

Relevanța conținutului web și a comportamentului utilizatorilor în analiza traficului

Student doctorand: Diana - Florina HALIȚĂ
Coordonator științific: Prof. Dr. Florian Mircea BOIAN
Universitatea "Babeș-Bolyai" Cluj-Napoca
Proiect de cercetare

10 Iunie 2015



CUPRINS

- 1 Impactul similarității web asupra traficului
 - MIGRAREA UNUI SITE WEB
 - SIMILARITATE ȘI BOUNCERATE
 - SCRAPER SITE
- 2 Interpretarea logurilor unei platforme de e-learning folosind FCA
 - ANALIZA FORMALĂ CONCEPTUALĂ
 - WEB USAGE MINING
 - INTERPRETAREA REZULTATELOR FOLOSIND CIRCOS



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



DE CE?

De ce să alegem un CMS?

- Întreținerea facilă - web;
- Independența conținutului de prezentare;
- Update-uri periodice de securitate;
- Roluri multiple pentru utilizatori;

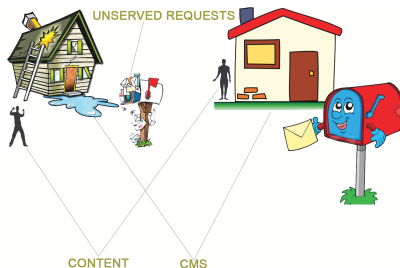
De ce este necesară migrarea?

- CMS-ul vechi este considerat învechit;
- CMS open source;
- Scalabilitatea CMS-ului nou;
- Fructificarea noilor tehnologii web: CSS3, Ajax, HTML5.



PROVOCĂRI

- expunerea conținutului site-ului web la un URL nou
- creșterea numărului de vizitatori care sunt induși în eroare
- pierderea diverselor beneficii câștigate în timp



SOLUȚII

Soluții posibile:

- migrarea automată sau manuală;
- estimarea timpului necesar migrării;
- evaluarea procesului de migrare.

Soluții propuse de alții:

- plugin-uri care țin cont de comportamentul utilizatorilor în timp.

Soluția propusă:

- potrivirea perechilor de URL-uri;
- **AVANTAJ:**
algoritmul se utilizează înaintea lansării site-ului.

REZULTATE

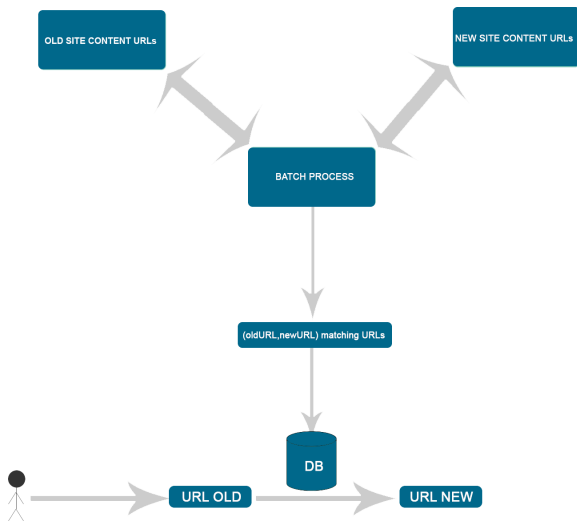


Figure: Batch processing

REZULTATE

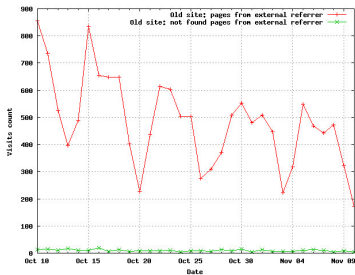


Figure: Vechiul site: numărul de pagini accesate de către un referrer extern și numărul de pagini care au generat eroarea 404 și au fost accesate de către un referrer extern

