

UNIVERSITY POLITEHNICA OF BUCHAREST
FACULTY OF AUTOMATIC CONTROL AND COMPUTERS
COMPUTER SCIENCE DEPARTMENT



Data mining & Data warehousing

-Friend suggestion on social networks-

Web Mining: SALSA

Enache Andrei-George

MT11A

BUCHAREST

2022

Cuprins

1. Importanta si aplicatiile practice ale algoritmului	3
2. Prezentarea generala a algoritmului	3
3. Rezultate si probleme cunoscute	Error! Bookmark not defined.
4. Seturi de date utilizate (nume si esantioane)	7
5. Rezultate si evaluarea acestora	8
6. Bibliografie.....	11

1. Importanta si aplicatiile practice ale algoritmului

Internetul reprezintă o colecție de informații nestructurate sau semi-structurate, în continua creștere, ce conține diverse surse și legături între pagini. În ziua de azi, atunci când o persoană caută o informație pe internet, aceasta face o interogare în cadrul unui browser, fiindu-i returnate o serie de link-uri și pagini cu conținut asemănător, cel mai probabil, celui căutat. În cazul interogărilor pe subiecte generale pot rezulta multe pagini irelevante pentru utilizator, îngreunând astfel procesul de găsire a informațiilor utile [2].

Câteva dintre problemele întâlnite în cadrul căutărilor pe motoarele de căutare web sunt reprezentate de [7]:

- Interogările restranse - informația este foarte puțină și greu de găsit, neputându-se stabili concret dacă este relevantă sau nu pentru utilizator;
- Interogări pe subiecte mari - există foarte multe pagini ce dezbate un anumit subiect, neputându-se stabili cu ușurință relevanța informațiilor afișate, utilizatorii având nevoie doar de câteva site-uri pentru a-și satisface nevoia de căutare;
- Sinonimia – căutarea anumitor informații poate duce la găsirea unor pagini web cu tema similară, dar diferită de cea pe care și-a dorit utilizatorul;
- Polisemia – rezultate din domenii diferite ale aceluiași termen, putând fi catalogat și ca o căutare ambiguă;
- Stiluri diferite ale autorilor – oameni din zone diferite ce folosesc un vocabular diferit pentru elemente din același domeniu, putând induce în eroare utilizatorul.

Din fericire, una dintre cele mai importante caracteristici ale site-urilor web o reprezintă existența legăturilor către alte site-uri, oferind astfel o structură generală asemănătoare între diferite site-uri, putând astfel face generalizări cu scopul de a îmbunătăți calitatea elementelor afișate în urma căutărilor în browser, utilizând diverse metode de web mining.

Web mining-ul reprezintă procesul de aplicare a algoritmilor și tehnicilor de căutare a informațiilor pentru a extrage datele direct din paginile web. Printre principalele categorii de web mining se Numără:

- Web content mining – procesul de găsire a informațiilor utile din conținutul text și media a paginilor web;
- Web structure mining – procesul de analizare a structurii paginilor web, a nodurilor de legătură și a conexiunilor între diferite pagini;
- Web usage mining – procesul de extragere a informațiilor despre modul cum utilizatorii folosesc un anumit site.

În cadrul acestei lucrări a fost implementat algoritmul de web structure mining SALSA (Stochastic Approach for Link-Structure Analysis) cu scopul de a face o analiză a importanței anumitor pagini web în funcție de structurile hyperlink-urilor existente. Acest algoritm este utilizat pentru a extrage modele din hyperlink-urile web și se focusează pe acelea care unesc pagini diferite (Inter-Document Hyperlink), nu și acelea care trimit utilizatorul către aceeași pagină (Intra-Document Dyperlink). Scopul implementării acestui algoritm este de a spori relevanța rezultatelor căutărilor în motoarele de căutare, dar poate fi utilizat și în alte scopuri, cum ar fi sugerarea elementelor ce au legătură cu utilizatorul (ex: sugestia prietenilor pe facebook).

