

UNIVERSITATEA "LUCIAN BLAGA" DIN SIBIU

Contract UEFISCDI nr. 27 / 04.08.2010

Programul: Resurse umane

Tipul proiectului: Proiect de cercetare post-doctorală

Cod proiect: PD_670 / 28.07.2010

SINTEZA DE CERCETARE

**SISTEM DE CLASIFICARE AUTOMATA A DATELOR NESTRUCTURATE
FOLOSIND METACLASIFICATOARE BAZATE PE METODE DE TIP SUPPORT
VECTOR MACHINE ȘI NAIVE BAYES**

**Etapa unica 2012
Anul III**

DIRECTOR PROIECT,

Sef lucrări dr. Ing. Daniel MORARIU

iunie, 2012

1 Rezumatul proiectului

Acest proiect se încadrează în domeniul clasificării automate a documentelor text și își propune îmbunătățirea rezultatelor clasificării prin realizare unor strategii de combinare a rezultatelor metodelor de clasificare folosite.

Un număr impresionant de documente se găsesc în format electronic, iar utilizatorul are nevoie de componente de clasificare automată a acestora pentru a le putea gestiona. Gestiunea lor a devenit o problema foarte importantă în ultima perioadă. Devine esențială existența unor programe inteligente de organizare automată a documentelor în categorii pentru a facilita analiza și prelucrarea acestor documente. Datorită domeniului foarte vast în care ar trebui să lucreze acestea, devine dificil de realizat un singur clasificator cu performanțe foarte bune. Abordarea actuală este de a utiliza mai mulți clasificatori de diferite tipuri combinați într-un meta-clasificator, sau realizarea unei clasificări hibride care se bazează pe predicția clasificatorului cel mai bun pentru o problemă particulară folosind caracteristicile vectorilor de intrare ale documentelor și istoria clasificărilor. Având mai mulți clasificatori de bază, de tipuri diferite (SVM - Support Vector Machine, Bayes, rețele neurale, etc), idea este de a învăța un meta-clasificator care prezice gradul de corectitudine pentru fiecare dintre clasificatorii de bază. Meta-etichetarea unei instanțe indică încrederea în clasificarea făcută de acesta, dacă instanța este clasificată corect de către acel clasificator dintre toți ceilalți clasificatori utilizați. Regula de clasificare a meta-clasificatorului este ca fiecare clasificator de bază să atribuie o clasă la instanța curentă și apoi meta-clasificatorul să decidă dacă clasificarea este demnă de încredere sau nu. Pe lângă creșterea acurateții de clasificare, prin exploatarea sinergismului mai multor clasificatoare, un alt avantaj al meta-clasificării constă în posibilitatea de exploatare a paralelismelor funcționale (multiprocesor).

2 Obiectivele proiectului pentru anul 2012

Etapa	Obiective	Activități	Categorii de buget (Valoare lei)
Unică	1. Implementarea unui sistem de răspuns automat la întrebări.	1.1 Realizarea conexiunii la motoare de căutare web. Realizare practica.	19000,00
		1.2 Clasificarea rezultatelor utilizând meta-clasificatoarele anterior antrenate și extragerea informațiilor relevante în sensul unui răspuns adecvat. Realizare practica	19500,00
Total etapa: 38500,00			

3 Sinteza activităților de cercetare realizate în 2012

În această etapă s-a realizat un sistem online de răspuns automat la întrebări care se găsește disponibil la adresa <http://193.226.29.26/ACUD/Default.aspx>. Acest sistem permite introducerea a unei întrebări formulate în limbaj natural în limba engleză. Aplicația realizată recunoaște 5 tipuri de întrebări: Person, Location, Time, Description și Reason. Ca și rezultat deocamdată întoarce o listă de snippet-uri. Această listă este sortată în funcție de încrederea pe care o dă analiza lexicală efectuată asupra snippet-ului. Încrederea este calculată relativ la faptul dacă în acel snippet se găsește sau nu răspunsul așteptat. Aplicația realizată analizează din punct de vedere sintactic cuvintele din interogare pedicționând pentru fiecare cuvânt tipul de vorbire al acestuia în funcție de contextul în care apare acesta în întrebare. În această aplicație contextul considerat este de maxim 3 cuvinte. Cuvintele considerate esențiale pentru interogare sunt transmise motorului de căutare Bing iar de la acesta se așteaptă o listă de maxim 20 de răspunsuri pe care motorul de căutare le consideră ca fiind cele mai relevante pentru cuvintele din interogarea specificată. Fiecare snippet este analizat din punct de vedere sintactic și în funcție de tipul întrebării se caută anumite părți de vorbire în vecinătatea cuvintelor din interogare. Snippetul este etichetat cu o valoare mai mare dacă în preajm cuvintelor specificate în interogare se găsesc cuvinte care sunt etichetate cu tipul de vorbire așteptat.

