

TEZĂ DE ABILITARE
REZUMAT

FLORENTINA HRISTEA

UNIVERSITATEA DIN BUCUREȘTI
FACULTATEA DE MATEMATICĂ ȘI INFORMATICĂ
DEPARTAMENTUL DE INFORMATICĂ

2016

Rezumat

Teza de față se referă la întreaga noastră activitate de cercetare postdoctorală. Ea se bazează însă, în mod special, pe acea parte a cercetării noastre care s-a concentrat asupra rolului dezambiguizării automate a sensului cuvintelor în îmbunătățirea rezultatelor din regăsirea informației. Întrebarea la care ne-am propus să răspundem este următoarea: poate *dezambiguizarea sensului cuvintelor* (“word sense disambiguation” - *WSD*) să îmbunătățească rezultatele obținute în *regăsirea informației* (“information retrieval” - *IR*)? Această întrebare a captat, cu predilecție, interesul nostru științific în ultimii ani și a generat o cercetare de amploare și de natură interdisciplinară ale cărei rezultate originale sunt prezentate în capitolele tezei de față.

Capitolul final al tezei de abilitare va arăta modul în care această temă de cercetare își găsește locul în cadrul activității noastre de cercetare postdoctorală globale.

Problema menționată a fost îndelung studiată și dezbătută în literatura domeniului IR, dar cu rezultate controversate. Și cu o largă majoritate de răspunsuri negative la întrebarea enunțată anterior. WSD a fost unanim recunoscută ca fiind utilă în multiple tipuri de aplicații ale diverselor subdomenii din inteligența artificială. WSD în IR a fost pusă sub semnul întrebării și declarată ca fiind ineficientă. Sperăm că această cercetare va duce la schimbarea perspectivei sceptice care există, în momentul de față, asupra acestei probleme.

Performanța sistemelor de IR este cunoscută ca fiind puternic dependentă de interogări. În mediul motoarelor de căutare, utilizatorii trimit interogări (“queries”) în concordanță cu nevoile lor informaționale. Ca răspuns la o interogare, sistemul regăsește și afișează o listă de documente pe care le consideră relevante și de interes pentru utilizator, din perspectiva interogării efectuate. Utilizatorul analizează rezultatele regăsite și decide care dintre ele sunt utile nevoilor sale informaționale. Interogările pentru care motorul de căutare nu este capabil să furnizeze informații relevante sunt considerate dificile (“difficult queries”).

Dificultatea interogărilor poate fi legată de un număr de cauze total diferite: ambiguitate, formulare neclară sau prea vagă, lipsa contextului, natura și structura colecției de documente etc. Ambiguitatea termenilor a fost identificată ca reprezentând o cauză majoră a dificultății unei interogări. Cercetarea noastră s-a concentrat asupra acelor interogări considerate dificile datorită ambiguității cuvintelor care intervin în structura lor.

În pofida faptului că procesul de IR se bazează, în general, pe suprapunerea (coincidența) termenilor din interiorul interogărilor și respectiv al documentelor, fără a fi luat în considerație sensul cuvintelor, abordarea noastră s-a bazat pe semantica indusă atât de interogări, cât și de documentele regăsite. Termenii care intervin nu au mai fost considerați ca fiind independenți, ci ca fiind, mai degrabă, înrudiți. Principala presupunere pe care s-a bazat abordarea noastră este, prin urmare, aceea conform căreia *contextul* poate îmbunătăți performanța sistemelor de IR.

În mod intuitiv, o procesare a textelor mai bogată, care ține seama de multiple aspecte lingvistice, ar trebui să conducă la o îmbunătățire a rezultatelor în IR, atât în cazul interogărilor dificile, în general, cât și în cel al interogărilor ambigue, în particular. De-a lungul timpului, tehnicile specifice domeniului procesării limbajului natural au fost utilizate, în diferite moduri, în IR. Cu toate acestea, în pofida eforturilor depuse, nu au fost constatate îmbunătățiri semnificative ale eficienței în IR, mai ales în cazul evaluărilor la scară largă. Suntem de părere că motivul acestui eșec îl constituie faptul că tehnicile procesării limbajului natural nu ar trebui folosite în cazul tuturor interogărilor, ci numai în cazul acelor interogări care necesită, în mod efectiv, o tratare (prelucrare) de acest fel. Propriul nostru studiu și-a propus să identifice cuvintele ambigue din structura unei interogări și să decidă ce tehnică de tip WSD trebuie aplicată în cazul punctual al respectivei interogări.

