

Dragoș VESPAN

**EXTRAGEREA CUNOȘTINȚELOR DIN DOCUMENTELE ELECTRONICE
PRIN TEXT MINING**

Dragoș VESPAN

**EXTRAGEREA CUNOȘTINȚELOR
DIN DOCUMENTELE ELECTRONICE
PRIN TEXT MINING**



Copyright © 2014, **Editura Pro Universitaria**

Toate drepturile asupra prezentei ediții aparțin

Editurii Pro Universitaria

Nicio parte din acest volum nu poate fi copiată fără acordul scris al

Editurii Pro Universitaria

Descrierea CIP a Bibliotecii Naționale a României

VESPAN, DRAGOȘ MARCEL

**Extragerea cunoștințelor din documentele electronice prin
text mining** / Dragoș Vespan. - București : Pro Universitaria, 2014

Bibliogr.

ISBN 978-606-647-968-4

004

INTRODUCERE

Evoluția internetului ca mijloc de transmitere a informațiilor a dus atât la creșterea volumului resurselor de cunoștințe disponibile on-line cât și la diversificarea formelor și formatelor de stocare și transmitere a acestora: text, date, video, audio. Deși restricțiile hardware în ceea ce privește spațiul de stocare și viteza de transmitere a datelor nu mai reprezintă o problemă, textul rămâne cea mai eficientă formă de prezentare a cunoștințelor pe internet în comparație cu diversele formate audio, video sau multimedia.

Webul a fost proiectat ca un spațiu informațional cu scopul de a depăși barierele mijloacelor clasice de comunicație și de a permite mașinilor să ajute utilizatorii să comunice unii cu alții. Cel mai important obstacol în realizarea acestui deziderat este reprezentat de faptul că informațiile și cunoștințele existente pe internet sunt destinate exclusiv consumului uman. Acest obstacol poate fi depășit prin utilizarea unor metode și tehnici de reprezentare a cunoștințelor existente în documentele text astfel încât acestea să poată fi automat achiziționate și procesate de către mașini.

Capitolul „Text mining – caracteristici și domenii de aplicabilitate” prezintă evoluția internetului ca mijloc de transmitere a informațiilor, evidențiind eterogenitatea și diversitatea documentelor text existente pe web. Este prezentat stadiul cunoașterii în text mining, punându-se accentul pe caracterul multidisciplinar, abordările existente și domeniile de aplicabilitate ale acestui domeniu și sunt descrise mai multe tipuri de arhitecturi funcționale ale sistemelor de text mining.

Capitolul „Reprezentarea documentelor” analizează documentul ca element de bază al achiziției cunoștințelor prin text mining. Sunt prezentate mai multe tehnici de reprezentare descriptivă a documentelor, cu accent pe utilizarea tabelor ca modalitate de organizare a informațiilor. Se definesc caracteristicile documentelor și se analizează tehnici de reprezentare a documentelor în spațiul vectorial.

Capitolul „Clasificarea automată a documentelor” abordează problematica achiziției cunoștințelor prin utilizarea algoritmilor de clasificare. Sunt identificați și analizați principalii algoritmi utilizați în clasificarea textelor. Se definesc

măsurile de evaluare a eficienței clasificatorilor de texte și sunt comparate performanțele algoritmilor de clasificare.

Capitolul „Webul Semantic – Caracteristici și limbaje” abordează problematica reprezentării cunoștințelor în contextul Webului Semantic. Sunt comparate modelele de distribuire a informațiilor pe internet și se utilizează cadrul de descriere a resurselor pentru a reprezenta datele preluate din documentele web ale ASE. Sunt descrise modele de procese pentru servicii oferite în cadrul ASE și prezentate modalitățile de dezvoltare a ontologiilor prin text mining.

