

Impactul similarității documentelor web asupra traficului

Absolvent: Diana - Florina HALIȚĂ

Coordonator științific: Lect. Dr. Darius - Vasile BUFNEA
Universitatea "Babeș-Bolyai" Cluj-Napoca

1 Iulie 2014



CUPRINS

- 1 MIGRAREA UNUI SITE WEB
 - De ce?
 - Provocări
 - Soluții
 - Rezultate
- 2 SIMILARITATE ȘI BOUNCERATE
 - De ce?
 - Soluții
 - Rezultate
- 3 SCRAPER SITE
 - De ce?
 - Soluții
 - Rezultate
- 4 CONCLUZII



DE CE?

De ce să alegem un CMS?

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

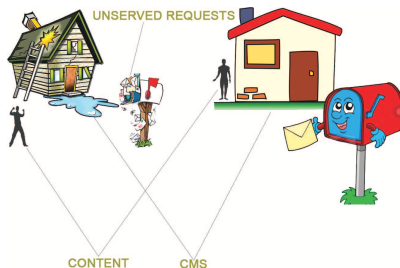
De ce este necesară migrarea?

- 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

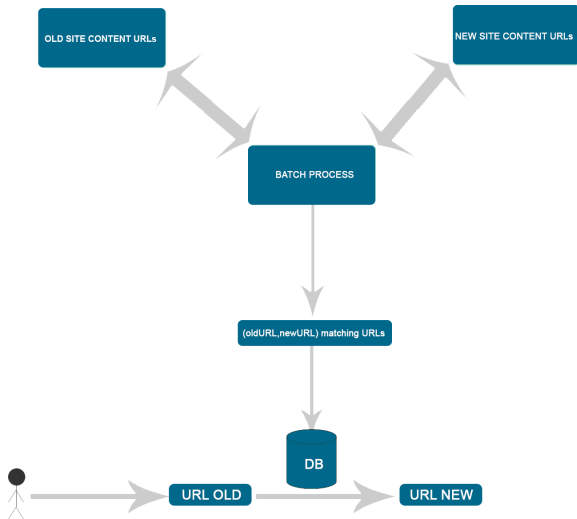


Figure : Batch processing



REZULTATE

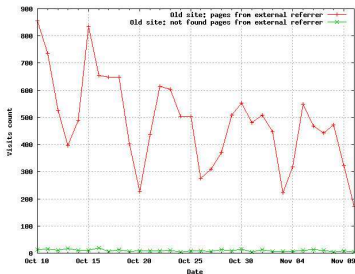


Figure : Vechiul site: numărul de pagini accesate de către un referrer extern și numărul de pagini care au generat eroarea 404 și au fost accesate de către un referrer extern

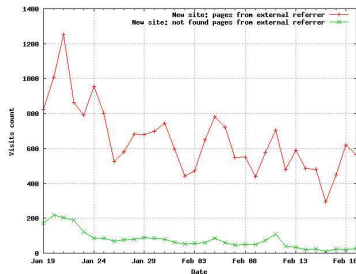


Figure : Noul site: numărul de pagini accesate de către un referrer extern și numărul de pagini care au generat eroarea 404 și au fost accesate de către un referrer extern

REZULTATE

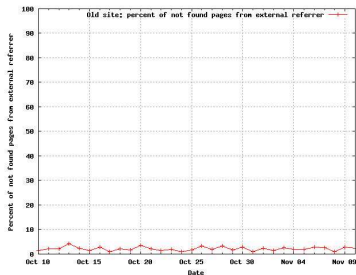


Figure : Vechiul site: procentul paginilor ce au generat eroare 404 și au fost accesate de către un referer extern

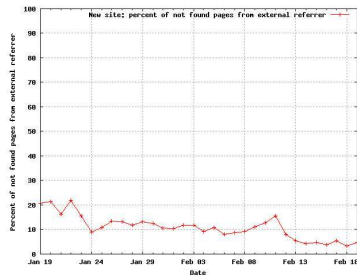


Figure : Noul site: procentul paginilor ce au generat eroare 404 și au fost accesate de către un referer extern

REZULTATE

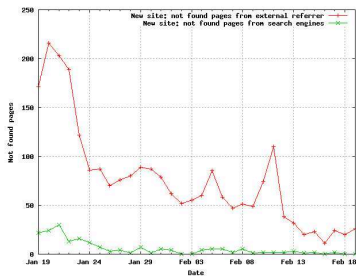


Figure : Noul site: adaptarea motoarelor de căutare la noua structură a site-ului

