

Metode utilizate la elaborarea unui sistem întrebare – răspuns (QA)

Victoria MAXIM

Technical University of Moldova
maxivica@yahoo.com

Abstract — In information retrieval and natural language processing (NLP), question answering (QA) is the task of automatically answering a question posed in natural language. To find the answer to a question, a QA computer program may use either a pre-structured database or a collection of natural language documents (a text corpus such as the World Wide Web or some local collection). QA research attempts to deal with a wide range of question types including: fact, list, definition, How, Why, hypothetical, semantically constrained, and cross-lingual questions. Search collections vary from small local document collections, to internal organization documents, to compiled newswire reports, to the World Wide Web. QA is regarded as requiring more complex natural language processing (NLP) techniques than other types of information retrieval such as document retrieval, thus natural language search engines are sometimes regarded as the next step beyond current search engines. [3]

Index Terms — Answer Extraction, Information Retrieval, Typical Question Answering System, Question Answering, Question Classification.

I. INTRODUCERE

Sistemele de întrebare - răspuns (în engleză "question answering systems", sau **sisteme QA**) sunt considerate ca fiind următorul pas în evoluția motoarelor de căutare a informației. Sistemele de tip QA sunt caracterizate prin faptul că primesc un set de întrebări în limbaj natural și, pe baza unei colecții de documente, trebuie să extragă răspunsul sau răspunsurile. Această colecție poate varia de la o simplă colecție locală până la întregul World Wide Web.

II. GENERALITĂȚI

La ora actuală Internetul reprezintă, fără îndoială, cea mai mare bază de cunoștințe, aflată într-o continuă extindere și actualizare. El este, în același timp, una dintre cele mai accesibile locații în care aceste cunoștințe pot fi consultate. Dar gradul de dezvoltare a Internetului are și aspecte negative: datorită multitudinii de informații disponibile, găsirea informației necesare la un moment dat poate fi dificilă sau și nesigură.

Cele mai eficiente metode de descoperire și de achiziție a informației o reprezintă, în prezent, **motoarele de căutare**. Scopul acestora este de a oferi utilizatorului un set de articole sau pagini web în care acesta să poată găsi informația care îi este necesară. De multe ori articolele oferite de motoarele de căutare nu îndeplinesc dezideratul utilizatorului de a obține un răspuns satisfactor.

Deasemenea, ele nu oferă răspunsul concret la problema utilizatorului, ci doar un set de pagini web, din care utilizatorul trebuie să extragă singur informația căutată.

Pasul următor în domeniul achiziției informației e constituit de dezvoltarea sistemelor capabile să răspundă la întrebări formulate de utilizator în limbaj natural. Dezideratul principal al unui astfel de sistem este să asigure un răspuns la întrebarea utilizatorului care să îndeplinească urmatoarele trei condiții:

- ✚ să fie corect,
- ✚ să fie formulat tot în limbaj natural și
- ✚ să fie suficient de succint.

Un sistem de răspuns la întrebări necesită o procesare a limbajului natural mult mai complexă decât sistemele de achiziție de documente.

În teorie, procesarea limbajului natural e un subiect foarte atractiv, datorită aplicabilității sale în domenii ca cel al interacțiunii om-mașină. În practică se constată, însă, o serie de dificultăți majore, datorate mai ales modului diferit în care o afirmație în limbaj natural poate fi interpretată și a multitudinii de sensuri pe care cuvintele constituente le pot lua. Sistemele de răspuns la întrebări, văzute ca un subdomeniu al procesării limbajului natural, moștenesc problemele acestora.

III. METODE UTILIZATE LA ELABORAREA UNUI SISTEM QA

Pentru construirea unui sistem de răspuns la întrebări există două variante:

✚ **Abordare de tip shallow, bazată pe cuvinte cheie.** În această metodă se folosesc cuvinte cheie pentru a găsi pasaje și propoziții în text care ar putea reprezenta răspunsuri valide la întrebări. Aceste potențiale răspunsuri urmează să fie analizate apoi mai în profunzime pentru a se stabili dacă sunt răspunsuri reale sau nu. Această metodă poate fi folosită cu succes în cazul întrebărilor scurte, factuale, când se caută nume, date, locații, cantități.

