

ETAPELE DE DEZVOLTARE A SISTEMULUI DE TIP „ÎNTREBARE-RĂSPUNS”

Victoria BOBICEV

Universitatea Tehnică a Moldovei

Abstract. *The paper presents a question-answering system created in the project „Research in the field of Information Retrieval for question answering system creation”. The system consist of question analysis module and answer retrieval and extraction module. At the first stage the system is working on the base of the documents provided by IDSI (Institutul de Dezvoltare a Societății Informaționale - Information Society Development Institute). At present the system is installed on-line with the aim to collect questions from different users and its improvement using collected information.*

Cuvinte-cheie: *procesarea automată a textului, sistem întrebare-răspuns, analiza automată a întrebărilor, căutarea informației, extragerea informației.*

1. Introducere

În prezent, internetul este cea mai mare sursă de cunoștințe ce se extinde și se actualizează incontinuu. Internetul mai este și una din cele mai accesibile locații în care aceste cunoștințe pot fi consultate. Însă dezvoltarea rapidă a internetului nu are loc fără aspecte negative, și anume: din cauza volumului mare de informații disponibile, găsirea informației necesare poate fi, uneori, dificilă sau nesigură.

Cele mai eficiente metode de găsire și de acumulare a informației o reprezintă, la ora actuală, motoarele de căutare, scopul cărora este de a oferi utilizatorului un set de articole sau pagini web în care acesta să poată găsi informația necesară. Deseori articolele propuse de motoarele de căutare nu îndeplinesc dezideratul utilizatorului de a căpăta un răspuns satisfăcător. În plus, ele nu oferă răspunsul concret la întrebarea utilizatorului, dar numai un set de pagini web, și utilizatorul e nevoit să-și extragă informația necesară.

Etapă următoare în domeniul achiziției informației constă în dezvoltarea sistemelor capabile să răspundă la întrebările formulate de utilizator în limbajul natural. Un sistem de răspunsuri la întrebări necesită o procesare a limbajului natural mult mai complexă decât sistemele de achiziție de documente [3][5].

2. Sistemul de întrebare-răspuns

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 [6].

Crearea unui sistem de întrebare-răspuns reprezintă un proces complex care implică multe componente și resurse importante. Un sistem de întrebare-răspuns bazat pe o colecție de documente, de obicei, este format din trei componente principale [5]:

1) Modul de analiză a întrebării - convertește întrebările din limbaj natural în interogări pentru motorul de achiziție de documente.

2) Modul de achiziție de documente - caută în totalitatea de documente articolele relevante pentru interogarea făcută de utilizator în baza datelor returnate de modulul de analiză a întrebării.

3) Modul de extragere a răspunsului - din colecția de articole returnate de modulul de achiziție de documente, extrage un răspuns succint și care constituie răspunsul în limbaj natural la întrebarea utilizatorului.

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) ca, de exemplu, „Codul cu privire la știință și inovare al Republicii Moldova”, „Acordul de parteneriat între Guvern și Academia de Științe a Moldovei pentru anii 2009-2012”, „Lege cu privire la parcurile științifico-tehnologice și incubatoarele de inovare”, și altele de același tip.

La momentul dat răspunsul la întrebarea pusă de utilizator se caută doar într-un document, deci, utilizatorul trebuie să aleagă documentul în care se va căuta răspunsul la întrebarea pusă din lista documentelor prestate de IDSI.

3. Analiza întrebărilor

Metodologia clasică de analiza semantică a întrebărilor propusă în [1] presupune împărțirea propoziției interogative în următoarele elemente: pronume interogativ (de exemplu, cine, care, unde, când) sau grup prepozițional interogativ (de exemplu, în ce, pe cine, de unde, la care) și propoziția inversată. Propoziția inversată începe cu grupul verbal după care urmează grupul nominal. La sfârșitul propoziției se presupune prezența așa numitului „gap” – loc liber prevăzut pentru răspunsul care va fi găsit. Astfel de metodologie este dezvoltată pentru întrebările formulate în limba engleză însă ea poate fi adaptată la analiza propozițiilor în limba română. Exemple de împărțire a propozițiilor engleze sunt prezentate în figura 1. Tabelul 1 conține exemplele propozițiilor în limba engleză și română.

