din prepoziție „de” și grup nominal „încetare al statutului”. În cazul dat regula este unică. Mai dificilă este problema unirii grupului prepozițional la textul

ce se află de partea stînga a prepoziției, la textul precedent. Un exemplu de atașare problematică este deja prezentat în articolul dat.

## V. CONCLUZIE

Sisteme de întrebare-răspuns (Question-Answering systems) au nevoie de procesare intensivă a textului, etapele de preprocesare fiind analiza morfologică și sintactică. În lucrarea dată este descris procesul de analiză a textelor în sistemul de întrebare-răspuns ce este creat în cadrul proiectului „Cercetarea în domeniul de Regăsire a Informației în scopul creării sistemului electronic de informare publică”. Prima etapă descrisă este marcarea morfologică în baza dicționarului, a doua este dezambiguizarea morfoogică utilizînd metoda statistică. Metoda statisitică propusă rezultă în text marcat cu rata aproximativă de erori 3-5%, ce este la nivel cu alte sisteme pentru limba română. A treia etapă descrisă este analiza sintactică parțială, numită „chunking” ce definește grupurile sintactice în propoziții.

## REFERINȚE

- [1] R. Baeza-Yates, B. Ribeiro-Nieto, Modern Information Retrieval, Addison Wesley, 2000.
- [2] Bobicev V. Metoda semiautomată de creare a gramaticii pentru analiza sintactică a grupurilor nominale în textele române. International Conference Trends in the Development of the Information and Communication Technologies in Education and Management, ASE, Moldova, 2003, pp. 301-304.
- [3] Bobicev V., Popescu A. Marcarea morfo-sintactică automată folosind modelul PPM. Conferința jubiliară Tehnico-științifică a colaboratorilor, doctoranzilor și studenților consacrată celei de-a 40-a aniversări a doctoranturii U.T.M., Moldova, 2006, pp. 155-158
- [4] Tomaz Erjavec, Nancy Ide, Dan Tufis: Development and Assessment of Common Lexical Specifications for Six Central and Eastern European Languages, ALLC-ACH '98 Conference, 1998.
- [5] Moffat, A., 1990. Implementing the PPM data compression scheme. IEEE Transactions on Communications, Vol. 38, No. 11, pp. 1917-1921.
- [6] Dan Tufiş, Oliver Mason: Tagging Romanian texts: A Case Study for QTAG, a Language Independent probabilistic tagger, First International Conference on Language Resources and Evaluation, 1998.
- [7] Dan Tufiş, Dan Ștefănescu, Radu Ion, and Alexandru Ceaușu. RACAI's Question Answering System at QA@CLEF 2007. In Alessandro Nardi and Carol Peters, editors, Working Notes for the CLEF 2007 Workshop, pages 15–21, 2007.
- [8] Dan Tufiş, Octav Popescu, “A Knowledge-Based Approach to Morpho-lexical Processing of Natural Language”, in Proceedings of the International Conference for Young Computer Scientists, Beijing, 1991.