✚ **Abordarea de tip deep, ce implică o analiză mai sofisticată, o procesare sintactică, semantică și contextuală.** Există o serie de metode ce pot fi încadrate în această categorie: abduction, named-entity recognition, relation detection etc.

Alegerea unuia dintre cele două modele depinde de complexitatea întrebărilor ce vor fi formulate și de gradul de performanță dorit de la sistem. Este clar că sistemele din cea de-a doua categorie sunt superioare primelor.

Arhitectura generală a unui sistem de tip întrebare-răspuns

Dacă la începuturile inteligenței artificiale, în anii 1960, cercetătorii erau fascinați de ideea de a putea construi sisteme capabile să răspundă la întrebări aparținând unor domenii restrânse (closed domains), în prezent dezvoltarea Internetului și pașii făcuți în ceea ce privește achiziției informației (information retrieval - IR) și a tehnicilor de procesare a limbajului natural (natural language processing - NLP), precum și cererea pentru acces facil la informație, a dus la creșterea interesului pentru sisteme care să ofere răspunsuri din domenii largi (open domains).

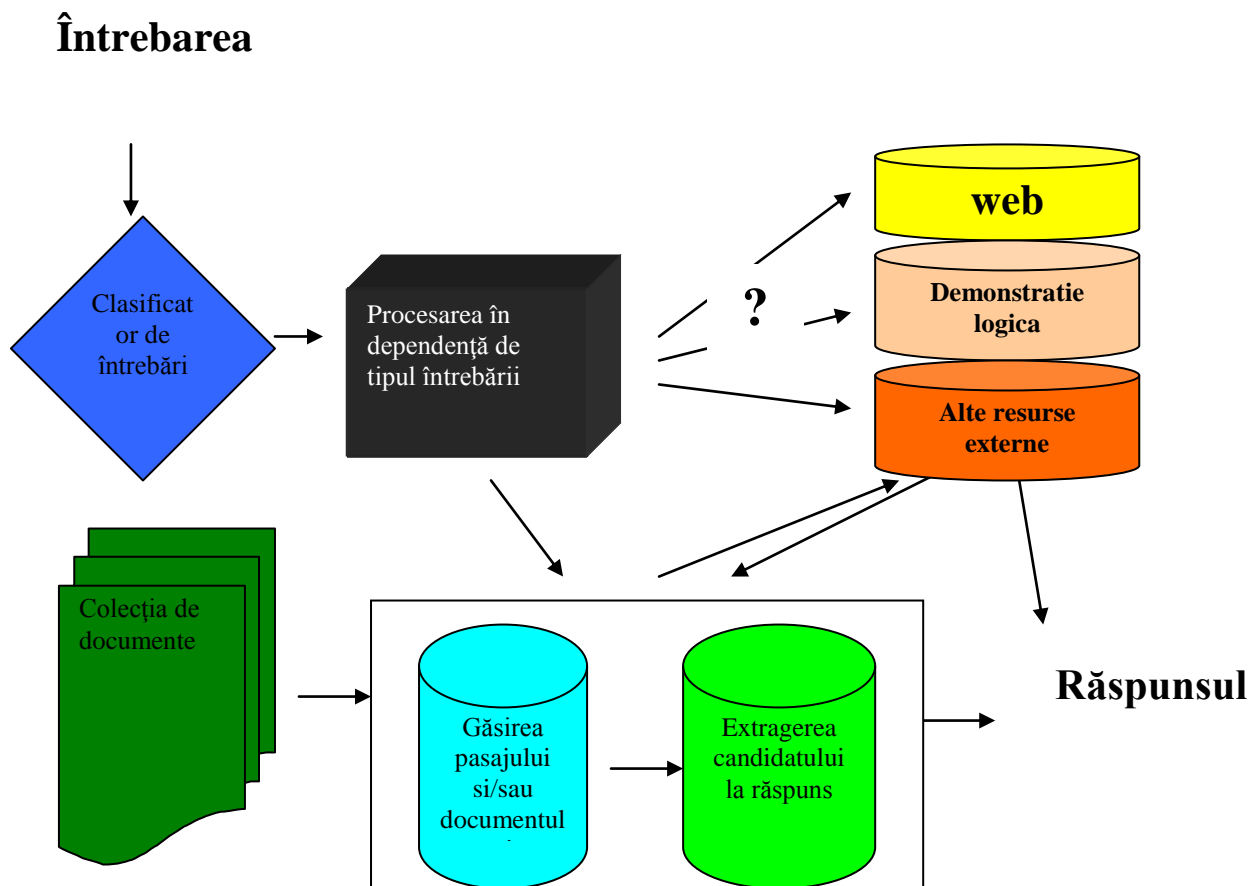


Figura 1. Un sistem QA tipic [2]

