

Îmbunatatirea clasificarii bayesiene prin algoritmi genetici, cu aplicatii în text mining

Flavian VASILE, masterand
Managementul Informatizat al Proiectelor, A.S.E. Bucuresti

This paper studies the possibility of increasing the accuracy of classification and obtaining knowledge about the quantitative correlations between the representative features for the classification. The algorithm starts from a naive Bayesian classifier and incrementally increase complexity with a state-of-the-art genetic algorithm: N.E.A.T.

The paper is based on the results obtained using the application "Mentat"[2], developed and presented in the graduation thesis "Evolutive Text Mining" by the same author [1]. It also contains a survey about the types of software developed for text mining and future development.

Keywords: text mining, Bayesian classifier, genetic algorithm, correlation, representative features, optimising, networks.

Introducere

Data mining-ul a aparut ca o continuare a metodelor traditionale de analiza a depozitelor de date (pe volume mari). Datorita costurilor de implementare fidicate a aparut necesitatea identificarii unor solutii de extragere a esentialului dintr-un volum foarte mare de date, la costuri mult mai reduse, solutii cunoscute astazi sub denumirea de data mining. Odata cu expansiunea Internet-ului si a informatiei de tip text în format electronic, a aparut necesitatea extragerii automate de cunostinte si din text si astfel data mining-ul a cunoscut o noua specializare: **text mining-ul**.

Spre deosebire de data mining, text mining-ul presupune un software care se adreseaza publicului larg consumator de servicii în retea, motivele pentru aceasta fiind universalitatea cererii de achizitionare de informatie în timp real si costurile mici (pretul conexiunii) de achizitionare a informatiei, comparativ cu data mining-ul.

Text mining-ul are drept obiectiv principal extragerea automata de cunostinte, ce trebuie sa îndeplineasca la rândul lor cerintele de: *noutate, validitate, operationalitate*.

Aplicatii de text mining

Aplicatii de text mining se pot clasifica în aplicatii online si aplicatii offline. În pre-

zent, cele mai cunoscute aplicatii ale text mining sunt:

a) Aplicatii online pentru:

- *cautarea inteligenta pe Internet*, care presupune *analiza de continut* (cu ajutorul tehnicilor de text mining documentele regasite de un motor de cautare sunt filtrate fiind pastrate doar rezultatele cu sensul cautat), *dezvoltarea unui profil al utilizatorului* (documentele sunt aduse automat fara o solicitare expresa din partea utilizatorului, plecând de la un profil al preferintelor conform carora programul cauta independent documente relevante).

- *regasirea stirilor interesante*: selectarea de stiri de pe Usenet este o provocare care apare în viata de zi cu zi a multor utilizatori de Internet.

- *regasirea de raspunsuri la întrebări*: exista întrebări frecvente la care altcineva poate raspunde sau a raspuns deja. Acestea se numesc *frequently asked questions* (F.A.Q.) si se gasesc împreuna cu raspunsurile lor în baze de date.

- *filtrarea postei*: partitionarea multimii scrisorilor electronice în grade de prioritate în functie de reactiile anterioare ale utilizatorului si emiterea de recomandari de stergere a unor mesaje pasibile de a fi comerciale sau neinteresante.

b) Aplicatii offline pentru:

- *clasificarea documentelor preluate de pe Internet.*
- *rezumarea documentelor* (obtinând astfel documente condensate si rezumate, abstracturi ale lucrarilor depozitate)
- *concatenarea documentelor* (de exemplu crearea de relatari ale unor stiri pe baza informatiilor preluate din mai multe surse)

Aplicatiile mentionate nu respecta întru totul obiectivul de extragere a cunostintelor, aceste sisteme realizând cel mult extragere de informatii. Extragerea cunostintelor reclama mai mult de la metodele de text mining si anume înțelegerea, macar aproximativa, a documentelor si crearea unui sistem de generare într-o forma inteligibila de cunostinte noi. În prezent, interpretarea rezultatelor este realizata în continuare de utilizator pe un numar însa mai mic de informatii mai exacte.

O directie de dezvoltare a text mining-ului.

Una dintre directiile de dezvoltare cu un impact potential deosebit de amplu este extragerea automata interdependentelor dintre termenii relevanti ai domeniului supus analizei.

