

Rolul bibliotecilor in organizarea sistemului informational al universitatii

Rodica Volovici

Universitatea "Lucian Blaga" din Sibiu

Bd. Victoriei Nr. 10

e-mail: rodica.volovici@ulbsibiu.ro

Abstract: Information technology and communication technology have fundamentally changed our life. They have not missed libraries. Since the discovery of printing it gathers and contains written and printed paper and makes it accessible under the form of books.

În cadrul relatiei biblioteca-cercetare, legaturile simple pâna nu demult, capata în noua societate informationala o complexitate si o dinamica neasteptate, care le modifica termenii si desfasurarile. Se stie ca bibliotecile au fost si sunt principala baza de referinta a tuturor dezbaterilor, conjuncturilor si risipelor gândului. Aceasta motiveaza cu o temeinicie greu de contestat validitatea relatiei biblioteca-cercetare. Si totusi, pentru ca biblioteca sa devina cu adevarat o componenta sine qua non a cercetarii de calitate în conditiile exploziei de cunoastere determinata strâns de dezvoltarea tehnicii informatiei, se impune o reorganizare a lucrului cu documentele si a gestiunii bibliotecilor.

Cu ajutorul noilor tehnici ale prelucrării informationale, activitatea de informare-documentare poate fi adusa mai aproape si mai repede catre beneficiarii bibliotecilor, apropiind fazele finale ale muncii creatoare. Calitatea informării documentare realizata de o biblioteca poate sa creasca substantial prin îmbunatarea unor indicatori de referinta între care figureaza cu necesitate:

- numarul de directii si teme prioritare acoperite cu informatii,
- cantitatea de informatie relevanta oferita beneficiarilor,
- distanta în timp între aparitia documentelor si furnizarea informatiei.

Abordarea problemelor referitoare la eficientizarea activitatii bibliotecilor impune oricarui analist o relatie strânsa cu teoria biblioteconomica si cu stiinta informării. Biblioteca poate fi privita ca un ansamblu cu functii si activitati dependente de structuri, relatii, dotari, metode, de organizarea si calitatea acestor componente si, nu în ultimul rând, de competenta. Gasim aici o colectie de entitati care sunt în relatie unele cu altele si cu mediul, astfel încât ele formeaza un întreg, un sistem complex.

Un concept cheie care rezulta din definitia sistemului este ca elementele componente ale acestuia luate fiecare separat, nu prezinta si nu reprezinta esenta sistemului. Este necesar sa consideram sistemul ca întreg pentru a-l înțelege. Deci întregul este mai mare decât suma partilor sale; aceasta proprietate este cunoscuta sub numele de sinergie. Aceasta este o proprietate evidenta a sistemelor organizationale în general, si în particular a sistemelor informationale.

O prima schita a sistemului informational universitar, nu trebuie sa prezinte toate conexiunile în complexitatea lor, ci numai elementele constitutive cele mai importante si cum interactioneaza aceste elemente. De asemenea, trebuie avut în vedere si cum interactioneaza sistemul cu alte sisteme ("furnizorii" si "clientii").

Optiunea pentru modernizare, pentru introducerea celor mai noi tehnici de achizitie, prelucrare, stocare si dimensionare a documentelor si informatiilor este si aici, ca în orice domeniu de activitate, o conditie hotarâtoare a cresterii eficientei. Transferul de informatie trebuie sa se realizeze în conditii si cu mijloace competitive cu cele ale structurilor care beneficiaza de acest proces, daca nu chiar în situatii privilegiate, dat fiind rolul propulsiv al informatiei stiintifice în toate sferile de cunoastere si actiune umana. Daca unul din factorii de baza ai eficientei este viteza de regasire a informatiilor, nici un sistem traditional nu poate oferi adrese si combinatii într-un timp mai scurt. Dar impactul informatizarii bibliotecilor asupra vietii economice si sociale este mult mai puternic si mai adânc.

Oricare ar fi determinarile materiale ale eficientei în activitatea bibliotecilor, principala conditie a unor rezultate remarcabile o constituie calitatea specialistilor antrenati. Însusirea obiectivelor economico-sociale, prin traducerea cerintelor în limbajul si tehnicile informarii documentare, organizarea achizitiei, prelucrarii, stocarii si diseminarii informatiei si documentelor, sustinerea unui dialog fertil cu beneficiarii, nu pot fi decât rodul unei participari profesionale, politice si cetatenesti marcata de competenta si clarviziune.

