

## Rezumat

În timp ce tot mai multe informații în format text devin disponibile, recunoașterea automată a acestora este dificilă fără o bună indexare și un rezumat bun al conținutului documentelor. Catalogarea textului este o soluție a acestei probleme și constă în clasificarea documentelor într-un set de categorii predefinite. În ultimi ani au fost aplicate un număr crescând de metode de clasificare (clustering) și de tehnici de învățare automată.

Documentele sunt de obicei reprezentate prin vectori rari modelați într-un spațiu vectorial, unde fiecărui cuvânt din vocabular îi corespunde o axă de coordonate și numărul de apariții ale cuvântului în document reprezintă valoarea componentei din vector în reprezentarea documentului. Dimensiunea mare a spațiului trăsăturilor este o problemă majoră în catalogarea textului. Spațiul trăsăturilor constă din termeni unici care apar în documente. Când modelăm o colecție de documente spațiul trăsăturilor poate fi de ordinul zecilor sau sutelor de mii de termeni. Antrenarea clasificatorilor pe o astfel de colecție mare de documente necesită mult timp și memorie. De aceea se încearcă aplicarea de metode diverse pentru a reduce spațiul de reprezentare și timpul de răspuns. Așa cum vom vedea rezultatele sunt mai bune când lucrăm cu o dimensiune mult redusă a vectorilor. Odată cu creșterea dimensiunii nu va crește neapărat și acuratețea clasificării, urmând ca aceasta să scadă substanțial când se utilizează o dimensiune foarte mare sau întreaga dimensiune a spațiului.

În acest raport voi prezenta un studiu comparativ între câteva metode de selecție al trăsăturilor caracteristice (Information Gain, Mutual Information și Support Vector Machine) precum și tipul de reprezentare al datelor de intrare în învățarea folosită pentru clasificarea documentelor. De asemenea voi prezenta o metodă folosită cu mult succes în ultimi ani în cadrul problemelor de clasificare pe date de intrare neliniar separabile. Voi începe prin a prezenta programele dezvoltate pentru partea de procesare a documentelor și crearea vectorilor caracteristici și voi continua cu cele dezvoltate pentru partea de clasificare (clusterare) a acestora utilizând tehnicile bazate pe vectori suport și nuclee.

Am considerat partea de minerit al textului ca o aplicație care folosește tehnicile de minerit al datelor pentru a extrage semnătura fiecărui document (a crea vectorul de trăsături). Ca și bază de date am utilizat baza Reuters care este des utilizată în articolele de procesare și clasificare a documentelor. Baza de date Reuters conține articolele de presă publicate de agenția Reuters. Pornind cu o mulțime de  $d$  documente și  $t$  termeni (cuvintele din documente) putem modela fiecare document ca un vector  $v$  într-un spațiu cu  $t$  dimensiuni  $\mathcal{R}^t$ . În faza de clasificare, am utilizat tehnica bazată pe învățarea cu nuclee și vectori suport. Aceasta este o tehnică puternică folosită cu mult succes în problemele de clasificare neliniar separabile. Avantajul ei este că poate fi aplicată pe mulțimi mari de date. Astfel putem testa ușor influența numărului de trăsături caracteristice asupra acurateții clasificării. Am implementat clasificatorul pentru 2 tipuri de nuclee, nucleul polinomial și nucleul Gaussian. Voi prezenta rezultatele obținute atât pentru clasificarea la două clase cât și pentru clasificarea în mai multe clase. Pentru clasificarea în două clase am luat în considerare doar documentele care aparțin la o clasă comparativ cu restul documentelor. Pentru clasificarea la mai multe clase am repetat clasificarea la două clase pentru fiecare topic (categorie în care este grupat documentul în faza de învățare) obținând mai multe funcții de decizie (funcții care fac distincția între documentele dintr-un topic și restul documentelor). În pasul de minerit al textului am utilizat această tehnică pentru selecția trăsăturilor caracteristice. O să prezint și o vizualizare grafică a rezultatului clasificării (clusterării) utilizând această tehnică.