Deoarece în text mining se pleaca de la ideea ca un sistem artificial de tratare a limbajului are obiective limitate si nu trebuie sa obtina o interpretare corecta a înțelesului textelor tratate, regulile sintactice folosite în limbaj nu trebuie neaparat implementate. Totusi, tratarea cuvintelor ca entitati independente si clasificarea informatiilor doar dupa indici derivati din frecvente si raportari la nivel de cuvânt are unele limite deoarece *reprezentarea datelor nu este o problema colaterala, ea exprima în fapt ipotezele care se fac asupra proprietatilor sau mai bine zis metaproprietatilor limbajului si a formei sale concrete, textul.*

În prezent, cea mai mare parte din metodele de text mining se bazeaza pe metode inductive (de construire a unor rețele neuronale sau a unor arbori de decizie), care de obicei fac o cautare locala, de tip greedy, în spatiul legilor.

Algoritmii genetici au caracteristici interesante pentru aplicatiile de data si text mining, cum ar fi:

- capacitatea de a explora mai în profunzime spatiul de solutii,
 - capacitatea de a înlatura problema interactiunii atributelor,
 - capacitatea de a opera o cautare cu final deschis pentru descoperirea de pattern-uri,
 - capacitatea de a descoperi cunostinte de nivel calitativ înalt (reguli if-then, etc.).
- Ca dezavantaj, algoritmii genetici au timpuri mari de prelucrare în cazul volumelor mari de date.

Implementarea algoritmilor genetici pe rețele neuronale. Algoritmul N.E.A.T.

Deoarece algoritmii genetici au nevoie de o solutie de start, aceasta solutie poate fi una din structurile traditionale de tratare a documentelor.

În acest scop am introdus în cadrul algoritmului un clasificator Bayesian naiv ca genom de start. Optimizarea structurii rețelei ar fi dus la surprinderea interdependentelor pornind de la presupunerea ca o structura ce reprezinta din ce în ce mai bine domeniul documentelor considerate drept interesante de catre utilizator se apropie asimptotic de structura reala a conceptelor din cadrul aceluia domeniului. Astfel o clasificare corecta a documentelor înseamna în cele din urma o "învatare" a regulilor de clasificare folosite de catre utilizator, iar interdependentele ce se stabilesc între termeni în cadrul structurii de implementare a regulii sunt valide si operationale.

Metoda de optimizare a structurii rețelei clasificatorului naiv propusa în aceasta lucrare se bazeaza pe *algoritmul NEAT* [3] (NeuroEvolution of Augmenting Topologies) prezentat în lucrarea "Evolving Neural Networks through Augmenting Topologies", autori: Kenneth O. Stanley si Risto Mukklainen, Department of Computer Sciences, The University of Texas at Austin, un raport tehnic datat 28 Iunie

2001, ceea ce face din acest algoritm unul din cei mai noi algoritmi genetici.

Algoritmul genetic dezvoltat de cei doi cercetatori americani evolueaza în paralel atât structura cât și ponderile rețelei, lucru rar întâlnit până acum în literatura de specialitate și de multe ori soldat cu eșecuri. Inovațiile pe care algoritmul le aduce cresc valoarea perspectivei genetice din care domeniul se hrănește cu idei. Astfel, *cele trei puncte cheie* pe care se bazează evoluția sunt:

1. codificarea genetică cu markere genetice ce identifică în mod unic genele și folosirea în cadrul reproducerii a informațiilor despre gene;
2. protejarea inovației genetice prin speciere;
3. creșterea incrementală de la o structură minimală.

1. *Codificarea genetică.* Una din cele mai cunoscute probleme din cadrul algoritmilor genetici este problema permutațiilor sau a competiției convențiilor. Aceasta înseamnă că există mai multe modalități de a reprezenta o soluție la o problemă de optimizare a ponderilor într-o rețea. Când genomurile ce reprezintă aceeași soluție nu au aceeași formă, reproducerea nu este împiedicată iar urmașul este foarte probabil nefezabil (ca și în cazul reproducerii umane cosangvine!).

Algoritmul rezolvă această problemă introducând o notatie unică a genelor, fiecare dintre acestea având un identificator dat de ordinea apariției, numit "historic marking" ce are omolog în natură: genele se interschimbează după ce genoamele participante sunt supuse unui test de asemănare și genele cu același stramos sunt aliniate.

