

Universitatea Babeş - Bolyai
Facultatea de Matematică și Informatică

Relevanța conținutului web și a comportamentului utilizatorilor în analiza traficului

Proiect de cercetare

Student doctorand: Diana-Florina HALIȚĂ
Coordonator științific: Prof. Dr. Florian Mircea BOIAN

Iunie 2015

Cuprins

Introducere	2
1 Impactul similarității web asupra traficului	6
1.1 Migrarea transparentă a unui site web între două sisteme de management de conținut	6
1.1.1 Provocările procesului de migrare	6
1.1.2 Tipuri de migrare a conținutului	7
1.1.3 Implementare și algoritm	9
1.1.4 Rezultatele obținute și evaluarea acestora	11
1.2 Analizarea legăturii dintre similaritatea documentelor web și a ratei de respingere generate de link-urile dintre aceste documente	12
1.2.1 Rezultate anterioare	14
1.2.2 Prezicerea ratei de respingere folosind similaritatea conținutului dintre sursă și destinație	14
1.3 Măsurarea și vizualizarea nivelului de <i>scrapping</i> al unui site web	17
1.4 Concluzii și direcții de cercetare	19
2 Interpretarea logurilor unei platforme de e-learning folosind Analiza Conceptuală Formală	21
2.1 Analiza Conceptuală Formală	21
2.2 Analiza comportamentului utilizatorilor web folosind Analiza Conceptuală Formală	23
2.3 Vizualizarea datelor triadice	25
2.4 Concluzii și direcții de cercetare	31
Lista figurilor	32
Bibliografie	36

Introducere

Navigarea pe Internet a devenit un aspect esențial din viața de zi cu zi a fiecăruia dintre noi. Totodată aceasta începe să aibă din ce în ce mai mult o dimensiune socială, devenind astfel un mecanism eficace prin intermediul căruia se pot dobândi cunoștințe. Dezvoltarea rapidă a World Wide Web, precum și varietatea resurselor care sunt disponibile în Internet, impun necesitatea existenței unor instrumente care pot extrage informații importante din resursele existente sau care pot duce la îmbunătățirea calității experiențelor utilizatorilor în timpul explorării acestora.

Astfel de instrumente joacă un rol important atât pentru administratorii site-urilor web, cât și pentru utilizatorii obișnuiți. Din perspectiva unor experți interesul major ar fi capturarea și analizarea traficului web, cu scopul de a deservi interesul utilizatorilor.

În acest context, lucrarea de față se concentrează pe prezentarea a două abordări diferite. Prima abordare se referă la descoperirea de informații importante disponibile pe web. Cea de-a doua se referă la urmărirea și analiza comportamentului utilizatorilor și a modalităților în care acesta este influențat de conținutul care se regăsește online. De ce aceste două abordări?

Analiza conținutului, ca tehnică de cercetare, este obiectivă, sistematică și poate defini documentul atât din punct de vedere cantitativ cât și din punct de vedere calitativ.

Privind prima abordare urmărită în cercetarea noastră, luând în considerare analiza unui document care este interesant sau util, se poate pune problema găsirii de documente similare care să prezinte un interes la fel de mare pentru studiul considerat.

În literatura de specialitate, tema abordată pentru realizarea acestei lucrări se află la confluența mai multor clase: analiza legăturilor dintre documente și analiza consecințelor ce pot apărea în urma acestor legături.

Privitor la cea de-a doua abordare, scopul acestei lucrări este acela de a găsi noi modalități de analizare a comportamentului utilizatorilor unui site, astfel încât acesta să poată fi adaptat permanent nevoilor acestora (atât din punct de vedere al design-ului, cât și din punct de vedere al tipului de conținut prezentat). Majoritatea cercetătorilor propun utilizarea tehnicilor din *Web Usage Mining* pentru a putea studia astfel de comportamente.

Motivația alegerii temei

Alegerea temei a fost influențată de câteva probleme concrete cu care autorul acestei lucrări împreună cu membrii grupurilor de cercetare din care face parte s-a confruntat în ultimul timp.

Obiective:

- migrarea transparentă a unui site web atât din perspectiva utilizatorilor, cât și din perspectiva motoarelor de căutare, asistarea motoarelor de căutare în redirecționarea corectă a vizitatorilor, precum și redirecționarea consecventă a vizitatorilor care ajung pe site, prin intermediul terților *referreri*;
- estimarea ratei de respingere pentru terțe site-uri cu posibile aplicații în ierarhizarea site-urilor web, precum și a *spam linking*-ului și a *spamdexing*-ului;
- identificarea unei măsurii numerice de cuantificare a nivelului de *scraping* a unui site precum și vizualizarea caracteristicilor unor astfel de site-uri;
- identificarea unor traiectorii în diverse seturi de date asociate unor site-uri cu scopul de a folosi informațiile identificate în reorganizarea acestora din punct de vedere structural sau vizual;
- detectarea comportamentului utilizatorilor într-o platformă de *e-learning* folosind Analiza Conceptuală Formală;
- identificarea de structuri triadice în logurile web ale platformei de *e-learning* și utilizarea acestora în determinarea dinamicii temporale a utilizatorilor, folosind Analiza Conceptuală Formală Triadică;
- identificarea *trendsetter*ilor care creează diverse traiectorii prin site și a grupurilor de utilizatori care urmează traiectorii similare, folosind Analiza Conceptuală Triadică și Analiza Conceptuală Formală Temporală;
- urmărirea interesului utilizatorilor față de informația prezentată pe platformă, considerând scopul pentru care aceasta a fost construită.

Structura

Lucrarea de față este structurată în 2 capitole, după cum urmează.

Capitolul I - Impactul similarității web asupra traficului prezintă prin intermediul celor trei secțiuni principale, modalitatea în care similaritatea documentelor regăsite online își pune amprenta asupra traficului web.

Secțiunea I - Migrarea transparentă a unui site web între două sisteme de management de conținut prezintă pașii care ar trebui urmați astfel încât să se poată face migrarea transparentă a unui site web de la un Sistem de Management bazat pe Conținut (CMS) la altul. Astfel, interesul cade asupra a ceea ce se poate face în privința vizitatorilor site-ului care ar putea fi direcționați greșit de motoarele de căutare sau de terți *referreri* la *URL*-uri vechi, inexistente.

Secțiunea II - Estimarea *bounce-rate*-ului pentru terțe site-uri cu posibile aplicații în ierarhizarea site-urilor, precum și a *spam linking*-ului analizează posibilitatea existenței unei legături între similaritatea documentelor web și a ratei de respingere generate de *link*-urile dintre aceste documente și identificarea *link*-urilor abuzive prin compararea conținuturilor aflate în pagina care conține *link*-ul și pagina spre care duce *link*-ul.

Secțiunea III - Concluzii și direcții de cercetare prezintă ideile principale de continuare a cercetării, printre care și definirea unei metode care ajută la identificarea

scraper site-urilor, adică a site-urilor care copiază întreg conținutul altor site-uri web cu scopul de a manipula motoarele de căutare și de a atrage reclame.

Capitolul II - Interpretarea logurilor unei platforme de e-learning folosind Analiza Conceptuală Formală prezintă în cele patru secțiuni principale cum se poate observa comportamentul utilizatorilor în navigarea lor pe Internet cu ajutorul log-urilor unor servere web și cum poate el determina performanța acestora.

Secțiunea I - Analiza Conceptuală Formală prezintă aspectele teoretice din Analiza Conceptuală Formală (FCA), precum și două ramuri ale acesteia: Analiza Conceptuală Formală Triadică (3FCA) și Analiza Conceptuală Formală Temporală (TCA). De asemenea, tot aici sunt prezentate principalele instrumente utilizate de comunitatea cercetătorilor FCA.

Secțiunea II - Analiza comportamentului utilizatorilor web folosind Analiza Conceptuală Formală prezintă o nouă modalitate în care aceasta poate fi aplicată în *Web Usage Mining* pentru a determina dinamica utilizatorilor în portalul de *e-learning* studiat.

Secțiunea III - Vizualizarea datelor triadice prezintă CIRCOS, un instrument de vizualizare grafică a datelor obținute folosind 3FCA. Scopul folosirii acestuia a fost de a evidenția legăturile importante între conceptele triadice.

Secțiunea IV - Concluzii și direcții de cercetare prezintă ideile principale de continuare a cercetării, printre care și construirea unui sistem adaptativ care se mulează pe descoperirea de noi informații în seturi de date folosind 3FCA și TCA. Prin intermediul 3FCA se determina grupurile de utilizatori care tind să urmeze același tipar în vizitarea site-ului web, iar prin TCA se determină efectiv lanțurile de pagini care fac parte din traiectoriile cele mai des utilizate de către vizitatorii site-ului.

Lucrarea se încheie cu concluzii, o anexă ce cuprinde lista figurilor introduse în lucrare precum și cu prezentarea bibliografiei utilizate pe parcursul redactării lucrării.

Cercetări anterioare în domeniu

Data fiind actualitatea temei de cercetare propusă în această lucrare, literatura de specialitate prezintă o serie de publicații sau aplicații care o abordează.

Totodată, dinamica subiectului necesită o abordare continuă și de substanță.

Soluțiile adoptate în cercetările anterioare referitoare la subtemele propuse în această lucrare sunt:

- migrarea transparentă a unui site web:
 - existența unor *plugin*-uri implementate pentru majoritatea CMS-urilor care se bazează pe date statistice stocate de-a lungul timpului cu ajutorul cărora s-a stabilit un comportament al utilizatorilor, și în funcție de acesta se determină tiparul pe care e posibil să vrea să îl urmeze noul vizitator;
 - trimiterea unui utilizator nou către o pagină care dă eroare 404 sau la o pagină care are implementat un modul de căutare în interiorul site-ului web ([49], [17], [45], [20]);
- *spam linking* și *spamdexing*:
 - dezvoltarea de tehnici anti *spam* ([38],[47]);

- utilizarea tehnicilor de clasificare supervizată și nesupervizată în detectarea paginilor care sunt *spam*-uri
- utilizarea metodelor învățării automate ([47]);
- *Web Usage Mining*
 - folosind analiza statistică a datelor ([28]), graf-uri, instrumente care fac analiza logurilor web (i.e., Google Analytics, [16]) , business intelligence ([37], [39])

Rezultate și contribuții proprii

Dintre elementele originale propuse de autor, menționăm următoarele:

- construirea *layer*-ului prin intermediul căruia se poate face migrarea transparentă de la un CMS la altul ([9]);
- alegerea unor măsuri de similaritate adecvate pentru algoritmul de regăsire a perechilor de URL-uri care prezintă conținuturi similare;
- găsirea unei legături între rata de respingere a unei pagini dintr-un site web și similaritatea acesteia cu *referrerii* care indică spre ea ([24]);
- aplicarea pentru prima dată a 3FCA în explorarea log-urilor web;
- utilizarea întregii suite FCA (3FCA și TCA) pentru explorarea log-urilor web ([12]);
- utilizarea Circos, instrumentul de vizualizare grafică a datelor obținute în urma folosirii 3FCA;
- determinarea dinamicii utilizatorilor folosind FCA și 3FCA ([14]);

Capitolul 1

Impactul similarității web asupra traficului

1.1 Migrarea transparentă a unui site web între două sisteme de management de conținut

Sistemele de management a conținuturilor au început să fie din ce în ce mai utilizate în ultimul timp. Acest lucru reiese din top-ul realizat de site-ul alexa.com, în care 22,5% dintr-un milion de site-uri web sunt construite folosind un astfel de CMS, iar peste 50% dintre cele care folosesc un CMS utilizează Wordpress [33].

Avantajele utilizării unui astfel de sistem sunt: întreținerea facilă a paginii web, independența conținutului de prezentare, migrarea ușoară de la un design la altul, update-uri periodice de securitate, control total asupra optimizării motoarelor de căutare, administrare web, funcționalități suplimentare oferite de terți.

Argumentele aduse mai sus reprezintă o motivație clară a faptului că CMS-urile sunt din ce în ce mai folosite și totodată aduc un plus migrării conținutului site-urilor vechi la un astfel de sistem, în special în cazul în care CMS-ul curent nu mai poate îndeplini scopurile site-ului web, nu suportă noile obiective sau nu poate implementa toate funcționalitățile, respectiv caracteristicile necesare îndeplinirii acestora.

1.1.1 Provocările procesului de migrare

Migrarea de la un CMS la altul poate fi făcută fie manual, fie automat. În cazul migrării de la un site static, procesul de migrare implică aproape în totalitate operații manuale. Pe de altă parte, migrarea de la un CMS la altul va implica operații automate.

Operațiile care se fac în timpul procesului de migrare implică migrarea conținuturilor paginilor web, dar mai mult decât atât cea mai dificilă parte a procesului de migrare este maparea vechiului conținut celui nou corespunzător. De asemenea este importantă și corelarea structurii organizatorice a vechiului site cu cea nouă dată de noul CMS.

Acest lucru poate avea o rezolvare graduală, urmărind trei pași: împărțirea conținutului în categorii, estimarea timpului necesar migrării, reevaluarea bazată pe ghidarea migrării.

Primul pas al procesului de migrare este clasificarea diverselor conținuturi, în funcție de rezultatul analizei acestuia și de tipurile obținute (prin tipuri se înțelege totalitatea categoriilor în care a fost subdivizat conținutul). Astfel se nasc două tipuri de reguli

care trebuie identificate: acele reguli care se referă concret la conținutul ales a fi necesar și după migrare, precum și regulile care se referă la ceea ce va putea și la ceea ce nu va putea fi migrat automat. Aceste reguli care se definesc sunt utile în stabilirea priorităților și în dozarea efortului care trebuie depus. În acest moment se poate decide ceea ce se poate migra automat și ce nu. În general se dorește să se automatizeze cât mai mult din ceea ce trebuie migrat.

Al doilea pas se referă la ghidarea și estimarea timpului necesar migrării, adică compararea timpului necesar migrării automate, respectiv migrării manuale.