Voi prezenta rezultate pe diferite tipuri de nuclee și pentru diferite tipuri de reprezentare a datelor de intrare încercând să găsc o corespondență optimă pentru îmbunătățirea acurateței clasificării. Am utilizat o formă cât mai simplificată a nucleelor, fără a reduce performanțele, uneori chiar îmbunătățindu-le, astfel parametrii utilizați de acestea să fie cât mai ușor de specificat. Pentru algoritmi de clasificare și clusterare implementați și prezentați în acest raport am folosit o tehnică nouă utilizată cu mult succes în ultimi ani mai ales pentru lucrul cu mulțimi de date neliniar separabile - Support Vector Machine.

Datele de intrare sunt reprezentate în diferite formate și am analizat influența acestora asupra tipului nucleului. Am utilizat trei tipuri de reprezentări. Reprezentarea *binară* în care atributele iau doar valoarea „0” sau „1” („0” dacă cuvântul respectiv nu apare în document și „1” dacă apare fără să ne intereseze numărul de apariții al acestuia). Reprezentarea *nominală* unde atributul păstrează numărul de apariții al cuvântului în vectorul de frecvențe normalizat folosind norma normală. În reprezentarea *Connel Smart* atributul își păstrează numărul de apariții al cuvântului în vectorul de frecvențe dar este normalizat utilizând o altă formulă. Rezultatele experimentale sugerează că alegerea unui procent de 4 – 10% din setul inițial de trăsături conduce la cele mai bune rezultate. Deci după cum se poate observa din experimente nu este nevoie să utilizăm multe trăsături caracteristice în procesul de învățare. Utilizarea prea multor trăsături poate fi chiar dăunătoare ducând și la scăderea performanțelor învățării.

În secțiunea 2 sunt prezentate metodele de selecție a trăsăturilor și detaliile de construcție a mulțimilor de test și de antrenament care vor fi utilizate în acest raport. Voi prezenta un scurt rezumat al metodelor de selecție a trăsăturilor în contextul strategiei de antrenare propusă. În secțiunea 3 voi descrie algoritmi de clasificare și de clusterare bazați pe tehnica de învățare utilizând vectori suport și nuclee. Aici voi prezenta și detaliile de implementare a acestor algoritmi. În secțiunea 4 voi prezenta modul de utilizare al aplicațiilor realizate pentru a crește acuratețea clasificării. Tot aici voi descrie experimentele realizate și voi compara rezultatele obținute cu rezultatele obținute cu altă aplicație (LibSvm). În secțiunea finală voi prezenta concluziile rezultate precum și stadiul lucrării și posibile dezvoltări viitoare.

## Mulțumiri

Pe lângă părinții mei o mulțime de oameni merită recunoștință, nu-i pot aminti pe toți, dar le mulțumesc tuturor. Țin totuși să le mulțumesc câtorva care mi-au fost dascăli și care m-au îndrumat în acest proiect. În primul rând țin să mulțumesc conducătorului meu de doctorat prof. univ. dr. ing. Lucian VINȚAN, căruia îi sunt recunoscător în primul rând pentru responsabilitatea cu care dânsul își privește meseria, apoi pentru discuțiile extrem de stimulative pe care le-am avut și pentru 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.

Aș vrea să mulțumesc firmei SIEMENS AG, Corporate Technology, Information & Communication Division, MUNCHEN Germania, și în special vicepreședintelui Dr.h.c. Hartmut RAFFLER, pentru sponsorizarea acordată acestui program de doctorat. Deosebit de utile mi-au fost sugestiile profesionale precum și suportul material oferit. De asemenea vreau să-i mulțumesc tutorelui meu de la firma SIEMENS, dr. Volker TRESP, Senior Principal Research Scientist in Neural Computation, pentru sprijinul oferit și pentru îndrumarea în acest imens și dificil domeniu al cercetării. Aș vrea de asemenea să-i mulțumesc d-lui Dr. Kai YU de la aceeași companie pentru informațiile furnizate în dezvoltarea ideilor.