2. Prezentarea generală a algoritmului

SALSA (Stochastic Approach for Link-Structure Analysis) reprezintă un algoritm de clasificare a paginilor web ce le atribuie acestora scoruri de butuci (hubs) și de autorități (authorities), pe baza hyperlink-urilor și legăturilor dintre acestea [1].

Această soluție a fost inspirată din alți 2 algoritmi de web structure mining, HITS și PageRank. La fel ca și HITS, SALSA conferă paginilor web scoruri în funcție de relevanța acestora în funcție de un anumit subiect. Scorurile pe care le poate primi un site sunt denumite Buturugi (Hubs) și Autorități

(Authorities). O autoritate mai mare reprezintă o relevanță mai mare a unui site asupra unui anumit subiect, având mai multe hyperlink-uri ce ținesc acea pagină. Cu cât valoarea hub-ului unei pagini este mai mare, relevanța acelei pagini scade, fiind considerate doar o pagină ce trimite utilizatorul către un site mai important. Asemănarea dintre acest algoritm și PageRank este modalitatea de calculare a scorurilor hubs și authorities prin intermediul unui drum la întâmplare printr-un lanț Markov ce reprezintă graful paginilor web.

În prima etapă a algoritmului se alege un subiect care este căutat prin intermediul unui motor de căutare. Din rezultatele găsite se va alege un număr de pagini care vor fi verificate ulterior, luându-se de asemenea hyperlink-urile prezente pe paginile alese, dar și cautându-se hyperlink-uri ce duc către acele site-uri.

Algoritmul pornește de la ipoteza că paginile cu nivel de autoritate mai mare asupra unui subiect "t" vor fi vizibile de pe mai multe pagini din graful "C" (format din noduri ce reprezintă site-urile, marginile fiind considerate că fiind link-urile dintre pagini). Utilizând graful direcționat format din paginile rădăcină C se formează un graf bipartit nedirecționat "G", având în componență subgraful "H" ce va conține nodurile hub (care conțin hyperlink-uri ce redirectionează utilizatorul către alte pagini) și subgraful "A" ce va conține nodurile autoritate (pagini către care utilizatorul este trimis prin apăsarea link-urilor de pe paginile hub)(Fig.1, Fig.2)[1].

- $V_h = \{s_h \mid s \in \mathcal{C} \text{ and } out-degree(s) > 0\}$ (the *hub side* of \tilde{G}).
- $V_a = \{s_a \mid s \in \mathcal{C} \text{ and } in-degree(s) > 0\}$ (the *authority side* of \tilde{G}).
- $E = \{(s_h, r_a) \mid s \rightarrow r \text{ in } \mathcal{C}\}$.

Fig.1. Împărțirea nodurilor inițiale în subgrafurile H și A

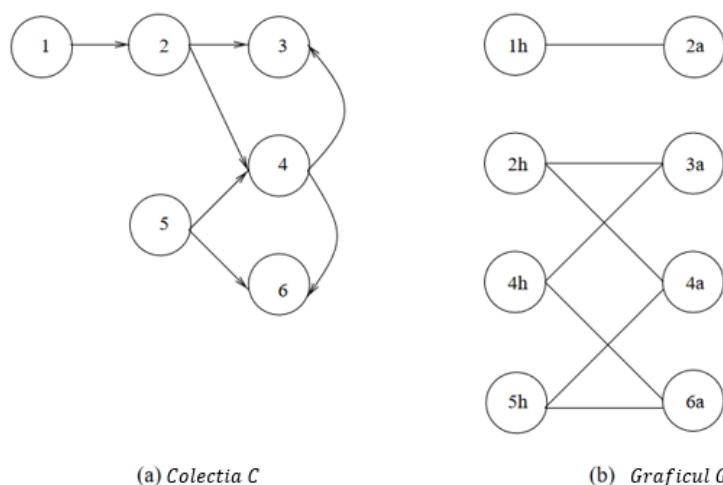


Fig.2. Transformarea grafului direcționat "C" în graful bipartit "G"[1]