În capitolul „Reprezentarea cunoștințelor din documentelor web în sistemul OntoDev”, sunt prezentate contribuțiile personale ale doctorandului. Este elaborat un algoritm de reprezentare a documentelor web cu păstrarea informațiilor legate de structura tabelară a acestora. Se analizează modul în care conceptele existente în cadrul ontologiilor sunt reprezentate în resursele web. Sunt prezentate tehnici de analiză a logurilor unei organizații prin care se identifică modele de comportament ale utilizatorilor de internet.

1. TEXT MINING – CARACTERISTICI ȘI DOMENII DE APLICABILITATE

Evoluția internetului și creșterea numărului de utilizatori a impulsionat foarte mult dezvoltarea motoarelor de căutare și identificarea unor soluții cât mai eficiente de regăsire a informațiilor pe internet. La începutul evoluției motoarelor de căutare, cea mai potrivită soluție pentru analizarea paginilor web și clasificarea lor în categorii, era utilizarea experților umani. Creșterea exponențială a numărului de site-uri web, de la 10.000 la începutul anului 1995, la peste 650.000 la sfârșitul anului 1996 (Gray, 1996) ajungându-se la 1 trilion de pagini indexate de Google la 25.08.2008 (Google, 2008), concomitent cu creșterea numărului de utilizatori de la 16 milioane în 2005 la 147 milioane în 1998 și aproape 1,5 miliarde în iunie 2008 (Stats, 2008), a făcut ca această soluție să devină complet ineficientă. Lansarea motorului de căutare Google, care automatiza procesul de căutare, a fost o adevărată revoluție, dovadă fiind faptul că astăzi peste 50% din căutările pe internet se fac prin utilizarea acestuia, în august 2007 fiind efectuate 31 de miliarde de căutări prin motorul Google, dintr-un total de 61 miliarde căutări la nivel mondial (Comscore, 2008).

Creșterea rapidă și continuă a numărului de site-uri dar și a numărului de utilizatori de internet va impune dezvoltarea motoarelor de căutare dincolo de analiza linkurilor și a rankurilor de pagină, spre înțelegerea semanticii paginilor analizate. Algoritmii de text mining sunt esențiali pentru dezvoltarea motoarelor de căutare în această direcție, pornind de la regăsirea unor pagini care să corespundă unui șir de cuvinte cheie, până la regăsirea unor cunoștințe care să corespundă nevoilor de informare ale utilizatorilor. O direcție de dezvoltare importantă a motoarelor de căutare este aceea de a oferi suport și pentru alte limbi decât cele de circulație internațională. Google depune eforturi intense în acest sens, oferind utilizatorilor atât posibilitatea restricționării căutărilor la o singură limbă (din peste 40 disponibile) dar și traducerea automată a anumitor pagini în limba selecționată, inclusiv română.

Motoarele curente de căutare pot analiza doar o mică parte a webului reprezentată de partea „vizibilă” a acestuia. O cantitate uriașă de valoroase

informații științifice sau de altă natură se regăsește în partea „invizibilă” a webului, reprezentată de paginile care nu au putut fi indexate de motoarele de căutare. Această parte a webului cuprinde documentele al căror format nu este recunoscut și nu poate fi procesat de către motoarele de căutare, dar și baze de date existente pe internet care nu pot fi indexate de acestea din motive de securitate sau de reprezentare a datelor. Mărimea webului „invizibil” nu poate fi precis estimată: Marcus Zillman estima la sfârșitul anului 2007 mărimea acestuia ca fiind de aproximativ 900 miliarde pagini, în comparație cu 20 de miliarde de pagini care erau indexate de motoarele de căutare la momentul respectiv (Zillman, 2007). În acest context, dezvoltatorii motoarelor de căutare trebuie să găsească soluții și metode adecvate pentru străpungerea barierelor ce separă webul „invizibil” de cel „vizibil”.

