

# Portal web de știri autonom bazat pe prelucrarea limbajului natural

**Traian Rebedea**

Universitatea „Politehnica”  
Bucuresti, Facultatea de  
Automatică și Calculatoare  
Splaiul Independenței, Nr.  
313, București, România  
trebedea@gmail.com

**Costin-Gabriel Chiru**

Universitatea „Politehnica”  
Bucuresti, Facultatea de  
Automatică și Calculatoare  
Splaiul Independenței, Nr.  
313, București, România  
chirucos@gmail.com

**Ștefan Trăușan-Matu**

Universitatea „Politehnica”  
Bucuresti, Facultatea de  
Automatică și Calculatoare  
Splaiul Independenței, Nr.  
313, București, România  
Institutul de Cercetare în  
Inteligența Artificială  
al Academiei Române,  
Calea 13 Septembrie, Nr 13,  
București, România  
trausan@cs.pub.ro

## REZUMAT

În această lucrare este prezentat un portal web de știri pentru limba română care încorporează un modul autonom de grupare și de clasificare a știrilor. Pentru a asigura funcționarea complet autonomă a portalului sunt folosite tehnici de prelucrare a limbajului natural. Știrile sunt colectate automat de la un număr însemnat de surse folosind sindicalizarea web. Apoi, sunt utilizate tehnici de învățare automată pentru clasificarea știrilor. În primul rând, acestea sunt grupate folosind un algoritm aglomerativ, iar cele mai numeroase grupuri corespund celor mai importante subiecte ale momentului. În acest fel se asigură și colectarea mai multor informații despre aceste subiecte, întrucât sunt utilizate știri din diverse surse. În al doilea rând, portalul folosește un algoritm de clasificare automată pentru a încadra fiecare grup de știri într-un număr predeterminat de categorii. Atât pentru testarea, cât și pentru evaluarea performanțelor clasificatoarelor au fost folosite peste o mie de știri pre-etichetate. Lucrarea conține și o comparație completă a rezultatelor obținute pentru algoritmii de clasificare.

## Cuvinte cheie

Portal de știri, Agent inteligent, Grupare, Clasificare, Prelucrarea limbajului natural

## Clasificare ACM

H5.2. Information interfaces and presentation (e.g., HCI): Miscellaneous.

## INTRODUCERE

În ultimii 10 ani am asistat la o adevărată explozie a web-ului, atât ca număr de domenii web existente, cât și ca trafic și număr de utilizatori. Avansul tehnologic din această perioadă a facilitat accesul la Internet al unui număr mare de utilizatori, oferindu-le acestora viteze de transfer din ce în ce mai mari. Studiile Internet World Stats [5] relevă faptul că în ultimii 5 ani numărul de utilizatori de Internet la nivel mondial s-a dublat, ajungând la 20% din populația lumii – adică peste 1 miliard de utilizatori. Creșterea numărului de utilizatori (consumatori) a determinat intrarea pe această piață a unui număr din ce în ce mai mare de ofertanți. Astfel, conform datelor oferite de Netcraft [7] – vezi figura 1, tot în ultimii 5 ani, numărul total al siturilor web s-a triplat, ajungându-

se în februarie 2008 la aproximativ 158 milioane de situri web, dintre care aproape jumătate sunt active – evoluția siturilor active fiind similară în această perioadă cu cea a siturilor totale. În aceste condiții, nu este o surpriză că volumul de informație disponibil pe web a crescut exponențial în această perioadă, ajungând la cifre impresionante. Yahoo a făcut public în august 2005 [14] că indexează peste 20 de miliarde de articole, dintre care 19,2 miliarde de documente web, 1,6 miliarde de imagini și peste 50 de milioane de fișiere audio și video. Numărul mare de documente aflate pe web, precum și starea de continuă schimbare a acestuia pot să convertească web-ul într-un sistem haotic. Astfel, ca o consecință directă a popularității de care se bucură ca sistem global de informare, web-ul este bombardat în mod constant de cantități mari de date și informație, care fac din găsirea informației utile pe web o experiență dificilă și frustrantă. Soluția pentru găsirea informației pe web sunt motoarele de căutare. Odată cu redefinirea importanței unei pagini web de către Brin și Page [8], motoarele de căutare au intrat într-o nouă etapă, oferind rezultate satisfăcătoare pentru majoritatea căutărilor. Totuși, odată cu dezvoltarea web-ului, rezultatele unei căutări, chiar dacă este bine formulată, oferă câteva sute de rezultate utilizatorului.

Volumul mare de informații nu este singura problemă indusă de dezvoltarea web-ului. O altă problemă importantă este redundanța informațională întrucât multe dintre aceste informații se repetă în surse diferite. Pentru a oferi rezultate mai bune utilizatorului se pot crea programe care să exploateze exact aceste caracteristici – prin procesarea cantității mari de informație existentă, precum și prin determinarea și prelucrarea datelor redundante se pot obține rezultate folositoare.

Ideea principală a acestei lucrări este de a folosi redundanța informațională din cadrul fluxurilor de știri preluate sub formă de sindicalizare web. Pentru aceasta se folosesc metode avansate de procesare a textului, incluzând gruparea și clasificarea, cu scopul de a obține un portal de știri autonom. Lucrarea continuă cu o secțiune în care este introdusă ideea de prelucrare inteligentă a știrilor și a conceptelor sale fundamentale. A treia secțiune descrie cele mai importante tehnici folosite pentru gruparea și clasificarea textelor. Secțiunea

următoare conține descrierea portalului web de știri pentru limba română ce folosește prelucrări inteligente, precum și rezultatele obținute. Lucrarea se încheie cu concluzii și referințe.