3.1 Motoare de căutare

3.1.1 Caracteristici generale

O dată cu dezvoltarea sa exponențială, World Wide Web-ul a devenit o sursă imensă de informații, fie ele text, imagini sau multimedia. Lipsa unei structuri centrale și libertatea de a nu urma o sintaxă strictă a făcut însă ca această cantitate vastă de informații să fie disponibilă pe Web, dar regăsirea ei să nu fie atât de ușoară. Pentru accesul eficient la aceste informații utile, sunt absolut necesare abordări și interfețe corespunzătoare pentru căutarea și navigarea prin aceste colecții imense, care să fie eventual capabile să întoarcă doar răspunsurile la căutările utilizatorului.

Motoarele de căutare (sau serviciile de căutare pe Web) sunt cea mai utilizată metodă folosită în zile noastre pentru a accesa informațiile de pe Web. Utilizarea motoarelor de căutare se bazează pe o paradigmă simplă: pentru o interogare (query) formată din unul sau mai multe cuvinte cheie, motorul de căutare răspunde cu o listă de rezultate ale căutării. În cazul general această listă reprezintă o serie de pagini web pe care s-au găsit cuvintele specificate în interogare. În această listă fiecare pagină este descrisă printr-un scurt rezumat numit snippet. Snippet-ul este compus dintr-un titlu, o adresa a documentului (URL-ul acestuia) și un text scurt care evidențiază conținutul documentului referit.

Motoarele de căutare parcurg periodic paginile web și indexează cuvintele din descrierile acestor pagini astfel în momentul căutării vor întoarce acele pagini care conțin cuvinte comune cu interogarea. Primul motor de căutare propriu-zis a fost Archie lansat în 1990 și era un serviciu de căutare a fișierelor FTP. Primul motor de căutare care a încercat să adune conținutul WEB-ului într-o bază de date a fost Wander lansat în 1993. Tot atunci apare și Aliweb care indexa paginile de pe un server astfel încât să devină mai simplă căutarea anumitor informații pe acel server.

Doar din 1995, putem vorbi deja de prima generație de motoare de căutare (1995-1997) care se bazau doar pe informația text din cadrul paginii (frecvența cuvintelor, limba). Exemple de astfel de motoare de căutare care au fost și folosite pe scară largă sunt: Lycos și Excite.

A doua generație a motoarelor de căutare apărute în anii 1998 este reprezentată de apariția sistemelor care folosesc și informații externe paginii propriu-zise. Aceste sisteme analizează link-urile utilizatorului, informațiile despre rezultatele alese de către utilizatori (click-through data) precum și referințele la pagina respectivă și referințele de la pagina respectivă la alte pagini. Cu cât pagina avea mai multe referințe din exterior și mai multe referințe către alte pagini importante era considerată mai relevantă. Cel mai reprezentativ motor de căutare a fost și este în continuare Google-ul.

A treia generație de motoare de căutare este cea actuală, încă la nivel experimental și în plină dezvoltare. Principalele îmbunătățiri ale acestora sunt legate de analiza semantică a textelor indexate, determinarea contextului căutării (locația utilizatorului, profilul acestuia și istoria căutărilor realizate de acesta) precum și îndreptarea atenției spre nevoile utilizatorului nu doar pe căutarea în sine. Motoarele de căutare din această generație sugerează utilizatorului anumite contexte de căutare, corectează greșelile gramaticale și prezintă utilizatorului interfețe cât mai sugestive pentru prezentarea rezultatelor.

Un motor de căutare conține mai multe componente dintre care cele mai importante sunt prezentate în figura 1.