Ultimul pas presupune evaluarea procesului de migrare automată. Cea mai mare problemă care se pune atunci când se dorește migrarea automată, este structura și regularitatea conținutului. Pe de altă parte, migrarea manuală spre un alt CMS presupune vizualizarea, editarea și mutarea manuală a conținutului, ceea ce probabil se va solda cu o risipire a resurselor de timp.

Procesul de migrare nu implică numai transformarea conținutului vechi în conținut nou, ci migrarea, fie ea automată sau manuală, ar trebui să ia în considerare comportamentul transparent al site-ului web, atât din punctul de vedere al vizitatorului sau al unui crawler al unui motor de căutare cât și al unui browser web. O problemă comună care este indusă de o astfel de migrare este dată de expunerea conținutului unui site web la un URL nou și diferit față de cel vechi. După migrare, acest lucru se va reflecta în creșterea numărului de vizitatori care vor fi induși în eroare de către link-urile apărute în *SERP* (*Search Engine Results Page*) sau de către terți *referreri*. Pe lângă faptul că vizitatorilor nu li se afișează conținutul dorit, site-ul ce tocmai a fost migrat ar putea pierde diversele beneficii câștigate în timp, pe care le avea datorită link-urilor postate pe diverse rețele de socializare sau a *back link*-urilor.

Astfel, a fost propus în [9] un *layer middleware* care să redirecționeze transparent *request*-urile către URL-uri din vechiul site la URL-uri care aparțin noului site (atunci când maparea nu se poate face automat), pe baza, nu a conținutului comun, ci pe baza semanticii conținutului descris de cele două URL-uri și pe baza semanticii unei informații adiționale, cum ar fi query-ul dat de motoarele de căutare sau conținutul de la URL-ul *referrer*-ului. Câteva dintre avantajele mecanismului propus, ar fi: reducerea numărului de *dead link*-uri venite dinspre *referreri 3rd party*, asistarea motoarelor de căutare în direcționarea corectă a utilizatorilor, precum și conservarea *pageranking*-ului și a rezultatelor din *SERP*. Spre deosebire de soluțiile propuse în cercetările anterioare ([49], [17], [45]), sistemul propus în [9] are avantajul de a fi implementat înainte de lansarea noului site, nedepinzând de date statistice strânse de-a lungul unei perioade de timp.

1.1.2 Tipuri de migrare a conținutului

În continuare, vom propune un mecanism de formalizare a notațiilor, adaptat CMS-urilor implicate în procesul de migrare (phpWebSite, Wordpress).

Conținutul site-ului web este prezentat în moduri diferite în vechiul, respectiv în noul site. Astfel, apar două tipuri de conținut ce se poate regăsi în structura site-ului: conținut static (adică un fișier ce rezidă pe un server web, de exemplu o resursă .pdf sau .jpg) sau conținut dinamic, preluat dintr-o bază de date.

Migrarea conținutului static în conținut static.

Chiar dacă URL-ul se schimbă numele fișierului rămâne la fel:

vechiul URL: `http://vechiulsite/vecheaaddressa/numefisier`

noul URL: `http://vechiulsite/nouaaddressa/numefisier`

Putem spune că similaritatea perfectă se obține atunci când similaritatea dintre URLvechi și URLnou este 1.

Migrarea conținutului dinamic în conținut dinamic.

Un astfel de conținut se poate regăsi atât la un URL vechi, cât și la unul nou, și cuprinde două părți: template-ul paginii web și conținutul stocat în baza de date.

Ceea ce face obiectul migrării este conținutul absolut al paginii web, deoarece în calculul similarității dintre conținuturile aflate în paginile de la url-ul vechi, respectiv nou, trebuie evitate informațiile care țin de template-ul site-ului și care nu au nici o legătură cu conținutul prezentat într-o anumită pagină.

Exemple:

- În Wordpress conținutul absolut se poate extrage direct din baza de date, sau direct de la URL-ul asociat. Acesta poate fi diferențiat față de template-ul paginii, datorită poziționării într-un div care este unic determinat de id-ul `content`.

Se vor nota în continuare paginile din vechiul site cu OP și cele din noul site cu NP . Conținutului absolut din site-urile construite în Wordpress sunt reprezentate de: pagini(atemporale), NP_{-T} , postări(temporale), NP_T , categorii, sau reuniune de postări și submulțime de postări dintr-o categorie anume.

- În phpWebSite conținutul absolut se regăsește în: pagini cu o singură secțiune, OP_{SS} (migrate într-o pagină NP_{-T}) sau pagini cu mai multe secțiuni $OP_{MS} = \bigcup content(section_i)$. O astfel de pagină este migrată fie într-o singură pagină NP_{-T} , dacă secțiunile din pagina de la vechiul URL sunt atemporale, sau fiecare secțiune devine o postare NP_T ce va aparține unei anumite categorii. Migrarea postărilor temporale în Word-press au ridicat anumite probleme, dat fiind faptul că în phpWebSite nu se rețin informații temporale legate de data creării sau modificării secțiunii din pagină.

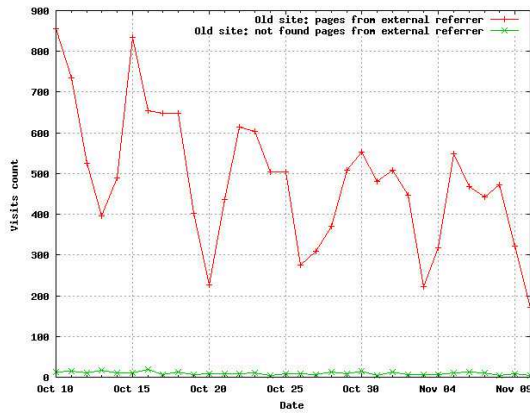
S-au întâlnit dificultăți la migrarea paginilor în postări temporale, datorită faptului că în phpWebSite nu sunt memorate date temporale.

Atunci când am comparat similaritatea conținuturilor dintr-o pagină cu secțiuni multiple prezentată în vechiul site a fost nevoie să se facă compararea cu conținutul absolut aflat într-o submulțime a postărilor dintr-o anumită categorie. Această necesitate va fi prezentată într-o secțiune ulterioară.

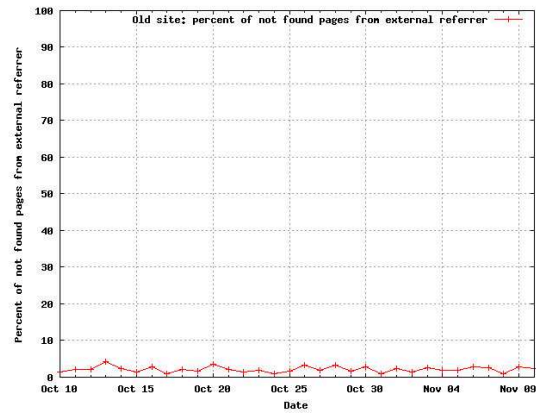
Migrarea conținutului static în conținut dinamic.

Această problemă apare atunci când e necesar să integrăm conținuturile fișierelor HTML în structura generală a site-ului. Astfel, în acest caz, similaritatea se va calcula folosind conținutul absolut din ambele site-uri.

Metoda propusă în cazul migrării unui site între două CMS-uri se bazează întotdeauna pe compararea conținuturilor absolute, deoarece:



(a) Vechiul site: număr



(b) Vechiul site: procent

Figura 1.1: Pagini ce au generat eroare 404 și au fost accesate de către un *referrer* extern

- este o metodă generală care poate fi luată în considerare atunci când se migrează și de la un site static. Ea nu se bazează pe compararea efectivă a conținuturilor din baza de date, ceea ce îi garantează caracterul general;
- alte metode nu acoperă cazurile enumerate anterior;
- este o metodă care nu necesită date de conectare la baza de date, datorită faptului că se bazează pe ceea ce este prezentat la un anumit URL. Acest lucru duce la ideea că nu este necesar să existe drepturi de administrator pentru că scopul este de a accesa cât mai ușor informația ce se regăsește la un anumit URL;
- procesul prin care se regăsesc perechile corespunzătoare se face din perspectiva clientului.

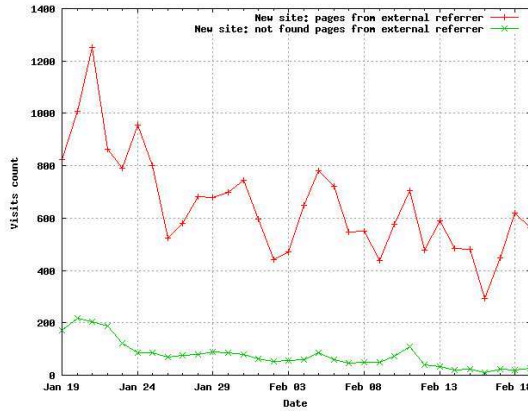
1.1.3 Implementare și algoritm

Voi prezenta în continuare ideile principale ale algoritmului utilizat pentru a potrivi cât mai bine perechile de URL-uri din vechiul și noul site. Pentru a găsi perechile potrivite, algoritmul utilizează o funcție de similaritate a două șiruri de caractere și anume similaritatea Cosinus. Algoritmul nu este dependent de funcția de similaritate aleasă, însă utilizând această funcție s-au obținut cele mai bune rezultate în cel mai scurt timp [36].

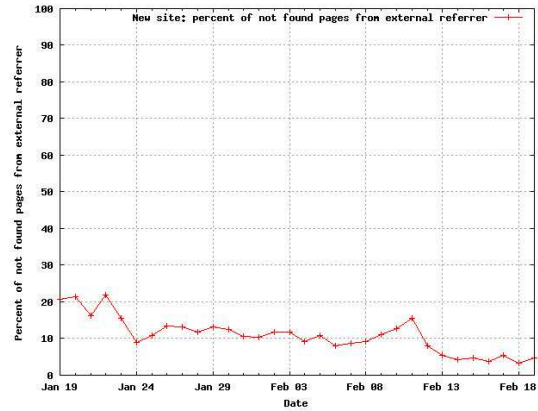
Procesarea informației nu s-a făcut în timp real. Soluția aleasă a fost aceea de a rula un program anterior accesării paginilor, astfel încât să se găsească cele mai bune potriviri între URL-uri. Am preferat acest lucru datorită faptului că atât în site-ul vechi, cât și în cel nou, sunt mii de url-uri care trebuie luate în considerare, astfel, complexitatea algoritmului fiind egală cu cardinalul produsului cartezian al celor două mulțimi de URL-uri. Acest fapt ar duce la deservirea răspunsului către client cu întârziere. În plus, algoritmul propus este complet independent atât de platforma utilizată cât și de limbajul de programare folosit în implementarea CMS-ului.

Prezentarea formală a algoritmului este următoarea:

Pentru fiecare URL din vechiul site la care se regăsește conținut static
 Identificarea URL-ului din noul site care indică spre același conținut



(a) Noul site: număr



(b) Noul site: procent

Figura 1.2: Pagini ce au generat eroare 404 și au fost accesate de către un *referrer* extern

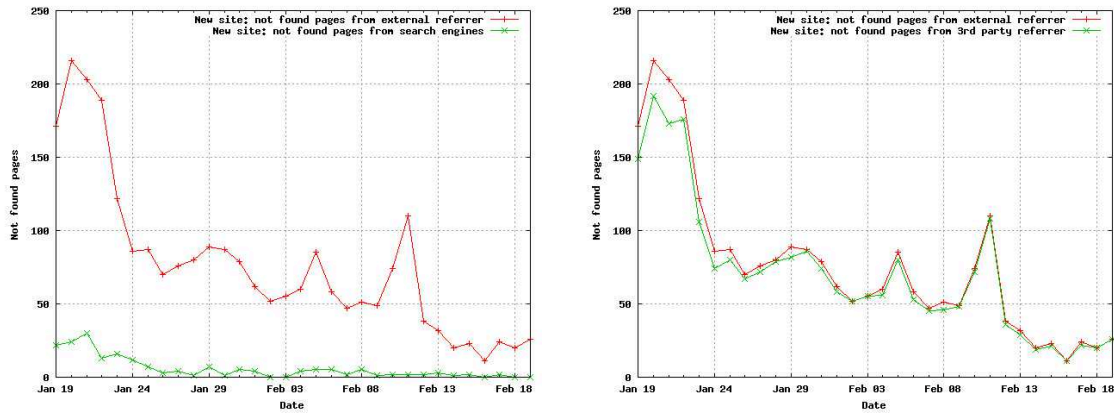
```

    acest conținut se bazează pe numele fișierului
    Dacă URL-ul din noul site a fost identificat
        elimină vechiul URL din lista de URL-uri care trebuie să fie procesate
    SfDacă
SfPentru
Pentru toate URL-urile neprocesate
    Dacă la acest URL se regăsește conținut static
        conținutul_absolut = conținutul efectiv
    altfel
        conținutul_absolut = conținutul(vechiulURL) - conținut (template-ul
            paginii)
    SfDacă
    Identifică URL-urile din noul site astfel încât conținutul absolut
        are cea mai bună similaritate cu conținutul absolut al paginii
        de la vechiul URL
    Perechea se inserează în baza de date, împreună cu:
        data, similaritatea obținută la rularea algoritmului și cea mai bună
        similaritate obținută pe parcursul tuturor rulărilor
SfPentru

```

Conform algoritmului de mai sus s-au salvat în baza de date mai multe date referitoare la similaritatea perechilor de URL-uri, printre care și data și similaritatea obținută la ultima rulare, precum și cea mai bună similaritate obținută de-a lungul timpului. Motivarea acestei alegeri vine de la faptul că postările dintr-o categorie evoluează în timp datorită noilor anunțuri ce pot să intervină și astfel similaritatea obținută la diferite rulări descrește. Deci, vom avea similaritate maximă doar dacă comparăm conținutul dintr-o anumită pagină cu secțiuni multiple (din vechiul site), cu o submulțime formată din mai multe anunțuri mai vechi (acestea vor corespunde cu cele mai noi anunțuri din noul site).