### PROCESAREA INTELIGENTĂ A ȘTIRILOR

Cele mai importante surse de știri din lume au introdus sindicalizarea web ca o modalitate nouă de furnizare a știrilor către cititori sau către alte situri web sau aplicații care doresc să le folosească. Formatele bazate pe XML folosite pentru sindicalizarea web au fost proiectate astfel încât conținutul lor să fie ușor de folosit de către programe numite cititoare de flux de știri sau agregatoare. Acestea automatizează procesul de colectare periodică a fluxurilor de știri și de prezentare a acestora utilizatorului într-o manieră ușor de urmărit. Însă nu rezolvă cele două probleme esențiale: cantitatea mare de informație și redundanța știrilor. Într-o fază inițială a prelucrării fluxurilor de știri, s-au dezvoltat portaluri web de știri care au folosit numărul de cititori al unei știri drept criteriu de calcul al importanței știrii. Deși prezintă unele avantaje clare – printre care determinarea importanței în funcție de preferințele utilizatorilor, această metodă are și multe dezavantaje: nu rezolvă nici problema cantității mari de informație, nici pe cea a redundanței știrilor, și nu oferă un criteriu automatizat pentru calculul importanței a priori a unei știri – înainte de a fi citită de cineva. În plus, fiecare știre este legată de sursa din care a fost preluată, fără a se oferi surse alternative în legătură cu acel subiect. Alternativa propusă în această lucrare dorește rezolvarea tuturor acestor probleme.

În ultima perioadă, mai multe echipe de cercetare, precum și companii de renume și-au axat eforturile spre prelucrarea inteligentă a fluxurilor de știri. De ce numesc aceste prelucrări inteligente? În primul rând, rezultatul acestor prelucrări determină automat subiectele zilei și oferă o clasificare a acestor subiecte în funcție de numărul de știri ce le alcătuiesc, prezentând în acest fel utilizatorului știrile cu adevărat importante. În acest caz, determinarea importanței unui subiect are la bază nu opțiunea subiectivă a cititorilor, ci pe cea “obiectivă” a specialiștilor, a agențiilor de știri și a ziarelor ce relatează diversele subiecte. În al doilea rând, pentru realizarea acestor prelucrări se utilizează o serie de tehnici de învățare automată.

Știrile preluate din cadrul diverselor fluxuri pot fi folosite pentru două prelucrări esențiale: prima este *gruparea* („*clustering*”) știrilor și aflarea subiectelor celor mai importante, iar a doua o constituie clasificarea automată a știrilor în diverse categorii de știri predefinite. În plus, unele lucrări din domeniul extragerii de informații din știri [11, 12] propun o prelucrare mai amănunțită a acestora, cu accentul pe extragerea de cunoștințe specifice – precum nume de persoane, de țări, de organizații și de companii – și folosirea acestor cunoștințe pentru gruparea și clasificarea știrilor. Desigur că există o multitudine de abordări în domeniul prelucrării inteligente a știrilor, majoritatea aflându-se încă în faza de cercetare. O idee interesantă o reprezintă asignarea de etichete (cuvânt sau grup de cuvinte cheie, foarte semnificative) fiecărui grup de știri rezultat din grupare. În acest fel, utilizatorul poate urmări mult mai ușor lista de subiecte relativ la o anumită

etichetă. Problema esențială a acestei abordări o reprezintă calitatea etichetării, întrucât nu s-au găsit algoritmi de etichetare care să ofere rezultate satisfăcătoare [11]. O metodă nouă este oferită de sistemul *NewsJunkie* [4] ce propune determinarea noutății informației unei știri – astfel având un flux de știri despre un subiect comun relatat de-a lungul unei perioade mai lungi de timp, sistemul își propune găsirea informațiilor cu adevărat noi și filtrarea articolelor ce relatează evenimente deja prezentate. O idee similară este propusă și în *Ranking a Stream of News* [2].

Gruparea știrilor exploatează redundanța informației obținute prin colectarea unui număr mare de știri din diverse surse, obținându-se grupuri de știri cu subiect comun. Redundanța informațiilor dintr-o colecție de știri poate fi utilizată spre atingerea a două scopuri. În primul rând, poate fi folosită drept un criteriu de calcul al importanței unei știri/subiect pornind de la ideea că numărul de apariții ale unui subiect în surse diferite este proporțional cu importanța acestuia. În al doilea rând, știrile pot fi grupate în jurul acestor subiecte, reducând cantitatea de informație prezentată utilizatorului. În plus, această modalitate oferă avantajul de a prezenta fiecare subiect relatat din perspectiva fiecărei surse, permițându-i cititorului urmărirea subiectului din sursa dorită sau compararea relatărilor din mai multe surse. Importanța unei știri este determinată de către comunitatea media a agențiilor de presă – această măsură putând fi considerată profesionistă și obiectivă.

Clasificarea automată a știrilor poate fi folosită cu succes împreună cu gruparea știrilor, descrisă anterior, pentru crearea unui portal de știri care să funcționeze complet automatizat și să ofere utilizatorului toate facilitățile unui portal de știri obișnuit. Toate sursele de știri oferă o structurare a știrilor în categorii pentru a furniza utilizatorilor informațiile necesare în funcție de interes. Chiar dacă standardele folosite pentru sindicalizarea web oferă posibilitatea definirii categoriei unui articol, faptul că nu există o taxonomie de știri valabilă universal, face ca această facilitate să fie practic inutilă în cadrul unui portal care culege știri de la un număr mare de surse. Singura opțiune este de a se folosi algoritmi de clasificare pentru asignarea fiecărei știri la o anumită categorie. Astfel, într-o primă etapă, se vor aduna știri care vor fi asignate manual sau pseudo-automat fiecărei categorii definite de portal. Întrucât există sindicări web ce conțin știri dintr-un singur domeniu, procesul de colectare a știrilor necesare pentru antrenare poate fi automatizat. Rezultatele obținute folosind diverși algoritmi de clasificare sunt promițătoare, permițând folosirea acestei tehnici cu rezultate foarte bune [3]. Deoarece știrile sunt grupate în funcție de subiect, algoritmul de clasificare poate fi folosit pentru grupuri și nu pentru fiecare știre individual – fiecare știre din cadrul unui grup aparținând categoriei în care este încadrat grupul respectiv.