2. *Specierea.* Inovația apare în cadrul evoluției rețelelor neuronale atunci când o nouă substructură este adăugată prin mutație. Această adăugare de noi noduri și/sau conexiuni prin intermediul mutațiilor nu poate să ducă din prima generație la un fitness crescut ci la unul mai scăzut. Până când acestor noi structuri vor dobândi o funcție și vor fi conectate la întreaga rețea, scorul genomului scade vertiginos și poate

fi eliminat din competiția reproducerii înainte de a putea să-și dovedească capacitățile. Dacă în schimb toate structurile noi ar fi păstrate artificial pentru un timp oarecare ar apărea probleme de calcul și de stabilire a acestei perioade. Autorii au folosit încă odată o soluție folosită încă și de natură: specierea. Astfel, fiecare individ este nevoit să concureze în interiorul speciei sale iar pe ansamblu speciile concurează între ele prin intermediul randamentului lor mediu.

Specierea are nevoie însă de o funcție de asemănare pentru a putea grupa indivizii în specii. Aceasta face uz tocmai de identificatorii genelor pentru a putea recunoaște indivizii asemănători. Populația se segmentează astfel pe specii, fiecare într-o anumită nișă variabilă. Funcția de asemănare este de forma:

$$d = \frac{c_1 E}{N} + \frac{c_2 E}{N} + c_3 \bar{W}$$

unde:

E = numărul de gene în exces (al caror identificator este mai mare de cel mai mare identificator comun, aparțin individului mai "tânăr" genetic)

D = numărul de gene diferite (al caror identificator este mai mic de cel mai mare identificator comun)

N = numărul de gene din genomul mai lung

W_{med} = diferența medie a ponderilor genelor comune

c_1, c_2, c_3 = ponderea fiecărui factor

Fitnessul unei specii este media fitnessurilor indivizilor (aritmetică, geometrică sau armonică).

3. *Creșterea incrementală de la o structură minimală.* În majoritatea cazurilor algoritmi genetici pornesc de la o colecție de genomuri generate aleator, ceea ce generează o multitudine de probleme: mărirea artificială a spațiului de căutare, numărul crescut de calcule, probleme de codificare și generare aleatoare, obținerea unei soluții prea complexe. Soluția este de obicei introducerea unui factor de penalizare a creșterii complexității în cadrul funcției de evaluare, dar algoritmul tot pleacă de la un

spatiu prea complex si ireal prin raport cu realitatea.

În realitate, un moment de start este determinat de evolutia incrementală de până atunci și situația respectivă este modelată tot de aceleași legi și nu este haotică. De aceea, algoritmul NEAT pleacă de la o populație de indivizi identici cu o structură minimală (rețea monostrat complet conectată – fără strat ascuns).

Evoluția va fi protejată prin speciere, lucru imposibil până acum. În acest mod, orice genom produs și perpetuat are o justificare de a exista, este o inovație validă și care a depășit soluția de start.

Elemente novative aduse de modificarea algoritmului pentru problemele de text mining.

Algoritmul și elementele novative au fost testate prin dezvoltarea aplicației Mentat. În cazul concret al aplicației „Mentat”, a fost folosită ca soluție de start un clasificator clasic pentru text mining, o rețea Bayesiană monostrat antrenată, ce opera deja clasificări performante. Folosirea rețelei antrenate este posibilă tocmai datorită ipotezei naive a independenței atributelor ce duc la forma monostrat a rețelei. Evoluția ulterioară a rețelei avea să producă două efecte în cazul unui succes:

- să îmbunătățească rezultatele clasificatorului naiv, considerat cel mai potrivit pentru aplicațiile de text mining;
- să ofere soluției o justificare teoretică prin eliminarea ipotezei simplificatoare și introducerea de dependențe prin evoluarea treptată a unui strat ascuns.

Acest strat ascuns reprezintă evoluția genetică a structurii și are două surse de apariție:

1. mutația de adăugare a unui nod pe un arc ce leagă direct stratul de intrare și cel de ieșire
2. mutația de adăugare a unui arc între două noduri neconectate până acum (deoarece toate nodurile de intrare și ieșire erau complet conectate, aceste arce nu apar decât între nodurile din stratul ascuns și celelalte noduri)

Noile rețele sunt separate pe specii conform algoritmului iar conectarea noilor noduri este “premiată” pentru fiecare conexiune nouă introdusă în sistem. Această premieră se păstrează doar în cazul în care ipoteza introdusă prin noua conexiune este verificată la noua evaluare a generației: dacă mai multe noduri din stratul de intrare se conectează la același nod din stratul ascuns înseamnă că sunt interdependente și determină împreună rezultatul din stratul de ieșire cu o mai mare putere decât fiecare din ele datorită efectului de sinergie. Acest principiu este modelat matematic prin introducerea premierii.