Crawler-ul numit și spider are rol de a descoperii automat noi pagini web pentru ca acestea să poată fi indexate și regăsite ușor ulterior. Această descoperire de pagini se face prin urmărirea hyperlink-urilor din paginile deja vizitate.

Indexarea documentelor se ocupă cu memorarea anumitor informații utile dintr-o pagină web pentru ca aceasta să poată fi regăsită în cazul căutării după anumite cuvinte cheie.

Cache-ul de documente se ocupă cu păstrarea resurselor indexate în forma originală sau păstrarea link-urilor la acele resurse după indexarea acestora

Sistem de ordonare a documentelor – de obicei pentru o căutare sunt întoarse ca rezultate milioane de astfel de linkuri. Este nevoie de metode și tehnici care să ordoneze link-urile întoarse de motorul de căutare astfel încât cele mai relevante și cele mai interesante pentru utilizator să apară primele în listă. Primele motoare de căutare ca Excite sau Infoseek utilizau algoritmi care ordonau

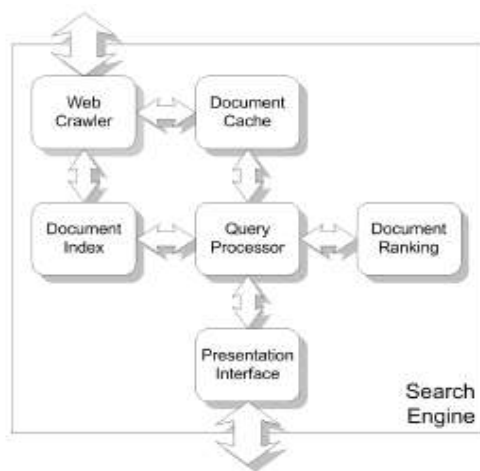


Figura 1. Componentele unui motor de căutare standard

documentele pe baza conținutului acestora. Această abordare avea deficiente majore pentru că se putea profita ușor pentru a-ți aduce propria pagina cât mai sus în ierarhie prin popularea excesivă a paginii cu cuvinte nu neapărat legate de conținutul real al paginii (proces de *spamming*). Motoarele de căutare din a doua generație implementează algoritmi de ordonare bazați pe analiza hyperlink-urilor dintre documentele web. Premisa este că numărul de hyperlink-uri care duc la o pagină, este direct proporțional cu popularitatea și calitatea acelei pagini. Prin urmare, cu cât o pagina are mai multe linkuri către ea, cu atât va apărea mai sus în ierarhie. În plus, este luată în considerare și poziția în ierarhie a paginilor care duc la pagina în discuție, deci cu cât sunt mai importante paginile de proveniență a linkurilor, cu atât îi crește poziția în ierarhie a paginii în discuție. Dezavantajul acestui tip de algoritm ar fi importanța prea mare pe care o acordă paginilor cu o vechime mare (care sigur că au multe legături din exterior) în detrimentul paginilor noi care pot fi uneori mai relevante. Cel mai cunoscut algoritm de acest tip este cel al motorului Google, PageRank.

Unitatea de procesat interogări are rolul de a coordona execuția căutării termenilor introduși de către utilizator. Executantul de interogări comunică cu index-ul de documente, cu cache-ul și cu sistemul pentru ordonarea documentelor pentru a obține o listă a documentelor relevante cu căutarea curentă. Această listă este apoi trimisă interfeței de prezentare a rezultatelor.

Interfața de prezentare a rezultatelor (Presentation Interface) se ocupă cu prezentarea într-un mod cât mai intuitiv și simplu a rezultatelor interogării curente astfel încât utilizatorul să găsească cât mai ușor paginile căutate. Motorul de căutare google folosește de exemplu o listă de snippet-uri ordonate, dar sunt motoare de căutare cum ar fi Yippy [yippy] care încercă pe baza unor algoritmi de învățare nesupervizată o grupare a paginilor întoarse de un motor de căutare standard în categorii, în funcție de cuvintele din snippet, altele decât cele din căutare. Astfel utilizatorul poate alege o categorie mai generală care să îl conducă mai repede la informațiile căutate.