În acest sens prezint un exemplu: din categoria **Anunțuri Recente** din noul site, un vizitator este foarte posibil să își dorească să vizualizeze cele mai noi articole postate. Astfel, redirecționarea se va face spre acestea, și nu spre un anunț care se potrivește



(a) Adaptarea motoarelor de căutare la noua structură a site-ului (b) Majoritatea erorilor 404 sunt generate de către utilizatorii care vin de la referreri *3rd party*

Figura 1.3: Adaptarea rapidă a motoarelor de căutare la noua structură a site-ului

perfect, din punct de vedere al similarității, atunci când utilizatorul vine de la un referrer care face trimitere spre pagina cu anunțuri recente din vechiul site.

De asemenea, pentru a se grăbi puțin procesul de regăsire a perechilor, se poate determina un prag experimental, astfel încât orice pereche de tipul (URL vechi, URL nou) ce va avea similaritatea mai mare decât pragul experimental considerat, să se excludă din lista URL-urilor ce intră în procesul de comparare.

1.1.4 Rezultatele obținute și evaluarea acestora

Pentru a demonstra avantajele metodei propuse vom prezenta rezultatele unor experimente, efectuate pe parcursul a 30 de zile, care evidențiază din punct de vedere numeric comportamentul noului și respectiv al vechiului site atunci când vizitatorul dorește să vizualizeze o pagină și o accesează prin intermediul unui referrer extern (o anumită pagină este accesată prin intermediul motoarelor de căutare sau a altor site-uri care conțin link-uri spre site-ul migrat, și nu este accesată direct prin URL sau bookmark).

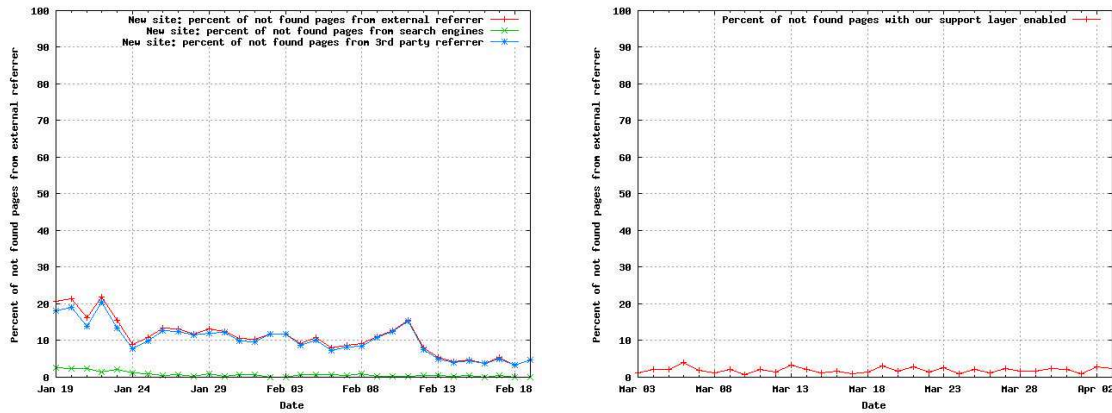
Astfel, au fost interesante rezultatele obținute privind numărul paginilor negăsite, care au dat eroare 404 la cererea unei anumite resurse, precum și procentul acestora relativ la numărul total de cereri HTTP care au fost efectuate.

Numărul paginilor care au generat eroare 404 înainte de migrare a fost relativ mic (1,9% Fig. 1.1(b)), acestea provenind în special de la referreri *3rd party*.

După migrare, numărul paginilor ce au generat eroare 404 și au fost accesate prin intermediul unui *referrer* extern a crescut în mod natural la aproximativ 11,2% (Fig. 1.2(b)) pe parcursul celor 30 de zile considerate.

S-a observat că motoarele de căutare s-au adaptat repede la noua structură a site-ului, astfel, ele redirectând utilizatorii corect după o săptămână de la migrare (Fig. 1.3(a), Fig. 1.4(a)).

Procentajul paginilor ce generau eroarea 404, și care aveau ca referreri motoarele de căutare, observate în prima săptămână de după migrare era în jur de 1,9%, în timp ce după patru săptămâni, acest procent a scăzut la 0,23%. Același lucru s-a întâmplat și în cazul paginilor care aveau referreri *3rd party*, procentul scăzând de la 14,7% la 4,25% în patru săptămâni.



(a) Pagini ce generează eroarea 404 și pagini ce (b) Pagini ce generează eroarea 404 având un re-
 generează eroarea 404 și vin de la motoarele de *ferrer* extern, odată ce a fost activat *layer*-ul su-
 căutare sau de la *referreri 3rd party* port propus

Figura 1.4: Adaptarea rapidă a motoarelor de căutare la noua structură a site-ului

Figura 1.4(b) prezintă rezultatele obținute odată cu punerea în funcțiune a *layer*-ului suport propus în lucrarea de față, observându-se pentru site-ul ce tocmai a fost migrat, o îmbunătățire a procentului de pagini ce genereau eroarea 404 la 1.6%.

În algoritm s-a observat existența perechilor de URL-uri care aveau o similaritate de peste 70% în peste 93% din cazuri.

În concluzie, rezultatele experimentale prezentate mai sus au desăvârșit motivația conform căreia un astfel de *layer*, prin intermediul căruia se poate face migrarea transparentă de la un CMS la altul, este absolut necesar. Metoda propusă acoperă atât aspectele practice cât și cele teoretice necesare mapeării URL-urilor din vechiul site, celor corespunzătoare din noul site, pe baza similarității dintre conținuturile prezentate la acele URL-uri. Totuși, scopul introducerii acestui *layer* rămâne acela de a redirecționa vizitatorul site-ului cu succes, indiferent de tipul referrer-ului de la care acesta vine.

1.2 Analizarea legăturii dintre similaritatea documentelor web și a ratei de respingere generate de link-urile dintre aceste documente

Unul dintre principalele scopuri ale unui link este să ofere vizitatorului mai multe informații înrudite semantic cu informația din documentul considerat. Când vine vorba de link-uri externe de foarte multe ori în internet acestea sunt folosite abuziv, doar pentru a crește *pagerank*-ul paginii sau al domeniului destinație și nu pentru a ghida utilizatorul spre una sau mai multe pagini care să îi ofere în continuare informația de care să fie efectiv interesat.

În majoritatea cazurilor asemenea link-uri abuzive sunt fie *sitewide* (cel mai ușor de detectat), dar pot să fie plasate automat (în cadrul conținutului absolut al unei pagini) de diferite module sau add-uri integrate în CMS-ul site-ului sursă.

Astfel, se pune problema existenței unei posibile legături între similaritatea conținutului documentului sursă cu conținutul documentul indicat printr-un link extern prezent în sursă și rata de respingere generată de link-ul respectiv în cadrul site-ului sau al domeniului destinație.

Rata de respingere asociată unei pagini web sau unui site reprezintă procentajul utilizatorilor care au vizitat o pagină din site și au preferat să părăsească site-ul în loc să continue cu vizitarea altor pagini prezentate în cadrul aceluiași site. În general acest comportament se poate declanșa în urma a două posibile scenarii.

Primul motiv este acela că utilizatorul a găsit prezentat la acea pagină exact informația pe care o căuta. De exemplu, cineva poate căuta definiția expresiei *rata de respingere*. Această informație este prezentată foarte bine chiar pe prima pagină care apare în SERP. El este mulțumit cu definiția pe care a găsit-o și părăsește site-ul fără să mai acceseze și alte pagini din interiorul acelui site.

Totuși, în majoritatea cazurilor lucrurile se întâmplă diferit, utilizatorul nefiind mulțumit cu informațiile oferite pe o pagină particulară din SERP, astfel părăsind-o pentru a putea accesa următoarele pagini din SERP.

Este bine cunoscut faptul că o rată de respingere mare, de obicei cauzată de al doilea caz prezentat anterior, are o semnificație negativă, iar aceasta este de obicei asociată cu calitatea conținutului prezentat.

Cu cât documentele prezente pe domeniul extern oferă un conținut mai similar cu documentul sursă, vizitatorului i se garantează accesul la tot mai multă informație care este de mare interes pentru el. Odată cu creșterea numărului de documente accesate ce aparțin site-ului destinație rata de respingere a acestuia scade.

Exemple:

- **forum dedicat iubitorilor de animale**

Un astfel de forum poate conține în unul din *topic*-urile deschise un link către un site ce se adresează crescătorilor de câini. Un astfel de link generează o rată de respingere mică, deoarece informația din pagina de destinație oferă un plus de conținut de calitate, de care utilizatorul să fie interesat.

- **site-ul admiterii la Universitatea Babeș-Bolyai**

Acest site prezintă informații detaliate despre procesul de admitere. Link-urile externe din orice pagină a acestui site reprezintă legături spre mai multe informații legate de admitere. Aceste informații sunt prezentate mai în detaliu pe site-ul corespunzător fiecărei facultăți, site ce este găzduit pe un domeniu diferit. Astfel de link-uri oferă informații legate din punct de vedere semantic care sunt mult mai detaliate decât cele oferite pe site-urile admiterii, astfel generându-se o rată de respingere mai mică pentru site-ul facultăților.

Contraexemple:

- **site pentru care majoritatea link-urilor externe sunt reclame**

Un link abuziv (spamlink, adlink, reclamă) în foarte multe cazuri duce spre site-uri cu un conținut total diferit față de site-ul sursă, iar prin urmarea unui astfel de link se generează o rată de respingere mai mare.

Pentru a analiza această posibilă legătură se vor folosi mai multe funcții de similaritate, precum: similaritatea Cosinus, similaritatea Jaccard, similaritatea Jaro-Winkler și similaritatea Sorensen.

1.2.1 Rezultate anterioare

Odată cu creșterea exponențială a Internetului și a posibilității navigării pe web, expunerea link-urilor abuzive a devenit un fenomen negativ ce afectează calitatea rezultatelor returnate de către motoarele de căutare.

Din ce în ce mai multe companii din industria optimizării motoarelor de căutare, precum și diverși cercetători și-au reunit forțele cu scopul de a identifica potențiale soluții referitoare la această problemă și de a limita efectul negativ al acestui fenomen.

Încă de când analiza link-urilor a fost folosită pentru optimizarea motoarelor de căutare, s-a încercat introducerea de diverse tehnici de spammare [47], obținându-se astfel multiple efecte negative ce au dus la noi provocări în ariile de cercetare.

Un studiu făcut asupra principalelor tehnici de detectare a link-urilor abuzive pe internet [38] clasifică aceste tehnici în trei clase: metode bazate pe analiza conținutului, metode bazate pe analiza link-urilor și metode bazate pe analiza datelor netradiționale (sesiuni HTTP, comportamentul utilizatorului pe web). Prin aceste tehnici se detectează până la 80% din paginile abuzive [4], iar recomandarea autorilor este să se folosească tehnicile enumerate mai sus împreună, pentru a crește procentul de pagini abuzive detectate (este recomandat să se folosească împreună cel puțin metodele ce analizează link-urile și metodele care analizează conținutul [3]).

Studiile anterioare care au propus diverse tehnici de detectare a *spam*-urilor web indică o posibilă clasificare automată a acestor tehnici prin metode de clasificare supervizate sau nesupervizate. De asemenea s-a mai propus și utilizarea de algoritmi probabilistici, algoritmi de detectare a coliziunilor sau matrici de confuzie pentru anumiți clasificatori [31]. Aceste propuneri au ridicat o mare provocare printre cercetători și anume ridicarea *spamdexing*-ului la nivel de problemă de învățare automată [47].

1.2.2 Prezicerea ratei de respingere folosind similaritatea conținutului dintre sursă și destinație

Similaritatea ideală

Pe baza celor spuse mai sus, paginile care au o similaritate mare din punctul de vedere al conținutului prezentat au șanse mai mici să genereze o rată de respingere mică, și reciproc, paginile care au o similaritate mai mică a conținutului generează o rată de respingere mai mare.

Astfel, problema care se pune este dacă este posibil să existe o legătură între rata de respingere și rezultatele obținute prin folosirea mai multor funcții de similaritate.

Metodă propusă în [24] pentru a găsi o funcție de similaritate ideală, ar putea să estimeze rata de respingere generată de link-uri de pe terțe site-uri spre terțe site-uri, fără a avea acces la diferite drepturi de administrare a domeniilor destinație (cum ar fi tool-urile oferite de Google Analytics).

O astfel de metodă ar permite identificarea pentru un site sursă a link-urilor abuzive, *add*-uri, impropriu plasate, care generează o rată de respingere mare, lucru care poate fi util în penalizarea site-ului sursă în cadrul motoarelor de căutare.

Totodată, o astfel de funcție de similaritate ideală ar trebui să se comporte ca și în graficul 1.5, adică se poate face o interpolare liniară la funcția de gradul I $f(x) = 100 - x$.

Corelarea perfectă dintre rata de respingere și similaritatea a două documente legate între ele (așa cum este prezentat în cazul ideal, în figura 1.5), nu este întotdeauna găsită, în special luând în considerare funcțiile uzuale de similaritate. Astfel, am luat

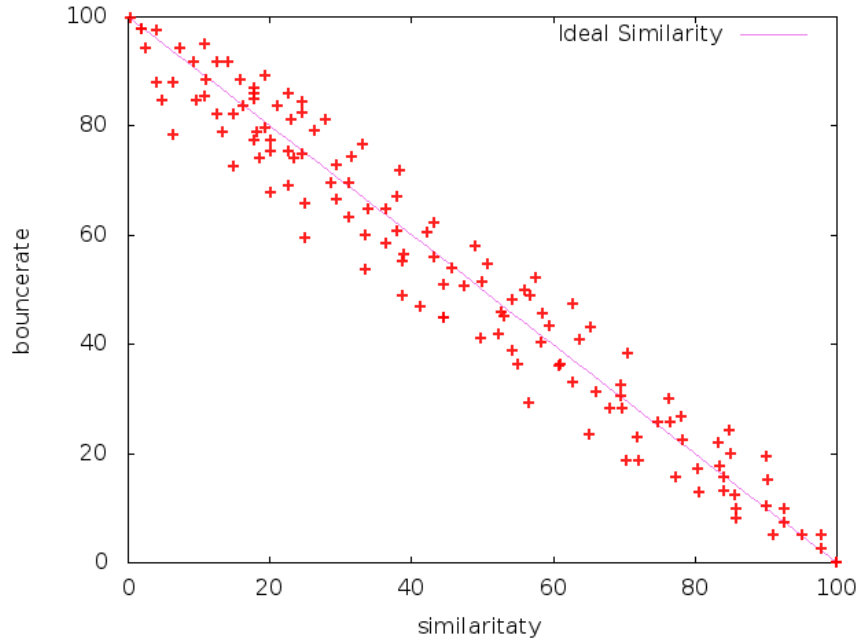


Figura 1.5: Similaritate ideală

în considerare mai multe astfel de funcții de similaritate pentru a determina care dintre ele se potrivesc cel mai bine cu similaritatea ideală descrisă mai sus.