Unde:

- G = Graful bipartit
- S = pagini web neizolate aparținând lui "C", reprezentate de 2 noduri ale lui "G" ("sh" și "sa")
- V_h = marginile care ies din hub și redirectionează utilizatorul către o autoritate
- V_a = marginile care ținesc spre autoritate
- E = marginile nedirecționate

a) Metoda de implementare 1

După organizarea subgrafurilor este necesară realizarea a două drumuri aleatorii distincte, fiecare dintre acestea pornind dintr-o latură a grafului, luându-se astfel în considerare un lanț de deplasare al hub-urilor și unul al autorităților. Pentru a realiza stările de tranziție ale fiecărui lanț, este necesară

parcurea a doua margini ale grafului, prima în fata (urmărind un link ce pornește din nodul de unde s-a pornit drumul) și a doua înapoi către partea inițială din graf [6]. Ponderile autorităților sunt calculate ca fiind distribuția lanțului ce verifica mai întâi un link direcționat către nodul de plecare, apoi urmandu-se un link direcționat către celelalte autorități ce se afla în legătură cu nodul respective. Ponderele hub-urilor este calculate similar, pornindu-se dintr-un nod aflat pe partea stangă a grafului și ajungând într-o autoritate de pe partea dreapta, de unde apoi se va continua drumul pe marginile ce țintesc către autoritatea pe care s-a staționat în acel moment.

Astfel, se ia în considerare un lant Markov al parcurgerii nodurilor autoritate ale lui "G" si unul al parcurgerii nodurilor hub ale acestuia în scopul crearii scorului de autoritate "A(u)" si a scorului de hub "H(u)" (Fig.3) [3]:

SALSA-Authority-Scores:

1. Let B^A be $\{u \in B : in(u) > 0\}$.

2. For all $u \in B$:

$$A(u) := \begin{cases} \frac{1}{|B^A|} & \text{if } u \in B^A \\ 0 & \text{otherwise} \end{cases}$$

3. Repeat until A converges:

(a) For all $u \in B^A$:

$$A'(u) := \sum_{(v,u) \in N} \sum_{(v,w) \in N} \frac{A(w)}{out(v)in(w)}$$

(b) For all $u \in B^A : A(u) := A'(u)$

SALSA-Hub-Scores:

1. Let B^H be $\{u \in B : out(u) > 0\}$.

2. For all $u \in B$:

$$H(u) := \begin{cases} \frac{1}{|B^H|} & \text{if } u \in B^H \\ 0 & \text{otherwise} \end{cases}$$

3. Repeat until H converges:

(a) For all $u \in B^H$:

$$H'(u) := \sum_{(u,v) \in N} \sum_{(w,v) \in N} \frac{H(w)}{in(v)out(w)}$$

(b) For all $u \in B^H : H(u) := H'(u)$

Fig. 3. Calculul scorului de autoritate si de hub al fiecarui nod din subgraful "A" si "H"

Calculul realizat în Fig. 3. sunt reprezentate grafic în Fig.4:

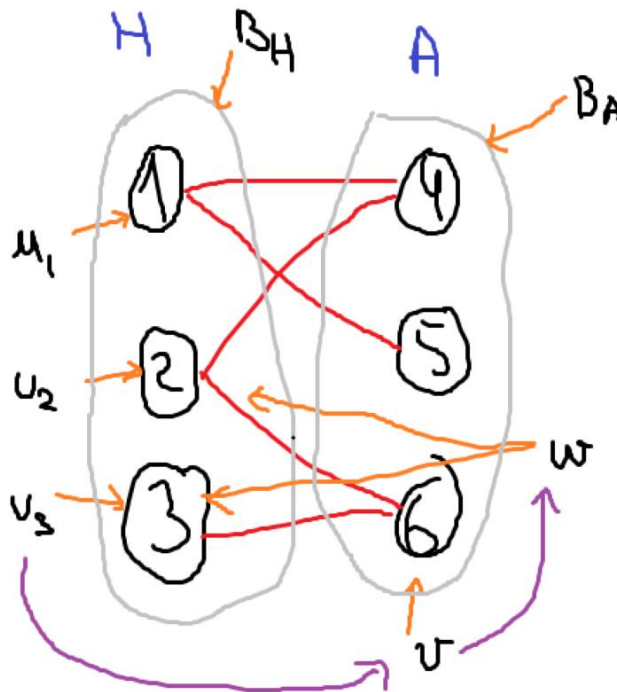


Fig. 4. Reprezentare grafica a drumurilor aleatorii realizate pentru "H", respectiv "A"

Pentru calculul scorului de autoritate se vor alege întâi toate nodurile din graf care au margini directionate de la "H" către "A". În cazul în care un nod nu va avea o margine directionată către el, atunci scorul de autoritate al acestuia va fi 0. În caz contrar, acestuia i se va atribui o valoare egală cu raportul nodului față de numărul total de noduri din subgraf.

În următoarea etapă, pentru toate nodurile aflate în "A" (notate cu "u"), se va verifica fiecare nod din subgraful "H" care țintește către cel autoritate de la pasul curent (vertex-urile aflate în "H" se vor nota cu "v"), acest calcul reprezentând prima parte de drum aleatoriu. Al doilea pas este cel de verificare a nodurilor din "A" de care nodul "v" este legat (aceste noduri din "A" fiind notate cu "w"). Odată aflate toate traseele posibile, noua valoare de autoritate va fi egală cu suma rapoartelor dintre valoarea autoritate a nodurilor "w" (la care un nod "u" poate ajunge prin trecerea către un nod "v" din subgraful "H" și înapoi printr-o legătură directă în partea "A") și produsul dintre numărul marginilor ce ies din nodurile "v" cu cel al marginilor ce intra în nodurile "w".

Această etapă de calcul se va repeta similar și pentru valoarea nodurilor hub, dar în sens invers (un drum începând din subgraful "H", ce va ajunge în "A" și se va întoarce înapoi în zona inițială). Procesul de calcul pentru fiecare nod din graf se va repeta de "k" ori, până când valoarea fiecărui nod se va normaliza, ajungând că după mai multe iterații să nu se mai observe nicio diferență după recalculare.

b) Metoda de implementare 2

Că alternativă la calculul prezentat anterior, se poate crea inițial o matrice de adiacență "W" a grafului directionat "C". Prin intermediul acesteia, se vor crea încă 2 matrice, unde "W_r" va fi matricea rezultată prin împărțirea fiecărei intrări diferite de 0 la suma elementelor de pe rânduri, respective "W_c", o matrice rezultată din împărțirea tuturor elementelor nenule ale celei inițiale la suma elementelor de pe coloana.