Pe 28 iulie 2008 a fost lansat un nou motor de căutare denumit „Cuil” care promite o abordare diferită a modului de realizarea a căutărilor pe internet. Dacă motoarele de căutare clasice se bazează pe tehnici de analiză a legăturilor și pe statistici de trafic, Cuil se orientează spre analiza contextului fiecărei pagini și a conceptelor existente în fiecare interogare, după care organizează căutările similare în grupuri, sortându-le pe categorii. Cuil prezintă utilizatorilor rezultatele pe baza conținutului documentelor și nu a popularității acestora.

80% din informațiile companiilor sunt stocate în documente text (Tan, Text Mining: The state of art and the challenges, 1999). Aceste informații pot fi structurate, reprezentate prin tabele, atribute, numere generate de anumite aplicații, sau nestructurate, așa cum se regăsesc în e-mailuri, rapoarte sau diverse documente. Informațiile structurate au un caracter repetitiv fiind, în esență, numerice și, deci, mai ușor de identificat și prelucrat. Ele sunt rezultatul unor tranzacții cu caracter repetitiv (completarea unui ciclu de producție, vânzarea unui produs, încasarea unui ordin de plată) care generează date structurate despre activitățile la care se referă. Singurele elemente care diferă de la o tranzacție la alta sunt seturile de valori luate de tipurile de date.

În contrast, informațiile nestructurate nu au format, structură sau repetabilitate strict definite fiind, în general, scrise în limbaj natural liber. Scrierea unui e-mail nu este condiționată din punctul de vedere al formatului sau conținutului. Un e-mail poate fi scris de oricine, în orice format, utilizând orice limbă și poate conține orice subiect. La fel, există și alte tipuri de documente text care conțin foarte multe informații și cunoștințe dar nu sunt condiționate de nici un fel de formatare: articole, cărți, pagini web, rapoarte.

1.1. CARACTERISTICILE DOCUMENTELOR UNEI ORGANIZAȚII

Într-o organizație complexă cum ar fi, de exemplu, o universitate, se regăsește o mare varietate de documente nestructurate ale căror caracteristici sunt neuniforme: cursuri, e-mailuri, contracte, rapoarte financiare.

Cursurile sunt documente de lungime variabilă ce pot conține terminologie specifică domeniului în care se încadrează. Sunt scrise într-o singură limbă și sunt identificate prin titularul de curs, specializarea și anul la care acestea sunt predate.

E-mailurile sunt relativ scurte în comparație cu alte documente și pot conține atât informații personale cât și informații legate de mediul de afaceri. De obicei sunt scrise într-o singură limbă, dar pot conține și texte scrise în diferite limbi. Sunt identificate prin adresa de e-mail a destinatarului, adresa de e-mail a expeditorului, data și ora trimiterii.

Contractele abundă în terminologie juridică și pot varia foarte mult în mărime. Sunt scrise într-un stil formal, putând fi traduse în mai multe limbi și se identifică prin titulari de contract, obiectul contractului, data semnării, perioadă de valabilitate și instituția care îl legiferează.

Rapoartele financiare conțin informații legate de mediul de afaceri care includ datele financiare ale organizației pe o anumită perioadă de timp, date referitoare la manageri și acționari precum și la relaționările organizației cu alte organizații. Aceste rapoarte pot include și grafice și indicatori financiari dar și previziuni referitoare la evoluția financiară a organizației. Sunt identificate prin organizația la care fac referire, perioada pentru care au fost făcute, organizația care le-a realizat și data la care au fost realizate.

Acestea reprezintă doar câteva tipuri de documente text nestructurate existente într-o organizație educațională, fiecare având propriile caracteristici.

Inmon și Nesavich au propus o serie de caracteristici care evidențiază lipsa de uniformitate a documentelor text nestructurate din interiorul unei organizații (Inmon & Nesavich, 2007): corelarea cu profilul de activitate al organizației, modalitatea de redactare, terminologia utilizată, gradul de repetitivitate și volumul datelor conținute, numărul documentelor, perioada de păstrare, mediul de stocare, posibilitatea de actualizare.