Intuitiv, se consideră că o funcție de similaritate este mai bună decât o altă funcție de similaritate, dacă perechile de tipul (*similaritate, rata de respingere*), ce se obțin din link-urile ce leagă două documente, sunt mai aproape de diagonala reprezentată în cazul ideal în figura 1.5.

Pentru reprezentarea grafică a rezultatelor s-a considerat un reper cartezian xOy , unde similaritatea reprezintă abscisa punctului, iar rata de respingere reprezintă ordonata punctului reprezentat în sistemul considerat.

Luând în considerare dreapta paralelă cu a doua bisectoare ce trece prin punctele de coordonate $(100,0)$ și $(0,100)$, o funcție de similaritate este cea mai bună dacă suma tuturor distanțelor de la reprezentările grafice ale punctelor la dreapta mai sus menționată este minimă, adică:

$$\sum \frac{|x_i + y_i - 100|}{\sqrt{2}}$$

este minimă, unde x_i și y_i sunt coordonatele punctelor reprezentate grafic în sistemul cartezian de axe.

Rezultate experimentale

Pentru a evalua rezultatele experimentale obținute prin metoda propusă, se vor lua în considerare: similaritatea Cosinus, similaritatea Jaccard, similaritatea Jaro-Winkler și similaritatea Sorensen.

Pentru a obține rezultate cât mai reale și pentru a avea o acuratețe cât mai bună, am testat toate funcțiile de similaritate pe conținutul absolut al paginilor web, ignorând conținutul din șablonul paginii web, adică header-ul, footer-ul, respectiv meniul.

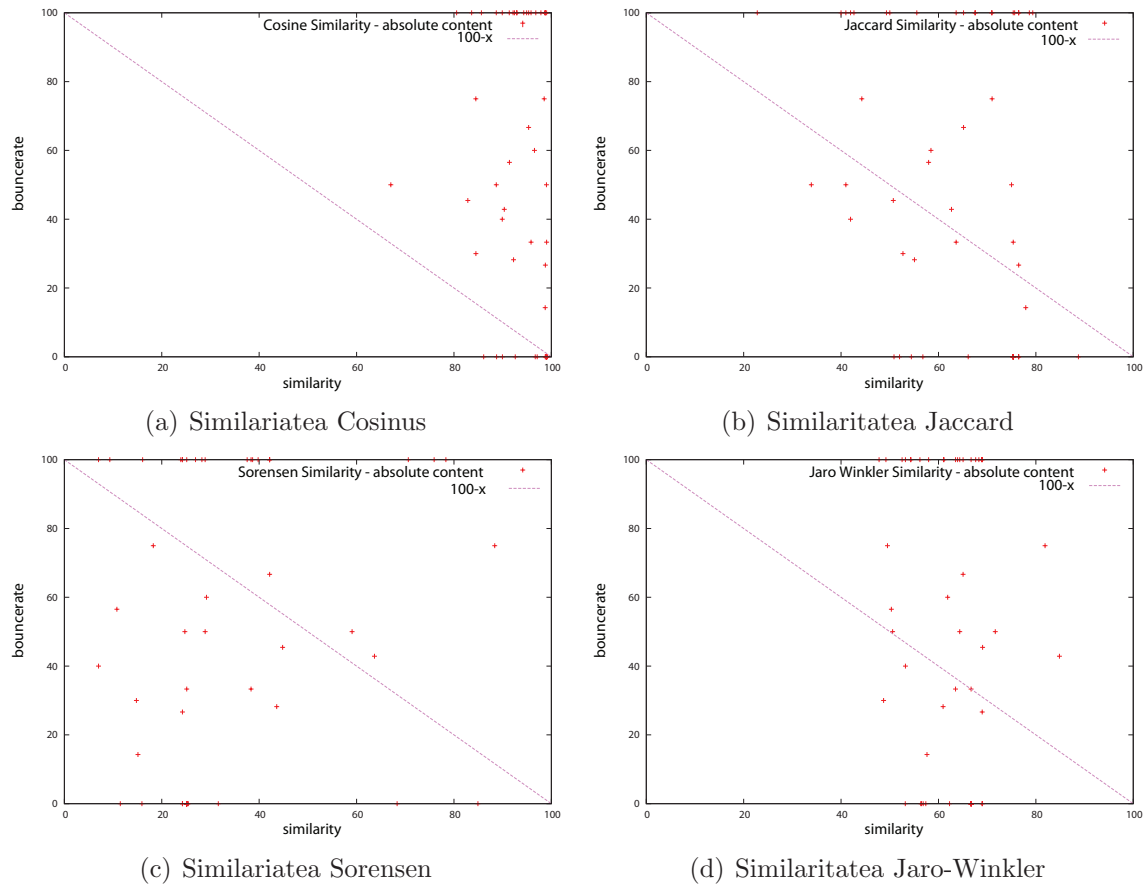


Figura 1.6: Similaritati

Conținutul absolut al paginii web a putut fi complet determinat folosind o librărie Java, numită *boilerpipe*. Această librărie, oferită de Apache License 2.0, furnizează algoritmi care detectează și elimină tot ceea ce ține de template-ul site-ului web, păstrând doar conținutul absolut al acestuia [1].

Pentru a demonstra ceea ce mi-am propus prin lucrarea de față, am luat în considerare site-ul Facultății de Matematică și Informatică al Universității Babeș-Bolyai (<http://www.cs.ubbcluj.ro>) și am generat tripletele (pagină accesată, referrer, rată de respingere), prin două metode:

- folosirea unui tool client-side oferit de Google (Google Analytics) ;
- folosirea unui tool server-side, integrat în template-ul site-ului, dezvoltat de autorii [24].

Reușita experimentului se bazează pe faptul că am avut drepturi de administrare asupra site-ului mai sus menționat, ceea ce a permis măsurarea corectă a ratei de respingere pentru fiecare link extern referrer al paginilor din acest site.

După cum se observă din figurile prezentate, cea mai bună funcție de similaritate care oferă rezultatele dorite, este similaritatea Jaccard. Matematic, faptul că similaritatea Jaccard este cea mai bună funcție de similaritate rezultă din calcularea sumei distanțelor de la fiecare punct ce e reprezentat grafic la dreapta de ecuație $y = x - 100$.

Funcție de similaritate	Metodă	Valoarea sumei
Cosinus	conținut absolut	1804.476000579978
Jaccard	conținut absolut	1414.03085464388
Sorensen	conținut absolut	1769.3699189543242
Jaro-Winkler	conținut absolut	1528.5359097516346

Tabela 1.1: Suma distanțelor de la toate punctele de pe grafic la dreapta de ecuație $y = x - 100$

1.3 Măsurarea și vizualizarea nivelului de *scrapping* al unui site web

Odată cu dezvoltarea Internetului numărul site-urilor a crescut considerabil, și astfel se poate explica ușor abundența informațiilor care există pe web. Din păcate în ultimul timp atunci când se dorește accesarea unei informații prin intermediul unui motor de căutare, atât user-ul cât și motorul de căutare sunt puși în situația unei probleme: prezența unor site-uri web în *SERP* care direcționează greșit utilizatorii, sau îi direcționează spre un site care preia complet informațiile postate pe alt site. Acestea sunt numite *scraper sites*, iar de cele mai multe ori sunt considerate pagini web legitime.

În această categorie intră următoarele tipuri de site-uri:

- site-uri care publică întregul conținut preluat de la un alt site, fără a-i adăuga nici o informație originală care să pună în valoare noua pagină;
- site-uri care copiază conținuturile altor site-uri și le postează într-o formă nouă, modificând automat anumite cuvinte din conținut (un exemplu în acest sens ar fi folosirea sinonimelor);
- site-uri care reproduc informațiile preluate din *feed*-urile *RSS* (*Rich Site Summary*) fără să adauge informații care sunt importante pentru vizitatorii site-ului
- site-uri care preiau fișiere multimedia (imagini, filme sau orice alt tip de media) și le prezintă utilizatorilor fără a oferi un plus de informație vizitatorilor site-ului

Așadar prin *scraper site* se înțelege un amalgam de conținuturi ce au fost preluate din alte surse, de cele mai multe ori fără permisiune. Astfel de site-uri web sunt în general pline de reclame, iar scopul lor este să se interpună între utilizator și site-ul care chiar dispune de informația pe care utilizatorul o caută. Astfel, *scraper*-ul își atribuie de fapt rolul motorului de căutare, scopul său fiind acela de a crește *pagerank*-ul *scraper*-ului în detrimentul site-ului original. Se dorește atingerea acestui scop prin prezentarea frecventă a unui conținut relevant și unic preluat de la site-uri sursă care sunt bine punctate de motoarele de căutare.

Rezultatele mixte obținute în *SERP* scad atât din performanța motorului de căutare cât și din mulțumirea utilizatorului referitoare la informația găsită. Dat fiind acest fapt, motoarele de căutare dezaproabă existența acestor tipuri de site-uri, tocmai din cauza faptului că acestea se interpun între user, motorul de căutare și site-ul destinație.

Identificarea site-urilor de acest tip s-a încercat să se facă utilizând diverși algoritmi, însă nu s-au obținut rezultatele dorite. În consecință politica pe care unele motoare de căutare au anunțat-o la începutul anului 2014 a fost de a identifica *scraper*-ele prin intermediul utilizatorilor și a *feedback*-ului pe care aceștia îl dau. Aceste site-uri odată catalogate nu sunt depunctate, ci clasificarea făcută de vizitatorii site-ului se utilizează în testarea algoritmilor prin intermediul cărora se detectează dacă un anumit site este *scraper* sau nu.

Totuși această dezaprobare nu este întotdeauna una fermă, ci mai mult una declarativă, dat fiind faptul că unele motoare de căutare permit postarea de reclame pe care le furnizează pe *scraper site*-uri, astfel aducând un câștig în plus atât site-ului (deoarece apare în SERP), cât și motorului de căutare (deoarece este un intermediar între cei care doresc să își facă reclamă prin *adlink*-uri și cei care doresc să publice pe site astfel de reclame).

Pentru a îmbunătăți calitatea rezultatelor oferite de motoarele de căutare, trebuie să se găsească o metodă de identificare a site-urilor care sunt *scraper*-e. Scopul identificării acestor site-uri este de a reduce, pe cât posibil, apariția lor în *SERP*. Astfel, se dorește în continuare studierea existenței unei metode prin intermediul căreia aceste site-uri să poată fi identificate și odată identificate acestea să fie penalizate de către motoarele de căutare astfel încât prezența lor în *SERP* să fie redusă cât mai mult posibil.

Datorită asistenței aproape permanente de care un utilizator are nevoie atunci când își dorește localizarea unor anumite informații pe web, calitatea acestor informații mai poate fi asigurată doar prin procesul de analizare a conținuturilor prezentate pe web.

Pornind de la această idee, se va considera ca o direcție de cercetare viitoare posibilitatea identificării unor astfel de pagini pe baza similarității dintre conținutul aflat la pagina care face parte dintr-un *scraper site* și pagina sursă de la care a fost preluat conținutul. De obicei, pentru a nu încălca anumite reguli de confidențialitate, *scraper site*-urile precizează pe fiecare pagină sursa de la care a preluat conținutul, ceea ce face ca acestea să fie foarte ușor recunoscute ”cu ochiul liber”.

Astfel, în procesul de determinare a unui *scraper site* se vor lua în considerare mai multe site-uri, urmând etapele:

- se vor compara conținuturile aflate în fiecare pagină din site-ul care se dorește a fi testat cu conținutul ce poate fi găsit urmând link-urile externe din pagina corespunzătoare,
- odată cu calcularea similarității dintre cele două conținuturi, se va observa că pentru un *scraper site* se va obține o similaritate mai mare în cazul în care link-ul extern se află în conținutul absolut și nu în template-ul paginii.

Pentru a studia fezabilitatea acestei idei am luat în considerare mai multe site-uri, care sunt actualizate zilnic. Cu ajutorul unui *crawler web* am indexat toate paginile site-urilor enumerate mai sus, după care, utilizând colecția de date obținută, am identificat în fiecare pagină conținutul absolut prezentat precum și sursa de la care a fost preluat conținutul.

Aplicând etapele algoritmului descris mai sus, au fost generate tripletele (*link intern*, *link sursa*, *similaritate*), pe baza cărora se va face identificarea *scraper site*-urilor.

Pentru a putea observa cât mai bine comportamentul unui astfel de site, pe baza tripletelor generate se vor reprezenta grafic punctele obținute într-un sistem cartezian după următorul model:

- pe axa absciselor se va reprezenta similaritatea dintre cele două link-uri; în calculul acesteia s-a folosit măsura de similaritate a Cosinusului [1].
- pe axa ordonatelor se va reprezenta numărul de apariții; prin acest număr de apariții se înțelege numărul de perechi (*link intern*, *link sursă*) care au o similaritate ce se află în intervalul de încredere al similarității ce reprezintă abscisa punctului considerat. Printr-un interval de încredere al unui număr se va înțelege un interval de forma $[s - \epsilon, s + \epsilon]$, unde s este numărul căruia i s-a construit intervalul de încredere și $\epsilon = 0.05$ este un prag determinat experimental.

