

Universitatea „Lucian Blaga” din Sibiu
Facultatea de inginerie “Hermann Oberth”
Catedra de calculatoare și automatizări



Extragerea unor trăsături relevante din masive de date ne-structurate semantic

Referat de doctorat nr. 3

Titlul tezei: „Contribuții la extragerea automată de
cunoștințe din masive de date”

Autor:

asist. drd. ing. Daniel MORARIU

Coordonator științific:

prof. univ. dr. ing. Lucian N. VINTAN

SIBIU, 2006

1 Rezumat

În ultimii ani, datorită impactului calculatoarelor în toate domeniile, cantitatea de informații în format text devine tot mai mare. Recunoașterea și clasificarea automată a acestora devine tot mai necesară. Acestea sunt considerate ca fiind date semi-structurate deoarece ele sunt mai degrabă nestructurate decât complet structurate. Ele conțin o mică parte de organizare a informațiilor, care în majoritatea cazurilor nu este completată de către autorul documentului.

Tehnicile tradiționale pentru recunoașterea informațiilor devin inadecvate pentru aplicațiile de căutare în depozite de date. De obicei, doar o mică parte din toate documentele disponibile sunt relevante pentru utilizator la un moment dat. Fără a ști ce conține documentul, este dificil să formulăm interogări pentru analiza și extragerea informațiilor necesare. Pentru a putea extrage doar acele documente relevante, utilizatorii au nevoie de componente care să îi ajute să poată compara documentele, cum ar fi eficiența și relevanța lor în raport cu o anumită interogare, sau pentru a putea găsi șabloane pentru a le putea indexa și regăsi mai ușor.

Un număr impresionant de documente se găsesc pe WEB, iar utilizatorul are nevoie de componente de clasificare automată a acestora pentru a le putea gestiona. Gestiunea lor a devenit o problemă importantă în ultima perioadă. Devine esențială existența unor programe de organizare automată a documentelor în clase pentru a facilita analiza și recunoașterea acestora. O posibilă procedură generală pentru rezolvarea acestei probleme este alegerea unei mulțimi de documente pre-clasificate (mult mai mică în comparație cu documentele existente) și considerarea acestora ca fiind documente de antrenament. Mulțimea de antrenament este apoi analizată pentru a putea extrage o schemă de clasificare. Această schemă este finisată utilizând o mulțime de documente de test. După aceasta, schema astfel obținută poate fi utilizată pentru clasificarea celorlalte documente existente. Analiza clasificării decide care mulțime de perechi de atribut-valoare are o putere de discriminare mai mare în determinarea claselor. O metodă eficientă pentru clasificarea documentelor este de a exploata clasificările bazate pe asocieri, unde documentele sunt clasificate pe baza unei mulțimi de asocieri și a frecvenței de șabloane. Metodele de clasificare pe baza asocierilor respectă următoarele etape: (1) cuvintele cheie și termenii pot fi extrași prin tehnici de recunoaștere a informațiilor și tehnici de analiză a asocierilor simple; (2) ierarhiile de concepte de cuvinte cheie sau termeni pot fi obținute utilizând termenii claselor disponibile, încrederea în cunoștințele expertului sau unele sisteme de clasificare pe bază de cuvinte cheie. Documentele din mulțimea de antrenament pot fi de asemenea clasificate în ierarhi de clase. Metodele de minerit pe baza asocierii termenilor pot fi apoi aplicate pentru a descoperi mulțimi de termeni de asociere care pot fi utilizați pentru a maximaliza diferențele dintre două clase de documente. Acestea produc o mulțime de reguli de asociere pentru fiecare clasă de documente. Astfel de reguli de asociere pot fi ordonate pe baza frecvenței lor de apariție și a puterii de discriminare, utilizate apoi pentru clasificarea de noi documente.

Clasificarea documentelor text este un proces foarte general incluzând o mulțime de cerințe care trebuie realizate pentru rezolvarea problemei. Fiecare dintre aceste cerințe are o influență foarte mare asupra rezultatului final al clasificării. Acest proces poate fi privit ca un flux de date în mai multe etape. În care, în fiecare etapă, văzută ca o cutie neagră, se primesc informații, se procesează și se transferă mai departe. Fiecare etapă din acest flux de date poate avea unul sau mai mulți algoritmi atașați. La un moment dat, pentru fiecare etapă putem alege unul dintre algoritmi atașați

cu diferiți parametri de intrare. În primul meu referat am prezentat câteva părți ale fluxului; în acest referat am îmbunătățit anumite părți deja existente și am finalizat fluxul de date.