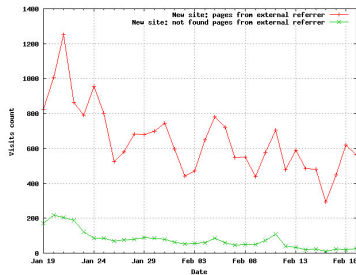


Figure: Noul site: numărul de pagini accesate de către un referrer extern și numărul de pagini care au generat eroarea 404 și au fost accesate de către un referrer extern



REZULTATE

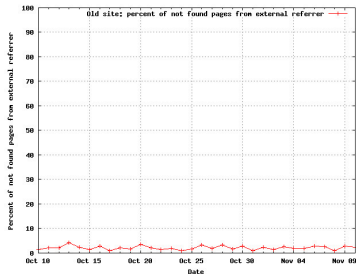


Figure: Vechiul site: procentul paginilor ce au generat eroare 404 și au fost accesate de către un referrer extern

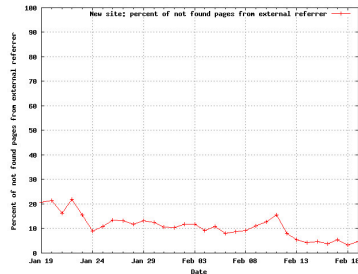


Figure: Noul site: procentul paginilor ce au generat eroare 404 și au fost accesate de către un referrer extern



REZULTATE

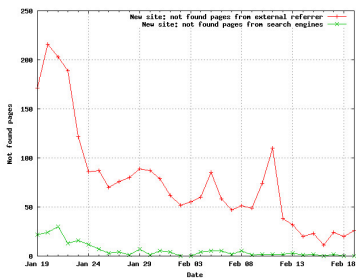


Figure: Noul site: adaptarea motoarelor de căutare la noua structură a site-ului

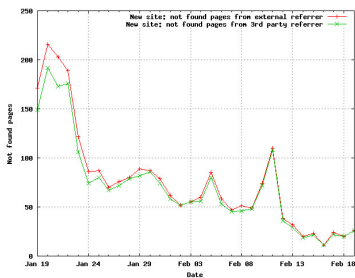


Figure: Noul site: majoritatea erorilor 404 sunt generate de către utilizatorii care vin de la refereri 3rd party

REZULTATE

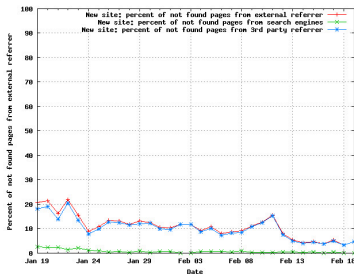


Figure: Noul site: procentul paginilor ce generează eroarea 404 și procentul paginilor ce generează eroarea 404 și vin de la motoarele de căutare sau de la referreri 3rd party

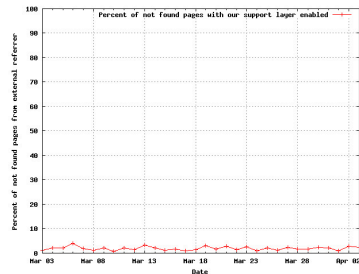


Figure: Procentul paginilor ce generează eroarea 404 având un referer extern, odată ce a fost activat layer-ul suport propus



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



DE CE?

Link-urile sunt folosite abuziv:

Scop:

- creșterea page rank-ului domeniului destinație;

Rezultat:

- utilizatorului nu îi este prezentată o informație de care să fie interesat;

Locație:

- sunt localizate fie sitewide, fie în cadrul conținutului absolut.



SOLUȚII

Tehnici de detectare a link-urilor abuzive:

- analiza conținutului
- analiza link-urilor
- analiza comportamentului utilizatorilor
- metode de clasificare automată, supervizată sau nesupervizată

Acuratețe:

- 80% din paginile care contin link-uri abuzive sunt detectate

Recomandare:

- combinarea tehnicilor ⇒ crește procentul de pagini abuzive detectate



REZULTATE

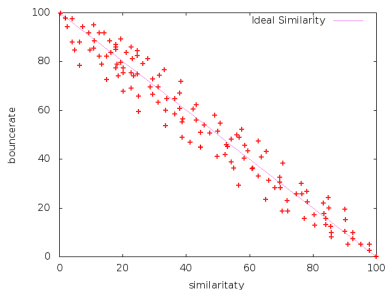


Figure: Similaritate ideală

- similaritate mare \Rightarrow rata de respingere mică
- similaritate mică \Rightarrow rata de respingere mare
- cea mai bună funcție de similaritate:
condiții: suma

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

este minimă



REZULTATE

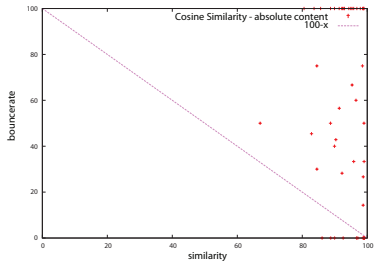


Figure: Similaritatea Cosinus

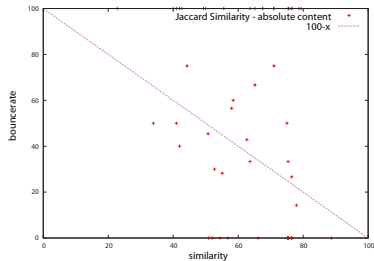


Figure: Similaritatea Jaccard

REZULTATE

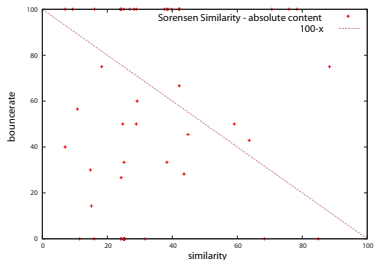


Figure: Similaritatea Sorensen

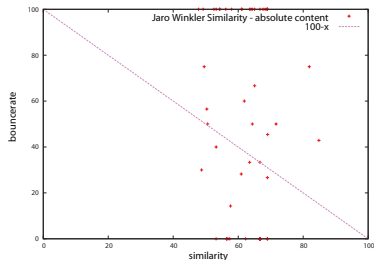


Figure: Similaritatea Jaro Winkler

REZULTATE

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

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



Măsurarea și vizualizarea nivelului de scrapping al unui site web



DE CE?

Probleme:

Prezența în SERP a unor site-uri care:

- direcționează greșit utilizatorii;
- direcționează utilizatorii spre un scraper site.

Consecințe:

- scăderea performanței motoarelor de căutare;
- nemulțumirea utilizatorilor referitoare la informația găsită



SOLUȚII

Identificarea scraper site-urilor:

- algoritmi automați;
- pe baza feedback-ului utilizatorilor - programe de învățare automată;
- 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.



REZULTATE

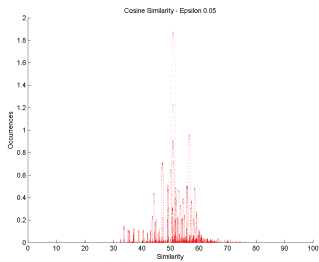


Figure: Non-Scraper site:
Similaritatea Cosinus - $\epsilon = 0.05$

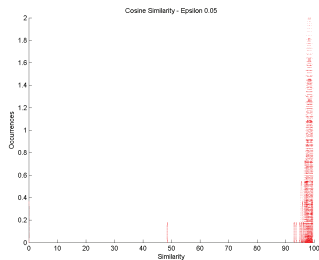


Figure: Scraper site:
Similaritatea Cosinus - $\epsilon = 0.05$

CONCLUZII

Migrarea unui site web

- ponderarea proprietăților conținutului prezentat la un anumit URL;
- redirecționarea vizitatorului spre o pagină similară cu cea a referrer-ului.

Similaritate și bouncerate

- analiza similarității conținuturilor
⇒ (referrer, \forall link intern);
- funcții de similaritate semantică;
- ponderarea funcțiilor de similaritate;
- ponderarea unor proprietăți specifice ale conținutului.