Grupul nominal în cazul dat este considerat fraza-cheie pentru căutarea răspunsului. Respectiv, întrebarea se reformulează pentru obținerea răspunsului. Grupul verbal și cel nominal în propoziția inversată se schimbă cu locul și dacă întrebarea începe cu o prepoziție aceasta se atașează la sfârșit. Locul liber pentru răspuns urmează după elementele acestea.

Pronumele interogativ are și el rolul său specific, el definește **tipul** răspunsului: se întreabă despre o persoană (*cine*), despre un obiect (*ce*), despre localitate (*unde*), despre timp (*când*). Sunt posibile și întrebări mai complicate ce necesită răspunsuri desfășurate, de exemplu: *Care sunt argumentele pro și contra ...* [2]. Răspunsuri la astfel de întrebări necesită crearea rezumatului din documente multiple [4].

4. Problemele rezolvate și nerezolvate la etapa dată

- A fost rezolvată problema semnelor diacritice. Cu scopul procesării corecte a textului român cu litere speciale cu semne diacritice toate documentele sistemului au fost recodificate în UTF-8. Codificarea dată păstrează literele speciale și asigură vizualizarea lor corectă în diferite medii cum ar fi editoare de text sau browsere.

Documentele html se codifică în UTF-8 cu ajutorul editoarelor în care sunt create. Ca regulă, editoarele de texte au opțiunea de a schimba codificarea textului în cadrul procesului de salvare. De exemplu, Notepad oferă opțiunea dată în fereastra de salvare, ultima casetă de selectare Encoding oferă o listă de codificări posibile ce include și opțiunea UTF-8. Fragmentul de fereastră este prezentat în figura 1.

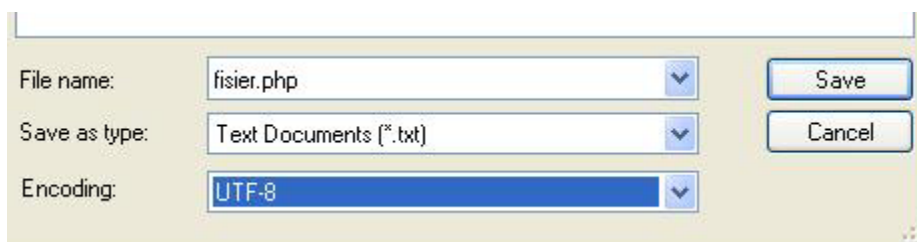


Figura 1. Fragmentul de fereastră de salvare a fișierelor în Notepad cu opțiunea de codificare a documentului.

În Notepad++ codificarea se indică cu ajutorul punctului meniului Encoding și apoi se salvează în mod normal. Pentru vizualizarea corectă a documentelor html în browsere este necesară adăugarea tag-ului meta în secțiunea head al documentului:

```
<meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
```

În cazul editării documentelor în Dreamweaver se selectează din Modify->Page Properties->Title/Encoding->Encoding->Unicode(UTF-8), iar la Unicode Normalization Form se selectează None (secvența meta din head va fi automat scrisă).

Transmiterea datelor introduse de utilizator în formular se efectuează prin CGI – Common Gateway Protocol care ca regulă recodifică textul în mod specific. Codificarea dată depinde de proprietățile documentului html, de atributele tag-ului form care conține elementele de introducere a textului, de setările serverului și de alte condiții. De exemplu, CGI permite transmiterea directă fără schimbări doar literelor alfabetului englez și unui număr limitat de caractere (. , ; ! ? @ # \$ și altelor). Restul caracterelor se codifică cu așa-numit percent-encoding. De exemplu, litera „ț” codificată în UTF-8 se va transmite recodificată în

%C4%83 unde semnele procentelor precedă fiecare din cele două octeți care formează litera dată în UTF-8. În afară de aceasta caractere de spații se înlocuiesc cu semnele „+”. Astfel programul care va analiza întrebarea ce a fost introdusă de utilizator în primul rând trebuie să decodifice textul obținut.

