

Dan Ștefănescu, Radu Ion, Alexandru Ceașu, and Dan Tufiș. Sistem întrebare-răspuns antrenabil pentru limba română.

In Adrian Iftene, Horia-Nicolai Teodorescu, Dan Cristea, and Dan Tufiș (eds.) Resurse Lingvistice și Instrumente pentru Prelucrarea Limbii Române, pp. 153-164, Iași, Romania, septembrie 2010.

Universitatea "Al.I. Cuza" Iași, Editura Universității "Al.I. Cuza" Iași. ISSN 1843-911X.

SISTEM ÎNTREBARE-RĂSPUNS ANTRENABIL PENTRU LIMBA ROMÂNĂ

DAN ȘTEFĂNESCU, RADU ION, ALEXANDRU CEAȘU, DAN TUFIȘ

Institutul de Cercetări pentru Inteligență Artificială, Academia Română

{danstef, radu, aceausu, tufis}@racai.ro

Rezumat

Lucrarea prezintă un sistem întrebare-răspuns dezvoltat la *Institutul de Cercetări pentru Inteligență Artificială* – ICIA în cadrul unui proiect național și evaluat independent în contextul competiției europene CLEF. Evaluarea a fost realizată în cadrul exercițiului *ResPubliQA* pentru limba română. Este descris modul de combinare a diferiților factori de relevanță pe baza cărora sistemul identifică paragrafele cele mai relevante ca răspunsuri la întrebările formulate în limbaj natural. Sistemul este disponibil on-line pe pagina de servicii web a ICIA. El este însă complet antrenabil, funcționalitatea sa fiind independentă de registrul lingvistic ce caracterizează datele de antrenare.

1. Introducere

Cercetările privind prelucrarea automată a limbajului natural (PLN), domeniu central al inteligenței artificiale, produc rezultate cu impact din ce în ce mai mare în societatea globalizată de fenomenul Internet. Sistemele de întrebare-răspuns (ÎR) în limbaj natural, în vogă în perioada anilor '70-80, au revenit în centrul cercetărilor PLN dar, de data aceasta având ca obiectiv identificarea răspunsurilor la întrebări arbitrare în spații de căutare incomparabil mai mari, la limită întregul web. Conținutul informațional al acestui spațiu virtual este atât de mare încât se consideră că orice solicitare rațională de informație își poate găsi măcar un răspuns pe Internet. Evaluarea calității răspunsurilor și respectiv asigurarea găsirii lor sunt însă probleme de cercetare, pentru care abordările tradiționale au devenit insuficiente. Astfel, implementările motoarelor de căutare moderne recurg din ce în ce mai mult la tehnici PLN, acestea fiind utilizate în toate etapele fluxului de prelucrare, începând de la nivelul specificării întrebării și până la extragerea fragmentului de text relevant dintr-unul sau mai multe documente. Odată cu utilizarea tehnicilor PLN au apărut și campaniile de evaluare în domeniul regăsirii inteligente a informației. Acestea constituie astăzi priorități ale cercetării de avangardă dedicată spațiului digital al cunoașterii. Ele au fost organizate mai întâi în SUA (*MUC-Message Understanding Conference*, *TREC-Text Retrieval Conference*, *DUC-Document Understanding Conference*, devenită *TAC-Text Analysis Conference*). În Europa, manifestarea similară este CLEF (Cross Language Evaluation Forum), ajunsă în anul 2009 la a 10-a ediție. Având ca subiect al analizei în primul rând limbile Uniunii Europene, începând cu anul 2006 limbile europene “cu resurse electronice limitate” (română, bulgară, cehă, greacă, portugheză etc.) au devenit “subiecte” de concurs.

În acest context, proiectul național SIR-RESDEC (lansat în 2007) a răspuns unei priorități europene, propunându-și realizarea unui sistem de ÎR în limbaj natural la nivelul celor mai avansate sisteme ale cercetării internaționale. Consorțiul SIR-RESDEC format din cercetători de la ICIA, UAIC și ICI și-a concentrat eforturile în

direcția realizării unor sisteme de ÎR pentru limbile română și engleză, în domeniul legislației Uniunii Europene și respectiv în domeniul geneticii umane. Grupul de cercetare al ICIA și-a concentrat eforturile în direcția realizării sistemului de ÎR pentru limba română în domeniul legislativ, având la dispoziție corpusul JRC-Acquis (Steinberger et al., 2006). În restul articolului de față vom descrie acest sistem.

2. *Considerente preliminare*

Pentru a testa performanțele sistemului, ICIA s-a înscris în anul 2009 în competiția QA@CLEF la secțiunea *ResPubliQA*¹ (Peñas et al., 2009), urmând tradiția participărilor la competițiile CLEF încă din 2006. Sarcina sistemelor dezvoltate de echipele înscrise la *ResPubliQA* a fost să identifice automat paragrafe relevante pentru răspunsuri la întrebări formulate în limbaj natural, în domeniul juridic acoperit de corpusul paralel de lucru al competiției. Pentru prima oară, evaluările sistemelor de întrebare-răspuns în limbaj natural au putut fi comparate interlingual, întrucât întrebările de test (500) au fost aceleași în 8 limbi (bască, bulgară, engleză, franceză, germană, italiană, română și spaniolă) răspunsurile trebuind a fi căutate în corpusul paralel (JRC-Acquis) al legislației europene „Acquis Communautaire” disponibil în toate limbile Uniunii Europene. Alinierea la nivel de paragraf a corpusului pentru toate limbile implicate a oferit posibilitatea evaluării răspunsurilor sistemelor indiferent de limba de interogare. În plus, tot în premieră, organizatorii ResPubliQA au calculat, pentru fiecare limbă, performanțele unui sistem de regăsire documentară (RD) de ultimă generație, dar fără componente de PLN. S-a urmărit în acest mod evaluarea cantitativă a rolului tehnologiilor de prelucrare a limbajului natural față de tehnicile standard utilizate în regăsirea informațiilor. Notând cu $A_{\text{ÎR}}$ acuratețea sistemului de ÎR (v. Secțiunea 6) și cu A_{RD} acuratețea sistemului de RD, atunci raportul $M = \frac{A_{\text{ÎR}}}{A_{\text{RD}}}$ cuantifică meritul tehnicilor de prelucrare a limbajului natural. O cifră de merit supraunitară semnifică faptul că prelucrarea limbajului natural îmbunătățește performanța unui sistem de regăsire documentară. Deși pare intuitiv ca pentru orice sistem de ÎR cifra de merit M să fie supraunitară, organizatorii ResPubliQA au constatat că dintre 28 de sisteme evaluate doar 14 au avut o cifră de merit supraunitară (v. Table 10 în (Peñas et al., 2009)).