FCA

Definiție

Un *context formal* este un triplet (G, M, I) , unde:

- G reprezintă mulțimea obiectelor;
- M reprezintă mulțimea atributelor;
- $I \subset G \times M$ reprezintă o relație binară între mulțimea obiectelor și mulțimea atributelor, numită relație de incidență

Observație

gIm se citește obiectul g are atributul m .



Definiție

Operatorul de derivare in FCA este o conexiune Galois între mulțimile G și M :

- $A' = \{m \in M \mid \forall g \in A, gIm\}, A \subseteq G;$
- $B' = \{g \in G \mid \forall m \in B, gIm\}, B \subseteq M.$



3FCA

Definiție

Un *triconcept* al unui *tricontext* (K_1, K_2, K_3, Y) este un *triplet maximal* (A_1, A_2, A_3) cu $A_i \subseteq K_i$.

Observație

Pentru un *triconcept* (A_1, A_2, A_3) :

- A_1 se numește *extent*-ul *triconcept*-ului;
- A_2 se numește *intent*-ul *triconcept*-ului;
- A_3 se numește *modus*-ul *triconcept*-ului.



WEB USAGE MINING

- tehnicile din Data Mining precum și cele din Statistică au fost utilizate pentru a extrage informații utile din log-urile web;
- tipologia site-urilor:
 - site-uri de e-commerce - **SCOP**: - vânzarea de produse;
 - site-uri de e-learning - **SCOP**: - oferirea de informații.
- o vizită pe un site educațional nu se aplica euristicilor utilizate de majoritatea instrumentelor de analiză.



REZULTATE

- **platforma de e-learning:** PULSE;
- **perioada de colectare a datelor:** Februarie - Iulie 2013;
- **ce investigăm:** tipare de comportament al utilizatorilor care folosesc PULSE;
- **date folosite:** clasele URL-urilor accesate, clase de referreri, timestamp-urile accesărilor username-ul studenților;
- ELBA, TOSCANA & Toscana2TRIAS.



CLASELE URL-URILOR ACCESATE (AF_CLASS)

- URL-urile accesate au fost împărțite în 9 clase.

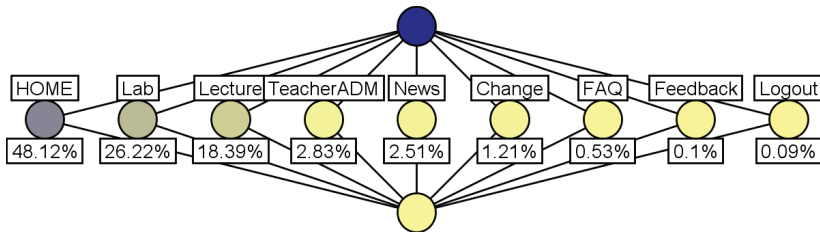


Figure: Clasele URL-urilor accesate

CLASE DE REFERRERI (R_CLASS)

- URL-urile referrer-ilor reprezintă paginile dinspre care utilizatorii vizitează o pagină dintr-un site;
- referrer-ii externi platformei PULSE nu au fost considerați în acest studiu.

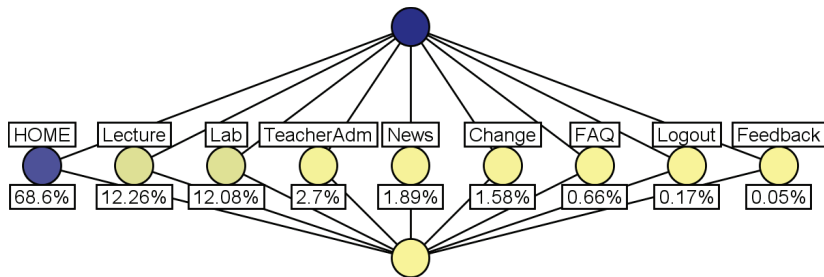


Figure: Clase de referreri

WEB USAGE MINING FOLOSIND 3FCA

- Toscana2Trias permite selectarea datelor triadice pornind de la un set de scale preprocesate în ToscanaJ;
- Date de intrare pentru Trias:
 - attribute: perechea (R_CLASS - AF_CLASS);
 - conditii: timestamp;
 - obiecte: studenti.