Programul de analiza a întrebărilor a fost scris în Perl care la fel are probleme cu codificarea UTF-8 însă problemele aceste pot fi rezolvate. În primul rând textul programului la fel a fost codificat în UTF-8. ‘use utf8;’ la începutul programului indică codificarea și permite utilizarea caracterelor cu diacritice în textul programului. use URI::Escape; indică că va fi utilizat modulul dat pentru recodificarea datelor obținute. use Encode ‘decode_utf8;’ introduce modulul pentru recodificarea datelor în UTF-8. și în final binmode(STDOUT, ‘:utf8’); permite output programului în UTF-8.

După obținerea datelor din formular textul este recodificat respectiv cu:

\$textul = uri_unescape (\$textul); și \$textul = decode_utf8 (\$textul); după ce textul arată absolut corect, este codificat în UTF-8 și poate fi procesat de program în orice mod necesar.

Trebuie de mențona că toate fișierele textuale cu care lucrează programul sunt la fel salvate cu codificarea UTF-8 și respectiv în program se indică faptul acesta la deschiderea lor: open(FISIER,<:encoding(UTF8), "\$filename"); pentru citirea datelor și respectiv open(RASPUNSURI, >:encoding(UTF8), "raspunsuri.html"); pentru salvarea datelor.

Toate procedeele date permit procesarea corectă a textului român scris corect cu semnele diacritice necesare.

▪ Există două tipuri de litere „ș” și „ț” cu coduri diferite: „ș”, „ș”, „ț”, „ț”. Verificând toate variantele posibile am observat și perechea de litere „ș”, „ț” ce arată similar care am adăugat în lista literelor posibile. Respectiv cu scopul procesării corecte tuturor literelor indiferent de codul lor am introdus modulul care înlocuiește în întrebarea introdusă literele „ș”, „ț” și „ș”, „ț” cu „ș”, „ț” care sunt utilizate în dicționarul și în documentele sistemului nostru. Astfel sistemul nu întâmpină dificultăți la căutarea cuvintelor-cheie extrase din întrebarea pusă în dicționar și în textele utilizate pentru căutarea răspunsului.

În cazul dacă utilizatorul intenționează să introducă întrebarea fără semne diacritice, înlocuind „ă”, „î”, „â”, „ș”, „ț” cu „a”, „i”, „a”, „s”, „t” respectiv textul întrebării nu va fi procesat corect. Sistemul presupune că întrebarea este scrisă cu semne diacritice și la compararea cuvintelor întrebării cu textul nu se va găsi coincidența. De exemplu, pentru întrebarea Ce este știința? introdusă fără semne diacritice ca Ce este stiinta? se va selecta cuvântul cheie stiinta care nu va fi găsit în text. Soluția posibilă este de a cerut utilizatorul de a indica faptul că semnele diacritice sunt intenționate omise. În acest caz răspunsul se va căuta în texte cu semnele diacritice eliminate la fel ca și în întrebare ce va permite coincidența cuvintelor și găsirea răspunsului.

Altă problemă poate fi diferența între modurile de scriere textului: cu litere „î” sau cu litere „â” la mijlocului cuvântului. Luând în considerație faptele că:

- în limba română sunt acceptabile ambele metode de scriere;
- scrierea cu „â” la mijlocul cuvântului este considerată mai modernă în România;
- limba de stat în Moldova se scrie cu „î” la mijlocul cuvântului;
- ambele tipuri de scris sunt utilizate intens inclusiv în Moldova;

putem presupune că întrebările pot fi introduse utilizând ambele tipuri de scriere. În acest caz putem utiliza modulul creat în cadrul cercetărilor precedente dedicate transformării modurilor de scriere [7].

▪ În sistemul creat a fost adăugat dicționarul morfologic ce permite căutarea și analiza tuturor formelor morfologice ale cuvintelor. De exemplu, dacă este introdusă întrebarea „Ce este cercetarea?”, sistemul va găsi răspunsul din cauza că în document este scris „cercetare”. Motoarele de căutare moderne deja lucrează asupra problemei formelor morfologice a cuvintelor. Noi dispunem de un dicționar morfologic a limbii române cu aproximativ 90000 de forme morfologice a cuvintelor care poate fi utilizat pentru obținerea tuturor formelor a cuvintelor-cheie. Problema constă în faptul că în majoritatea cazurilor sunt căutate nu dor cuvinte dar fraze-cheie. De exemplu, dacă avem fraza „Cercetare fundamentală”, formele ei sunt: „Cercetării fundamentale”, „Cercetarea fundamentală”, „Cercetărilor fundamentale”, și altele. Aceasta înseamnă că trebuie de schimbat forma tuturor cuvintelor în frază respectând acordul morfologic între ele. Problema aceasta poate fi rezolvată însă necesită intervenția lingviștilor.