### TEHNICI DE GRUPARE ȘI CLASIFICARE A TEXTELOR

Metodele de grupare și clasificare sunt utilizate des în prelucrarea limbajului natural. O știre poate fi reprezentată ca un vector, ale cărui elemente sunt asociate cuvintelor cheie din text. Valorile elementelor pot fi și variabile booleene reprezentând apariția sau lipsa unui cuvânt cheie

din textul asociat știrii, și anume titlul și descrierea preluate din sindicalizarea web. O îmbunătățire a acestei reprezentări este de a folosi frecvența fiecărui termen sau măsura TF-IDF [1, p. 56-57]. Atât algoritmi de grupare, cât și cei de clasificare necesită calcularea similarității între două elemente, în acest caz, între documente text. Spațiul caracteristicilor folosit pentru reprezentarea documentelor are dimensiunea egală cu numărul de termeni distincți din toate documentele, indiferent de metodele prezentate anterior. Transformarea distanțelor metrice în măsuri de similaritate nu oferă rezultate bune pentru texte. De aceea, pentru calculul similarității între documentele text sunt folosite alte măsuri [10], iar una dintre cele mai utilizate este cosinusul unghiului format între vectorii ce corespund documentelor care sunt comparate.

Gruparea datelor [6, p. 495-528] este o tehnică de analiză statistică a datelor, folosită cu succes în diverse subdomenii ale inteligenței artificiale, precum învățarea automată, mineritul datelor, recunoașterea formelor și analiza imaginilor. Gruparea este o modalitate de partiționare a datelor în subseturi – grupuri, *cluster* – astfel încât datele din fiecare grup să aibă o caracteristică comună. De obicei, gruparea datelor se face în funcție de proximitatea elementelor de grupat, pentru aceasta folosindu-se o funcție distanță. Gruparea datelor poate fi privită ca o tehnică de generalizare, întrucât datele de intrare sunt împărțite în clase de echivalență folosind anumite criterii. Toate elementele ce formează o clasă de echivalență moștenesc atributele grupului din care fac parte.

Algoritmi de grupare pot fi împărțiți folosind mai multe criterii [10]. O prima clasificare este în algoritmi divizivi („*top-down*”) și algoritmi aglomerativi („*bottom-up*”). Procesarea top-down pornește cu toate datele într-un singur grup, pe care apoi îl rafinează în subgrupuri. Algoritmi aglomerativi, consideră fiecare element ca un grup separat, iar acestea sunt grupate ulterior. În funcție de rezultatul grupării, avem de a face cu algoritmi ierarhici și algoritmi de tip centroid („*flat*”). Algoritmi ierarhici oferă ca rezultat o structură arborescentă ce poate fi vizualizată ca o dendrogramă. De asemenea, gruparea poate fi tare („*hard*”) sau slabă („*soft*”). Gruparea tare va produce cluster ce nu se întrepătrund, adică fiecare element aparține unui singur grup. Gruparea slabă asignează fiecărui obiect o probabilitate de a fi membrul unui grup.

Algoritmi ierarhici folosesc aproape întotdeauna o grupare tare, numai cei de tip centroid abordând și gruparea slabă. Algoritmi de grupare ierarhici sunt algoritmi greedy, ce folosesc pentru construirea arborelui rezultat ori o abordare top-down, ori una bottom-up. Algoritmi divizivi pornesc cu toate elementele considerate într-un singur cluster și la fiecare pas, se determină grupul cel mai puțin coerent, acesta urmând ales să fie divizat. Grupurile ce conțin elemente mai similare au o coerență mai mare decât grupurile ce conțin elemente mai puțin similare. Algoritmi aglomerativi încep cu fiecare element ca un cluster separat și, la fiecare pas, se alege cele mai similare două cluster spre a fi unite. În cadrul algoritmilor ierarhici aglomerativi există trei metode de determinare a similarității dintre două grupuri –

pentru a afla ce grupuri trebuie unite la fiecare pas. Prima abordare, definește similaritatea între două grupuri ca similaritatea între cei mai similari doi membri, câte unul din fiecare grup. Această variantă se numește legătură simplă („*single link*”) și corespunde unui algoritm de determinare a arborelui minim de acoperire pentru setul de puncte determinate de elementele ce trebuie grupate. A doua abordare, denumită legătură completă („*complete link*”), definește similaritatea a două grupuri ca similaritatea celor mai puțin similare două elemente ce aparțin unuia fiecărui grup. Această abordare produce rezultate mult mai bune decât cea anterioară, întrucât folosește calitatea globală a unui grup și nu pe cea locală – prima abordare are tendința de a produce cluster „alongate”. Dezavantajul metodei legătură completă îl reprezintă complexitatea mai mare decât cea a metodei legătură simplă. Pentru a păstra performanțele metodei legătură completă și complexitatea redusă a metodei legătură simplă, se folosește o abordare ce definește similaritatea a două grupuri ca fiind similaritatea medie între elementele fiecărui grup în parte.