În această situație, utilizatorii ar putea beneficia de o interfața de căutare în care să poată introduce nevoia informațională sub forma unei întrebări în limbaj natural, iar apoi să primească răspunsul la acea întrebare, fie sub forma unui simplu cuvânt, fie un document, site, fișier audio sau video. Spre exemplu, dacă un utilizator dorește să afle în ce an s-a născut compozitorul austriac Mozart, să fie de ajuns să scrie într-un motor de căutare întrebarea (“When was Mozart born?”) și să primească înapoi răspunsul corect (“1756”), eventual alături de un document (URL) care conține acel răspuns pentru lămuriri suplimentare.

O motivație suplimentară pentru care este necesară aplicarea și integrarea algoritmilor de Question Answering în motoarele de căutare este ca 15-20% din log-urile cu căutări ale motoarelor de căutare conțin de fapt căutări sub forma de întrebare. Deci utilizatorul are o nevoie informațională tot mai mare, știe ca pe Web poate găsi orice și este deci tot mai tentat să caute pe Web așa cum știe el, adică întrebând motorul de căutare în limbaj natural.

3.1.2 Conectarea la API-urile motoarelor de căutare

Cele mai importante motoare de căutare oferă ca facilitate conectarea directă din aplicație la acestea, prin intermediul API-ului propriu. Prin funcții parametrizabile specifice se apelează sistemul de căutare al motorului care va returna sub forma de fișier XML sau structuri de date rezultatele căutării pentru a fi procesate de către aplicația utilizatorului. Comunicarea se realizează prin protocoale SOAP sau REST.

Pentru această aplicație au fost accesate interfețe către următoarele servicii oferite de motoarele de căutare: Google API și Bing API. Aceste servicii necesită doar o înregistrare prealabilă (momentan google nu mai oferă gratuit acest serviciu) pentru a obține o cheie (Application Key) pentru a ține evidența utilizării serviciilor proprii în diverse sisteme software. Cheia permite și descărcarea de pe Web a unei biblioteci care conține funcțiile ce trebuie apelate pentru accesarea serviciului de căutare. În continuare o să prezint conectarea la motorul de căutare BING.

Chiar dacă serviciul Bing API este gratuit, este nevoie de un ID de aplicație pentru a utiliza serviciul. Acest ID se poate crea accesând adresa <http://www.bing.com/developers/createapp.aspx>. De la acea adresă după înregistrare se obține un AppID. După crearea proiectului trebuie adăugată o referință web pentru adresa: <http://api.search.live.net/search.wSDL?AppID=YourAppId>, unde YourAppID este ID-ul obținut. În momentul achiziției ID-ului se preia și biblioteca cu funcțiile API necesare conectării al motorul de căutare urmând ca în continuare să se utilizeze funcțiile din acele biblioteci pentru transmiterea interogării și parametrizarea numărului de răspunsuri returnate de motorul de căutare. Astfel într-un obiect de tip *SearchRequest* se poate specifica întrebarea adresată motorului, numărul maxim de rezultate pe care le dorim de la motorul de căutare precum și ofsetul adică de la ce index din lista proprie în continuare dorim să obținem rezultatele. La apelul metodei *Search* se realizează conectarea efectivă la motorul de căutare și preluarea rezultatelor care vor fi întoarse utilizatorului într-un obiect de tip *SearchResponse* de unde ulterior aplicația le poate procesa.

3.2 Detectarea părților de vorbire

Pentru a obține de la motorul de căutare rezultate cât mai relevante cu interogare specificată trebuie specificat acestuia cuvintele de interogare esențiale care se stabilesc pe baza analizei lexicale. De asemenea răspunsul este eticheta în funcție de analiza sintactică a acestuia. Astfel problema majoră abordată în momentul de față este detectarea automată și cât mai corectă a tipului întrebării (pentru a ști tipul de răspuns) și părților de vorbire din cadrul unei propoziții în limba engleză pentru a putea interpreta corect semantica propoziției. Domeniul care se ocupă cu „mineritul” sensurilor unui cuvânt se numește WSD (Word Sense Disambiguation) [Ste03] și este un domeniu foarte apropiat de lexicografie care se ocupă de descoperirea și descrierea sensurilor cuvintelor. În 1994 Hanks vorbește despre două feluri de extragere a înțelesurilor „normală și exploratorie”. Un cuvânt are un sens comun și în majoritatea timpului acest sens este cel utilizat de către interlocutor. Acest sens constituie sensul normal dar o limbă este întotdeauna deschisă interpretărilor, combinațiilor și astfel apar noi „setări” ale cuvintelor.