După cum se știe, deși procesul de WSD este, în general, ușor pentru subiecții umani, el reprezintă o adevărată încercare pentru soft. Problema devine și mai greu de rezolvat în mod automat atunci când un cuvânt ambiguu intervine în cadrul unei “bucăți” scurte de text, cum este cazul interogărilor dintr-un sistem de IR.

Multe studii argumentează eșecul îmbunătățirii performanței în IR prin folosirea tehnicilor de WSD ca fiind datorat ineficienței algoritmilor de dezambiguizare existenți, din nou o problemă care se accentuează în cazul fragmentelor scurte de text reprezentând o interogare.

În literatură există numeroase încercări de a se folosi WSD bazată pe cunoștințe în IR. Acele studii care au raportat rezultate pozitive se caracterizează, însă, prin testări efectuate pe eșantioane mici de date. Atunci când evaluarea a fost adusă la scara unei colecții de testare ample, îmbunătățirile erau raportate relativ la un nivel de referință (“baseline”) slab, așa cum se comentează, în detaliu, în Cap. 7 al tezei de față. În general, ori de câte ori au fost raportate concluzii și rezultate pozitive, cuantumul de îmbunătățiri obținute a fost mic (a se vedea Cap. 7).

Puținii autori care, mai recent, au exprimat o încredere crescândă în beneficiile aduse de WSD pentru IR, folosesc, cu toții, o tehnică de WSD supervizată (a se vedea Cap.7 pentru detalii). Este, însă, cunoscut faptul că WSD supervizată nu poate fi folosită pe scară largă în practică, datorită absenței corpusurilor adnotate / paralele pe care le necesită.

Contrar opiniei tuturor acestor autori, cercetarea noastră propune și investighează utilizarea unei tehnici WSD nesupervizate în IR. În situația în care principala noastră presupunere este aceea conform căreia contextul poate îmbunătăți performanța unui sistem de IR, am prezentat în lucrările noastre o abordare care dorește să identifice cluster-e de contexte similare, cu fiecare cluster desemnând un cuvânt polisemantic folosit printr-un anumit sens al său. Suntem de părere că IR reprezintă o aplicație pentru care acest tip de analiză este folositor și, mai mult, suficient pentru mărirea performanței. Prin urmare, abordarea noastră nu își propune o dezambiguizare directă a sensului, ci o discriminare între sensurile posibile ale unui cuvânt ambiguu care intervine într-o interogare, cu impact în IR.

În vederea pregătirii utilizării tehnicilor de WSD în IR, studiul nostru s-a concentrat asupra a două aspecte fundamentale: tipul de metodă de clustering care ar trebui utilizat pentru WSD nesupervizat și tipul de cunoștințe care ar trebui furnizate algoritmului de WSD (lipsit de cunoștințe) rezultat. Studiile noastre preliminare ne-au încurajat să urmăm ambele linii de investigație menționate, cu scopul testării referitor la interogările ambigue în contextul unei aplicații de IR.

Pentru a studia rezultatele dezambiguizării nesupervizate, am luat în considerație două tehnici de clustering complet diferite, care au fost testate cu și fără selecție de caracteristici. În prezentarea din această teză, parametrii modelului matematic utilizat pentru WSD vor fi întotdeauna formulați prin tehnici nesupervizate. Vom compara performanțele algoritmului lipsit de cunoștințe rezultat în prezența și respectiv în absența efectuării unei selecții de caracteristici bazate pe cunoștințe. Cercetarea noastră a investigat câteva surse de cunoștințe de naturi total diferite, considerate adecvate pentru a se realiza selecția caracteristicilor (a se vedea Cap.3-6), precum și posibilitatea combinării unora dintre ele (a se vedea Cap. 5).

În fine, ca urmare a acestor demersuri, am propus (2015) o nouă metodă nesupervizată care utilizează WSD (prin discriminare între sensuri) în IR (a se vedea Cap. 7). Metoda pe care am dezvoltat-o se bazează pe spectral clustering și reordonează o listă de documente regăsite inițial prin “împingerea” la începutul listei a acelor documente care sunt semantic similare cu interogarea. Testele efectuate de noi corespunzător câtorva colecții TREC ad-hoc au demonstrat (a se vedea Cap. 7) că metoda noastră este utilă în cazul interogărilor ce contin cuvinte ambigue.