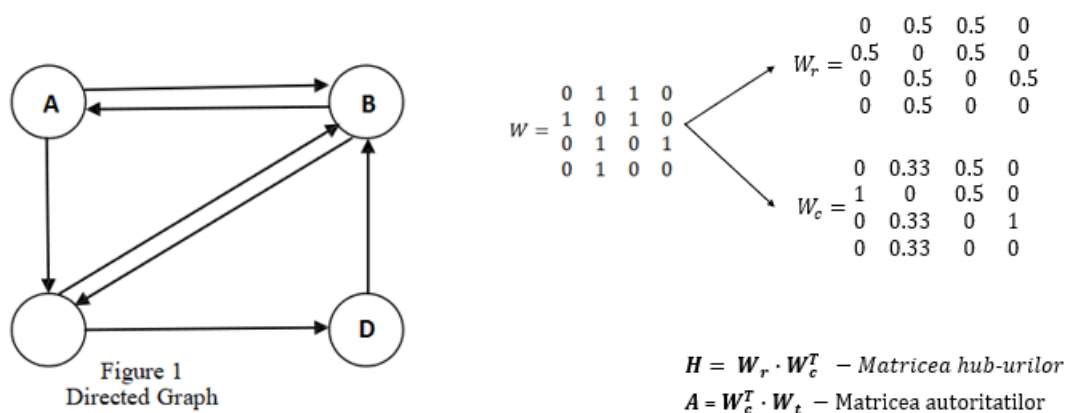


Fig.5. Crearea matricei de adiacență în funcție de graful inițial

După calculul acestor matrice, principalele comunități de autorități și hub-uri vor fi constituite din primele "k" pagini având ca cele mai mari intrări în vectorul principal al A și H.

În cadrul implementării acestui algoritm, s-a mers pe prima metodă de implementare, fiind mai rapid de calculat prin intermediul bibliotecii "networkx" din Python.

3. Rezultate și probleme cunoscute

Unul dintre avantajele acestui algoritm în comparație cu celelalte metode utilizate în web structure mining este acela al lipsei sensibilității acute la efectul TKC (Tightly Knit Community), o problemă majoră pentru algoritmi de web mining ce nu pot să ofere notații de importanță precise grupurilor strânse, unde mai multe elemente legate de același subiect se leagă între ele, formând un

amalgam de trimiteri ce își consolidează reciproc importanța, fiecare pagină având multe link-uri tip hub și multe referințe tip autoritate).

Un alt avantaj al acestui algoritm este acela al rezultatelor remarcabile atunci când este supus la căutări generale, având rezistența la hub-urile cu o valoare foarte mică, astfel crescând șansa găsirii rezultatelor relevante. De asemenea, un avantaj similar algoritmului PageRank este acela că ponderile nodurilor pot fi puse explicit, putând astfel influența într-o anumită măsură căutările inițiale, în funcție de scopul dorit (căutări generale, căutări de informații restrânse, etc.). În plus, spre deosebire de celelalte metode de web structure mining, SALSA nu are topic drift, astfel prezentând un avantaj în cazul problemelor de polisemie.

Printre dezavantajele cunoscute ale acestui algoritm este acela al favorizării în unele cazurilor a comunităților mai mici, fără, poate, prea mare legătură cu subiectul (ce cauzează automat mai puține noduri, ajungând că rezultatele să se normalizeze mult mai repede, data de cazurile cu eșantioane mari de date).

Un alt dezavantaj al acestui algoritm este reprezentat de rezultate mai slabe decât HITS sau PageRank în cazul căutărilor specifice, neavând capacitatea de a normaliza rezultatele în urma mai multor iterații pe un număr mic de noduri. În plus, acest algoritm este dependent de interogări, astfel putând avea rezultate mai puțin bune în funcție de căutările inițiale (graficul inițial "C" având un impact major asupra rezultatelor extrapolate din graficul "G").

3. Seturi de date utilizare (nume și eșantioane)

În cadrul acestui proiect au fost utilizate 4 eșantioane de date privind relațiile dintre persoane de pe diferite rețele sociale, obținute de pe site-ul <https://snap.stanford.edu/data/index.html>.

Seturile de date sunt:

- Twitter ("Social circles: Twitter"): Dataset ce conține cercuri de prieteni de pe Twitter luate din surse publice. Dataset-ul formează un graf direcționat constituit din 81.306 noduri cu 1.768.149 muchii;
- Google Plus ("Social circles: Google+"): Dataset ce conține cercuri de prieteni de pe Google+, informațiile fiind colectate de la oameni care și-au distribuit relațiile cu ceilalți membri. Dataset-ul formează un graf direcționat constituit din 107.614 noduri și 13.673.453 muchii;
- Facebook ("Social circles: Facebook"): Dataset ce conține listele de prieteni de pe Facebook, informațiile fiind colectate de la participanții unor studii pe "Facebook app". Dataset-ul formează un graf nedirecționat constituit din 4.039 noduri și 88.234 muchii;
- GitHub ("GitHub Social Network"): Dataset ce conține oameni activi pe GitHub care se urmăresc între ei. Dataset-ul formează un graf nedirecționat constituit din 37.700 noduri și 289.003 muchii.