Cu experiența dobândită în competițiile CLEF precedente (Pușcașu et al., 2007; Tufiș et al., 2008c; Ion et al., 2009a) și cerințele specifice ale competiției din 2009, obiectivele tehnice principale au fost perfecționarea modulului de regăsire a paragrafelor relevante și respectiv a celui de reordonare a acestora, pe baza unei analize complexe a relevanței paragrafelor candidat. Un impact semnificativ l-a constituit implementarea unei metode similare metodei de optimizare MERT (Och, 2003), în cadrul etapei de reordonare a paragrafelor. Am păstrat abordarea din anii precedenți în ceea ce privește construcția sistemului, însă diferitele module care îl alcătuiesc au fost implementate ca servicii și/sau aplicații Web (Tufiș et al., 2008b):

- *serviciul de analiză a întrebării*² cu ajutorul căruia se clasifică fiecare întrebare atașându-i-se o etichetă ce indică tipul de răspuns pe care acea întrebare îl cere;

¹ <http://celct.isti.cnr.it/ResPubliQA/>

² <http://shadow.racai.ro/JRCACQCWebService/Service.aspx?WSDL>

- *serviciile de generare a cererilor*³ cu ajutorul căruia o întrebare în limbaj natural este transformată în interogări în limbaj formal compatibile cu motorul de căutare;
- *serviciul de interogare a motorului de căutare*⁴ se ocupă de partea de regăsire a paragrafelor relevante pentru o interogare în limbaj formal furnizată la intrare;
- *modulul de reordonare a paragrafelor* preia rezultatele furnizate de motorul de căutare sub forma unei liste de paragrafe, calculează scoruri suplimentare de relevanță pentru fiecare paragraf și, în funcție de o interpolare liniară (obținută aplicând o optimizare de tip MERT) a acestora, asignează scoruri paragrafelor. Paragraful sau paragrafele cu scorurile cele mai ridicate sunt întoarse utilizatorilor.

În faza de indexare a corpusului JRC-Acquis, am luat în considerare doar textul propriu-zis al documentelor, acesta fiind în prealabil preprocesat cu ajutorul TTL (Ion, 2007) dezvoltat la ICIA. Textul a fost segmentat la nivel de unitate lexicală în funcție de terminologia Eurovoc, adnotat la parte de vorbire, și lematizat.

3. *Identificarea terminologiei*

Având în vedere caracterul juridic specializat al corpusului de lucru, o etapă importantă a fost identificarea și tratarea ca unități lexicale a unor anumite expresii sau termeni multi-cuvânt. Acest lucru a fost realizat cu ajutorul tezaurului multilingv Eurovoc⁵ (Ștefănescu and Tufiș, 2006). Descriptorii Eurovoc sunt termeni tehnici care trebuie să apară cu consecvență în toate documentele juridice pentru toate limbile implicate. Recunoașterea acestora atât în faza de analiză a întrebărilor cât și faza de preprocesare a documentelor (ce trebuie ulterior indexate) devine esențială pentru performanțele oricărui sistem QA pe acest corpus.

În consecință am realizat un modul care, după preprocesarea corpusului, recunoaște termenii Eurovoc și generează unități lexicale corespunzătoare (în acest fel, unui termen multi-cuvânt îi corespunde o singură unitate lexicală). Etapa de identificare a terminologiei se desfășoară de-a lungul a 6 etape: (i) termenii din Eurovoc sunt identificați în corpus, în forma lor de teaur; (ii) pentru fiecare formă ocurență identificată la pasul 1 se extrage secvența de leme implicată și se adaugă unui inventar; (iii) se identifică în corpus toate ocurențele secvențelor de leme din inventarul construit la pasul 2; (iv) termenii astfel identificați sunt aduși la forma ocurență din corpus; (v) fiecărui termen identificat i se asignează un descriptor morfo-sintactic – cum termenii sunt de fapt grupuri nominale, descriptorul asignat termenului este același cu descriptorul centrului grupului nominal; (vi) fiecărei ocurențe a unui termen i se asignează ca leme descriptorul corespunzător Eurovoc.

De exemplu, termenul *adunare parlamentară* apare în corpus cu formele flexionate: *adunarea parlamentară*, *adunările parlamentare*, *adunărilor parlamentare*. Toate aceste unități lexicale primesc ca leme descriptorul Eurovoc *adunare parlamentară*, iar ca descriptor morfosintactic, descriptorul corespunzător centrului grupului (i.e., *adunare*, *adunarea*, *adunările*, *adunărilor*).