Folosindu-ne de cele două coordonate ale unui punct de pe grafic, cu scopul de a evidenția numărul de perechi pentru care s-a obținut similaritate foarte mare între conținutul regăsit la link-ul intern și cel regăsit la link-ul sursă, pe reprezentarea grafică se regăsesc desenate cercuri pline cu centrul în punctul considerat și cu raza direct proporțională cu numărul de apariții.

În urma testelor efectuate se deduce faptul că pentru un *scraper site* graficul ce se obține este mult *shiftat* în dreapta (Fig:1.7(a)). Aceste rezultate pot fi dovedite și matematic, luând în considerare modul de construcție a site-urilor de acest tip (similaritatea a două pagini (destinație, sursă) este direct proporțională cu numărul de apariții).

Tocmai acesta este și motivul pentru care s-a ales ca în reprezentarea grafică să fie considerate și cercurile de raze diferite. Dacă s-ar fi renunțat la această reprezentare, atunci multe dintre punctele aflate pe grafic ar fi fost suprapuse, datorită diferențelor foarte mici dintre similaritățile corespunzătoare, iar efectul produs de *shiftare* la dreapta ar fi fost redus considerabil.

Totuși, efectuând aceleași teste asupra unui site care citează pentru fiecare articol sursa și care prezintă o informație modificată și adaptată tipului de utilizatori căruia i se adresează, se observă că rezultatele obținute sunt total diferite față de cele anterioare (Fig:1.7(b)). Graficul obținut în cel de-al doilea caz indică un comportament normal și chiar urmează curba lui Gauss, majoritatea punctelor de pe grafic aglomerându-se în centrul acestuia.

O analiză interesantă ar fi calcularea similarităților pentru mai multe astfel de site-uri și sursele lor, observând astfel dacă rezultatul ce se va obține se va putea generaliza pentru orice fel de site, iar în urma analizei rezultatului se va putea decide categoria în care acesta poate fi plasat.

1.4 Concluzii și direcții de cercetare

Prezentarea direcțiilor viitoare de cercetare începe cu prezentarea unor idei, care vin ca o continuare naturală a aspectelor prezentate în secțiunile anterioare.

O eventuală completare ce ar putea fi adusă *layer*-ului de migrare transparentă a conținuturilor poate urmări mai multe aspecte:

- evaluarea diferitelor funcții de similaritate, cu scopul de a descrește timpul de regăsire a perechilor de URL-uri care prezintă conținuturi similare;

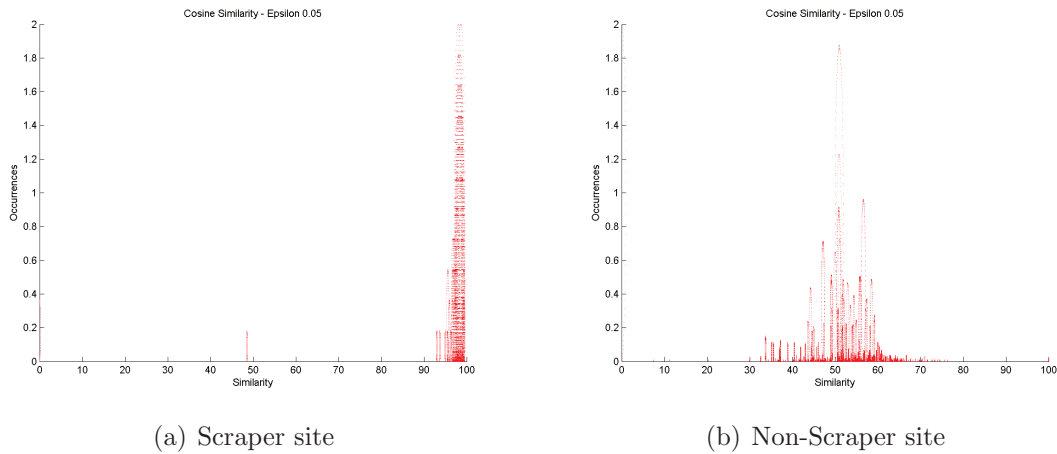


Figura 1.7: Similaritatea Cosinus - $\epsilon = 0.05$

- ponderarea unor anumite proprietăți ale conținutului prezentat la un anumit URL (cum ar fi: URL, *heading*-urile din pagină web: $h1, h2, \dots, h6$, titlul paginii web, cuvintele cheie corespunzătoare care provin în urma accesării unui link din rezultatele afișate în SERP);
- redirectionarea vizitatorului spre o pagină similară cu cea a *referrer*-ului.

În plu, este posibil să se obțină rezultate mai bune dacă se va considera analiza similarității conținuturilor de la *referrer* cu conținutul prezentat la orice link intern găsit în pagina accesată.

De asemenea, o alte direcție de cercetare ar fi: utilizarea funcțiilor de similaritate care să interpreteze asemănarea dintre conținuturi din punct de vedere semantic, ponderarea funcțiilor de similaritate sau alegerea unei funcții de similaritate și ponderarea unor proprietăți specifice ale conținutului.

Capitolul 2

Interpretarea logurilor unei platforme de e-learning folosind Analiza Conceptuală Formală

2.1 Analiza Conceptuală Formală

Analiza Conceptuală Formală (FCA) a fost dezvoltată constant pe parcursul ultimilor 30 de ani ([43]), pornind de la restructurarea teoriei laticilor. FCA cuprinde nu numai importante aspecte teoretice despre latici, despre proprietățile lor și despre structuri asemănătoare, ci și algoritmi cu aplicații în *knowledge discovery* și *knowledge representation*. De asemenea de-a lungul timpului au apărut și diverse extinderi ale acestei teorii, precum: Fuzzy FCA, FCA Temporal, FCA Triadic, FCA relațional.

În cazul datelor tridimensionale este necesară o abordare triadica.

Analiza Conceptuală Formală este un instrument matematic prin intermediul căruia se pot procesa conceptele.

Structura de bază este un *context formal*, din care se pot extrage și procesa cunoștințe folosind o conexiune Galois, numită operator de derivare.

Definiție 2.1.1 Un *context formal* $\mathbb{K} := (G, M, I)$ este format din două mulțimi G și M și o relație binară I între G și M . Elementele din G se numesc **obiecte** și elementele din M se numesc **atribute**. Relația I se numește relație de incidență a contextului formal. Uneori se folosește notația: gIm în locul $(g, m) \in I$, pentru a exprima că obiectul g are atributul m .

Definiție 2.1.2 Fie $\mathbb{K} := (G, M, I)$ un context formal. Pentru mulțimea de obiecte $A \subseteq G$ se definește:

$$A' := \{m \in M \mid gIm \forall g \in A\}$$

mulțimea tuturor atributelor comune obiectelor din A . Analog, pentru mulțimea atributelor $B \subseteq M$ se definește:

$$B' := \{g \in G \mid gIm \forall m \in B\}$$

mulțimea tuturor obiectelor comune atributelor din B .

Conceptele sunt considerate unitatea de bază din care se pot desprinde cunoștințe, iar ele sunt extrase din contextul formal utilizând operatorii de derivare.

Definiție 2.1.3 *Un concept formal este definit ca o pereche (A, B) , unde $A \subseteq G$, $B \subseteq M$ și $A' = B, B' = A$. Mulțimea A se numește **extent** și conține toate obiectele care au legătură cu conceptul, iar B se numește **intent**, și reprezintă mulțimea tuturor atributelor comune obiectelor din A .*

Mulțimea tuturor conceptelor dintr-un context dat (G, M, I) se notează $\mathfrak{B}(G, M, I)$. Se poate defini o ierarhie a conceptelor pe $\mathfrak{B}(G, M, I)$ folosind relația dintre subconcept-superconcept.

Definiție 2.1.4 *Fie (G, M, I) un context. Relația subconcept-superconcept pe $\mathfrak{B}(G, M, I)$ este definită astfel: $(A, B) \leq (C, D)$ dacă și numai dacă $A \subseteq C$ (sau echivalent, $D \subseteq B$), unde $(A, B), (C, D)$ sunt două concepte formale. Conceptul (A, B) se numește subconcept (sau specializare) pentru (C, D) , în timp ce (C, D) se numește superconcept (sau generalizare) pentru (A, B) .*

Relația subconcept-superconcept este o relație de ordine pe $\mathfrak{B}(G, M, I)$. În plus, $\mathfrak{B}(G, M, I)$ este o latice completă, numită **ierarhie de concepte**.

În astfel de diagrame de ordine, fiecare concept poate fi reprezentat într-un nod, iar relația de tip subconcept-superconcept este reprezentată sub forma unor muchii ce interconectează astfel de noduri.

Analiza Conceptuală Formală Triadică a fost introdusă în [29, 44].

Definiție 2.1.5 *Un context triadic (sau: tricontext) este un cvadruplu (K_1, K_2, K_3, Y) , unde K_1, K_2 și K_3 sunt mulțimi și $Y \subseteq K_1 \times K_2 \times K_3$ este o relație ternară între ele. Elementele din K_1, K_2, K_3 sunt numite obiecte, atribute și respectiv condiții. Expresia $(g, m, b) \in Y$ se citește obiectul g are atributul m luând în considerare condiția b .*

Următoarea definiție arată cum contextele diadice pot fi obținute dintr-un context triadic.

Definiție 2.1.6 (Contexte derivate) *Pentru orice context triadic (K_1, K_2, K_3, Y) se pot obține următoarele contexte diadice prin proiectare pe una din componente:*

$$\mathbb{K}^{(1)} := (K_1, K_2 \times K_3, Y^{(1)}) \text{ cu } gY^{(1)}(m, b) := \Leftrightarrow (g, m, b) \in Y,$$

$$\mathbb{K}^{(2)} := (K_2, K_1 \times K_3, Y^{(2)}) \text{ cu } mY^{(2)}(g, b) := \Leftrightarrow (g, m, b) \in Y,$$

$$\mathbb{K}^{(3)} := (K_3, K_1 \times K_2, Y^{(3)}) \text{ cu } bY^{(3)}(g, m) := \Leftrightarrow (g, m, b) \in Y.$$

Pentru $\{i, j, k\} = \{1, 2, 3\}$ și $A_k \subseteq K_k$, se definesc $\mathbb{K}_{A_k}^{(ij)} := (K_i, K_j, Y_{A_k}^{(ij)})$, unde $(a_i, a_j) \in Y_{A_k}^{(ij)}$ dacă și numai dacă $(a_i, a_j, a_k) \in Y$ pentru orice $a_k \in A_k$.

Operatorii de derivare în cazul triadic se definesc folosind operatorii de derivare diadici pentru contextele diadice obținute după proiectare.

Definiție 2.1.7 (Operatorii de derivare - (i)) *Pentru $\{i, j, k\} = \{1, 2, 3\}$ cu $j < k$ și pentru $X \subseteq K_i$ și $Z \subseteq K_j \times K_k$ Operatorii de derivare-(i) se definesc prin:*

$$X \mapsto X^{(i)} := \{(a_j, a_k) \in K_j \times K_k \mid (a_i, a_j, a_k) \in Y \text{ for all } a_i \in X\}.$$

$$Z \mapsto Z^{(i)} := \{a_i \in K_i \mid (a_i, a_j, a_k) \in Y \text{ for all } (a_j, a_k) \in Z\}.$$

Bineînțeles, acești operatori de derivare corespund cu operatorii de derivare din contextul diadic $\mathbb{K}^{(i)}, i \in \{1, 2, 3\}$.

Definiție 2.1.8 (Operatorii de derivare - (i, j, X_k)) Pentru $\{i, j, k\} = \{1, 2, 3\}$ și $X_i \subseteq K_i, X_j \subseteq K_j, X_k \subseteq K_k$, Operatorii de derivare (i, j, X_k) se definesc prin:

$$X_i \mapsto X_i^{(i,j,X_k)} := \{a_j \in K_j \mid (a_i, a_j, a_k) \in Y \text{ for all } (a_i, a_k) \in X_i \times X_k\}$$

$$X_j \mapsto X_j^{(i,j,X_k)} := \{a_i \in K_i \mid (a_i, a_j, a_k) \in Y \text{ for all } (a_j, a_k) \in X_i \times X_k\}.$$

Operatorii de derivare - (i, j, X_k) corespund operatorilor de derivare din contextul diadic: $(K_i, K_j, Y_{X_k}^{(ij)})$.

Similar noțiunii de cocept formal din contextul diadic, se pot introduce conceptele triadice.

Definiție 2.1.9 Un concept triadic (sau: **triconcept**) al $\mathbb{K} := (K_1, K_2, K_3, Y)$ este un triplet (A_1, A_2, A_3) cu $A_i \subseteq K_i$ pentru $i \in \{1, 2, 3\}$ și $A_i = (A_j \times A_k)^{(i)}$ pentru orice $\{i, j, k\} = \{1, 2, 3\}$ cu $j < k$. Mulțimile A_1, A_2 , și A_3 se numesc **extent**, **intent**, și **modus** al conceptului triadic. Mulțimea $\mathfrak{T}(\mathbb{K})$ reprezintă mulțimea tuturor triconceptelor mulțimii \mathbb{K} .

Propoziție 2.1.1 Triconceptele contextului triadic (K_1, K_2, K_3, Y) sunt tripletele maxime $(A_1, A_2, A_3) \in \mathfrak{P}(K_1) \times \mathfrak{P}(K_2) \times \mathfrak{P}(K_3)$ cu $A_1 \times A_2 \times A_3 \subseteq Y$, în raport cu relația de incluziune.

2.2 Analiza comportamentului utilizatorilor web folosind Analiza Conceptuală Formală

Navigarea pe Internet are din ce în ce mai mult o dimensiune socială și a devenit în ultimul timp un mecanism eficient de dobândire de cunoștințe. Așadar se impune necesitatea îmbunătățirii atât a conținutului prezentat utilizatorilor unui site web, cât și a design-ului acestuia în special datorită faptului că au fost dezvoltate din ce în ce mai multe instrumente de analiză a comportamentului online al utilizatorilor.

Cercetările efectuate în ultimii ani relevă modul în care utilizatorii interacționează pe Internet, impactul acestuia asupra vieții de zi cu zi și asupra culturii [41],[22],[7].