Aceasta este pentru moment nimic mai mult decât o ipoteză: nu a fost verificată pe o instanță reală. Dacă instanța probează ipoteza, premiera rămâne. Dacă nu, pentru fiecare conexiune neverificată rețeaua primește o penalizare mai mare decât premiera inițială. Astfel sunt promovate structurile cu adevărat inovative.

De ce este nevoie de speciere? Deoarece testarea se face la fiecare generație din motive de calcul doar pentru un singur document pe clasă, ceea ce înseamnă că probarea unei ipoteze în întregime să se face de-a lungul mai multor generații.

Rezultate

Conform rezultatelor obținute, perspectivele algoritmului sunt promitatoare. Clasificarea, în cazul unei antrenări medii (10 documente/domeniu), se face foarte aproape de o clasificare umană, iar valoarea extragerii cunostintelor exprimată prin frazele-cheie extrase la nivelul claselor deși are o interpretare subiectivă, a fost probată prin căutări multiple ale căror rezultate au fost de fiecare dată adăugarea la aceeași clasă cu un scor crescut.

Concluzii

Versiunea actuală a aplicației a urmarit în primul rând probarea unor ipoteze din cadrul algoritmului și adâncirea cunoașterii informației structurată textual. Deoarece rezultatele au fost satisfăcătoare, o versi-

une ulterioara, ar putea cuprinde ca îmbunătățiri:

- *integrarea pe motor de cautare/agent inteligent de cautare*: cautarile dupa frazele-cheie ale unui domeniu s-ar putea face automat iar rezultatele s-ar constitui ca noi cunostinte ale domeniului, astfel întreaga etapă de achiziționare de cunostinte ar putea deveni automată;
- *depistarea unor noi domenii de interes* nu numai din bazinul documentelor neclasificate ci și din adresele cautate;
- *integrarea altor surse de informații* în format electronic: ziare electronice, newsgroup-uri, discussion groups, e-mail;
- *aducerea optională și a imaginilor* ce însoțesc documentația, pentru a obține o bază multimedia de cunostinte;
- *explorarea și mai detaliată a fiecărui domeniu* cu ajutorul arborilor de decizie implementați pe grupurile de cuvinte interdependente identificate.

Bibliografie

1. Flavian Vasile – “Text Mining Evolutiv” – Lucrare de licență, ASE, București, 2002;
2. Flavian Vasile – “Mentat” – aplicație text mining pentru Internet 2002 www.idea.boys.ro;
3. Constanta-Nicoleta Bodea – “Inteligența artificială – Calcul neuronal”, Editura ASE 2002;
4. Kenneth O. Stanley, Risto Mukkainen – “Evolving Neural Networks through Augmenting Topologies”, Technical Report TR-AI-01-290, June 28, 2001;
5. Anders Holst – “The Use of a Bayesian Neural Network Model for Classification Tasks”, Dissertation - September 1997, STOCKOLMS UNIVERSITET;
6. Dunja Mladenic – „Machine Learning on non-homogeneous, distributed text data”, Ljubljana, 1998
7. Robert Armstrong, Dayne Freitag, Thorsten Joachims, Tom Mitchell – „Web Watcher – A Learning Apprentice for the World Wide Web”;
8. Dunja Mladenic, Marko Grobelnik – „Predicting content from hyperlinks”
9. Peter F Brown, Vincent J Della Pietra, Peter V de Jonza, Jenifer C Lai – „Class-based n-gram models of natural language”;
10. David D Lewis, William A Gale – „A Sequential Algorithm for Training Text Classifiers” SIGIR 1994;
11. Andrew McCallum, Kamal Nigam, Jason Rennie, Kristie Seymore – „A Machine Learning Approach to Building Domain-Specific Search Engines” www.cora.justresearch.com;
12. Jose M. Bernardo – „Bayesian Statistics” Departamento de Estadístico, Facultad de Matemáticas, 46100 – Burjassot, Valencia, Spain;
13. Kamal Paul Nigam – „Using Unlabeled Data to Improve Text Classification”, May 2001;
14. David Heckerman (heckerma@microsoft.com) - „A tutorial on Learning with Bayesian Networks” November 1996;
15. Jason D.M.Rennie - „Improving Multi-class Text Classification with Naive Bayes”;
16. McCallum, K. Nigam – „Employing EM and Pool-Based Active Learning for Text Classification”.