Astfel, în primul referat am prezentat unele tehnici de pre-procesare a documentelor text. În special am prezentat pre-procesarea bazei de date Reuters-2000. Am continuat cu o scurtă introducere a procesării paginilor web, axându-mă pe diferențele dintre acestea și documentele text. În această etapă datele de intrare sunt reprezentate prin documente text (fișiere text sau pagini web) și datele de ieșire sunt reprezentate de vectori de trăsături care caracterizează fiecare document din mulțimea de documente. Trăsăturile din acești vectori sunt de fapt frecvențele de apariție ale cuvintelor în documentul respectiv. Această etapă conține un modul de eliminare a cuvintelor irelevante din punct de vedere semantic (stop-words), un modul de extragere a rădăcinii cuvintelor pentru limba engleză și un modul de numărare a cuvintelor. Datorită dimensiunii impresionante a vectorilor rezultați această etapă continuă cu o etapă de selectare a trăsăturilor relevante din acești vectori. Astfel, în cel de-al doilea referat am prezentat trei metode de selecție a trăsăturilor relevante: selecția aleatoare, selecția bazată pe câștigul informațional și selecția bazată pe vectori suport. O nouă metodă de selecție a trăsăturilor caracteristice relevante bazată pe algoritmi genetici este prezentată și implementată în acest referat.

În cel de-al doilea referat am prezentat în detaliu algoritmul utilizat pentru clasificarea documentelor, bazat pe vectorii suport (SVM) și învățarea utilizând nuclee. Acolo m-am axat în general pe procesul de clasificare a documentelor. De asemenea am prezentat o adaptare a acestui algoritm pentru a putea fi folosit și în contextul învățării nesupervizate. Apoi am prezentat o metodă nouă de corelare a parametrilor nucleelor și îmbunătățirile aduse aplicației noastre în comparație cu o implementare a acestei tehnici des folosită în literatură, LIBSVM. În acest referat vom prezenta rezultatele care justifică corelarea parametrilor din nuclee precum și metodologia de alegere a acestor corelații.

În ultima etapă din acest flux prezentăm o tehnică de dezvoltare a unui meta-clasificator pentru a îmbunătăți acuratețea clasificării. Clasificatorii de bază din acest meta-clasificator utilizează tehnica vectorilor suport. Această metodă este prezentată în acest referat și încearcă să obțină rezultate mai bune decât oricare dintre clasificatorii de bază componenți.

Acest referat conține unele contribuții care îmbunătățesc fluxul clasificării făcându-l mult mai fiabil pentru lucrul cu mulțimi mari de date. Una dintre acestea este abilitatea aplicației noastre de a lucra cu mult mai multe documente. Rezultatele prezentate până acum sunt obținute pe o dimensiune mică a setului de date în comparație cu toate datele existente în baza de date. Utilizarea tuturor datelor existente ar face procesul de învățare inefficient din punct de vedere al timpului. În acest referat am prezentat o metodologie care face ca aplicația noastră să fie capabilă să lucreze cu o dimensiune impresionantă a datelor de intrare și cu pierderi cât mai mici din punct de vedere al acurateții clasificării. Această metodologie are două obiective majore antagonice: o dimensiune mare a datelor de intrare și un timp de învățare minim pentru o clasificare optimă.

Pentru a avea date cât mai realiste în cercetarea noastră, ar trebui să prezentăm media rezultatelor obținute pe mai multe perechi de mulțimi de date de antrenare-testare. Pentru fiecare pereche de mulțimi ar trebui calculată separat acuratețea clasificării. În toate rezultatele am prezentat acuratețea obținută pe o singură pereche de mulțimi de antrenare-testare. În ultimul capitol al acestui referat am prezentat o parte din rezultatele obținute și pe o altă pereche de mulțimi. Am făcut aceste experimente pentru a vedea dacă rezultatele noastre sunt sau nu apropiate de rezultatele ce s-ar obține în urma calculului mediei. Aceste rezultate ne vor furniza nivelul de încredere al aplicației noastre din punct de vedere al acurateții prin calcularea mediei acuratețiilor obținute pe mai multe perechi de mulțimi de antrenare-testare.