O atenție specială trebuie acordată platformelor educaționale. Acestea trebuie să vină cu o structură clară și cu un conținut organizat.

Pentru a determina comportamentul utilizatorilor față de astfel de platforme online, sunt necesare colectarea de date adiționale prin intermediul logurilor serverelor web.

Metodele folosite până acum care implică extragerea de informații utile din aceste loguri web sunt metodele statistice sau tehnicile de *data mining*.

Instrumentele de analiză a logurilor web utilizează frecvent câteva metrici de bază. Însă acestea oferă o perspectivă reală în special site-urilor comerciale. Pentru site-urile educaționale ele nu sunt eficiente, în special datorită faptului că procesul de învățare poate dura mai mult timp, iar în acest caz o vizită pe un site educațional nu se supune euristicilor folosite de majoritatea instrumentelor de analiză [16].

Analiza pe care ne-am propus să o facem folosește datele colectate de o platformă educațională numită PULSE [15], în semestrul II al anului universitar 2012-2013. Pentru perioada considerată, au existat 40768 accesări pe platformă. Informațiile colectate conțin informații despre: *URL*-ul complet al paginii accesate, *URL*-ul *referrer*-ului, *cookieID*, *login*, *timestamp*-ul accesării, adresa IP, *User Agent*, *Screen Resolution*. Informațiile înregistrate sunt referitoare la două cursuri: Sisteme de Operare (SO1 -

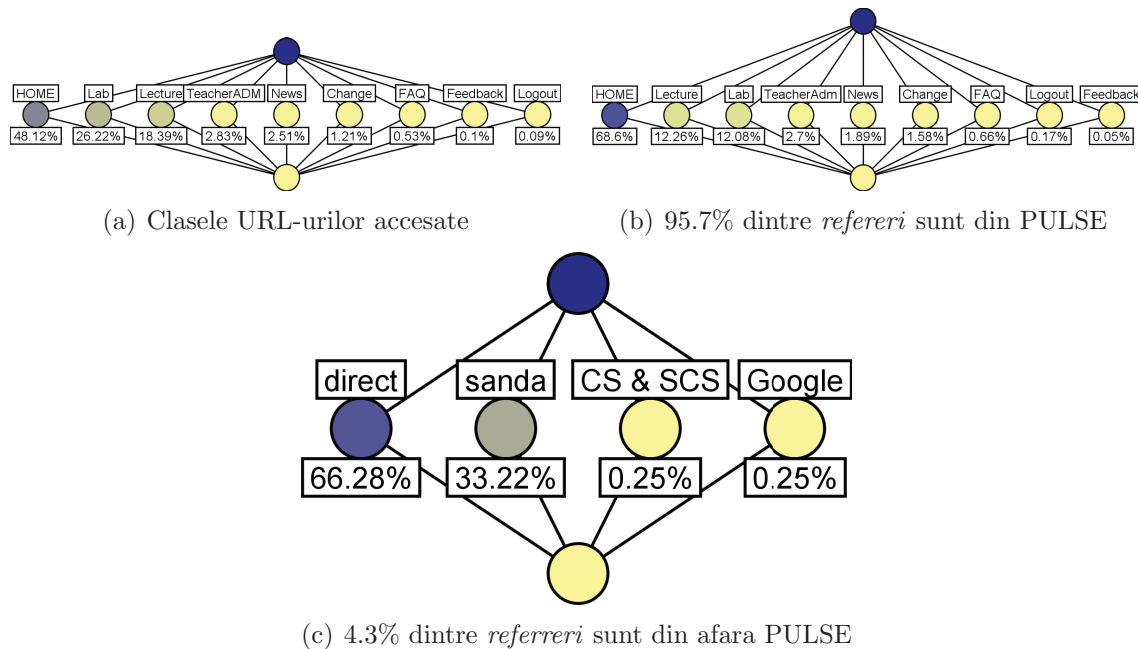


Figura 2.1: Scale nominale folosind ToscanaJ

curs obligatoriu) și Proiectare Web și Optimizare (WDO - curs opțional). Grupele de studenți care au participat la cursul SO1 sunt notate "ar" și "ri", iar grupele de studenți care au participat la cursul WDO sunt notate "ei" și "ie".

Am ales să nu folosim tehnicile clasice de data mining care ne asigură o perspectivă cantitativă, ci să avem o perspectivă calitativă prin utilizarea tehnicilor impuse de Analiza Conceptuală Formală, prin intermediul unui instrument numit ToscanaJ.

ToscanaJ [5] este un instrument utilizat îndeosebi pentru reprezentarea informațiilor dobândite prin analiza datelor, și oferă importante perspective asupra conexiunilor dintre datele analizate.

Pentru analizele efectuate am luat în considerare doar trei câmpuri din baza de date: *URL*-ul paginii accesate, *referrer*-ul și *cookieID*-ul. Deoarece datele ce urmează a fi analizate sunt destul de consistente, este necesară pregătirea unei etape de preprocesare a acestora (există 751 de *URL*-uri, 471 de *referreri* și 3472 *cookieID*-uri distincte).

Astfel, *URL*-urile accesate au fost clasificate în 9 clase disjuncte, care au fost vizualizate în ToscanaJ folosind o scală nominală (Figura 2.1(a)).

PULSE este o platformă care a fost construită în special pentru a fi folosită în timpul laboratoarelor, cu scopul de a consulta suportul teoretic oferit de către profesor, de a vizualiza problemele alocate sau de a vizualiza notele și prezențele. Fiecare vizitator PULSE se poate identifica în mod unic, deoarece este necesară trecerea printr-o etapă de autentificare înainte de a exista posibilitatea vizualizării oricărei informații enumerate anterior. Există două tipuri de utilizatori: student și profesor.

După autentificare, utilizatorul de tip student intră pe o pagină generală (din categoria *HOME*) care oferă informații generale despre: prezențe, problemele de laborator alocate, suporturile teoretice aferente, note și anunțuri curente.

Conform figurii 2.1(a) procentul celor mai multe accesări este atribuit clasei *HOME*. Acest lucru este de așteptat dat fiind faptul că această clasă este un liant între toate celelalte clase. Următoarele clase frecvent vizitate sunt *LAB* și *LECTURE*, care conțin materiale referitoare la laboratoare și cursuri, exemple, lucrări, rezultate la lucrări,

precum și soluțiile propuse pentru rezolvarea acestora.

Utilizatorului de tip profesor îi sunt atribuite drepturi administrative: atribuirea de probleme, acordarea de note, înregistrarea prezențelor, postarea materialelor necesare, adăugarea de noutăți. Toate aceste pagini intră în clasa TeacherADM.

Alte clase vizitate sunt NEWS (pagina pe care sunt publicate toate anunțurile), FAQ (pagina care conține cele mai frecvente întrebări legate de platformă și răspunsurile aferente), CHANGE (paginile cu materialele puse la dispoziția studenților de la alte cursuri sau din anii anteriori) și respectiv LOGOUT (pagina prin intermediul căreia se închide o sesiune pe PULSE).

Analog clasificării URL-urilor, s-a făcut o clasificare a *refererilor*, care reprezintă site-ul sau pagina de pe care au ajuns vizitatorii pe una din paginile din PULSE. *Refererii* pot fi atât pagini din PULSE cât și pagini din afara PULSE.

Folosind din nou ToscanaJ au fost vizualizate scalele nominale ale referer-ilor din PULSE și a celor din afara platformei (vezi figurile 2.1(b), 2.1(c)).

Refererii interni au fost clasificați în aceleași clase ca și URL-urile accesate. *Refererii* externi simbolizează fie accesările directe ale utilizatorilor (fie prin scrierea URL-ului, fie prin bookmark), fie accesările de pe rețele de socializare și motoare de căutare (cum ar fi facebook, google), fie de pe site-ul facultății și site-urile personale ale profesorilor.

2.3 Vizualizarea datelor triadice

Extensia Toscana2TRIAS permite selectarea datelor triadice, pornind de la o mulțime de scale preprocesate în ToscanaJ. Pentru a obține un set de date triadice, am ales ca și attribute clasele de *referreri*, ca și condiții clasele de URL-uri accesate, iar ca și obiecte perechea (login, IP). Astfel am generat toate triconceptele folosind TRIAS.

Problema vizualizării datelor triadice nu a fost încă rezolvată satisfăcător, singurele opțiuni disponibile fiind trilaticile sau graf-urile. Din aceste motive am ales să folosim CIRCOS, un instrument nou de vizualizare a datelor, care a fost construit tocmai pentru a ajuta la investigarea tiparelor în date.

CIRCOS este un instrument prin intermediul căruia datele și conexiunile dintre ele pot fi vizualizate într-un format circular [10]. Datele de intrare pentru CIRCOS au fost obținute din tricontext utilizând operatorii de derivare. A fost necesară implementarea unui algoritm care transformă XML-ul rezultat din TRIAS în formatul valid de intrare acceptat de CIRCOS.

XML-ul rezultat din TRIAS conține toate triconceptele care pot fi derivate din tricontextul definit. Fiecare triconcept este definit de către *extent*, *intent* și *modus*. Formatul datelor de intrare acceptat de CIRCOS este un tabel bidimensional $R \times C$, care conține valori numerice. Acesta a fost obținut după cum urmează:

Pornind de la tricontextul (G, M, B, Y) , am construit proiecția diadică $\mathbb{K}_{32} := (G, (B, M), I)$, unde $(g, (b, m)) \in I \Leftrightarrow (g, m, b) \in Y$. După aceea, pentru fiecare pereche (b, m) am evaluat conceptul corespunzător $\mu_{\mathbb{K}_{32}}$ și am determinat cardinalitatea *extent*-ului $(b, m)'$.

Mulțimea care definește coloanele tabelului, notată C , este mulțimea obținută prin proiectarea relației de incidență Y pe M , $\text{pr}_M(Y) := \{m \in M \mid \exists (g, b) \in G \times C. (g, m, b) \in Y\}$. Analog, se obține mulțimea care definește liniile tabelului, notată R , și se obține prin proiectarea relației de incidență Y pe mulțimea condițiilor B .

Algoritmul implementat, construiește un tabel care are ca indicatori de coloană și de linie mulțimile construite anterior, și calculează valorile numerice din tabel în modul în care urmează. Pentru fiecare pereche $(c, r) \in C \times R$, numărul de elemente din *extent* $(c, r)' \in \mathbb{K}_{32}$ este calculat direct din XML-ul rezultat din TRIAS. Cardinalul acestei mulțimi reprezintă valoarea numerică ce se va regăsi în tabel la intersecția coloanei c cu linia r .

Exemplu de XML generat din TRIAS

```

<Triconcepts>
  <Triconcept id="1">
    <Extent>
      <Object>e1</Object>
      <Object>e2</Object>
    </Extent>
    <Intent>
      <Attribute>i1</Attribute>
      <Attribute>i2</Attribute>
    </Intent>
    <Modus>
      <Condition>m1</Condition>
    </Modus>
  </Triconcept>
  <Triconcept id="2">
    <Extent>
      <Object>e1</Object>
    </Extent>
    <Intent>
      <Attribute>i2</Attribute>
    </Intent>
    <Modus>
      <Condition>m2</Condition>
      <Condition>m3</Condition>
    </Modus>
  </Triconcept>
</Triconcepts>

```

Exemplu de date de intrare formate pentru CIRCOS

```

-  i1 i2
m1 2  2
m2 0  1
m3 0  1

```

Figura 2.2 prezintă un exemplu de format circular generat de CIRCOS folosind o mulțime de triconcepte generate de TRIAS având ca și obiecte: *cookieID*, ca și atribute: clasele de *referreri* și condiții: clasele de URL-uri accesate. Deoarece clasele de *referreri* (R_CLASS) și clasele de URL-uri accesate (AF_CLASS) au elemente în comun, mulțimile C și R nu sunt disjuncte. În figura 2.2 segmentele sunt reprezentate circular, în sensul acelor de ceasornic, fiind ordonate descrescător, începând cu segmentul verde "HOME", segmentul mov "LECTURE", segmentul albastru "LAB" și segmentele "NEWS", "CHANGE", "FAQ", "ADM", "FEEDBACK" and "LOGOUT". Reprezentarea circulară obținută poate fi interpretată conform explicațiilor din Figure 2.2.

Fiecare bandă corespunde unei perechi (R_CLASS, AF_CLASS). Culoarea benzii este dată de culoarea segmentului *referrer*-ului. De exemplu banda verde de la "HOME" la "LECTURE" corespunde tuturor paginilor accesate din clasa "LECTURE" de la *referrer* din clasa "HOME". Pentru o înțelegere mai bună a conexiunilor, am reprezentat cu ajutorul unui graf orientat (Figura 2.3) conexiunile dintre clasele de *referrer* și de URL-uri accesate. Nodurile graf-ului au culorile corespunzătoare benzilor din reprezentarea grafică rezultată din CIRCOS. Deoarece mulțimile claselor de *Referreri* și de URL-uri accesate nu sunt disjuncte, vor exista atât muchii orientate în ambele direcții (înspre și dinspre același nod), precum și cicluri.