Un sistem de răspuns la întrebări bazat pe o colecție de documente are, în mod tipic, trei componente principale:

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

formulata de utilizator, pe 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 și care constituie răspunsul în limbaj natural uman la întrebarea utilizatorului. Dacă un astfel de răspuns nu există în colecția de documente considerată de modulul de achiziție

de articole, e de preferat ca sistemul să nu răspundă la întrebare, în loc de a întoarce un răspuns eronat.

Chiar dacă Internetul este un mediu plin de informații din toate domeniile, găsirea unui răspuns la o întrebare simplă poate fi uneori o sarcină dificilă. Unele din dificultățile ce pot să apară în dezvoltarea unui astfel de sistem:

Formularea corectă a interogărilor. Transformarea unei întrebări din limbaj natural într-o interogare pentru un motor de căutare este o sarcină dificilă. Dacă întrebarea e prea generală, va fi extras un număr prea mare de documente. Deasemenea, temele descrise în colecția de documente extrase e posibil să nu conțină tocmai răspunsul la întrebarea utilizatorului. Dacă sunt extrase prea multe documente, timpul de procesare va fi sporit. Dacă setul de cuvinte căutate e prea mic, e posibil să nu fie găsit articolul care răspunde la întrebare. De aceea se cere ca setul de cuvinte folosite în interogare să fie bine formulat, pentru a fi returnate, pe cât posibil, doar documente care să conțină informație utilă.

Chiar dacă e găsit setul de cuvinte corecte pentru a realiza o interogare care să aibă posibilitatea să întoarcă articole cu informație utilă, motorul de căutare poate întoarce un număr foarte mare de articole care să nu răspundă întrebării utilizatorului. De exemplu, în cazul întrebării „Cine a fost primul om în spațiu?”, modulul de achiziție de informații va primi, printre altele, și următoarele cuvinte cheie după care să realizeze căutarea: „primul om în spațiu”. Motorul de căutare e posibil să aibă indexate pagini despre „primul turist în spațiu” sau despre „cel mai bătrân om în spațiu” și să le returneze și pe acestea, împreună cu articolele care conțin răspunsul corect, și anume cele despre Iuri Gagarin, primul cosmonaut în spațiu.

Informații false. Chiar dacă întrebarea e bine formulată, motorul de achiziție de articole întoarce articole care corespund subiectului întrebării, e posibil ca unele dintre aceste articole să conțină informații greșite. Acesta este unul dintre cazurile cele mai nefavorabile, deoarece sistemul nu are cum să își dea seama care răspunsuri sunt corecte și care nu.

Resurse limitate. Când este construit un sistem de răspuns la întrebări trebuie să se țină cont de limitările impuse de lucrul cu cantități mari de informații. Este neindicat să se trimită sistemului un set prea mare șiruri de cuvinte pentru interogări. Cu toate ca motoarele de căutare actuale sunt suficient de rapide și întorc răspunsuri la obiect, căutarea în liste prea mare de articole consumă mult prea mult timp, și utilizatorul sistemului nu este dispus să aștepte minute pentru ca sistemul să îi ofere un răspuns.

IV. CLASIFICAREA ÎNTREBĂRILOR

Un sistem QA caută expresia sau enunțul (sau pasajul, documentul, setul de documente) care este răspunsul exact la o întrebare.

- întrebările au multe nuanțe;

- majoritatea căutărilor sunt axate pe întrebări factice;
- răspunsurile pot fi enunțuri sau expresii;
- răspunsurile complete ar trebui să fie găsite într-o sursă.

Tipuri de întrebări:

- ✓ Factice : Cine l-a omorât pe Martin Luther King?
- ✓ Sarcini : Cum pot să aplic pentru un pașaport?
- ✓ Opinii : Care a fost cel mai bun film anul acesta?
- ✓ Definiții : Cine este Jane Goodall?
- ✓ Liste : În ce filme s-a produs Jude Law?
- ✓ Explicații : Care a fost motivul războiului din Coreea?
- ✓ Da-Nu : Este legal să mergi la culoarea roșie a semaforului?