A vedea este aproape similar cu a citi. Este atunci firesc primul pas facut de programatori în directia aceasta, si anume, sa învete calculatorul sa recunoasca literele. Tentativa este cunoscuta sub numele de OCR, acronim pentru Optical Character Recognition, adica recunoasterea optica a caracterelor. Avantajul unui astfel de soft este evident atunci când se doreste stocarea documentelor într-o forma electronica (text ASCII). Marea majoritate a sistemelor OCR încep cu o imagine bit-map culeasa via scanner sau fax-modem si primul pas consta din identificarea pe pagina a blocurilor de text, pe baza unor caracteristici tipografice, cum ar fi alinierea la stânga sau la dreapta. Apoi aceste blocuri sunt descompuse în semne individuale, care deseori corespund literelor.



Algoritmul de recunoastere încearcă în următoarea etapă să "ghicească" cât mai bine semnificația fiecărui semn, după care este construită întreaga pagină cu formatul corespunzător. Cel mai bun sistem OCR poate realiza achiziții cu o acuratețe de peste 99% pentru paginile tiparite și scrise cu fonturi obișnuite. Cu toate acestea acest procent pare aproape perfect, nivelul de eroare este destul de ridicat dacă ne gândim că pe o pagină standard sunt aproape 1500 de caractere. În acest caz, chiar și un procent de corectitudine de 99,9% generează una sau două erori pe pagină necesitând intervenția omului pentru corectare și asigurarea fidelității la citire.

În viața de zi cu zi ne întâlnim rareori cu situații ideale: text "curat", font standard, litere tiparite complet, astfel ca un "o" să nu semene cu un "c". Imaginile murdare sunt o problemă deoarece chiar și o pată mică poate acoperi o parte importantă din litera făcând un "m" să semene cu un "n". Dacă documentul este neclar, dacă contururile literelor sunt marite sau lipite unele de altele, pot să apară erori de identificare a caracterelor, deoarece sistemele OCR consideră orice semn continuu ca un singur caracter. Probleme pot să apară chiar și dacă paginile sunt "curate". Dezvoltarea algoritmilor care pot recunoaște caracterele, în ciuda acestor probleme, trebuie să asigure un echilibru între flexibilitate și acuratețe. Dacă soft-ul nu este destul de flexibil va da eroare dacă întâlnește cea mai mică diferență între corpurile de litere și, pe de altă parte, o mare flexibilitate poate conduce la confundarea a două caractere asemănătoare cum ar fi "b" și "h".

Rezolvările tehnice găsite în ultima vreme sunt foarte bune dacă textul este "curat" și forma literelor identică în întregul text.

Softurile OCR lucrează, de regulă, pornind de la imaginea unei pagini "furnizate" de un scanner sau care este preluată de la un fax-modem, deci ele pot conversa direct cu scannerele, iar unele se pot atașa aplicațiilor astfel încât se poate translata o pagină de text direct în editorul de texte.

Primul pas al procesului OCR este "spargerea" imaginii paginii în blocuri de text, grafice, imagini, tabele. Există algoritmi specifici pentru identificarea naturii fiecărui bloc. În general rândurile sunt grupate în paragrafe utilizând criterii cum ar fi: alinierea, lungimea, corpul de literă, spațierea. Pentru a se face distincția între grafică și text, de regulă, se utilizează algoritmi bazati pe analiza diferențelor statistice între text și grafică, precum: ritm regulat la text față de lipsa unei periodicități la grafică.

După reconstruirea paginii, semnele individuale sunt translatare în text ASCII. Există mai multe modalități de a face acest lucru, dar în mare există două categorii: metode bazate pe șabloane (template-based) și metode bazate pe trăsături caracteristice (feature-based). Dintre sistemele OCR bazate pe șabloane existente acum pe piață amintim sistemul Expert Vision, care recunoaște pată de cerneală comparând-o cu șablonul cu care se aseamănă cel mai mult. Un exemplu elocvent de soft care utilizează cealaltă metodă de recunoaștere, cea bazată pe trăsături specifice, este OmniPageProfessional (Caere) care prezintă avantajul că nu trebuie "reglat" pentru anumite corpuri de literă. Softul conține aproximativ 100 de "sisteme expert", în esență 100 de algoritmi, pentru identificarea unor caractere de la "A" la "Z" și apoi de la "a" la "z", la care se adaugă numerele și semnele de punctuație. Acest tip de soft încearcă să surprindă esența fiecărei litere într-un algoritm.