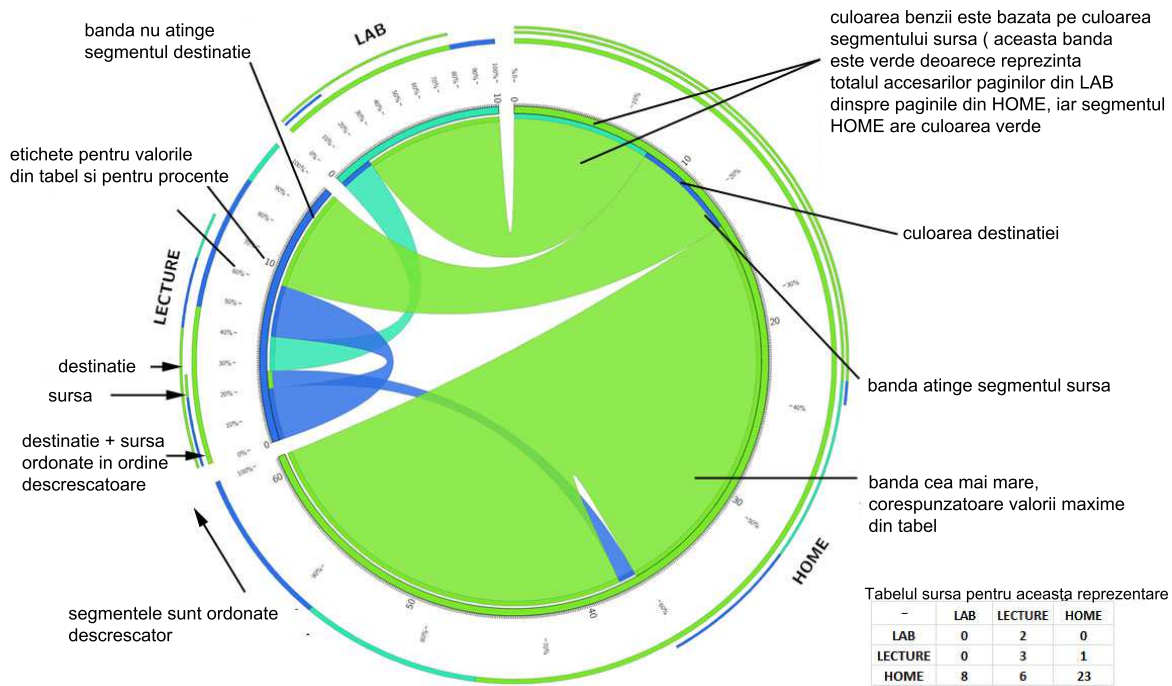


Figura 2.2: Objects Set: CookieId, Attribute Set: Referrer Class, Conditions Set: Access File Class

Luând în considerare datele logate de PULSE, am cautat diferite structuri triadice în acestea, pentru a putea fi interpretate folosind Toscana2TRIAS. Dintre toate testele efectuate, rezultate reprezentative pentru tiparuri în traiectoria studenților în navigarea lor prin site [12].

Testele au fost efectuate luând ca și obiecte: *cookieID*, atribute: *R_class* și condiții *AF_class*. Am considerat doar *referrerii* care sunt din interiorul PULSE (Figura 2.4(a)).

Se observă că cele mai multe accesări sunt făcute din clasa HOME spre clasa HOME. Acest lucru se explică prin tranzițiile de la pagina de login la pagina principală, prin selectarea anului de studiu și a cursului pentru care utilizatorul dorește să vizualizeze informațiile sau reîncărcarea paginii.

După ce sunt urmați acești pași, cele mai multe vizite sunt de la HOME la paginile din clasele LAB, LECTURE sau NEWS. Aceste rezultate dovedesc scopul principal al platformei educaționale PULSE și anume deservirea de suport teoretic pentru laboratoare, cursuri și postarea de anunțuri.

Accesarea paginilor din clasa LECTURE semnifică faptul că utilizatorii au vizitat atât slide-urile de la curs cât și informațiile despre lucrările care s-au dat pe parcursul

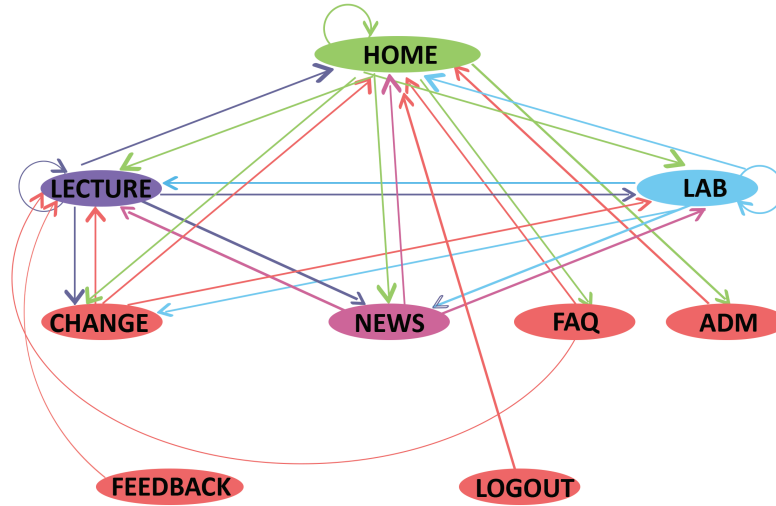
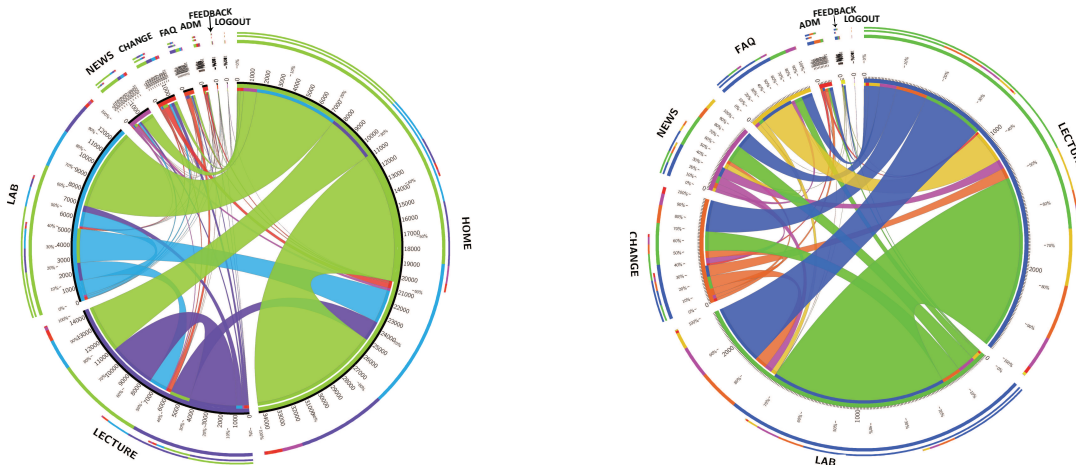


Figura 2.3: Vizualizarea cu ajutorul graf-ului orientat a conexiunilor dintre clasele de *referreri* și URL-uri accesate

semestrului sau explicații detaliate despre rezolvarea lor. Aceste tiparuri în navigare în aceeași clasă sunt naturale datorită modului de construcție al claselor (fiecare clasă conține mai multe pagini). Mai mult decât atât aceste tiparuri se regăsesc în toate clasele existente.

Cu scopul de a observa și alte pattern-uri am eliminat din datele de intrare pentru CIRCOS astfel de legături. De asemenea am eliminat toate interacțiunile dintre HOME și alte clase (Figurile 2.4(a), 2.4(b)).



(a) Interpretarea trend-urilor triadice deduse din triconcepte, considerând R_CLASS, AF_CLASS și cookieID

(b) O nouă interpretare a Figurii 2.4(a) după eliminarea interacțiunilor

Figura 2.4: Interpretarea trend-urilor triadice deduse din triconcepte

După eliminarea elementelor de mai sus, se pot observa și alte tipare în navigarea

studenților prin platformă: activitate intensă de la paginile cu laboratoare la paginile cu cursuri, sau reciproc de la cursuri la laboratoare, utilizarea facilității de CHANGE, care e cel mai adesea accesată de către studenții din paginile din clasele LAB și LECTURE, și reciproc, vizitarea paginilor din clasele LAB și LECTURE din paginile din clasa FAQ.

Alte detalii sunt dificil de observat luând în considerare aceste reprezentări. Așadar, pentru a vedea trend-urile cele mai importante între clasele principale din PULSE, am eliminat gradual valorile (numărul de vizite) mai mici decât 50 și respectiv 100, considerând că nu sunt foarte reprezentative dat fiind faptul că valoarea maximă este mai mare decât 1000. Rezultatele se pot observa în Figura 2.5.

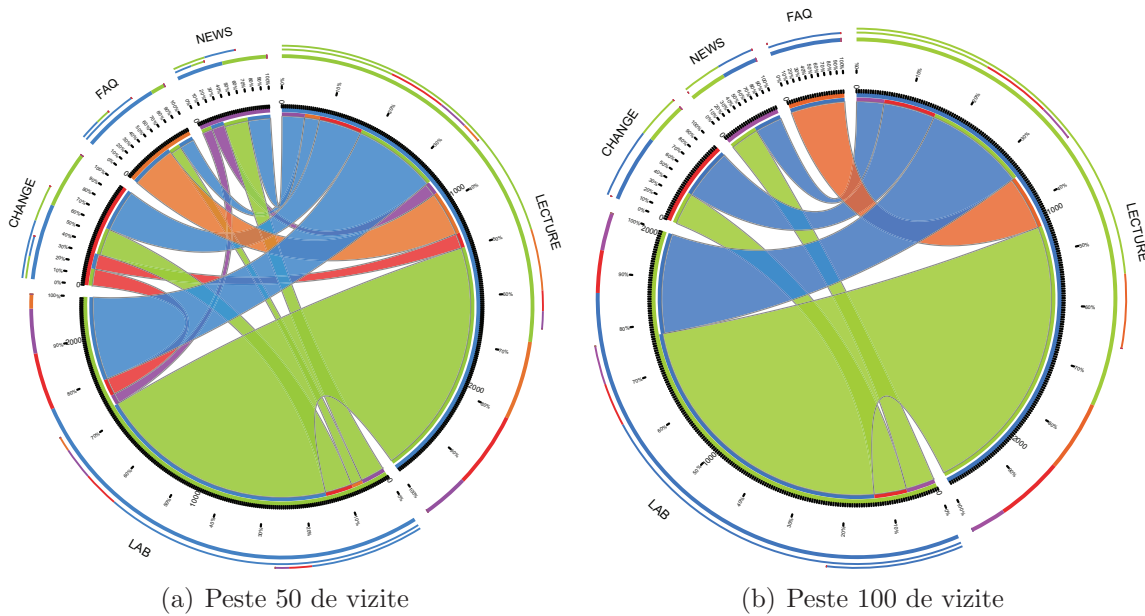


Figura 2.5: O nouă interpretare a Figurii 2.4(b) fără a lua în considerare clasele cu prea puține vizite

Principalele clase din PULSE sunt: LECTURE, LAB, NEWS, CHANGE și FAQ. Figura 2.5(a) arată că sunt mai mult de 50 de vizite dinspre LECTURE către orice altă clasă. Același lucru se poate spune și despre LAB. Pentru CHANGE și NEWS, interacțiunile sunt bidirecționale, dar doar cu clasele LAB și LECTURE, în timp ce înspre FAQ vizitatorii pleacă atât de la LAB cât și de la LECTURE, dar se întorc la FAQ doar de la LECTURE.

Pentru interacțiunile cu cel puțin 100 de vizite, direcția de vizitare la la LAB sau LECTURE este unidirecțională către CHANGE și NEWS.

Numărul mare de accesări de la FAQ la LECTURE a fost surprinzător. Investigând motivul pentru care studenții urmează acest tipar am ajuns la concluzia că pagini din cele două clase sunt consecutiv amplasate în meniul de navigare, prin urmare există posibilitatea ca utilizatorii care au ajuns pe pagina de FAQ de multe ori să fi ajuns din greșeală. Așadar se impune o modificare a design-ului paginii în acest sens.

Pentru a continua abordarea triadică am luat în considerare R_CLASS ca și obiecte, AF_CLASS ca atribute și timestamp ca și condiții. Pentru a investiga comportamentul temporal al studenților pe parcursul unui semestru am analizat datele pe intervale de timp [12].

Primele intervale de timp considerate au fost egale cu aproximativ o treime din

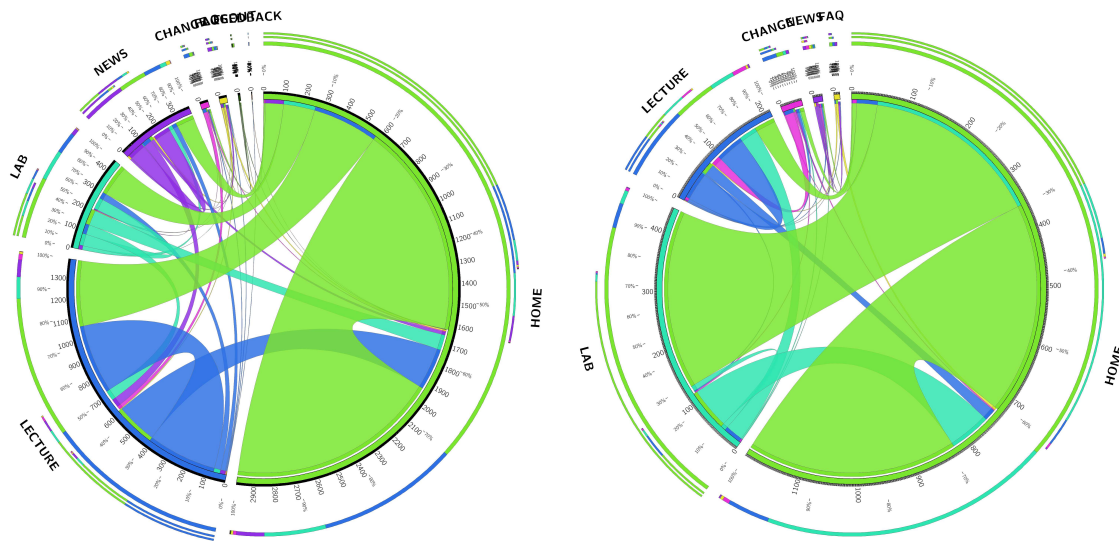
semestru, și au fost denumite *început*, *mijloc* și *sfârșit*. Astfel de intervale au conținut un volum prea mare de date, și astfel nu au generat trend-uri semnificative în datele considerate.

Așadar am considerat intervalele de timp săptămânile din semestru. Lista completă a rezultatelor obținute se regăsește pe <http://www.cs.ubbcluj.ro/~fca/tests-2013>. În această analiză observăm trei tendințe de comportament: relaxat, intens și normal.