³ <http://shadow.racai.ro/QADWebService/Service.aspx?WSDL>

⁴ <http://www.racai.ro/webservices/search.aspx?WSDL>

⁵ <http://en.wikipedia.org/wiki/Eurovoc>

4. Clasificarea automată a paragrafelor și întrebărilor

Specificațiile competiției au definit 5 tipuri de întrebări posibile: (i) *factoid* (factual) – întrebări care cer ca răspuns persoane, locații, instituții, momente în timp, etc.; (ii) *definition* (definiție) – întrebări care cer ca răspuns o definiție; (iii) *procedure* (procedură) – întrebări care cer ca răspuns o procedură juridică; (iv) *reason* (motiv) – întrebări care cer ca răspuns un motiv, o cauză; (v) *purpose* (scop) – întrebări care cer ca răspuns un scop, un obiectiv. Numărul redus de clase și faptul că răspunsul corect pentru o întrebare nu se poate întinde pe mai multe paragrafe au condus la ideea clasificării paragrafelor în funcție de probabilitatea lor de a răspunde la un tip de întrebare sau altul. Etichetarea unui paragraf cu tipul de întrebare la care acel paragraf ar putea răspunde cel mai bine, oferă, în mod evident, posibilitatea reducerii complexității etapei de identificare a documentelor relevante, pe care o vom numi etapă de *regăsire documentară*. Acolo unde tipul întrebării este corect identificat, răspunsul este căutat, în principal, în paragrafele de tip identic cu cel al întrebării. Acest lucru a presupus construcția unui modul de clasificare a întrebărilor care a fost optimizat și antrenat pentru tipurile de întrebări asociate corpusului de tip juridic așa cum vom arata mai jos.

Problema de clasificare a paragrafelor este similară cu cea a selecției propozițiilor dintr-un text în vederea generării automate a rezumatului aceluși text (Ion et al., 2009b). Se observă însă că cele două probleme diferă în ceea ce privește numărul de clase care trebuie considerate și tipul entităților (pe de o parte paragrafe, iar de cealaltă, fraze) ce trebuie supuse procesului de clasificare. În cazul nostru, clasele pe care le-am considerat diferă ușor de cele furnizate de organizatorii competiției în specificațiile date. Astfel, clasele *reason* și *purpose* au fost unite, clasa obținută având denumirea de *reason-purpose*. Motivul constă în dificultatea dezambiguizării automate între cele două clase. Pentru a îmbunătăți precizia în faza de regăsire a paragrafelor relevante, am adăugat o altă clasă, etichetată *delete*, cu scopul de a elimina astfel, încă din faza de căutare, acele paragrafe care nu ar fi putut conține răspunsuri corecte pentru vreo întrebare. În această categorie intră paragrafele ce conțin titluri (e.g., „Articolul 1”), părți de tabele ale căror formatare nu s-a mai păstrat, semnături, etc.

Pentru faza de antrenare a clasificatorului am utilizat o colecție de aproximativ 800 de paragrafe etichetate manual cu următoarele etichete: *factoid*, *definition*, *procedure*, *reason-purpose* și *delete*, iar pentru a testa precizia, am folosit doar 89 de paragrafe. Metoda de clasificare la care am recurs este aceeași pe care am folosit-o și pentru clasificarea întrebărilor în anii precedenți: *principiul maximizării entropiei* (Ratnaparkhi, 1998). Trăsăturile pe care le-am luat în considerare s-au bazat pe cuvinte cheie, descriptori morfo-sintactici, punctuație, lungimea propoziției. Primele 5 cuvinte au fost considerate trăsături, ele având un puternic caracter de discriminare pentru clasele alese. Ca trăsături de tip morfo-sintactic, am luat în considerare descriptorul morfo-sintactic al verbului principal, o altă trăsătură fiind existența sau absența în cadrul propoziției a unui substantiv propriu. Alte trăsături sunt numărul de virgule din propoziție, numărul de ghilimele (indicând numărul de citate), semnul de punctuație cu care se termină propoziția. Din punct de vedere ortografic, o trăsătură importantă este numărul cuvintelor din propoziție care încep cu majusculă. Trăsăturile legate de lungime includ numărul propozițiilor în paragraf și lungimea paragrafului în cuvinte.

Precizia clasificatorului pe datele de test a fost de 94% (Ion et al., 2009b; Ștefănescu, 2010). Deși gradul de încredere statistică a evaluării preciziei este redus datorită numărului mic de 89 de paragrafe conținute de datele de test, rezultatele finale obținute folosind clasele asignate paragrafelor au fost mult îmbunătățite. În condițiile în care doar primele 50 de paragrafe întoarse de motorul de căutare au fost luate în considerare în etapele următoare ale fluxului de prelucrare, am constatat că, folosind clasificarea paragrafelor, în majoritatea cazurilor răspunsurile corecte s-au aflat în printre acestea.