Metodele bazate pe cunoaștere reprezintă o categorie distinctă a WSD, alături de metodele bazate pe corpusuri. Performanțele acestora este întrecută de către cele bazate pe corpusuri, însă ele au o acoperire (în sensul de aplicabilitate) mai mare. De obicei aceste metode sunt aplicate dezambiguizării tuturor cuvintelor în texte nerestricționate pe când cele bazate pe corpusuri sunt aplicate, în principiu, adnotării cuvintelor în corpusuri. În [Cret12] am prezentat o serie de algoritmi utilizați în WSD pentru descoperirea de cunoștințe.

Pentru detectarea automată a părților de vorbire am folosit baza de date WordNet și am dezvoltat o serie de algoritmi de predicție a părților de vorbire. Această problemă este dificilă chiar și pentru limba engleză deoarece o parte semnificativă din cuvintele limbii engleze au mai multe forme sintactice și sunt cuvinte a căror sens depinde de context.

3.2.1 WordNet

Baza de date WordNet [WorNet] este des folosită în literatura de specialitate și este un proiect lansat în 1985 la Princeton University. Este o bază de date lexicală iar structura sa se dorește a fi o rețea semantică a sensurilor cuvintelor. Acest dicționar semantic a fost proiectat pentru a stabili conexiuni între patru tipuri de părți de vorbire: substantive, verbe, adjective și adverbe.

Cea mai mică unitate din WordNet este synset, care reprezintă o anumită semnificație a unui cuvânt. Acesta include cuvântul, explicația acesteia, și sinonimele sale. Înțelesul specific al unui cuvânt sub un singur tip de parte de vorbire este numit un sens. Fiecare sens al unui cuvânt este într-un synset diferit. Synsets sunt echivalente cu sensurile = structuri care conțin seturi de termeni cu sensuri sinonime. Fiecare synset are un gloss care definește conceptul care îl reprezintă. Cuvintele polisemantice pot aparține mai multor synset-uri.

3.2.2 Metode de detectare a părților de vorbire

Pentru evaluarea algoritmilor realizați pentru detectarea părților de vorbire am folosit benchmark-ul Brown Corpus [BrownCorp] care reprezintă o colecție de texte pentru cercetare în domeniul limbajelor naturale. Aceste benchmark-uri conțin după fiecare cuvânt sau semn de punctuație o etichetă care reprezintă tipul de vorbire iar fiecare tip de vorbire cu formele lui, existând aproximativ 80 de astfel de taguri disponibile. Am testat mai multe metode de etichetare iar pentru fiecare metodă am calculat acuratețea de etichetare ca și raport între numărul de etichete corecte (comparativ cu etichetele deja existente pentru fiecare cuvânt din benchmark, pe care le-am considerat perfecte) și numărul total de cuvinte de etichetat.

Folosind doar WordNet pentru a eticheta fiecare cuvânt din benchmark am obținut o limită maximă a acuratețe de 72.93%. În cazul introducerii în algoritm și a informațiilor din limba engleză referitoare la sufixele unor cuvinte și partea de vorbire am crescut această limită a acurateței de etichetare la 83.4%. Gramatica limbii engleze are o serie de reguli întotdeauna valabile, de exemplu întotdeauna după articolul „the” urmează un substantiv sau un adjectiv. Am introdus și aceste reguli în algoritm cu scopul de a reduce lista de părți de vorbire posibile identificate la pasul anterior. Acuratețea de clasificare nu s-a îmbunătățit în acest caz. Valorile maxime obținabile pentru fiecare parte de vorbire întoarsă de către WordNet sunt prezentate în figura 2.