Clasificarea automată este definită ca procesul de asignare a unor obiecte din universul problemei în două sau mai multe categorii sau clase. Procesul de clasificare statistică poate fi împărțit în două etape distincte: etapa de antrenare a clasificatorului și etapa de utilizare efectivă a acestuia. Pentru etapa de antrenare, datele de intrare sunt numărul de categorii distincte și elemente ce aparțin fiecărei categorii (un element poate aparține uneia sau mai multor categorii). Aceste obiecte formează setul de antrenament și poate fi reprezentat ca un set de  $(x, c)$  unde  $x$  este un vector  $m$ -dimensional ce codifică un obiect din setul de antrenament, iar  $c$  este clasa căreia obiectul îi aparține. Pentru fiecare algoritm se definește o procedură de antrenare al cărei rezultat este un model de date (clasificatorul) pe baza căruia se pot clasifica obiectele din universul problemei. Acest model este implicat în cadrul etapei de utilizare efectivă a clasificatorului. Algoritmi de clasificare diferă prin modul cum este definită procedura de antrenare, deci și prin modelul folosit pentru clasificator.

Există două categorii importante în care se împart algoritmi de clasificare. Prima variantă folosește o separare a universului problemei, care este de obicei un spațiu vectorial multidimensional, în regiuni diferite, fiecărei regiuni fiindu-i asignată câte o categorie. Acești algoritmi nu folosesc praguri de încredere sau clase de probabilitate și în această categorie se află algoritmul NN („*Nearest Neighbour*”) și algoritmul de selecție a caracteristicilor. A doua categorie privește problema clasificării ca pe o problemă de estimare a unei funcții de forma  $P(\text{clasa} | \bar{x}) = f(\bar{x}, \bar{\theta})$ , unde  $x$  este vectorul de intrare, iar funcția  $f$  este de obicei parametrizată cu vectorul  $\bar{\theta}$ . În acest caz, algoritmul de clasificare trebuie să determine din familia de funcții parametrizate funcția ce clasifică cel mai bine setul de antrenament – deci avem o problemă de optimizare. Acești algoritmi folosesc un model matematic mai complex, dar oferă rezultate mult mai bune. Printre aceștia se numără rețelele Bayes și mașinile vectorilor de suport („*SVMs – Support Vector Machines*”).

Categorizarea documentelor text [1, p. 124-169] este una dintre aplicațiile algoritmilor de clasificare în domeniul prelucrării limbajului natural, având scopul de a clasifica textele în funcție de subiect. Algoritmii de clasificare sunt folosiți de o perioadă lungă de timp în inteligența automată, recunoașterea formelor și mineritul (bazelor) de date; trebuie avut însă în vedere că în toate aceste domenii datele sunt mult mai „structurate” decât în cazul textului. Una dintre caracteristicile definatorii pentru texte este numărul extrem de mare de termeni, de ordinul miilor sau zecilor de mii, spre deosebire de mineritul datelor unde numărul de coloane poate fi de ordinul zecilor sau sutelor, în cel mai rău caz.

Clasificatoarele Nearest Neighbour (NN) [1, p. 133-138] sunt cei mai folosiți algoritmi ce separă spațiul caracteristicilor în regiuni corespunzătoare fiecărei categorii. Ideea de bază este că documentele similare ar trebui să aparțină aceleiași clase. Faza de antrenare presupune o simplă indexare a datelor de test în funcție de clasa din care fac parte – deci este foarte rapidă. Pentru clasificarea unui document  $d_q$ , se determină  $k$  documente cât mai similare cu  $d_q$ , unde  $k$  e o constantă. Clasa căreia îi aparțin cele mai multe dintre aceste documente reprezintă clasa în care este încadrat  $d_q$ . O variantă mai bună este de a utiliza un scor pentru fiecare categorie astfel: dacă  $d$  face parte din cele  $k$  documente cele mai apropiate de documentul de clasificat  $d_q$  și  $d$  este clasificat în categoria  $c_d$ , atunci scorul acestei categorii acumulează

pentru  $off$  și  $k$ , ce se obțin în urma antrenărilor succesive, clasificatorul poate deveni foarte puternic, oferind „o acuratețe comparabilă cu cea a celor mai bune clasificatoare de texte” [1, p. 135]. Avantajele clasificatoarelor NN sunt faptul că sunt ușor de antrenat și de folosit, nefiind nevoie de construcții matematice suplimentare. Pe lângă aceasta, oferă și rezultate bune, dacă  $k$  și deplasamentele categoriilor sunt bine alese. Dezavantajele lor țin de timpii mari de clasificare, precum și de cantitatea mare de memorie consumată în cadrul procesului de clasificare. Pentru a îmbunătăți performanțele clasificatoarelor NN se poate folosi un algoritm greedy de selecție a caracteristicilor esențiale pentru fiecare clasă [1, p. 136-143].

Algoritmii ce folosesc estimarea unei distribuții de probabilitate pentru fiecare clasă sunt puternici în clasificarea textelor, producând rezultate foarte bune. Dar pentru a obține aceste rezultate este necesar un număr mare de date în setul de antrenare, pentru a obține un estimator cât mai bun al fiecărei categorii. În condițiile unui set de antrenament suficient de mare și a unor categorii bine definite, categorizoarele de texte bazate pe rețele Bayes și SVM-uri au avut o acuratețe de peste 80% [3].

#### PORTAL WEB DE ȘTIRI INTELIGENT PENTRU LIMBA ROMÂNĂ

Prin combinarea facilităților oferite de tehnologiile web precum sindicalizarea și a avantajelor furnizate de către

The screenshot shows a news portal interface. At the top, there's a search bar and a navigation menu with categories like Romania, Politica, Economie, International, Sport, Mădăria, Cultura, and High Tech. The main content area displays a list of news articles under the 'International' section. The first article is titled 'Europa îl cere lui George W. Bush eliminarea vizelor' and discusses the summit in Vienna. Other articles include 'Bush, huiduit la Viena', 'Bush s-a grabit să ajungă în Noua Europă', and 'Bush găsește aliați la Viena'. Each article snippet includes a date, time, source, and a coefficient value.