Exemple de întrebări:

1. *Aspartame mai este cunoscut și ca?*
2. *La ce vârstă Rossini a încetat să mai scrie opere?*
3. *În ce zi este sărbătorită Ziua Boxului?*
4. *Definește Thalassemia.*
5. *Cât de mare este Galaxia noastră în diametru?*
6. *Cât de rece trebuie să fie un frigider?*
7. *CPR este acronimul la ce?*
8. *Cât durează gestația la om?*
9. *Ce lungime are râul Nistru?*

Exemple de întrebări factice:

Q: Când s-a născut Mozart?

A: 1756.

Q: Ce ste un nanometru?

A1: O miliardime de metru.

A2: O milionime de milimetru.

Q: Când a avut loc Marea Depresie?

A1: Anii 1930.

A2: 1931.

A3: 1932.

Q: Cine este Abesalom?

A1: lider afro-american, primul negru capitan de navă de vânătoare de balene.

A2: Fiul lui David (biblic), care la- trădat pe tatăl său.

Clasificarea întrebărilor factice după tipul răspunsului așteptat/scontat.

Există 2 abordări de bază:

1. Clasificarea

Avantaj – ușor de înțeles/implementat, sigur.

Dezavantaj – nu oferă o altă informație decât categoria.

Expresii uzuale

Avantaj – oferă informații suplimentare la categorie.

Dezavantaj – foarte incert/nesigur.

Categoriile de clasificare: etichete PDV, cuvinte, WordNet, cuvinte WH, arbori sintactici, taguri NE.

Expresii uzuale

- cuvinte simple WH
- compuse tipologia QA

precizie/acuratețe 90%.

2. Procesarea în dependență de tipul întrebării

Rescrierea întrebării pentru web

Transformarea întrebării într-o cerere, combinarea probelor multiple

Demonstrații logice

- Tentativa de a raționa despre întrebare
- Filtrarea răspunsului
- Generarea răspunsurilor ce par corecte, dar sunt imposibile.

- Analiza întrebării pentru modele (for patterns)
 - Modele din întrebare sugerează modelele din răspuns

Resurse externe

- web (problematic)
- rezumate Web
- Validarea Răspunsului
- Mărirea datelor experimentale
- etichetare PDV, extractoare de NE, separatoare de grupuri nominale
- Gazetteers, ontologii, tezaure
- WordNet, ConceptNet
- Demonstrații logice
- Întrebări soluționate anterior

V. EXTRAGEREA RĂSPUNSULUI

Extragerea Răspunsului (mai simplu)

- Extragera răspunsului care conține tipul corect numit al entității din enunțul superior.
- Extragera răspunsurilor din enunțurile superioare N, și votarea.
- Extragera răspunsurilor candidate din enunțurile superioare N, validarea pe Web.

Marcarea/etichetarea răspunsului

- Tratați etichetarea răspunsului drept marcarea unei entități numite.
- răspunsurile adesea nu sunt de tipul unei entități numite, de exemplu, celebrele Wh-întrebări.
- Răspunsurile nu sunt previzibile și nu au întotdeauna indicatori previzibili.

Caracteristicile răspunsurilor nu sunt direct secvențiale și sunt adesea prea lungi.

Caracteristicile unui tip de întrebări pot să nu se refere la alte tipuri de întrebări

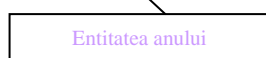
Etichetarea răspunsului (Easy)

- Determinați tipul de răspuns la întrebare
- Preluați un enunț corect
- Întoarceți entitatea numită corespunzătoare

Q: Când s-au născut gemenii lui Shakespeare ?



A: 2 ani mai târziu are loc nașterea gemenilor lui Shakespeare Judith și Hamnet, fată și băiat, botezați în anul 1585.