Pentru clasificarea întrebărilor cele mai multe din sistemele actuale QA folosesc un modul specializat pentru a determina ce tip de răspuns ar trebui căutat în corpusurile avute la dispoziție. Desigur, clasificarea se poate realiza în mai multe feluri, cele mai simple metode folosind reguli sub forma unor expresii regulate. În cazul nostru, am apelat din nou la clasificatorul bazat pe maximizarea entropiei. Acesta a fost antrenat pentru a identifica 8 clase: *reason-purpose*, *procedure*, *definition*, *location*, *name*, *numeric*, *temporal* și *factoid*. Clasa *location* cuprinde întrebările care necesită ca răspuns o locație; clasa *name* conține întrebările care cer ca răspuns un nume de persoană, organizație sau numele unei entități (e.g., comisie, țară, etc.); clasa *temporal* cuprinde acele întrebări care cer drept răspuns o dată calendaristică sau un interval de timp; clasa *numeric* conține întrebările care cer drept răspuns un număr („*Câți membri sunt în comisia de ...*”), iar clasa *factoid* conține toate întrebările de tip factoid care nu sunt în clasele tocmai descrise. Trăsăturile pe care le-am folosit sunt următoarele: primul cuvânt de tip *WH* din întrebare, primul verb principal din întrebare, primul substantiv din întrebare, descriptorii morfo-sintactici ai tuturor substantivelor, verbelor, adjectivelor, adverbilor și numeralelor din întrebare, ordinea de apariție a primului verb și a primului substantiv în cadrul întrebării analizate. Desigur, extragerea acestor trăsături survine abia după preprocesarea întrebării.

Pentru antrenare, am plecat de la exemplele furnizate de organizatori și am construit 200 de întrebări (incluzând și acele exemple) pentru care am atașat manual clasa din care fac parte (Ștefănescu, 2010). Modelul a fost construit considerând doar trăsăturile care apăreau de cel puțin două ori în datele de antrenare, în final acesta având o acuratețe de 99% pe aceste date. Deși evaluarea performanțelor folosind datele de antrenare („*biased evaluation*”) trebuie evitată pentru aplicațiile de învățarea automată, ea ne-a permis totuși să ajungem la o selecție mai bună de trăsături. După cum era de așteptat, acuratețea a scăzut pentru cele 500 de întrebări furnizate odată cu startul competiției, însă la numai 97.2%. Scorul nu este atât de surprinzător pe cât pare, datorită diferențelor mari dintre trăsăturile caracteristice claselor considerate. De altfel, am ales clasele astfel încât clasificatorul să nu aibă probleme în identificarea lor corectă. Ne interesează în principal să nu avem erori de clasificare, chiar dacă rămânem cu clase nu foarte rafinate. Cu alte cuvinte, încercăm să restrângem cât mai mult spațiul de căutare fără însă a îngădui prea multe erori de clasificare (Ștefănescu, 2010).

5. Sistemul QA

Sistemul nostru este implementat ca un flux de prelucrare realizat peste arhitectura serviciilor și aplicațiilor web dezvoltate la ICIA. Sistemul este optimizat pentru a maximiza un scor de relevanță global $S(p)$ folosit la identificarea paragrafului cel mai plauzibil a constitui răspunsul adecvat la o întrebare adresată sistemului. Scorul global

$S(p)$ se calculează ca o combinație liniară a unor scoruri asignate paragrafelor în funcție de criteriile de relevanță în raport cu o întrebare furnizată de un utilizator:

$$S(p) = \sum_i \lambda_i s_i, \quad \sum_i \lambda_i = 1 \quad (1)$$

unde s_i ($s_i \in [0,1]$) este unul din următoarele scoruri de relevanță:

- s_1 este 1 când clasificarea întrebării corespunde cu cea a paragrafului, în caz contrar, valoarea sa fiind 0. După cum am arătat, clasele modulului de clasificare a paragrafelor (5 la număr) diferă de cele ale modulului de clasificare a întrebărilor (8) dar între ele există o corespondență bine definită (Ion et al., 2009b);
- s_2 este un scor de similaritate lexicală bazat pe lanțuri lexicale, coeziune lexicală și identificarea perechii verb principal–argument; este scorul care încearcă să cuantifice gradul de similaritate lexicală dintre o întrebare și paragrafele întoarse de motorul de căutare ca fiind relevante pentru acea întrebare. Se calculează după următoarea metodă: pentru o întrebare Q și un paragraf candidat P, construim două liste conținând lemele ce corespund cuvintelor conținut din întrebare, respectiv din paragraf: LQ și LP. Numim lemele din lista LQ *cuvinte cheie*. Pe baza lor asignăm lui P un scor de relevanță. Astfel, s_2 se calculează ca un produs a trei alte scoruri ce caracterizează: (i) distanța semantică dintre lemele celor două liste (DS), (ii) coeziunea cuvintelor cheie în paragraf (CC) și (iii) identificarea în paragraf a unui posibil cuplu verb-argument extras din întrebare (VA). Aceste scoruri au fost descrise pe larg în (Ștefănescu, 2010) și nu vom mai insista asupra lor. Scorul final se calculează ca produsul celor trei scoruri:

$$s_2 = DS \times CC \times VA$$

- s_3 este un scor similar cu scorul BLEU (Papineni et al., 2002) care avantajează paragrafele în care cuvintele cheie ale întrebării apar în aceeași ordine ca în întrebare; ca și scorul precedent, se calculează considerând lemele cuvintelor conținut din întrebare, respectiv paragraf. Ideea implementării acestui scor (Ion et al., 2009b) are la bază observația că de foarte multe ori, în realitate, formularea răspunsului corespunzător unei întrebări conține o parte din acea întrebare. Principiul comparării n-gramelor a fost folosit în cazul de față pentru a evalua similaritatea dintre întrebare și paragrafele candidat.
- s_4 și s_5 sunt scorurile de relevanță pentru paragraf și document întoarse de motorul de căutare.

5.1. Prelucrarea întrebării și alegerea răspunsului

După ce sistemul primește la intrare o întrebare, ea este trimisă serviciului web TTL pentru a fi preprocesată. Este apelat apoi serviciul web care se ocupă de clasificarea întrebărilor pentru a obține tipul de răspuns care trebuie căutat. În următorul pas, folosind informația adnotată după preprocesare, întrebarea este transformată în interogări într-un limbaj formal înțeles de motorul de căutare. Folosim 2 algoritmi pentru a genera două interogări diferite, ambele conținând ca termen de căutare clasa întrebării. Trebuie menționat că în faza de indexare au fost indexate odată cu paragrafele și clasele corespunzătoare acestora, pentru ca reducerea spațiului de căutare să se facă direct din faza de regăsire documentară.