Figura 1 Grup de știri corespunzător unui subiect în cadrul portalului

valoarea similarității dintre  $d$  și  $d_q$ . Clasa ce va avea scorul cel mai mare, după ce sunt considerate toate cele  $k$  documente cele mai apropiate de  $d_q$ , va fi clasa în care acesta este încadrat. Pentru o acuratețe și mai bună, pentru fiecare clasă se poate folosi un deplasament,  $off$ , ce va fi adunat scorului categoriei. Considerând diverse valori

tehnicile de minerit ale textelor este posibilă dezvoltarea unui portal de știri care poate funcționa cu un minim de intervenție umană. Acesta oferă o alternativă viabilă pentru portalurile de știri tradiționale, datorită avantajelor oferite precum autonomia față de un administrator uman,

precum și metodologia folosită pentru a prezenta știrile pornind de la importanța unui subiect.

### Funcționarea și arhitectura portalului

Una dintre trăsăturile fundamentale ale aplicației este lipsa intervenției umane în timpul funcționării – portalul trebuie să funcționeze automat. Aceasta face ca portalul de știri construit să fie considerat un agent autonom.

Funcționarea de bază a portalului web poate fi descrisă de către următoarea secvența de pași:

Inspectarea periodică a siturilor agențiilor de știri și ziarelor ce oferă RSS și colectarea automată a sindicărilor noi;

Introducerea știrilor din cadrul fiecărui RSS nou într-o bază de date, pentru a oferi acces mai ușor la aceste informații;

Prelucrarea informației text a fiecărei știri noi, prin aplicarea diverselor tehnici lingvistice, pentru determinarea vectorului caracteristic – numărul de apariții ale fiecărui cuvânt;

Gruparea știrilor folosind un algoritm de clustering pentru texte, pornind de la reprezentarea știrilor în spațiul m-dimensional al cuvintelor;

Clasificarea fiecărui grup de știri în cadrul unei categorii predefinite, folosind un clasificator care este reantrenat în mod regulat;

Generarea automată a paginilor web corespunzătoare subiectelor celor mai importante dintr-o anumită perioadă, grupate în diverse modalități, inclusiv pe fiecare categorie de știri. Aceste pagini web constituie rezultatul final al funcționării portalului și sunt vizibile de către utilizatori;

Generarea automată a sindicărilor RSS și Atom oferite de portal, în conformitate cu subiectele folosite la pasul 6.

A acțiunile elementare descrise anterior, pot fi rulate atât

precizat, cât și individual, la momente de timp diferite. Modul de funcționare este determinat de cantitatea fluxurilor RSS preluate, de numărul știrilor noi într-o anumită perioadă, cât și de intervalul de timp la care situl este modificat automat.

Din modul de funcționare al portalului descris anterior se poate remarca că acesta este un agent reactiv (răspunde schimbărilor din mediul său), necomunicativ (poate comunica numai într-o mică măsură, prin sindicalizare web) și neadaptiv. Pe lângă acest comportament de agent autonom, portalul trebuie să implementeze și un modul web care să permită căutarea text în cadrul subiectelor, cât și o personalizare a portalului pe bază de utilizator, prin alegerea categoriilor, a ordinii acestora și a numărului de știri din fiecare categorie.

Intervenția unui administrator al portalului este necesară doar într-un singur punct și anume pentru administrarea listei de adrese web a sindicărilor colectate în prima etapă a funcționării automate. Legat de administrarea listei de sindicări folosite, administratorul va decide care dintre fluxurile RSS vor fi folosite pentru etapa de antrenare a clasificatorului de știri.

O etapă importantă în cadrul funcționării aplicației o constituie antrenarea inițială și, eventual, reantrenarea periodică, a clasificatorului. Antrenarea inițială trebuie efectuată în momentul când avem exemple suficiente pentru fiecare categorie considerată. Reantrenarea se poate face periodic, pe măsură ce numărul setului de antrenament va crește – în acest fel acuratețea clasificatorului ar trebui să se îmbunătățească.

După cum se poate deduce și din descrierea funcționării portalului, acesta este alcătuit din două module cu o funcționare relativ independentă: un modul ce conține partea de agent a portalului (partea inteligentă) și un modul web. Ambele module folosesc o bază de date în

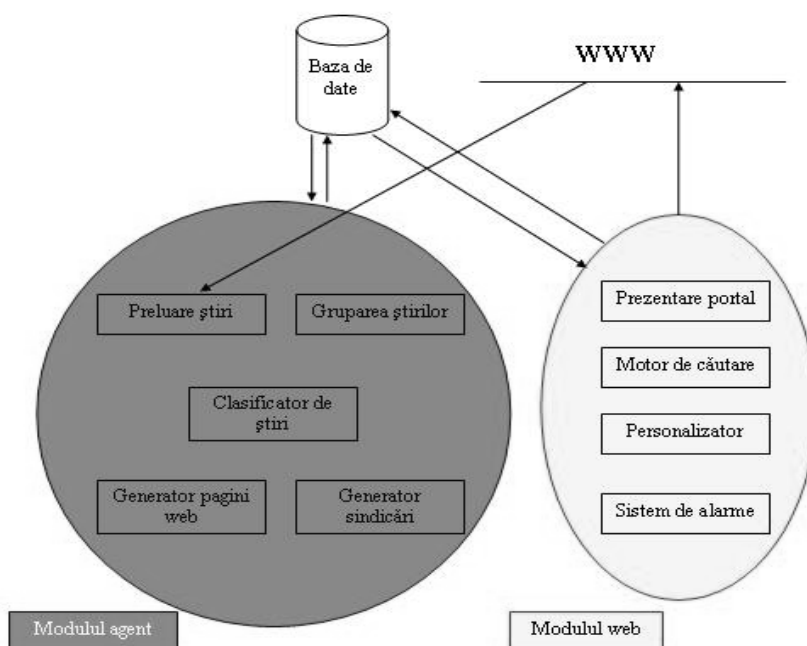


Figura 2 Arhitectura portalului web de știri – modulele agent și web

într-o singură etapă, în mod secvențial, urmând ordinea

care sunt salvate informațiile despre știri, după cum se poate observa și din figura 2.