Ca și metodă de etichetare utilizată pentru determinarea părții de vorbire corecte în contextul în care se folosește WordNet, reguli și sufixe este metoda care ia în considerare doar 3 cuvinte succesive din propoziție. Metoda încearcă să predicționeze partea de vorbire a unui cuvânt pe baza cuvântului anterior acestuia, pe baza cuvântului următor acestuia și pe baza combinațiilor între toate sensurile respectivelor cuvinte, alegând întotdeauna combinația care obține scorul cel mai mare. Acest scor se calculează în funcție de probabilitatea de apariție a unui tip de vorbire după și înaintea unui alt tip de vorbire relativ la WordNet și benchmark-ul Brown Corpus. Testat pe benchmark-ul Brown corpus această metodă obține o acuratețe de etichetare de 68.83% în medie și de 87.4% pentru etichetarea substantivelor.

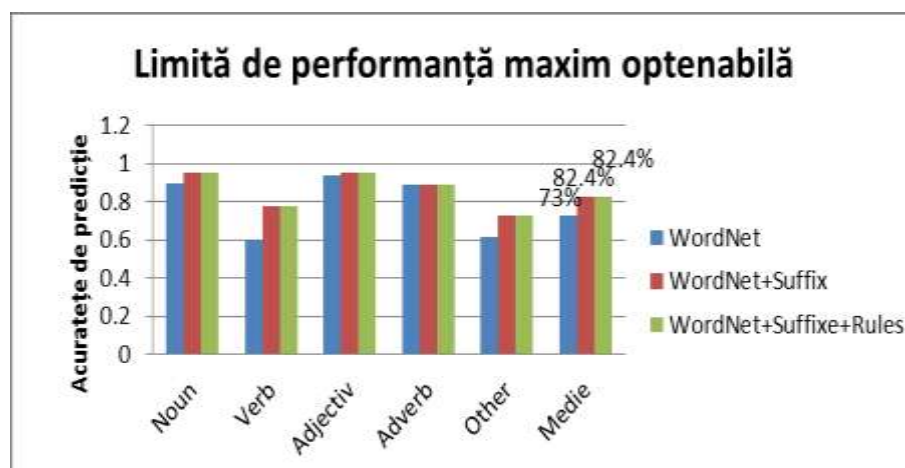


Figura 2. Limita maximă obținută pentru fiecare parte de vorbire.

3.3 Modulele aplicației de răspuns automat la întrebări

Aplicația de răspuns automat la întrebări conține mai multe module principale: modulul de analiză a întrebării, modulul de conectare la motorul de căutare și aducere snippet-uri, modulul de procesare a snippet-urilor și modulul de extragere a răspunsurilor.

3.3.1 Modulul de analiză a întrebării

Modulul de analiză a întrebării are mai multe etape. În prima etapă se identifică tipul răspunsului (Person, Location, Time, etc) în funcție de cuvântul cu care începe întrebarea și astfel după anumite cuvinte cheie se obține tipul de răspuns. O altă etapă este identificarea cuvintelor esențiale prin stabilirea părții de vorbire a fiecărui cuvânt din propoziție. În ultima etapă se formează interogarea din cuvintele esențiale. Cuvintele esențiale sunt acele cuvinte care definesc întrebarea utilizatorului și astfel vom acorda o încredere mai mare răspunsurilor în care apar cuvintele esențiale. Pentru identificarea cuvintelor esențiale am folosit WordNet și o serie de reguli și sufixe folosite în limba engleză pentru stabilirea părții de vorbire.

3.3.2 Modulele de aducere și procesare snippet-uri

Modulul de aducere a snippet-urilor - are rolul de a realiza conectarea cu motorul de căutare, precum și interogarea formată într-un query și returnează o listă de snippet-uri.

Modulul de procesare a snippet-urilor – are rolul de a împărți snippet-urile primite de la Modulul de Aducere a Snippet-urilor în propoziții. Acest modul se ocupă de extragerea paragrafelor din fiecare snippet. Acest pas este necesar deoarece un snippet este alcătuit din una sau mai multe fraze, delimitate de semne de punctuație. Un snippet poate conține mai multe propoziții delimitate de anumite caractere dar poate conține și mai multe părți din propoziții diferite.