Pentru fiecare din eșantioanele prezentate anterior s-au realizat verificări pe diferite eșantioane: 100, 1000, 10.000, 100.000 și întregul dataset.

```

214328887 34428380
17116707 28465635
380580781 18996905
221036078 153460275
107830991 17868918
151338729 222261763
19705747 34428380
222261763 88323281
19933035 149538028
158419434 17434613
149538028 153226312
364971269 153226312
100581193 279787626
113058991 69592091
151338729 187773078
406628822 262802533
460282402 88323281
280935165 437804658
222261763 27633075
285312927 151338729
279787626 131613362
158419434 17675120
394263193 100581193
254839786 88323281
204317520 21548772
67864340 172883064
270119528 297801196

```

Fig.6. Exemplu de eșantion din setul "Twitter"

4. Rezultate și evaluarea acestora

În urma rularii programului pe un set de date format din cercuri sociale de prieteni de pe platforma socială "Twitter", s-a observat că după 10 iterații ale algoritmului asupra grafului, atât valorile autoritatilor cât și ale grafulor s-au normalizat spre o singură valoare unitară. Din aceste calcule se poate observa că atenuarea valorilor autoritatilor și hub-urilor este realizată rapid din punct de vedere al iterațiilor, fiind foarte eficient în găsirea imediată a răspunsurilor dorite.

```

Node: 214328887 -- Auth: 0 -- Hub: 0.015151515151515152
Node: 34428380 -- Auth: 0.0136986301369863 -- Hub: 0.009493078937523384
Node: 17116707 -- Auth: 0.008563287671232878 -- Hub: 0.015151515151515152
Node: 28465635 -- Auth: 0.0136986301369863 -- Hub: 0
Node: 380580781 -- Auth: 0.01764733334794521 -- Hub: 0.015151515151515152
Node: 18996905 -- Auth: 0.0136986301369863 -- Hub: 0
Node: 221036078 -- Auth: 0 -- Hub: 0.026707504734848488
Node: 153460275 -- Auth: 0.007063356164383561 -- Hub: 0
Node: 107830991 -- Auth: 0.010959256440611083 -- Hub: 0.015151515151515152
Node: 17868918 -- Auth: 0.0136986301369863 -- Hub: 0
Node: 151338729 -- Auth: 0.010216003139269406 -- Hub: 0.03871391382965457
Node: 222261763 -- Auth: 0.010411452158478437 -- Hub: 0.05456320340057856
Node: 19705747 -- Auth: 0 -- Hub: 0.015151515151515152
Node: 88323281 -- Auth: 0.022405197631278538 -- Hub: 0
Node: 19933035 -- Auth: 0 -- Hub: 0.015151515151515152
Node: 149538028 -- Auth: 0.0136986301369863 -- Hub: 0.011379224342187307