*Modulul agent* este, la rândul său alcătuit din următoarele submodule:

*Submodul de preluare a știrilor* – responsabil cu citirea din baza de date a listei de adrese web ale sindicărilor web, preluarea RSS-urilor de pe Internet și salvarea acestora atât ca fișiere de sine stătătoare, pentru arhivă, cât și în baza de date. Rezultă că acest submodul este responsabil și de parsarea fișierelor RSS;

*Submodul de grupare a știrilor în subiecte* – se ocupă cu citirea celor mai recente știri din baza de date, aplicând apoi un algoritm de grupare pentru determinarea subiectelor. Grupurile sunt apoi inserate în baza de date;

*Clasificatorul de știri* – funcționează în două etape: prima etapă o constituie antrenarea folosind un set de date preluate din RSS-uri predefinite pentru fiecare categorie, iar a doua etapă constă în clasificarea subiectelor determinate la pasul anterior;

*Generatorul de pagini web* – creează paginile web cu știrile cele mai importante din fiecare categorie, paginile corespunzătoare fiecărui subiect, precum și alte pagini web.

*Generatorul de sindicări* – rulează odată cu submodulul anterior, și creează fișierele RSS și Atom folosite pentru sindicalizarea conținutului portalului.

*Modulul web* reprezintă interfața dintre portal și utilizatori și conține patru submodule:

*Interfața (prezentarea) portalului* – conține paginile web și sindicările generate de către submodulele 4 și 5 ale modulului agent, împreună cu alte pagini web statice, imagini și fișiere necesare funcționării portalului web;

*Motor de căutare* – permite căutarea în cadrul subiectelor și știrilor a unor cuvinte cheie;

*Submodul de personalizare* – permite personalizarea paginii de start a portalului în funcție de interesele utilizatorului;

*Sistem de alarme* – definirea de cuvinte cheie pentru trimiterea de alarme prin poșta electronică la întâlnirea acelor cuvinte în cadrul știrilor.

Cele două module pot funcționa independent, fiecare pe câte un server. De asemenea, baza de date poate să fie găzduită pe un server de sine stătător. De altfel, datorită volumului mare de date pe care portalul le va aduna și prelucra de-a lungul funcționării, o astfel de arhitectură este foarte utilă, chiar necesară. În plus, arhitectura are un grad mare de paralelism, întrucât modulul agent poate să fie la rândul său distribuit pe mai multe calculatoare. Submodulele constituente pot să funcționeze independent, necesitând doar acces la baza de date.

#### **Prelucrarea textului și gruparea știrilor**

Fluxurile de știri trebuie să fie procesate înainte de aplicarea metodelor de grupare și clasificare. Deoarece diacriticele nu sunt folosite de către toate sursele de știri, acestea au trebuit să fie eliminate, indiferent de schema de codificare utilizată de furnizor. De asemenea, sunt eliminate etichetele și entitățile HTML, precum și cuvintele de stop. Textul rezultat este descompus în

cuvinte și fiecare termen este adus la forma de *stem*, folosind un algoritm special pentru limba română, pentru a reduce numărul de forme ale cuvintelor. Implementarea unui algoritm de eliminare a sufixelor pentru limba română este extrem de dificilă deoarece regulile de flexionare sunt numeroase și destul de complicate, afectând și structura internă a cuvintelor. Mai mult, fiecare regulă de eliminare a sufixelor poate aduce dezavantaje [14], motiv pentru care s-a optat pentru folosirea unui set restrâns de reguli. Acestea reduc numărul de termeni cu aproximativ 20-25%, depinzând de numărul de știri procesate.

Nr. știri	Spațiu frecvență			Spațiu boolean		
	Nr. grup-uri	Dim. medie	Dim. max.	Nr. grup-uri	Dim. medie	Dim. max.
666	545	1.22	10	556	1.20	6
674	524	1.29	14	510	1.32	38
641	530	1.21	12	532	1.20	13
545	440	1.24	12	430	1.27	35
644	520	1.24	22	533	1.21	15
780	650	1.20	14	641	1.22	39
1024	828	1.22	18	845	1.20	10

**Tabelul 1 Rezultatele grupării**

După etapa de procesare inițială prezentată anterior, se extrage vectorul caracteristic pentru fiecare articol. Deoarece spațiul caracteristicilor este foarte mare, iar o știre conține doar un număr mic de termeni, acești vectori sunt foarte rari, motiv pentru care este utilă o reprezentare folosind un set ordonat după termeni.

Algoritmul de grupare folosește o strategie ierarhică aglomerativă, cu asignare tare și cu un calcul de tip legătură medie pentru similaritatea între grupuri. La fiecare pas, două grupuri sunt unite numai dacă similaritatea lor este mai mare decât un prag determinat în faza de testare a funcționării portalului. De fapt, algoritmul utilizează două praguri diferite, o valoare mai mare pentru crearea unui prim set de grupuri cu articole foarte similare și o valoare mai scăzută ce este utilizată numai când nu se mai găsesc grupuri care să aibă o similaritate mai mare decât prima valoare. Pragul mai mare evită unificarea, într-o primă fază, a articolelor care au subiecte asemănătoare, dar totuși diferite. În acest fel, fiecare grup format în prima fază va acoperi doar un singur subiect. Pentru îmbunătățirea timpului de procesare, a fost folosit un *buffer* pentru memorarea similarităților între grupuri.