3.3.3 Modulul de extragere răspunsuri

Modulul de extragere răspunsuri are rolul de a analiza propozițiile primite de la Modulul de procesare a snippet-urilor, caută răspunsurile probabile după tipul răspunsului și calculează un punctaj pentru fiecare răspuns probabil. Fiecare propoziție va fi analizată lexical, această analiză presupune identificarea cuvintelor din documentul de intrare. În principiu se va separa textul în cuvinte (token-uri), în funcție de anumiți separatori speciali. Se vor elimina cuvintele de legătură și în funcție de tipul de răspuns solicitat propozițiile vor fi procesate diferit. Deocamdată se detectează răspunsuri de tip Person, Location, Time, Description și Reason.

După ce a fost extrasă lista de răspunsuri posibile trebuie ca aceste răspunsuri să fie filtrate și la fiecare răspuns, în funcție de tipul de răspuns se stabilește o anumită încredere. Încrederea este calculată în funcție de nr. de cuvinte din interogare, nr. de propoziții în care apare răspunsul, încrederea în site-ul de unde vine răspunsul și o încredere bazată pe tipul de vorbire întors de WordNet pentru răspuns și tipul de vorbire așteptat pentru categoria respectivă de răspuns. De exemplu pentru o întrebare de tip Location ne așteptăm ca și cuvântul din răspuns să fie substantiv și acordăm o valoare calculată pe baza sensurilor din cuvânt pe baza WordNet și a regulilor.

4 Concluzii

Aplicația de răspuns automat la întrebări realizată practic se găsește online la adresa <http://193.226.29.26/ACUD/Default.aspx>. Permite introducerea unei întrebări în limbaj natural, și în această etapă deocamdată afișează o listă de snippet-uri în care se poate găsi răspunsul la întrebarea formulată. Această listă este ordonată în funcție de încrederea pe care a primit-o fiecare snippet în parte. Această încredere pentru fiecare snippet se calculează pe baza tipurilor de vorbire ale cuvintelor din vecinătatea cuvintelor din interogare raportat la tipurile de vorbire pe care le așteptăm în funcție de tipul întrebării. Deocamdată aplicația funcționează pentru 5 tipuri de întrebări (Person, Location, Time, Reason și Description).

Deoarece este dificil de calculat care este acuratețea predicției de etichetare pentru algoritmi implementați în cazul utilizării direct în aplicația de răspuns automat am efectuat câteva experimente pe baza unui benchmark gata etichetat. Valoarea maximă obținută de algoritmul

implementat este de 68.83% în medie pentru toate cele 4 tipuri de vorbire întoarse de WordNet și este de 87.4% pentru predicția etichetelor substantivelor. Pentru acest breanchmark am calculat de asemenea valoare maximă care am fi putut să o obținem care este de 82.4% ceea ce arată că algoritmul implementat este încă departe de performanțele maxime. Încercările avute folosind algoritmi de învățare (de exemplu Naive Bayes sau predictorul markov) pentru predicția părții de vorbire deocamdată nu a dat rezultatele asteptate și a crescut timpul de răspuns pentru o astfel de aplicație destul de mult. De exemplu predictorul Markov pentru breanchmark-ul folosit a obținut 68.53% dar a crescut timpul necesar așteptării pentru răspuns având în vedere că predicția părții de vorbire este folosită atât în partea de analiză a întrebării cât și în partea de analiză a snippet-urilor.

5 Referințe bibliografice

- [BrownCorp] Brown University Standard Corpus of Present-Day American English (Brown Corpus) <http://icame.uib.no/brown/bcm.html>.
- [Cret12] R. Cretulescu., D. Morariu, M. Breazu, L. Vintan – Word Sense Disambiguation for Text Mining, The third international conference in Romania of Information Science and Information Literacy, ISSN 2247-0255, April 2012
- [Han94] Hanks, P. "Linguistic Norms and Pragmatic Explanations, or Why Lexicographers need Prototype Theory and Vice Versa" in F. Kiefer, G. Kiss, and J.Pajzs (eds.), Papers in Computational Lexicography, 1994.
- [Ste03] M. Stevenson, Word Senseisambiguation, The case for Combinations of Knowledge Sources, CSLI Publications, 2003.
- [WordNet] Princeton University WordNet-A lexical database for English
[/http://wordnet.princeton.edu](http://wordnet.princeton.edu).
- [yippy] <http://search.yippy.com/>