```

Fig.7. Rezultatele după 5 iterații

```

Node: 214328887 -- Auth: 0.0032625926845726485 -- Hub: 0.009907267382273308
Node: 34428380 -- Auth: 0.023927425171368883 -- Hub: 0.003298167029871596
Node: 17116707 -- Auth: 0.0054424962683631475 -- Hub: 0.011005331854211757
Node: 28465635 -- Auth: 0.01304713533469059 -- Hub: 0.00879216591631387
Node: 380580781 -- Auth: 0.005443882739421821 -- Hub: 0.005503206755244298
Node: 18996905 -- Auth: 0.008738068610653441 -- Hub: 0.0010990630914830058
Node: 221036078 -- Auth: 0.004351250640615458 -- Hub: 0.009894061579780454
Node: 153460275 -- Auth: 0.004349272658765832 -- Hub: 0.0010996732841222158
Node: 107830991 -- Auth: 0.00870689760186844 -- Hub: 0.012111220921415441
Node: 17868918 -- Auth: 0.015239028871468838 -- Hub: 0
Node: 151338729 -- Auth: 0.014153069176043268 -- Hub: 0.014299110799584493
Node: 222261763 -- Auth: 0.007623495786210507 -- Hub: 0.012111612823149848
Node: 19705747 -- Auth: 0.003263496341432622 -- Hub: 0.0011000441583216257
Node: 88323281 -- Auth: 0.014211906470951923 -- Hub: 0.002198599936401756
Node: 19933035 -- Auth: 0.0021754891980913854 -- Hub: 0.005497673668119314
Node: 149538028 -- Auth: 0.005436843678458254 -- Hub: 0.012104998124978249

```

Fig.8. Rezultatele după 10 iterații


```

Node: 214328887 -- Auth: 0.0032809254730473475 -- Hub: 0.00993769784166706
Node: 34428380 -- Auth: 0.024060120136966935 -- Hub: 0.0033125659471625
Node: 17116707 -- Auth: 0.0054682091222487935 -- Hub: 0.011041886490774868
Node: 28465635 -- Auth: 0.013123701892575612 -- Hub: 0.008833509192444387
Node: 380580781 -- Auth: 0.005468209122258372 -- Hub: 0.0055209432452883464
Node: 18996905 -- Auth: 0.008749134595484757 -- Hub: 0.001104188649021634
Node: 221036078 -- Auth: 0.004374567297628145 -- Hub: 0.009937697841634018
Node: 153460275 -- Auth: 0.0043745672974517215 -- Hub: 0.0011041886490474123
Node: 107830991 -- Auth: 0.0087491345950065 -- Hub: 0.012146075140330736
Node: 17868918 -- Auth: 0.015310985542001837 -- Hub: 0
Node: 151338729 -- Auth: 0.014217343717315742 -- Hub: 0.014354452437919949
Node: 222261763 -- Auth: 0.007655492771273756 -- Hub: 0.012146075140222408
Node: 19705747 -- Auth: 0.0032809254729827134 -- Hub: 0.0011041886490797553
Node: 88323281 -- Auth: 0.014217343717990733 -- Hub: 0.002208377298002715
Node: 19933035 -- Auth: 0.002187283648712131 -- Hub: 0.005520943245220583

```

Fig.9. Rezultatele dupa 100 iteratii

In urma rularii, s-a constatat problema depistata la punctul 3, aceea ca algoritmul SALSA este dependent de interogari (si mai ales de numarul acestora), valorile multor noduri devenind 0, atat la autoritati cat si la hub-uri, din cauza lipsei unor valori catre care vectorii sa fie indreptati. In acest sens, cu cat marimea esantionului este mai mica, cu atat algoritmul este mai putin precis (lucru care se poate aplica, de asemenea, si la cautarile specifice unde informatia nu este atat de abundenta (Fig.9, Fig.10, Fig.11).

```

Node: 214328887 -- Auth: 0 -- Hub: 0.1
Node: 34428380 -- Auth: 0.1111111111111111 -- Hub: 0
Node: 17116707 -- Auth: 0 -- Hub: 0.1
Node: 28465635 -- Auth: 0.1111111111111111 -- Hub: 0
Node: 380580781 -- Auth: 0 -- Hub: 0.1
Node: 18996905 -- Auth: 0.1111111111111111 -- Hub: 0
Node: 221036078 -- Auth: 0 -- Hub: 0.1
Node: 153460275 -- Auth: 0.1111111111111111 -- Hub: 0
Node: 107830991 -- Auth: 0 -- Hub: 0.1
Node: 17868918 -- Auth: 0.1111111111111111 -- Hub: 0
Node: 151338729 -- Auth: 0 -- Hub: 0.1
Node: 222261763 -- Auth: 0.1111111111111111 -- Hub: 0.1
Node: 19705747 -- Auth: 0 -- Hub: 0.1