Pentru fiecare din cele două interogări generate, motorul de căutare întoarce două liste L_1 și L_2 conținând fiecare 50 de paragrafe sortate după scorul de relevanță descris de ecuația (1). Răspunsul întors de sistem este acel paragraf care se găsește atât în L_1 cât și în L_2 și care satisface:

$$\operatorname{argmin}_p(\operatorname{rang}_1(p) + \operatorname{rang}_2(p)), \quad \operatorname{rang}_1(p) \leq K, \operatorname{rang}_2(p) \leq K, K \leq 50 \quad (2)$$

unde $\operatorname{rang}_1(p)$ este poziția paragrafului p în L_1 , iar $\operatorname{rang}_2(p)$, poziția paragrafului p în L_2 . Dacă nu există un paragraf care să satisfacă condițiile din (2), atunci sistemul semnalizează că nu poate găsi un răspuns la întrebarea dată, întorcând șirul de caractere NOA (*no answer* – nu există răspuns). În urma experimentelor, cele mai bune rezultate au fost obținute pentru $K=3$.

Sistemul dezvoltat este antrenabil, ponderile λ_i folosite în interpolarea liniară a scorurilor de relevanță fiind obținute cu ajutorul unei tehnici de optimizare similare cu metoda MERT. Metoda de antrenare pentru găsirea celui mai probabil răspuns a fost folosită pe cele 200 de întrebări utilizate și în cazul clasificatorului pentru întrebări și are următoarele etape: (i) rularea sistemului pentru cele 200 de întrebări și păstrarea primelor 50 de paragrafe întoarse de motorul de căutare pentru fiecare întrebare, ordinea paragrafelor fiind dată DOAR de scorurile de paragraf întoarse de motor (s_4); (ii) calcularea scorurilor s_i , $i = \overline{1,5}$, pentru fiecare din paragrafele extrase; (iii) pentru fiecare combinație de ponderi λ , cu $\sum_{i=1}^5 \lambda_i = 1$ și un pas de incrementare de 10^{-2} , se calculează scorul *Mean Reciprocal Rank* (MRR) (Radev et al., 2002) pentru întreg setul de 200 de întrebări, lista paragrafelor întoarse fiind sortată după ecuația (1); (iv) reținerea combinației de ponderi λ care maximizează scorul MRR.

Cele două moduri de a genera cereri în limbaj formal conduc la rezultate diferite ale motorului de căutare. Putem considera că avem de a face cu două sisteme ale căror rezultate sunt combinate pentru a obține răspunsul final. Astfel, cele două sisteme au fost individual optimizate, în ceea ce privește ponderile λ , fără a avea în vedere, în vreun moment, posibilitatea ca un sistem să întoarcă răspunsuri NOA. După antrenarea celor două sisteme, observând că sistemul de evaluare al competiției *ResPubliQA* favorizează răspunsurile NOA față de cele incorecte, am luat decizia de a introduce ecuația (2) pentru a păstra doar acele răspunsuri pentru care avem un grad de încredere ridicat. La momentul transmiterii rezultatelor, setul de ponderi λ utilizat nu ținea cont de tipul întrebărilor furnizate la intrare. Vom arăta în secțiunea de evaluare că putem obține rezultate îmbunătățite antrenând setul de ponderi pentru fiecare clasă de întrebări.

5.2. Generarea de cereri formale pentru motorul de regăsire a documentelor

Etapa de regăsire documentară este similară cu cea descrisă în (Ion et al., 2009a), folosind și de această dată motorul de căutare Lucene (Hatcher and Gospodnetić, 2004). Cuvintele documentelor au fost filtrate în funcție de descrierea lor morfo-sintactică astfel încât să rămână pentru indexare doar cuvintele conținut. Mai mult, cuvintele au fost normalizate la forma lor lemă. Eventualele erori de dezambiguizare morfo-lexicală sau lematizare au fost luate în calcul și pentru a compensa astfel de erori (mai ales în cazul unităților lexicale absente din dicționarul sistemului TTL) am indexat și forma ocurență pe aceeași poziție cu lema. În acest fel, un termen poate fi căutat atât ca lema, cât și în forma sa ocurență. Am procedat la construirea a doi indecși pentru colecția de

documente: un index pentru paragrafe și unul pentru documente. Dată fiind o anumită cerere, motorul întoarce o listă alcătuită din paragrafele cele mai relevante în raport cu acea cerere. Astfel, avem deja două scoruri calculate pentru un paragraf: scorul de relevanță pentru paragraf în indexul construit pe paragrafe (s_4) și cel de relevanță pentru documentul din care face parte paragraful în indexul construit pe documente (s_5).

Abilitatea sistemului nostru de a întoarce pentru unele întrebări mesajul că nici un răspuns nu a fost găsit se datorează combinării rezultatelor diferite obținute folosind cei doi algoritmi de generare a interogărilor formale către motorul de căutare.

ALGORITMUL DE GENERARE A CERERILOR FORMALE – VARIANTA 1 ia în considerare toate cuvintele conținut ale întrebării (substantive, verbe, adjective și adverbe) cu care construiește o disjuncție de termeni, ce sunt practic lemele acelor cuvinte. Singura condiție impusă asupra includerii unui termen în cerere este ca scorul TF-IDF (Salton and Buckley, 1988) al aceluși termen să fie peste un anumit prag (Ion et al., 2009b; Ștefănescu, 2010). În urma experimentelor, am stabilit ca valoarea acestui prag să fie 9. Folosirea acestui scor vizează eliminarea din interogare a termenilor comuni (cu frecvență de ocurență mare în documente) și deci, cu capacitate discriminatorie mică.