Problemele nerezolvate:

▪ La momentul dat sistemul procesează corect întrebările scrise cu semnele diacritice, însă dacă utilizatorul introduce întrebarea fără semne diacritice sistemul va avea probleme analizând astfel de text și cel mai probabil nu va găsi răspunsul potrivit.

- Trebuie de evaluat lucrul sistemului în browsere diferite; ce puțin cele care sunt acum utilizate intens: Internet Explorer, FireFox, Google Chrome, Opera. De exemplu, parea AJAX care este încadrată în pagina cu răspunsuri nu funcționează în Internet Explorer din cauza problemei detectate în browser-ul dat: acesta nu execută AJAX corect în cazul indicării codificărilor locale cum ar fi Unicode (UTF-8). Noi am codificat situl nostru în UTF-8 cu scopul de a procesa corect textul cu semne diacritice, însă aceasta a dus la probleme cu AJAX.

- Dicționarul utilizat are unele neajunsuri: număr de cuvinte este destul de mare, ce duce la mărirea timpului de procesare; însă nu toate cuvinte care apar în text se găsesc în dicționar ce duce la probleme în căutarea răspunsului corect.

- Se procesează un număr limitat de tipuri de întrebări. La momentul de față sistemul procesează întrebările de tip definiție. Sistemul este instalat online cu scopul completării listei de întrebări de la utilizatori și memorizează toate adresările utilizatorilor, întrebările lor și răspunsurile găsite. Pe parcursul completării listei date modulul de analiză a întrebărilor va fi perfectat ca să fie capabil să proceseze toate tipurile întrebărilor de la utilizatori.

- Răspunsul se caută numai într-un document. Versiunea curentă a sistemului propune utilizatorului să aleagă documentul în care se va căuta răspunsul. IDSI a prestat un număr redus de documente și nu este problematic de încercat toate documentele pe rând. Pe parcursul utilizării sistemului numărul de documente va fi mărit și sistemul va fi modificat în așa mod că utilizatorul să poată căuta în câteva documente simultan sau în toate documente.

5. Concluzii

În lucrarea dată sunt discutate problemele sistemului de tip întrebare-răspuns creat în cadrul proiectului „Cercetarea în domeniul de Regăsire a Informației în scopul creării sistemului electronic de informare publică”. Modulul de analiză a întrebărilor și modulul de căutare și selectare a răspunsului au nevoie de rezolvare a problemelor codificării literelor cu semne diacritice. La prima etapă sistemul a funcționat fără semne diacritice; la momentul dat sistemul este codificat în UTF-8 cu scopul procesării corecte a literelor cu semnele diacritice. Altă modificare efectuată este adăugarea dicționarului morfologic și utilizarea tuturor formelor ale cuvintelor în procesul de căutare.

Bibliografie

1. R. Baeza-Yates, B. Ribeiro-Nieto, *Modern Information Retrieval*, Addison Wesley, 2000.
2. R. Botnaru, V. Bobicev *Studiul tipurilor de întrebări din sistemul de întrebare-răspuns*. ICMCS, Chișinău 2011.
3. L. Carcea. *Abordări în dezvoltarea sistemelor „întrebare-răspuns”*. ICMCS, Chișinău 2011.
4. V. Lazu., T. Prodan. *Metodologia reprezentării sensului întrebării în forma logică*. Conferința științifică a colaboratorilor doctoranzilor și masteranzilor UTM, 2011.
5. V. Maxim. *Metode utilizate pentru elaborarea unui sistem „întrebare-răspuns”*. ICMCS, Chișinău 2011.
6. Dan Tufiș, 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.
7. Bobicev V., Angheluș, V., *Corectarea automată a textului român*. ICTEI 2010, Chișinău, Moldova.