Interesul nostru s-a concentrat asupra îmbunătățirii nivelului de precizie după 5, 10 și 30 de documente regăsite ($P@5$, $P@10$ și respectiv $P@30$), acestea fiind situațiile cele mai utile aplicațiilor practice din lumea reală. Am arătat (a se vedea Cap. 7) că, în aceste situații, precizia poate fi ridicată cu cel puțin 7,9% deasupra pragurilor de referință state of the art.

Drumul care a condus la nașterea acestei metode ce utilizează WSD în IR nu a fost unul ușor¹. Investigația noastră a început (2008) cu utilizarea unui algoritm pentru WSD nesupervizat lipsit de cunoștințe și care se baza pe clasicul model Bayesian naiv ca tehnică de clustering. Alegerea a fost motivată de credința noastră că potențialul acestui model statistic, cu privire la WSD nesupervizat, rămăsese insuficient explorat la acea vreme. Cu alte cuvinte, am considerat că dezambiguizarea nesupervizată nu exploatase încă la maximum posibilitățile oferite de acest model statistic simplu, dar puternic. În egală măsură, am fost animați de credința că este necesar ca modelului Bayes naiv să-i fie furnizate cunoștințe pentru ca acesta să poată acționa în mod eficient ca tehnică de clustering pentru WSD nesupervizat. Această din urmă convingere ne-a determinat să plasăm procesul de dezambiguizare semantică la granița dintre tehnicile nesupervizate și cele bazate pe cunoștințe și, în cele din urmă, să desfășurăm întregul nostru demers științific la această graniță. Într-adevăr, studiile noastre (2008 – 2013) au arătat că, un

¹ Și a inclus 7 articole științifice de revistă, un articol de conferință și o carte, pentru care am fost autor sau coautor în perioada 2008-2015.

algoritm de dezambiguizare simplu, de bază și lipsit de cunoștințe, în cazul nostru reprezentat de modelul Bayes naiv, poate acționa destul de bine atunci când este alimentat cu cunoștințe în mod adecvat. Aceste cunoștințe pot fi furnizate în diferite moduri și pot fi de diverse tipuri. În lucrările noastre am examinat trei surse complet diferite de astfel de cunoștințe: WordNet (a se vedea Cap. 3), relații de dependență (a se vedea Cap. 4) și web N-gram (a se vedea Cap. 5).

Rezultatele pozitive raportate în literatura domeniului IR (2012) în urma folosirii selecției de caracteristici bazate pe WordNet, propusă de către noi pentru modelul Bayesian naiv, în cadrul unei aplicații de IR (a se vedea Cap. 3), ne-au încurajat să continuăm discuția cu privire la selecția caracteristicilor în cazul acestui model statistic. În mod succesiv am studiat și propus selecție de caracteristici sintactice bazate pe relații de dependență și selecție de caracteristici bazate pe N-gram culese de pe web, acestea din urmă fiind utilizate pentru prima oară, în WSD nesupervizat, de către noi (2012). Rezultatele noastre cu privire la selecția caracteristicilor pentru modelul Bayesian naiv, folosit în dezambiguizarea nesupervizată, au fost publicate de Editura Springer (Hristea 2012)².

Pasul următor al studiului nostru a fost acela de a lua în considerație o tehnică de clustering total diferită, “state of the art”, reprezentată prin spectral clustering (a se vedea Cap. 6). Alegerea noastră a fost determinată nu numai de faptul că aceasta este o tehnică de clustering modernă și cunoscută ca fiind puternică, performantă, ci și de faptul că ea nu necesită în mod esențial selecție de caracteristici. Din contră, testele noastre cu privire la spectral clustering în prezența și respectiv în absența selecției de caracteristici bazate pe WordNet (a se vedea Cap. 6) arată în mod clar faptul că această tehnică de clustering avansată lucrează mai bine atunci când își realizează propria ponderare a caracteristicilor. Spectral clustering este folosit de noi pentru prima oară în WSD și, relativ la seturile de date pentru care s-a făcut testarea, furnizează cele mai bune rezultate în dezambiguizare. Acest fapt ne-a determinat să implementăm spectral clustering (în prezența propriilor caracteristici) în cadrul unei aplicații de IR (2013 – 2015; a se vedea Cap. 7).