ALGORITMUL DE GENERARE A CERERILOR FORMALE – VARIANTA 2 folosește, asemeni variantei 1, informația obținută prin preprocesarea întrebării. Varianta precedentă a acestui algoritm, dezvoltată pentru căutările pe corpusul Wikipedia (Ion et al., 2009a) a fost extinsă și optimizată pentru noul corpus. Ca și în versiunea precedentă, algoritmul folosește verbele principale ale propoziției și grupurile nominale identificate în faza de preprocesare pentru a genera cererea formală către motorul de căutare. Noutatea constă în faptul că pentru fiecare grup nominal sunt construiți doi termeni ai cererii. (i) Primul este o expresie obținută prin concatenarea cu spațiu între ele a tuturor lemelor cuvintelor de tip conținut în ordinea apariției lor în grupul nominal. Setăm apoi 2 parametri pentru expresia formată. Primul, pe care îl vom numi *parametru de proximitate*, se referă în principal la numărul de cuvinte ce pot fi intercalate între termenii expresiei. Valoarea sa este egală cu 1 plus numărul cuvintelor funcționale ce se găsesc în grupul nominal. Al doilea parametru, pe care îl vom numi *parametru de amplificare*, definește importanța expresiei ca termen al interogării. Dacă acest parametru este setat la o valoare n , acest lucru înseamnă că identificarea ulterioară a termenului echivalează cu identificarea a n termeni obișnuși. Am setat acest parametru la o valoare egală cu numărul unităților lexicale din expresia formată, pentru a favoriza paragrafele care conțin grupuri nominale identice cu cele din întrebare. (ii) Al doilea termen este o expresie Booleană formată doar din conjuncția lemelor corespunzătoare cuvintelor conținut din grupul nominal.

Observăm că primul termen este mai restrictiv pentru că impune o ordine de apariție și o anumită poziționare în documente a cuvintelor din întrebare. În schimb, pentru al doilea termen condițiile sunt relaxate. Urmărim ca atunci când cuvintele se găsesc în texte, dar nu strict în condițiile impuse de primul termen, să putem extrage totuși acele fragmente care le conțin. Ca și în cazul versiunii precedente, pentru fiecare verb principal din întrebare (cu excepția verbelor difuze semantic⁶) se generează un termen corespunzând lemei aceluși verb. Operatorul boolean folosit implicit de Lucene este SAU (OR) și

⁶ Verbe cu putere discriminatorie redusă (a fi, a avea, a putea, etc.) întrucât utilizarea lor este comună în orice domeniu.

astfel, cererile sunt disjuncții logice. Pe de altă parte, există și situația în care nici un paragraf nu conține o anumită expresie. În acest caz, cu ajutorul expresiilor conjunctive se încearcă identificarea celor care conțin cât mai multe cuvinte conținut ce apar și în întrebare.

Cele două cereri furnizate către motorul de căutare conduc la obținerea a două seturi de paragrafe ca rezultate. Putem vorbi astfel de două sisteme, chiar dacă, în foarte multe din componentele ce le compun, aceste sisteme sunt similare.

Una din cele mai utilizate practici pentru a îmbunătăți performanțele unui sistem QA este aceea de a îmbogăți termenii mono-cuvânt din interogările formale cu sinonime extrase din lexicoane sau ontologii lexicale precum Pinceton WordNet (Fellbaum, 1998). Conținând aproximativ 60.000 de mulțimi sinonimice, varianta românească a WordNet-ului (Tufiş et al., 2008a) este o excelentă resursă ce ar putea fi utilizată în acest scop. Cu toate acestea, datorită limbajului juridic specializat, caracteristic corpusului JRC-Acquis, am ales să nu facem uz de ontologia lexicală la nivelul generării interogărilor, ci într-o etapă ulterioară în care a fost necesară calcularea unui scor de similaritate lexicală dintre întrebare și paragrafele având probabilitate ridicată de a conține răspunsul corect.

6. Evaluare și concluzii

Evaluarea acurateței sistemelor înscrise în competiția ResPubliQA a fost făcută de organizatori folosind un program ce implementează formula (3):

$$A = \frac{1}{n} (n_{RC} + n_{NOA} \times \frac{n_{RC}}{n}) \quad (3)$$

unde n reprezintă numărul total de întrebări (500), n_{RC} este numărul de răspunsuri corecte, iar n_{NOA} reprezintă numărul întrebărilor la care am răspuns cu NOA. După cum se observă, formula de evaluare punctează suplimentar (termenul al doilea al formulei) capacitatea unui sistem de a discerne situațiile în care nu are certitudinea necesară de a da un răspuns corect și decide să emită un răspuns de tip „nu știu” (NOA=„no answer”). Ipoteza organizatorilor a fost că acuratețea în aceste situații ar fi aceeași cu acuratețea răspunsurilor corecte ($\frac{n_{RC}}{n}$). Cu alte cuvinte, atunci când sistemul decide că „nu știe” răspunsul corect, el este la fel de precis ca în cazurile în care furnizează un răspuns. Din analiza efectuată după primirea rezultatelor competiției reiese că în cazul sistemului nostru, precizia cu care răspunsurile greșite au fost evitate prin soluții NOA a fost chiar mai mare decât precizia răspunsurilor corecte.