Interpretarea rezultatelor obținutelor în 3FCA folosind CIRCOS



INTERPRETAREA REZULTATELOR OBȚINUTELOR ÎN 3FCA FOLOSIND CIRCOS

Observație

- *XML-ul generat de TRIAS conține toate triconceptele derivate din tricontext;*
- *fiecare triconcept e definit de extent, intent și modus;*
- *datele de intrare pentru CIRCOS trebuie reprezentate într-un tabel bidimensional, $R \times C$.*



INTERPRETAREA REZULTATELOR OBȚINUTELOR ÎN 3FCA FOLOSIND CIRCOS

Observație

- $(G,M,B,Y) \Leftrightarrow (G,(B,M),I), (g, (b, m)) \in I \Leftrightarrow (g, m, b) \in Y;$
- $\forall (b,m) \Rightarrow$ numărul de elemente din extent este $(b,m)';$
- $C =$ mulțimea care definește indicatorii de coloană, obținută prin proiectarea relației Y pe $M;$
- $R =$ mulțimea care definește indicatorii de linie, obținută prin proiectarea relației Y pe $B;$
- valorile numerice ale tabelului se calculează astfel:
 $\forall (c, r) \in C \times R,$ numărul de elemente din extent $(c,r)'$ este calculat direct din XML-ul rezultat din Trias.



INTERPRETAREA REZULTATELOR OBȚINUTELOR ÎN 3FCA FOLOSIND CIRCOS

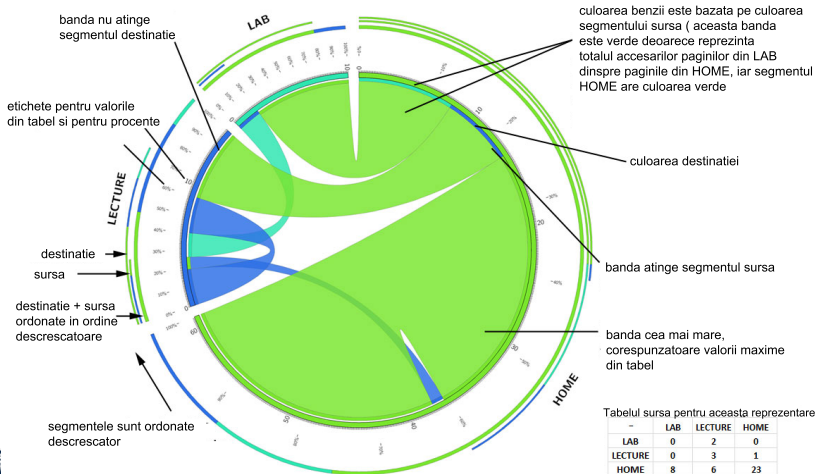


Figure: Grupa: "ar" - săptămână

REZULTATE

Test specification

- **durata:** un semestru;
- **date:** grupe de studenți;
- **triconcepte:**
 - **obiecte:** R_class;
 - **atribute:** AF_class;
 - **condiții:** timestamps;
- **comportament:**
 - relaxat;
 - intens;
 - normal.

Rezultatele obținute sunt publicate:

<http://www.cs.ubbcluj.ro/~fca/tests-2013/>



REZULTATE

■ comportamentul relaxat:

- apare în special în timpul vacanțelor;
- **observație:** număr redus de URL-uri accesate;
- **tipar:** HOME ⇒ LAB ⇒ LECTURE;
- **tipar:** HOME ⇒ LECTURE ⇒ HOME;



REZULTATE

- comportamentul normal:
 - apare în timpul semestrului, când nu sunt examene sau vacanțe;
 - **observație:** aproape toate clasele de URL-uri accesate sunt vizitate;
 - paginile din clasa LAB sunt cele mai vizitate.