Deoarece modalitatea de calcul a similarității este foarte importantă pentru gruparea ierarhică, au fost evaluate rezultatele folosind toate măsurile prezentate în [10], pentru a le găsi pe cele mai potrivite. Pentru spațiile vectoriale ale caracteristicilor bazate pe frecvență, a fost aleasă similaritatea cosinus, iar pentru cele booleene a fost aleasă o măsură definită de către autori și adaptată din teorema cosinusului. Rezultatele obținute folosind aceste metode sunt prezentate în tabelul 1.

### Clasificarea subiectelor de știri

Determinarea celor mai importante subiecte ale zilei îmbunătățește calitatea informației oferite consumatorului de știri. Dar sunt foarte mulți utilizatori care doresc să citească numai știrile dintr-un anumit domeniu de interes. Pentru aceasta, a fost implementat un clasificator automat pentru grupurile de știri, având următoarele opt categorii: *România, Politică, Economie, Cultură, Internațional, Sport, High Tech* și *Monden*. Criteriile determinante în alegerea categoriilor au fost: categoriile să fie cât mai bine definite (să nu se suprapună), să includă domeniile generale de interes și să existe fluxuri de știri specializate pentru fiecare categorie, care să poată fi utilizate în cadrul antrenării clasificatorului.

Trei variații ale clasificatoarelor de tip Nearest Neighbour au fost implementate și comparate: prima este un k-NN simplu, iar cea de a doua folosește un clasificator k-NN cu un scor pentru fiecare categorie, ce este calculat prin însumarea similarităților celor mai apropiate k documente ce aparțin aceleiași clase. A treia metodă seamănă cu un clasificator de tip NN, dar folosește o abordare un pic diferită. În loc de a calcula similaritatea între documentul ce trebuie să fie clasificat și cei mai apropiați vecini ai săi, algoritmul calculează similaritatea între acest document și centroidul fiecărei categorii, alegând categoria de al căru centroid este cel mai similar. Din această cauză, metoda aceasta se numește cel mai apropiat centru (*nearest centre* – NC) sau NN bazat pe centroid și funcționează foarte bine când obiectele dintr-o categorie nu sunt distribuite uniform în jurul centroidului său. Spre deosebire de clasificatoarele NN clasice, această metodă necesită o fază de antrenare simplă ce calculează vectorul caracteristic al fiecărei categorii prin însumarea vectorilor ce aparțin clasei respective din setul de antrenare.

Categorie	Dimensiune set de antrenare	Dimensiune set de validare
Romania	640	325
Politică	313	157
Economie	182	92
Cultură	88	45
Internațional	481	241
Sport	316	158
High-Tech	76	38
Moden (HL)	84	43
Total	2180	1099

**Tabelul 2 Dimensiunile seturilor de antrenare și de evaluare folosite pentru evaluarea clasificatoarelor**

Deși în alegerea categoriilor s-a încercat să se profite de existența unor fluxuri de știri specifice categoriilor, acest lucru nu a fost posibil pentru fiecare categorie, motiv pentru care a fost necesară o hartă a categoriilor furnizorilor de știri către cele ale portalului. În acest fel, a fost posibilă construirea semi-automată a datelor din setul de antrenare. Clasificatoarele au fost comparate folosind tehnica de validare încrucișată, pentru a determina metoda cu cea mai bună acuratețe. Astfel, datele din setul de antrenament, au fost împărțite în două subseturi: două

treimi au fost folosite pentru antrenarea efectivă, iar ultima treime a fost utilizată pentru evaluarea rezultatelor. Dimensiunile fiecărui subset, pentru fiecare categorie, sunt prezentate în tabelul 2.

Numărul de articole din fiecare categorie nu este foarte mare și acestea sunt distribuite inegal, existând o diferență mare între categoria cu cele mai puține articole în setul de antrenament și cea cu cele mai multe. Acestea au afectat procesul de antrenare, rezultatul reflectându-se astfel în performanța clasificatoarelor. Parametrii prezentați în continuare se pot îmbunătăți odată cu creșterea dimensiunii corpusului de antrenament.

Utilizând datele menționate mai sus, au fost testate următoarele clasificatoare: unul de tipul cel mai apropiat centru (NC), unul de tipul NN și două de tipul k-NN (pentru  $k = 3$  și  $k = 5$ ). Parametrii evaluați au fost: acuratețea clasificatorului definită ca numărul de articole clasificate corect împărțit la numărul total de știri, durata procesului de antrenare și durata procesului de clasificare efectivă. În plus, pentru fiecare clasificator a fost calculată matricea de confuzie, dar aceste date nu vor fi prezentate în această lucrare. Pentru fiecare clasificator, au fost utilizate ambele tipuri de criterii de similaritate.

După cum se poate observa și din tabelul 3, ce prezintă acuratețea clasificatoarelor, modelele bazate pe frecvență sunt mai precise decât cele binare, indiferent de algoritmul folosit. Cea mai bună acuratețe este oferită de către algoritmul cel mai apropiat centru (NC), care depășește algoritmiile NN și k-NN, chiar și variantele care folosesc însumarea similarității celor mai apropiați vecini din fiecare categorie. Deși faza de antrenare a algoritmului NC necesită mai mult timp decât cea a algoritmilor NN și k-NN, acesta compensează prin viteză de clasificare, fiind de câteva ori mai rapid decât algoritmiile de tip NN.