Etichetarea răspunsului (mai complicat)

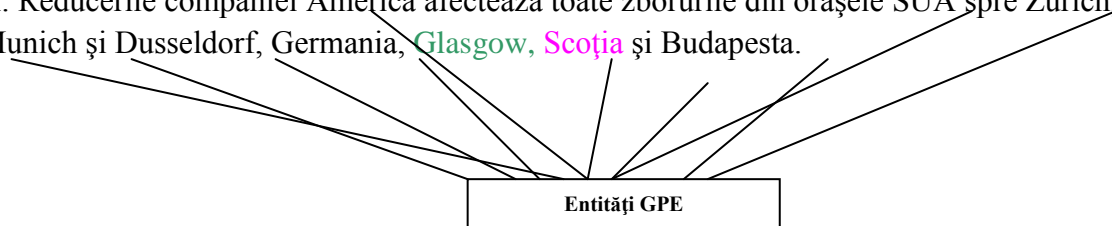
- Determinați tipul de răspuns la întrebare

- Preluați un enunț corect
- Întoarceți entitatea numită corespunzătoare

Q: Unde se află Glasgow ?

A: Recesiunea a atins Glasgow mai târziu, precum și restul Scoției.

A: Reducerile companiei America afectează toate zborurile din orașele SUA spre Zurich, Elveția, Munich și Dusseldorf, Germania, Glasgow, Scoția și Budapesta.



Modele de răspuns

- Modelele de răspuns se pot afla în textul de vecinătate imediată a unui răspuns la o întrebare faptică
- În funcție de tipul de întrebare
- Independent de o întrebare specifică

Exemple de modele de răspuns:

Modelul “inventator”

<NUME> , inventat de <RĂSPUNS>

<NUMELE>al cui <RĂSPUNS>

<NUME> a fost inventat de <RĂSPUNS>

<RĂPUNS> a inventat <NUME>

Exemplu:

Și demonstrația s-a întâmplat la 115 ani după ce *Edison a inventat becul electric.*

Modelul “anul de naștere”

Dodi Fayed s-a născut în 1956.

În viața sa scurtă și moartea tragică Jessica Dubroff (1988 – 1996) a devenit o metaforă pentru orice, pornind de la idealismul timpuriu, până la excesele Neew-Age.

VI. CONCLUZII

Sistemele de extragere a răspunsurilor la întrebări în limbajul natural uman se înscriu în categoria sistemelor de achiziție de informații.

Putem caracteriza drept cunoaștere starea unui sistem informațional cuplat cu un ambient. Prin sistem informațional se poate înțelege atât un individ (în sensul de persoană privită ca unitate distinctă față de alte persoane) care interacționează în mod direct și conștient cu ambientul și care își construiește prin mijloace proprii (observare, cogniție, etc.) universul de cunoaștere, dar și un sistem tehnic ce primește informația de la diverși agenți umani și o procesează sau o pune, mai departe, la dispoziția altor indivizi. Internetul poate fi descris ca un astfel de sistem informațional care, interacționând cu mediul prin intermediul agenților umani, își construiește o bază proprie de cunoaștere.

Proceedings of the Third International Conference on Language Resources and Evaluation (LREC 2002).

- [7] Burger, J., Cardie, C., Chaudhri, V., Gaizauskas, R., Harabagiu, S., Israel, D., Jacquemin, C., Lin, C-Y., Maiorano, S., Miller, G., Moldovan, D., Ogden, B., Prager, J., Riloff, E., Singhal, A., Shrihari, R., Strzalkowski, T., Voorhees, E., Weishedel, R. Issues, Tasks and Program Structures to Roadmap Research in Question Answering (QA).

REFERINTE BIBLIOGRAFICE

- [1] http://en.wikipedia.org/wiki/Question_answering
- [2] <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.14.143&rep=rep1&type=pdf>
- [3] Dragomir R. Radev, John Prager, and Valerie Samn. Ranking potential answers to natural language questions. In Proceedings of the 6th Conference on Applied Natural Language Processing, Seattle, WA, May 2000.
- [4] John Prager, Eric Brown, Anni Coden, and Dragomir Radev. Question-answering by predictive annotation. In Proceedings, 23rd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Athens, Greece, July 2000.
- [5] Hirschman, L. & Gaizauskas, R. (2001) Natural Language Question Answering. The View from Here. Natural Language Engineering (2001), 7:4:275-300 Cambridge University Press.
- [6] Lin, J. (2002). The Web as a Resource for Question Answering: Perspectives and Challenges. In