Pentru evaluare, ICIA a trimis organizatorilor competiției două seturi de rezultate. Primul set, ICIA091RORO, conține rezultatele obținute prin combinarea ieșirilor celor două variante ale sistemului fără a folosi clasele întrebărilor ca termeni ai interogărilor formale, în timp ce al doilea set, ICIA092RORO, folosește aceste clase ca termeni. Utilizarea claselor întrebărilor ca termeni în interogările generate a condus la rezultate semnificativ îmbunătățite datorită faptului că în faza de indexare a paragrafelor am adăugat un câmp care să conțină clasa acestora.

Faza de antrenare desfășurată pe cele 200 de întrebări (v. formula (1)) a condus la obținerea următoarelor ponderi λ :

Tabelul 1: Ponderile pentru cele două variante de interogare ale sistemului

	λ_1	λ_2	λ_3	λ_4	λ_5
varianta 1 (interogări generate în funcție de scorul TF-IDF)	0,22	0,1	0,1	0,19	0,39
varianta 2 (interogări generate folosind structura de grupuri propoziționale a întrebărilor)	0,32	0,14	0,17	0,25	0,12

Fiecare din variantele sistemului a primit la intrare cele 500 de întrebări ale competiției, pentru ca apoi să întoarcă o lista a celor mai relevante 50 de paragrafe pentru fiecare întrebare, folosind parametrii din tabelul de mai sus. Dintre acestea, în cazul în care un paragraf a satisfăcut ecuația (2) acesta a fost furnizat ca răspuns. În caz contrar, răspunsul generat a fost NOA. Tabelul 2 conține evaluările oficiale ale rezultatelor trimise de ICIA: ICIA091RORO și ICIA092RORO. Scorul de acuratețe (formula 3) pentru sistemul ICIA092RORO a fost cel mai bun dintre toate cele 44 de seturi de răspunsuri (28 trimise de participanți, plus încă 16 scoruri de referință calculate de organizatori). Scorul de acuratețe pentru sistemul ICIA092RORO a fost al 4-lea.

Tabelul 2: Rezultatele oficiale ale ICIA la competiția ResPubliQA pe limba română

	ICIA091RORO)	ICIA092RORO
Întrebări la care a fost dat un răspuns	393	344
Întrebări la care NU fost dat un răspuns	107	156
Întrebări la care răspunsul a fost CORECT	237	260
Întrebări la care răspunsul a fost INCORECT	156	84
Întrebări la care răspunsul a fost NOA	107	156
scorul c@1	0,58	0,68

Interesant de menționat faptul că evaluarea cifrei de merit $M = \frac{A_{IR}}{A_{RD}}$ (v. (Peñas, et al., 2009) și Secțiunea 2) a pus în evidență mai obiectiv performanța sistemelor înscrise în concurs în raport cu limba de întrebare și răspuns. Cele două sisteme ale ICIA au obținut cele mai mari cifre de merit (ICIA092RORO=1,55, respectiv ICIA091RORO=1,32). Următoarele cifre de merit au fost obținute de sistemele de ÎR pentru limba italiană (1,24) și limba spaniolă (1,175). În raport cu scorul de acuratețe c@1, aceste două sisteme s-au clasat pe locurile 8 și respectiv 12. Aceste rezultate sunt foarte interesante și ele demonstrează un lucru mai puțin evident: pentru limba engleză câștigul de acuratețe adus de tehnologiile de limbaj natural față de tehnologiile standard de regăsire a informației este relativ mic (15%); acest lucru se explică prin faptul că limba engleză are o morfologie foarte simplă dar și prin faptul că cele mai avansate tehnici de regăsire documentară au fost dezvoltate, testate și optimizate pe documente scrise în limba engleză. Explicația se confirmă și prin faptul că analiza scorurilor de acuratețe ale sistemelor pentru diferite limbi cu morfologie semnificativă (germană, franceză, spaniolă, italiană, bulgară) au fost mai mici decât scorurile de acuratețe pentru limba engleză, deși unele dintre sistemele pentru limba engleză au fost elaborate de aceleași echipe care au implementat sistemele pentru limbile lor naționale. De asemenea, ierarhizarea în funcție de cifra de merit M, plasează cel mai bun sistem de ÎR pentru limba engleză⁷ pe locul al 6-lea, în timp ce sistemul pentru limba germană, cu al 17-lea scor c@1, are cifra de merit 1,158 care-l plasează pe poziția a 5-a.

⁷ Sistemul uned092 a obținut al doilea scor de acuratețe (c@1) 0,61, în timp ce cifra sa de merit (M) a fost de 1,151.

SISTEM ÎNTREBARE-RĂSPUNS ANTRENABIL PENTRU LIMBA ROMÂNĂ

Am testat și ipoteza conform căreia antrenarea ponderilor λ în funcție de clasa întrebării va conduce la îmbunătățirea rezultatelor. Am folosit încă odată cele 200 de întrebări pentru antrenare, partiționate în funcție de clasă, și astfel am obținut rezultatele din tabelul 3. Scorurile $c@1$ și M s-au îmbunătățit cu peste 3% față de cel mai bun rezultat precedent.