Documentele pot fi direct sau indirect corelate cu profilul de activitate al organizației. În cazul unei organizații educaționale, cursurile, planurile de învățământ, contractele și rapoartele financiare pot fi considerate documente direct corelate cu profilul de activitate al acesteia. Rapoartele de evaluare a

resurselor umane, fișele angajaților sau e-mailurile sunt documente indirect corelate cu profilul de activitate al organizației.

Modalitatea de redactare a documentelor poate fi formală sau informală. Exemple de documente informale sunt scrisorile, e-mailurile sau paginile web. Exemple de documente formale sunt contractele sau rapoartele financiare.

Terminologia utilizată în cadrul unui document se referă la tipul de limbaj utilizat. Unele documente tind să abunde în terminologie specifică unui anumit domeniu putând fi corect înțelese și interpretate doar de experți care posedă cunoștințe anterior dobândite în domeniul respectiv. Rapoartele de cercetare tind să conțină foarte mulți termeni științifici iar contractele conțin, de regulă, foarte mulți termeni juridici. E-mailurile sau paginile web sunt, însă, mult mai generale.

Gradul de repetitivitate a datelor se referă la măsura în care se repetă atributele care reprezintă datele din cadrul documentelor. Paginile web ale site-urilor de comerț electronic au, de regulă, un grad de repetitivitate ridicat, în timp ce e-mailurile au un grad scăzut de repetitivitate.

Mediul de stocare reprezintă suportul de prezentare a documentelor. Paginile web și e-mailurile sunt, de regulă, stocate în format electronic în timp ce contractele sunt tipărite pe hârtie. În anumite circumstanțe, documentele pot fi tipărite pe hârtie din formatul electronic sau pot fi transpuse de pe hârtie în format electronic.

Volumul de date asociat documentelor poate varia de la câțiva kilobaiți la mărimi de ordinul zecilor de megabaiți în cazul e-mailurilor. În general, documentele tipărite pe hârtie au asociat un volum mic de date, în timp ce documentele în format electronic pot avea asociate volume mari de date.

Numărul de e-mailuri este, în general, mult mai mare decât numărul de contracte.

Perioada de păstrare se referă la cât timp vor fi păstrate documentele. Paginile web nu sunt păstrate, de obicei, pentru o perioadă foarte mare, în timp ce rapoartele anuale sunt păstrate pentru o perioadă foarte lungă de timp.

Posibilitatea de actualizare se referă la posibilitatea de schimbare a conținutului unui document original, odată ce acesta a fost creat. E-mailurile nu sunt niciodată modificate, în unele țări acest lucru fiind ilegal. Paginile web însă, sunt actualizate sau modificate cu regularitate.

Tipul de document	Corelare cu profilul de activitate	Modalitate de redactare	Terminologie	Grad de repetitivitate a datelor	Mediul de stocare	Volum de date asociat	Număr documente	Perioada de păstrare	Possibilitatea de actualizare
Pagini web	Direcți	Informală	Variază	Variază	Electronic	Mediu	Mare	Variază	Da, Frecventă
E-mail	Indirecți	Informală	Generală	Mic	Electronic	Poate fi foarte mare	Foarte mare	Variază	Nu
Rapoarte financiare	Direcți	Formală	Specifică	Ridicat	Hârtie, Electronic	Mic	Mic	Mare	Da, Uneori
Cursuri	Direcți	Informală	Variază	Mic	Hârtie, Electronic	Mic	Mare	Mare	Da
Rapoarte de cercetare	Direcți	Formală	Specifică	Mic	Hârtie, Electronic	Mic	Mare	Mare	Da
Contracte	Direcți	Formală	Specifică	Foarte mic	Hârtie	Mic	Moderat	Foarte mare	Da, Ocazional
Fișe angajați	Indirecți	Formală	Generală	Ridicat	Hârtie, Electronic	Mic	Mare	Mare	Da

Tabelul 1. Caracteristicile documentelor text existente într-o organizație educațională