Astfel, în finalul acestui amplu studiu, ne-am întors la aplicația de tip IR descrisă în Cap.3 și am utilizat WSD nesupervizat realizat cu spectral clustering relativ la aceeași problemă de IR. Succesul nostru (precizie îmbunătățită cu cel puțin 7,9% deasupra pragurilor de referință state of the art) deschide o întreagă linie de investigație viitoare și câteva noi și promițătoare perspective în această cercetare (așa cum se descrie în Postfața lucrării).

Încă de la începutul acestei cercetări, scopul nostru nu a fost acela de a câștiga o competiție de WSD, ci acela de a furniza o metodă (nouă) de dezambiguizare care ar putea ajuta la câștigarea unei competiții de IR. Din această cauză, testele noastre privitoare la modelul Bayesian naiv implică seturi de date relativ restrânse. În acel stadiu al studiului nostru eram interesați doar să

² Hristea, F.: The Naïve Bayes Model for Unsupervised Word Sense Disambiguation. Aspects Concerning Feature Selection. SpringerBriefs in Statistics Series, Springer (2012).

“simțim” reacția acestui model matematic la caracteristici de naturi complet diferite, precum și la posibile combinații ale lor.

Seturile de date implicate în testări au fost alese de noi, în egală măsură, pentru a facilita comparația cu *singura* abordare similară (selecția caracteristicilor pentru modelul Bayesian naiv antrenat cu algoritmul EM) existentă în literatura de specialitate³ (a se vedea Cap. 1). Același mediu de testare a fost apoi organizat pentru spectral clustering, tot în scopul facilitării comparațiilor. Din contră, testele pe care le-am organizat în cazul aplicației de tip IR finale utilizează colecțiile de referință pentru evaluare, care există și au fost create în mod special pentru task-uri de IR⁴. S-a folosit unul dintre cele mai importante seturi de colecții de testare (a se vedea Cap. 7), care este utilizat pe scară largă de mediile academice și care provine de la Text REtrieval Conference (TREC) - principalul cadru de evaluare în IR.

La momentul studiilor noastre inițiale, literatura domeniului IR afirmase deja (Stokoe, Oakes și Tait 2003)⁵ că singura problemă rămasă în dezbatere este dacă o acuratețe în dezambiguizare mai mică de 90% poate conduce la îmbunătățiri în eficacitatea regăsirii informației, un rezultat despre care suntem de părere că nu poate fi atins de către dezambiguizarea nesupervizată. Prin urmare, încă de la început, ne-am confruntat cu o adevărată provocare științifică: aceea de a accepta sau a nu accepta părerea predominantă existentă conform căreia WSD nu poate îmbunătăți performanța în IR. Am pus sub semnul întrebării, cu tărie, acest punct de vedere al specialiștilor și ne-am propus să demonstrăm că *tipul de WSD* (nesupervizat) folosit în IR este esențial, precum și *tehnica de clustering* pe care acesta o presupune.

Cu excepția primelor două capitole, introductive și tratând cadrul general al cercetării, fiecare capitol al acestei teze conduce la o concluzie științifică distinctă, de sine stătătoare și, din această cauză, este conceput ca fiind independent. Cu toate acestea, numai parcurgerea tuturor capitolelor, în ordinea lor de aici, poate sugera șirul evenimentelor științifice din cariera noastră care ne-au condus la succes în cazul aplicației de tip IR finale. Ultimul capitol al tezei arată modul în care această temă de cercetare își găsește locul în cadrul cercetării noastre postdoctorale globale.

³ Ne exprimăm gratitudinea față de Profesorul Ted Pedersen de la Universitatea din Minnesota, Duluth, S.U.A., care ne-a furnizat seturile de date necesare pentru efectuarea testelor și a comparațiilor prezentate, referitor la adjective și verbe.

⁴ Accesul la datele de testare specifice domeniului IR a fost asigurat de către partenerul nostru francez pentru acest studiu, IRIT (Institut de Recherche en Informatique de Toulouse).

⁵ Stokoe, C., Oakes, M.P., Tait, J.: Word sense disambiguation in information retrieval revisited. In: Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '03), pp. 159-166. Toronto, Canada (2003).