Comportamentul **relaxat** apare în principal în timpul vacanței (de exemplu săptămâna 10). Tiparul urmat poate fi vizualizat în Figura 2.2 și poate fi distins datorită numărului mic de URL-uri accesate (de obicei doar URL-urile din clasele HOME, LAB, LECTURE). Traiectoriile de navigare prin site observate pe parcursul acestei săptămâni sunt simple: vizitarea paginilor dinspre HOME spre LAB sau LECTURE; dinspre LAB spre LECTURE, dinspre LECTURE spre HOME. Pentru cursul opțional (WDO) rezultatele evidențiază traiectorii de navigare mult mai relaxate datorită faptului că acest curs implică cercetare. Așadar, materialul de curs nu este la fel de mult vizitat ca în cazul cursului obligatoriu (SO1). Acest tip de comportament apare și după examenele finale (săptămânile 18 și 20 pentru grupa "ar", 18 și 19 pentru grupa "ri" și săptămâna 14 pentru grupele "ei" și "ie").

Comportamentul **intens** apare în timpul perioadei de examinare (săptămânile 17, 19 și 20 pentru grupele "ar" și "ri" și săptămânile 7, 9 și 13 pentru grupele "ei" și "ie"). Traiectoriile utilizatorilor pot fi desprinse vizualizând Figura 2.6(a) care prezintă un număr ridicat de accesări. Aceste tiparuri pot apărea și în săptămânile de dinaintea examinării.

Comportamentul **normal** apare pe parcursul semestrului, atunci când nu este perioadă de examinare sau vacanță. Tiparul observat este evidențiat în Figura 2.6(b), adică pagini din aproape toate clasele au fost vizitate.



(a) Săptămâna 17 - intens

(b) Săptămâna 5 - normal

Figura 2.6: Compararea comportamentului grupei "ar": intens versus normal

Deși în Figura 2.6 cele două comportamente par similare, există diferențe semnificative, și anume faptul că în perioada cu comportament intens paginile din clasa LECTURE sunt cele mai vizitate, dat fiind faptul că studenții se pregătesc pentru examene, iar în perioada cu comportament normal paginile din clasa LAB sunt cele

mai vizitate, dat fiind faptul că studenții trebuie să își rezolve probleme de laborator pe care le-am primit. O altă diferență ar fi chiar numărul de accesări al paginilor din clasa HOME. În perioada cu comportament intens, clasa HOME pare a fi un liant între celelalte facilități oferite de PULSE.

Așadar, informațiile oferite de această interpretare, bazată pe triconcepte sunt calitative, spre deosebire de diverse instrumente statistice (cum ar fi histogramele) care oferă informații cantitative.

2.4 Concluzii și direcții de cercetare

Ca direcție de cercetare viitoare ne propunem ajustarea parametrilor în fișierele de configurare ale CIRCOS: ajustarea ordinii de afișare a benzilor, ajustarea transparenței benzilor în funcție de distribuția valorilor în datele de intrare, ascunderea sau ștergerea unor benzi care nu corespund unor condiții date, scalarea valorilor din celule.

În plus, plănuim să alegem și alte formate de reprezentare grafică în afară de cel circular.

O altă problemă pe care am vrea să o abordăm în cercetările viitoare este numărul mare de concepte generate de Trias. Pentru a rezolva această problemă propunem construirea unui algoritm bazat pe probabilități condiționale pentru a șterge conceptele cele mai puțin importante.

De asemenea, componenta temporală este foarte importantă în log-urile web. Ne propunem să folosim Analiza Conceptuală Formală Temporală pentru a studia comportamentul utilizatorilor, atât individual cât și pe grupuri.

Listă de figuri

1.1	Pagini ce au generat eroare 404 și au fost accesate de către un <i>referrer</i> extern	9
1.2	Pagini ce au generat eroare 404 și au fost accesate de către un <i>referrer</i> extern	10
1.3	Adaptarea rapidă a motoarelor de căutare la noua structură a site-ului	11
1.4	Adaptarea rapidă a motoarelor de căutare la noua structură a site-ului	12
1.5	Similaritate ideală	15
1.6	Similaritati	16
1.7	Similaritatea Cosinus - $\epsilon = 0.05$	20
2.1	Scale nominale folosind ToscanaJ	24
2.2	Objects Set: CookieId, Attribute Set: Referrer Class, Conditions Set: Access File Class	27
2.3	Vizualizarea cu ajutorului graf-ului orientat a conexiunilor dintre clasele de <i>referrer</i> și URL-uri accesate	28
2.4	Interpretarea trend-urilor triadice deduse din triconcepte	28
2.5	O noua interpretare a Figurii 2.4(b) fără a lua în considerare clasele cu prea puține vizite	29
2.6	Compararea comportamentului grupei ”ar”: intens versus normal . . .	30

Bibliografie

- [1] Apache Lucene, A high-performance, full-featured text search engine library, <http://lucene.apache.org/>
- [2] Apache HTTP Server, http://httpd.apache.org/docs/current/mod/mod_alias.html
- [3] L. Becchetti, C. Castillo, D. Donato, *Link-Based Characterization and Detection of web Spam*, 2nd International Workshop on Adversarial Information Retrieval on the web, AIRweb, Seattle, USA, August 2006, pp. 1-8
- [4] L. Becchetti, C. Castillo, D. Donato, S. Leonardi, R. Baeza-Yates, *Link Analysis for web Spam Detection: Link-based and Content Based Techniques*, ACM Transactions on the web (Tweb), Volume 2, Issue 1, New York, USA, February 2008, pp. 1-41
- [5] P. Becker, J. Hereth, G. Stumme, *ToscanaJ: An Open Source Tool for Qualitative Data Analysis*, *Advances in Formal Concept Analysis for Knowledge Discovery in Databases.*, page 1-2. Lyon, France, (July 2002).
- [6] B. Berelson, *Content analysis in communication research*, New York, Free Press, 1952
- [7] T. Berners-Lee , W. Hall, J. Hendler, N. Shadbolt, D.J. Weitzner, *Science* 313, 769 (2006).
- [8] G. Beydoun, R. Kultchitsky, G. Manasseh, *Evolving semantic web with social navigation*, *Expert Systems with Applications*, 32(2007), pp. 265–276.
- [9] D. Bufnea, **D. Haliță**, *A server-side support layer for transparent web content migration*, Proceedings of the International Conference on Knowledge Engineering, Principles and Techniques, KEPT 2013 Cluj-Napoca, Studia Universitatis Babeş-Bolyai Informatica, 3/2013, pp. 78-89
- [10] Circos, a circular visualization tool, www.circos.ca
- [11] A. Dieberger, *Supporting social navigation on the world wide web.*, International Journal of HumanComputer Studies, 46(6), 805825, 1997.
- [12] Sanda Dragoş, Diana Haliță, Christian Săcărea, Diana Troancă, *An FCA grounded study of user dynamics through log exploration*, Studia Universitatis Babeş-Bolyai Series Informatica, Volume LIX, no. 2, 2014, pp. 82-97

- [13] S. Dragoş, C. Săcărea, *Analysing the Usage of Pulse Portal with Formal Concept Analysis*, Studia Universitatis Babeş-Bolyai Series Informatica, LVII (2012), pp. 65–75.
- [14] Sanda Dragoş, Diana Haliţă, Christian Săcărea, Diana Troancă *Applying Triadic FCA in Studying Web Usage Behaviors*, Knowledge Science, Engineering and Management, 7th International Conference, Volume 8793, 2014, pp. 73-80
- [15] S. Dragoş, *PULSE Extended*, in *The Fourth International Conference on Internet and Web Applications and Services*, Venice/Mestre, Italy, May 2009, IEEE Computer Society, pp. 510–515.
- [16] S. Dragoş, *Why Google Analytics can not be used for educational web content*, 2011 International Conference on Next Generation Web Services Practices.
- [17] Dynamic404, Logical error-pages,
<http://www.yireo.com/software/joomla-extensions/dynamic404>
- [18] M. Eirinaki, M. Vazirgiannis, *Web mining for web personalization*, ACM Transactions on Internet Technology (TOIT), 3 (2003), pp. 1–27.
- [19] A. Farahat, M. Bailey, *How Effective is Targeted Advertising?*, Proceedings of the 21st World Wide web Conference 2012, Lyon, France, April 16-20, 2012, pp. 111-120
- [20] G. Gan, M. Chaoqun, W. Jianhong, *Data Clustering: Theory, Algorithms, and Applications*, ASA-SIAM Series on Statistics and Applied Probability, SIAM, Philadelphia, Alexandria, 2007
- [21] B. Ganter, R. Wille, *Formal Concept Analysis. Mathematical Foundations*. Springer, Berlin-Heidelberg-New York(1999).
- [22] B. Goncalves, J. Ramasco, *Human dynamics revealed through Web analytics*, Phys. Rev. E 78, 026123 (2008).
- [23] Z. Gyongy, H. Garcia-Molina, P. Berkhin, J. Pedersen, *Link Spam Detection Based on Mass Estimation*, 32nd International Conference in Very Large Data Bases (VLDB), Seoul, Korea, 2006, pp. 439-450
- [24] **D. Haliţă**, D. Bufnea: *A study regarding inter domain linked documents similarity and their consequent bounce rate*, Studia Universitatis Babeş-Bolyai Informatica, 1/2014, pp. 83-91
- [25] A. Huang, *Similarity Measures for Text Document Clustering*, Proceedings of the New Zealand Computer Science Research Student Conference, Hamilton, New Zealand, 2008, pp. 49-56
- [26] R. Jaeschke, A. Hotho, C. Schmitz, B.Ganter, G. Stumme, *Trias - An Algorithm for Mining Iceberg Trilattices*, Proceedings of the IEEE International Conference on Data Mining, pp. 907-911, Hong Kong, IEEE Computer Society, 2006.
- [27] The Joint Industry Committee for Web Standards (JICWEBS), *Reporting standards. website traffic*. Auditing Bureau of Circulations electronic (ABCe), Report 1, 2011.

- [28] R. Kosala, H. Blockeel, *Web mining research: A survey*, ACM Sigkdd Explorations Newsletter, 2 (2000), pp. 1–15.
- [29] F. Lehmann, R. Wille, *A Triadic Approach to Formal Concept Analysis*, Conceptual Structures: Applications, Implementation and Theory, vol. 954 of Lecture Notes in Artificial Intelligence, Springer Verlag, 1995.
- [30] J. Leskovec, A. Rajaraman, J. Ullman, *Mining of Massive Datasets*, Cambridge University Press, 2010
- [31] M. Najork, *Detecting Spam web Pages through Content Analysis*, International World Wide web Conference Committee, Edinburgh, Scotland, 2006, pp. 83-92
- [32] J.P. Norguet, B. Tshibusu-Kabeya, G. Bontempi, E. Zimanyi, *A Page-Classification Approach to Web Usage Semantic Analysis*, Engineering Letters, 14:1, EL 14 1 21.
- [33] L. Rao, *WordPress Now Powers 22 Percent Of New Active websites In The U.S.*, August, 2011, TechCrunch
- [34] C. Romero, P. G. Espejo, A. Zafra, J. R. Romero, S. Ventura, *Web usage mining for predicting final marks of students that use moodle courses*, Computer Applications in Engineering Education, 21 (2013), pp. 135–146.
- [35] C. Romero, S. Ventura, A. Zafra, P. D. Bra, *Applying web usage mining for personalizing hyperlinks in web-based adaptive educational systems*, Computers & Education, 53 (2009), pp. 828–840.
- [36] A. Singhal, *Modern Information Retrieval: A Brief Overview*, Bulletin of the IEEE Computer Society Technical Committee on Data Engineering, 2011, 24 (4): 35-43
- [37] M. Spiliopoulou, L. C. Faulstich, *Wum: a tool for web utilization analysis*, in The World Wide Web and Databases, Springer, 1999, pp. 184–203.
- [38] N. Spirin, J. Han, *Survey on web Spam Detection: Principles and Algorithms*, ACM SIGKDD Explorations Newsletter, Volume 13, Issue 2, December 2011, pp. 50-64
- [39] J. Srivastava, R. Cooley, M. Deshpande, T. Pang-Ning, *Web usage mining: Discovery and applications of usage patterns from web data.*, SIGKDD Explorations, 1(2), 2000.
- [40] R. Wille, *Conceptual landscapes of knowledge: a pragmatic paradigm for knowledge processing*, Proceedings of the International Symposium on Knowledge Representation, Use, and Storage Efficiency. Simon Fraser University, Vancouver 1997, 2-13.
- [41] D.J. Watts, *Nature 445*, 489 (2007).
- [42] Wille, R.: Conceptual Landscapes of Knowledge: a Pragmatic Paradigm for Knowledge Processing, In: Gaul, W.; Locarek-Junge H. (Eds.): Classification in the Information Age, Proceedings of the 22nd Annual Gfki Conference, Dresden, March 4-6, 1998, pp. 344–356.

- [43] R. Wille, *Methods of Conceptual Knowledge Processing, Formal Concept Analysis*, 4th International Conference ICFCA 2006, Dresden, Germany, LNAI 3874, Springer 2006, 1–29.
- [44] R. Wille, *The Basic Theorem of triadic concept analysis*, Order, no. 2, 1995, pp. 149-158
- [45] Wordpress, 404 Redirected Wordpress *plugin*,
<http://wordpress.org/extend/plugins/404-redirected/>
- [46] B. Zhou, S. C. Hui, K. Chang, *A formal concept analysis approach for web usage mining*, in Intelligent information processing II, Springer, 2005, pp. 437–441.
- [47] D. Zhou, C. Burges, T. Tao, *Transductive Link Spam Detection*, Proceedings of the 3rd International Workshop on Adversarial Information Retrieval on the web, AIRweb, New York, USA, ACM Press, 2007, pp. 21-28
- [48] B. Zhou, S. C. Hui, A. C. Fong, *Web usage mining for semantic web personalization*, in Workshop on Personalization on the Semantic Web, 2005, pp. 66–72.
- [49] Zyxxware Technologies, Search404: Automatically search for content when a 404 error occurs,
<http://drupal.org/project/search404>