Experimentele prezentate în acest referat sunt efectuate folosind colecția de documente Reuters 2000, care are o dimensiune de 948 Mb de știri în format comprimat. Colecția include un număr de 806791 știri publicate de agenția de presă Reuters în perioada 20.07.1996 - 19.07.1997. Statistic vorbind articolele au 9822391 paragrafe, 11522874 propoziții și au 310033 rădăcini distincte de cuvinte. Documentele sunt clasificate de Reuters după trei categorii distincte. După *Regiunea* la care se referă articolul, existând 366 regiuni distincte. După un *Cod Industrial* în funcție de ramura industrială la care se referă articolul, existând un număr de 870 coduri industriale. Și după *Categorie* care este propusă de către Reuters, existând 126 categorii distincte. Douăzeci și trei dintre aceste categorii nu au nici un document atașat. Datorită acestei dimensiuni imense a datelor de intrare am prezentat rezultatele doar pe o submulțime a acestora. Astfel din toate aceste documente am selectat doar acelea pentru care codul industrial este „System software”. Am obținut un număr de 7083 documente care pot fi reprezentate folosind 19038 trăsături distincte și 68 categorii. Am reprezentat documentele ca și vectori de frecvențe de termeni aplicând o filtrare a cuvintelor de legătură prin folosirea unei liste standard de 510 cuvinte (pentru limba engleză) și am extras rădăcina cuvintelor rămase. Din aceste 68 de categorii rezultate am eliminat acele categorii care sunt slab sau excesiv reprezentate. Astfel am eliminat acele categorii care conțin mai puțin de 1% din numărul total de documente (7083). De asemenea am renunțat la acele categorii care conțin mai mult de 99% din documente. Eliminarea a fost necesară deoarece păstrând aceste categorii există riscul ca algoritmul nostru de învățare să învețe doar o singură categorie și să clasifice toate documentele în acea categorie, ignorând celelalte categorii existente. După acest pas am obținut 24 de categorii diferite și 7053 documente care au fost împărțite aleator în setul de antrenament (4702) și cel de test (2531). În partea de extragere a trăsăturilor am luat în considerare atât întreg conținutul articolului cât și titlul articolului.

Capitolul doi conține experimentele care au condus la alegerea corelațiilor dintre parametrii nucleelor. Capitolul trei conține o nouă metodă de selecție a trăsăturilor caracteristice bazată pe algoritmi genetici utilizând o funcție de calcul a acurateței bazată pe tehnica vectorilor suport. Următorul capitol încheie procesul de clasificare automată a documentelor prin prezentarea câtorva metode pentru implementarea unui meta selector care să conducă la îmbunătățirea acurateței finale a clasificării. Capitolul cinci prezintă influența unei mari cantități de date de intrare asupra algoritmului nostru. Tot aici prezentăm o strategie care permite algoritmului nostru să lucreze cu date de intrare mari într-un timp mic și cu pierderi minore. În capitolul 6 prezentăm câteva rezultate obținute pe o altă împărțire a setului de date în date de test și date de antrenament. Ultimul capitol prezintă și dezbate cele mai importante rezultate obținute și propune câteva tendințe de viitor.

Mulțumiri

În încheierea acestei introduceri aș vrea să exprim sincerele mele mulțumiri către conducătorul meu de doctorat prof. dr. ing. Lucian VINȚAN pentru coordonarea științifică pe perioada de pregătire a acestui doctorat, pentru discuțiile extrem de stimulative pe care le-am avut și pentru întreg sprijinul acordat. De asemenea, rămân îndatorat pentru sprijinul profesional competent acordat mie cu generozitate, încă de la începuturile perioadei de pregătire prin doctorat, de către profesorii: mat. Ioana MOISIL, ing. Boldur BĂRBAT, ing. Daniel VOLOVICI, ing. Dorin SIMA și ing. Macarie BREAZU.

În același timp aș dori să mulțumesc firmei SIEMENS AG, CT IC MUNCHEN, Germany, în special vicepreședintelui Dr. h. c. mat. Hartmut RAFFLER, pentru sugestiile profesionale foarte folositoare și pentru suportul financiar acordat pe parcursul acestui program de doctorat. Aș dori să mulțumesc tutorelui meu de la firma SIEMENS, Dr. Volker TRESP, Senior Principal Research Scientist în Neural Computation, pentru sprijinul științific acordat și pentru îndrumarea în acest imens și dificil domeniu al cercetării. De asemenea aș dori să mulțumesc d-lui Dr. Kai Yu de la aceeași companie pentru informațiile furnizate în dezvoltarea ideilor mele. La sfârșit aș dori să mulțumesc tuturor acelorora care m-au ajutat în pregătirea acestui referat.