Tipul clasificatorului	Timp antrenare (s)	Timp clasificare (s)	Acuratețe (%)
NC – cosinus	279	17	64
NC – binar	280	18	63
NN – cosinus	1	60	61
3-NN – cosinus	1	59	57
5-NN – cosinus	1	59	54
3-NN – cosinus (sumă)	1	59	58
5-NN – cosinus (sumă)	1	59	56

**Tabelul 3 Acuratețea clasificatoarelor implementate**

Clasificatorul de tip NC, folosind vectorul caracteristic al frecvențelor și similaritatea cosinus, s-a dovedit cea mai bună opțiune atât din punctul de vedere al acurateții, cât și a duratei clasificării. În figura 3 este prezentată matricea de confuzie a acestui clasificator – ultima coloană conține precizia, iar ultima linie acoperirea, pentru fiecare categorie. În plus, parametrii săi descriptivi caracteristici sunt următorii: acoperirea medie = 59%, precizia medie = 62%, acuratețea = 64%, F1 (media armonică a preciziei și acoperirii) = 61%.

	Rom	Pol	Eco	Cul	Int	Spo	Teh	HL	Tot	Prec
Rom	177	28	32	11	37	12	10	18	<b>325</b>	0.54
Pol	15	117	10	0	12	2	1	0	<b>157</b>	0.75
Eco	16	2	57	2	10	1	2	1	<b>92</b>	0.63
Cul	11	1	0	21	3	3	2	3	<b>45</b>	0.48
Int	37	12	16	4	162	3	3	4	<b>241</b>	0.67
Spo	5	3	7	4	7	128	1	3	<b>158</b>	0.81
Teh	1	0	9	1	0	0	27	0	<b>38</b>	0.71
HL	6	0	5	10	4	0	2	16	<b>43</b>	0.37
<b>Tot</b>	<b>268</b>	<b>163</b>	<b>136</b>	<b>53</b>	<b>235</b>	<b>149</b>	<b>48</b>	<b>45</b>		
Rec	0.66	0.71	0.42	0.40	0.69	0.86	0.56	0.36		

Figura 3 Matricea confuziei pentru clasificatorul de tip NC

## CONCLUZII

Lucrarea prezintă o alternativă pentru portalurile de știri clasice, proiectată pentru a rezolva problemele cantității mari de știri și a redundanței informației, prin folosirea acestora din urmă ca un avantaj. În plus, portalul folosește sindicalizarea web și tehnici de prelucrare a limbajului natural pentru a atinge o funcționare autonomă, independentă de intervenția umană, ce poate fi numită inteligentă deoarece determină automat importanța unui subiect, precum și categoria sa.

Fluxurile de știri oferă o soluție simplă pentru preluarea știrilor de la un număr mare de surse și gruparea este folosită pentru a exploata știrile similare și pentru a le grupa în funcție de subiect. Utilizatorului îi sunt prezentate subiectele de știri, ordonate în funcție de numărul de știri ce alcătuiesc fiecare subiect. În plus, cititorii au posibilitatea de a alege sursa favorită pentru a citi un articol dintr-un anumit subiect.

Clasificarea automată a subiectelor de știri aduce avantaje importante față de clasificarea unei singure știri deoarece caracteristicile unui grup sunt mai consistente decât cele ale unui singur articol. Acest lucru este de o importanță crucială în aplicație deoarece sindicalizarea oferă numai o scurtă descriere pentru fiecare știre în parte.

Dezvoltarea ulterioară a portalului urmărește două direcții esențiale: îmbunătățirea tehnicilor de grupare și de clasificare, precum și încercarea de a face funcționarea sa independentă de limbă sau, cel puțin, multilingvă.

## PRECIZĂRI

Cercetările prezentate în această lucrare au fost finanțate parțial din proiectul european FP7 STREP LTfLL și din proiectul național CNCSIS K-Teams.

## REFERINȚE

- Chakrabarti, S.: Mining the Web: Discovering Knowledge from Hypertext Data. Morgan Kaufmann Publishers (2002)
- del Corso, G., Gulli, A., Romani, F.: Ranking a Stream of News. Proceedings of the 14th international conference on World Wide Web, Chiba (2005) 97–106
- Dumais, S., Platt, J., Heckerman, D., Sahami, M.: Inductive Learning Algorithms and Representations for Text Categorization. Proceedings of the Seventh International Conference on Information and Knowledge Management, Bethesda (1998) 148–155
- 
- 
- Gabrilovich, E., Dumais, S., Horvitz, E.: Newsjunkie: Providing Personalized Newsfeeds via Analysis of Information Novelty. Proceedings of the Thirteenth International World Wide Web Conference (2004)
- Internet World Stats: World Internet Usage and Population Statistics, disponibil online la adresa <http://www.internetworldstats.com/stats.htm> (2007)
- Manning, C., Schütze, H.: Foundations of Statistical Natural Language Processing. MIT Press, Cambridge, Massachusetts (2003)
- Netcraft: February 2008 Web Server Survey, disponibil online la adresa [http://news.netcraft.com/archives/web\\_server\\_survey.html](http://news.netcraft.com/archives/web_server_survey.html) (2008)
- Page, L., Brin, S., Motwani, R. and Winograd, T.: The page-rank citation ranking: Bring-ing order to the web. Technical Report, Stanford University (1998)
- Porter, M.F.: An algorithm for suffix stripping. Journal of the Society for Information Science, 3(14) (1980) 130–137
- Strehl, A.: Relationship-based Clustering and Cluster Ensembles for High-dimensional Data Mining. Doctoral dissertation, University of Texas at Austin (2002)
- Toda, H., Kataoka, R.: A Clustering Method for News Articles Retrieval System. Special interest tracks and posters of the 14th international conference on World Wide Web, Chiba (2005)
- Ueda, Y., Oka, M., Yamashita, A.: Evaluation of the Document Categorization in "Fixed-point Observatory". Proceedings of NTCIR-5 Workshop Meeting, Tokyo (2005)
- Ullman, J.: Data Mining Lecture Notes (2000)
- Yahoo! Search Blog: Our Blog is Growing Up And So Has Our Index, disponibil online la adresa <http://www.ysearchblog.com/archives/000172.html> (2005)