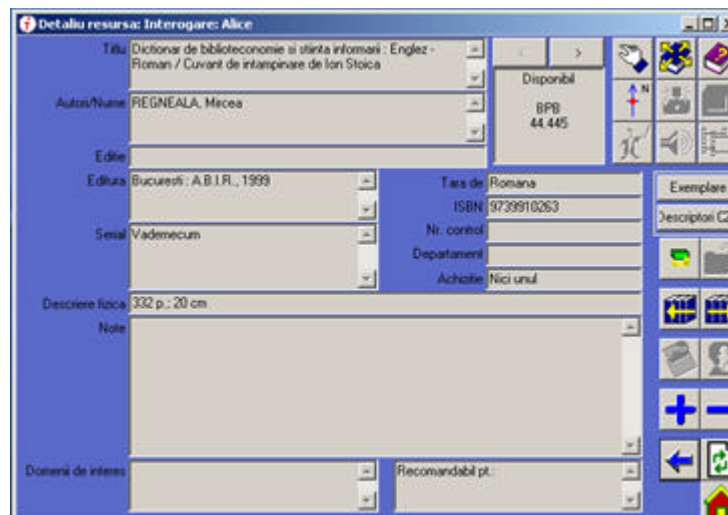
Unele dintre produsele software specializate în recunoașterea caracterelor sunt orientate pe context, ele conțin o serie de informații cu privire la regulile de scriere existente și la formulările cele mai uzuale dintr-o limbă sau un grup de limbi. Și asta, deoarece majoritatea softurilor OCR sunt dedicate lingvistic. Indiferent de tehnologia de recunoaștere utilizată, semnele diacritice, sau în alte cazuri, chiar forma diferită a caracterelor pot fi motive de reducere drastică a procentului de recunoaștere. Este un lucru obișnuit ca același soft OCR să fie utilizat pentru o singură limbă sau pentru un grup de limbi. De pildă Recognita

(Recognita Corporation) este implementat în mai multe versiuni: numai pentru engleza, pentru limbi vorbite pe continentele Americane (engleza, franceza, spaniola, portugheza) si o versiune internationala care ofera posibilitatea de a alege din 22 de limbi, una dintre ele fiind si româna.

Tehnologiile de ultima ora, cum sunt de pilda cele bazate pe modelul retelelor neuronale, ofera si posibilitatea autoinstruirii. Un exemplu este produsul Five-Star-Read (Epson). Softul posedea functii speciale pentru a învăta cum anume trebuie interpretat un anumit semn, dezvoltate din cercetarile privind recunoasterea scrisului de mâna. Daca textul care trebuie recunoscut este degradat si algoritmul de recunoastere da în prima instanta gres, majoritatea softurilor încearca un proces de îmbunatatire a imaginii caracterelor, de pilda adaugând un pixel lipsa când este posibila refacerea unei linii mai lungi. Dar, foarte frecvent si aceste procedee dau gres. Si atunci, cel mai bun OCR ramâne în continuare omul. Verificarea corectitudinii se dovedeste a fi cel mai mare consumator de timp în întregul proces OCR.

Echipea noastra si-a propus sa studieze acuratetea cu care sistemele OCR existente acum pe piata sunt capabile sa recunoasca caracterele de pe vechile fise de carte. Se vor analiza comparativ performantele obtinute în recunoasterea fiecarui caracter si se vor identifica situatiile în care sistemele existente sunt suficient de performante pentru a putea fi folosite pentru ducerea la bun sfârșit a obiectivului propus, si, se vor identifica si situatiile în care performantele sunt sub un nivel minim acceptabil. Pentru aceste ultime situatii, se vor dezvolta noi sisteme OCR implementate fie prin metode clasice fie folosind retele neuronale, astfel încât în final sa se obtina un sistem OCR deosebit de performant.

În concluzie, se va implementa un sistem OCR dedicat, deosebit de performant care sa poata fi folosit pentru conversia fiselor vechi de carte, existente în format dactilografiat, în fise electronice de carte care sa poata fi încorporate în baza de date a Bibliotecii Universitare Sibiu.



Bibliografie:

1. Thompson, Ronald; Cats-Baril, William: Information Technology and Management, 2nd edition, Boston: McGraw-Hill, 2003
2. Volovici, R. M.: Cercetari asupra cresterii fiabilitatii sistemului informatic si al informatiei, referat de doctorat, ULB Sibiu, Sibiu, 2003