REZULTATE

■ comportamentul intens:

- apare în special în perioadele cu examene;
- **observație:** număr mare de accesări;
- paginile din clasa LECTURE sunt cele mai vizitate;
- HOME reprezintă un liant între restul facilităților oferite de PULSE;



REZULTATE - GRUPA: "AR"



Figure: SO:
comportament relaxat

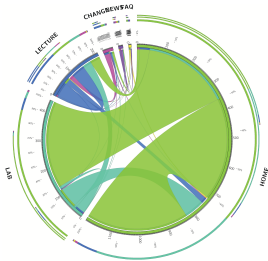


Figure: SO:
comportament normal

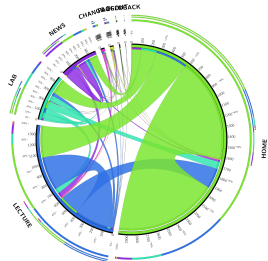


Figure: SO:
comportament intens



REZULTATE - GRUPA: "RI"

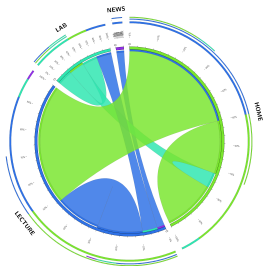


Figure: SO:
comportament relaxat

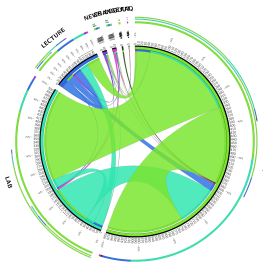


Figure: SO:
comportament normal

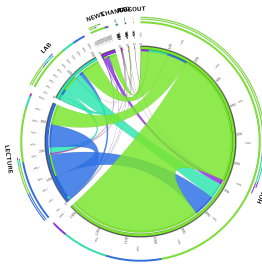


Figure: SO:
comportament intens



REZULTATE - GRUPA: "IE"

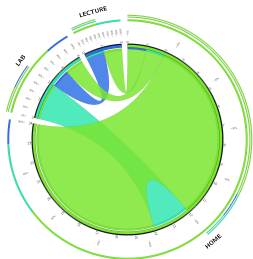


Figure: WDO: comportament relaxat

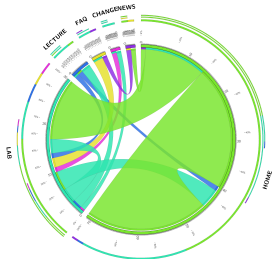


Figure: WDO: comportament normal

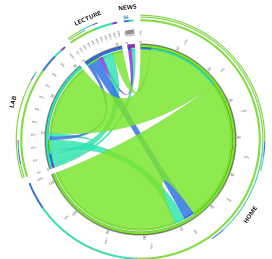


Figure: WDO: comportament intens



REZULTATE - GRUPA: "EI"

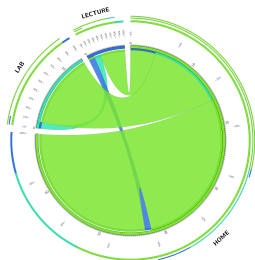


Figure: WDO:
comportament relaxat

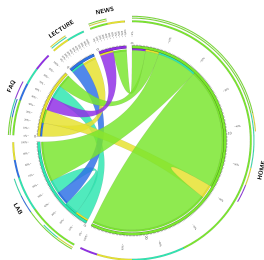


Figure: WDO:
comportament normal

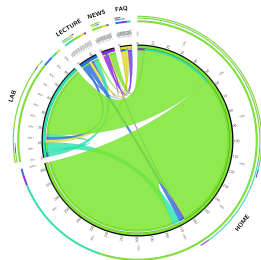


Figure: WDO:
comportament intens

CONCLUZII

- ajustarea parametrilor în fișierele de configurare ale CIRCOS, precum și afișarea datelor într-un alt format decât cel circular.
- construirea unui algoritm bazat pe probabilități condiționale pentru a șterge conceptele cele mai puțin importante.
- utilizarea TCA pentru a evidenția din punct de vedere temporal comportamentul utilizatorilor, atât individual cât și pe grupuri.



Mulțumesc!

Q & A