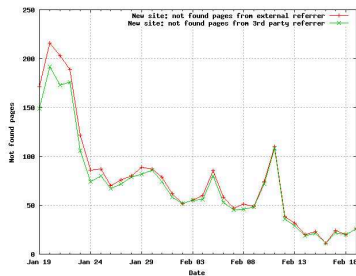


Figure : Noul site: majoritatea erorilor 404 sunt generate de către utilizatorii care vin de la refereri 3rd party

REZULTATE

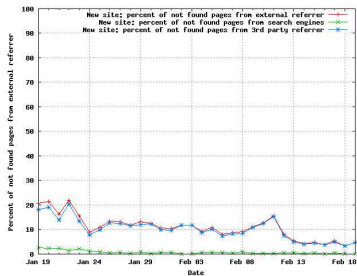


Figure : Noul site: procentul paginilor ce generează eroarea 404 și procentul paginilor ce generează eroarea 404 și vin de la motoarele de căutare sau de la referreri 3rd party

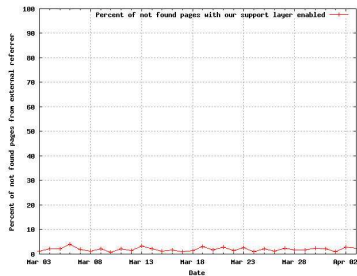


Figure : Procentul paginilor ce generează eroarea 404 având un referer extern, odată ce a fost activat layer-ul suport propus

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

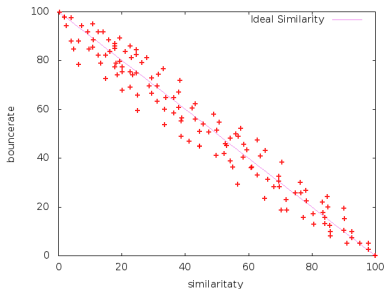


Figure : Similaritate ideală

- 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

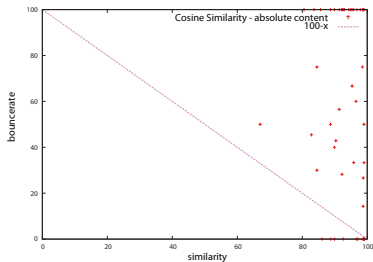


Figure : Similaritatea Cosinus

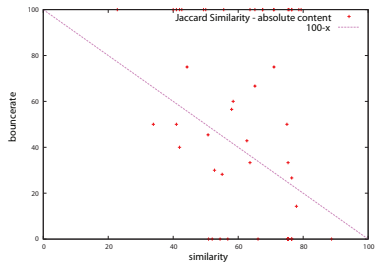


Figure : Similaritatea Jaccard

REZULTATE

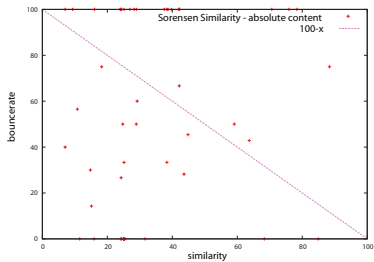


Figure : Similaritatea Sorensen

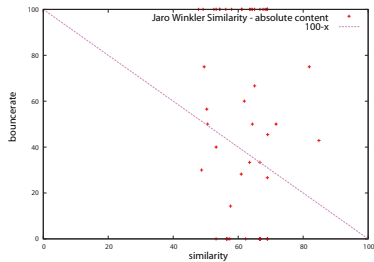


Figure : Similaritatea Jaro Winkler

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 $y = x - 100$



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

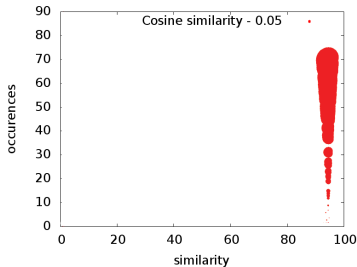


Figure : Scraper site:
Similaritatea Cosinus - $\epsilon = 0.05$

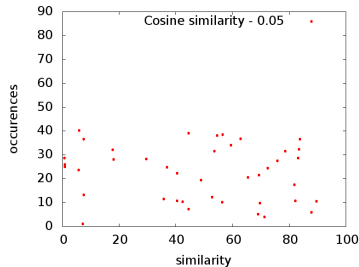


Figure : Non-Scraper site:
Similaritatea Cosinus - $\epsilon = 0.05$

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.



Mulțumesc!

Q & A