Tabelul 3: Ponderile rezultate după antrenarea pe fiecare clasă

		λ_1	λ_2	λ_3	λ_4	λ_5
varianta cu interogări generate în funcție de scorul TFIDF	Factoid	0,1	0	0,2	0,4	0,3
	Definition	0,2	0,15	0,05	0,15	0,45
	Reason-Purpose	0,1	0	0,15	0,3	0,45
	Procedure	0,1	0	0,15	0,15	0,6
varianta cu interogări generate folosind structura de grupuri propoziționale a întrebărilor	Factoid	0,15	0	0,3	0,3	0,25
	Definition	0,05	0,5	0,15	0,1	0,2
	Reason-Purpose	0,2	0	0,4	0,2	0,2
	Procedure	0,15	0,1	0,25	0,2	0,3

În final, trebuie să menționăm că sistemul nostru este disponibil on-line⁸, sub forma unei aplicații web ce apelează, pentru fiecare interogare, diferitele servicii web ce intră în arhitectura sistemului. În viitorul apropiat intenționăm să ne concentrăm eforturile în direcția dezvoltării unor capacități cros-linguale care să dea posibilitatea utilizatorilor să interogheze sistemul în limbaj natural fie în limba română, fie în engleză, iar sistemul să întoarcă rezultatele în oricare din aceste limbi. De altfel, secțiunea în limba engleză a corpusului JRC-Acquis a fost deja preprocesată și avem în vedere folosirea aplicațiilor de traducere automată dezvoltate la ICIA (Ceașu, 2009; Irimia, 2009) pentru a traduce fie întrebările în limbaj natural, fie cererile formale către motorul de căutare.

Mulțumiri. Activitatea de cercetare descrisă a fost sprijinită de CNMP – MECT prin proiectul național PNII SIR-RESDEC, D11007/18.09.2007.

Referințe bibliografice

- Ceașu, A. (2009) *Tehnici de traducere automată și aplicabilitatea lor limbii române ca limbă sursă*, București, România: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.
- Fellbaum, C. (1998) *WordNet: An Electronic Lexical Database*, MIT Press.
- Hatcher, E. and Gospodnetić, O. (2004) *Lucene in Action*.
- Ion, R. (2007) *Metode de dezambiguizare semantică automată. Aplicații pentru limbile engleză și română*, București: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.
- Ion, R., Ștefănescu, D., Ceașu, A. and Tufiș, D. (2009a) *RACAI's QA System at the Romanian-Romanian QA@CLEF2008 Main Task, Lecture Notes in Computer Science*, vol. 5706, September, p. 393–400.

⁸ <http://www2.racai.ro/sir-resdec/>

- Ion, R., Ștefănescu, D., Ceașu, A., Tufiș, D., Irimia, E. and Barbu-Mititelu, V. (2009b) *A Trainable Multi-factored QA System*, Proceedings of CLEF2009 Workshop, Corfu, Greece.
- Irimia, E. (2009) *Metode de traducere automată prin analogie. Aplicații pentru limbile română și engleză*, București: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.
- Och, F.J. (2003) *Minimum Error Rate Training in statistical machine translation*, Proceedings of the 41st Annual Meeting on Association for Computational Linguistics, Sapporo, Japan, 160-167.
- Papineni, K., Roukos, S., Ward, T. and Zhu, W.-J. (2002) *BLEU: a method for automatic evaluation of machine translation*, Proceedings of the ACL-2002: 40th Annual meeting of the Association for Computational Linguistics, Philadelphia, USA, 311-318.
- Peñas, A., Forner, P., Sutcliffe, R., Rodrigo, Á., Forascu, C., Alegria, I., Giampiccolo, D., Moreau, N. and Osenova, P. (2009) *Overview of ResPubliQA 2009: Question Answering Evaluation over European Legislation*, Proceedings of the QA@CLEF workshop, Sept. 30 - Oct. 3.
- Pușcașu, G., Iftene, A., Pistol, I., Trandabăț, D., Tufiș, D., Ceașu, A., Ștefănescu, D., Ion, R., Orășan, C., Dornescu, I., Moruz, A. and Cristea, D. (2007) *Cross-Lingual Romanian to English Question Answering at CLEF 2006*, Lecture Notes in Computer Science, Available: ISBN: 978-3-540-74998-1.
- Radev, D.R., Qi, H., Wu, H. and Fan, W. (2002) *Evaluating Web-based Question Answering Systems*, Demo section, LREC 2002, Las Palmas, Spain.
- Ratnaparkhi, A. (1998) *Maximum Entropy Models for Natural Language Ambiguity Resolution*, Philadelphia, PA, USA: PhD Thesis, University of Pennsylvania.
- Salton, G. and Buckley, C. (1988) *Term-weighting approaches in automatic text retrieval*, Information Processing and Management, pp. 513-523.
- Steinberger, R., Pouliquen, B., Widiger, A., Ignat, C., Erjavec, T., Tufiș, D. and Varga, D. (2006) *The JRC-Acquis: A Multilingual Aligned Parallel Corpus with 20+ Languages*.
- Ștefănescu, D. (2010) *Extragere inteligentă de informații din corpusuri multilingve*, București: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.
- Ștefănescu, D. and Tufiș, D. (2006) *Aligning Multilingual Thesauri*, Proceedings of The 5th Language Resources and Evaluation Conference (LREC), Genoa, Italy.
- Tufiș, D., Ion, R., Bozianu, L., Ceașu, A. and Ștefănescu, D. (2008a) *Romanian Wordnet: Current State, New Applications and Prospects*, Proceedings of the 4th Global WordNet Conference (GWC-2008), Szeged, Hungary, 441-452.
- Tufiș, D., Ion, R., Ceașu, A. and Ștefănescu, D. (2008b) *RACAI's Linguistic Web Services*, Proceedings of the 6th Language Resources and Evaluation Conference (LREC 2008), Marrakech, Morocco.
- Tufiș, D., Ștefănescu, D., Ion, R. and Ceașu, A. (2008c) *RACAI's Question Answering System at QA@CLEF 2007*, Lecture Notes in Computer Science: Advances in Multilingual and Multimodal Information Retrieval (CLEF 2007), pp. 3284-3291.