```

Fig.10. Cautari pe un esantion de 10 noduri din set

```

Node: 214328887 -- Auth: 0 -- Hub: 0.015151515151515152
Node: 34428380 -- Auth: 0.0136986301369863 -- Hub: 0.010771099244343144
Node: 17116707 -- Auth: 0.008771513876715483 -- Hub: 0.015151515151515152
Node: 28465635 -- Auth: 0.0136986301369863 -- Hub: 0
Node: 380580781 -- Auth: 0.017532606278677164 -- Hub: 0.015151515151515152
Node: 18996905 -- Auth: 0.0136986301369863 -- Hub: 0
Node: 221036078 -- Auth: 0 -- Hub: 0.021211787917171463
Node: 153460275 -- Auth: 0.008219479417351827 -- Hub: 0
Node: 107830991 -- Auth: 0.009125382970396011 -- Hub: 0.015151515151515152
Node: 17868918 -- Auth: 0.0136986301369863 -- Hub: 0
Node: 151338729 -- Auth: 0.011417754708904111 -- Hub: 0.032504402998127985
Node: 222261763 -- Auth: 0.00911846388420818 -- Hub: 0.043123057829303035
Node: 19705747 -- Auth: 0 -- Hub: 0.015151515151515152
Node: 88323281 -- Auth: 0.017164258632925197 -- Hub: 0
Node: 19933035 -- Auth: 0 -- Hub: 0.015151515151515152

```

Fig.11. Cautari pe un esantion de 100 noduri din set

```

Node: 21432887 -- Auth: 0.0032625926845726485 -- Hub: 0.009907267382273308
Node: 34428380 -- Auth: 0.023927425171368883 -- Hub: 0.003298167029871596
Node: 17116707 -- Auth: 0.0054424962683631475 -- Hub: 0.011005331854211757
Node: 28465635 -- Auth: 0.01304713533469059 -- Hub: 0.00879216591631387
Node: 380580781 -- Auth: 0.005443882739421821 -- Hub: 0.005503206755244298
Node: 18996905 -- Auth: 0.008738068610653441 -- Hub: 0.0010990630914830058
Node: 221036078 -- Auth: 0.004351250640615458 -- Hub: 0.009894061579780454
Node: 153460275 -- Auth: 0.004349272658765832 -- Hub: 0.0010996732841222158
Node: 107830991 -- Auth: 0.00870689760186844 -- Hub: 0.012111220921415441
Node: 17868918 -- Auth: 0.015239028871468838 -- Hub: 0
Node: 151338729 -- Auth: 0.014153069176043268 -- Hub: 0.014299110799584493
Node: 222261763 -- Auth: 0.007623495786210507 -- Hub: 0.012111612823149848
Node: 19705747 -- Auth: 0.003263496341432622 -- Hub: 0.0011000441583216257
Node: 88323281 -- Auth: 0.014211906470951923 -- Hub: 0.002198599936401756
Node: 19933035 -- Auth: 0.0021754891980913854 -- Hub: 0.005497673668119314

```

Fig.12. Cautari pe un esantion de 1000 noduri din set

O alta observatie in urma diferitelor rulari este aceea a implementarii tipurilor de grafuri, lucru ce influenteaza semnificativ rezultatele finale. In cadrul implementarii s-a considerat realizarea unui graf directionat pentru a putea obtine mai usor marginile care intra si ies din diferite noduri, insa la verificarea setului de date cu grafuri nedirectionate, numarul valorilor hub-urilor egale cu 0 au crescut considerabil. Acest lucru, insa, va putea fi remediat prin intermediul unei implementari diferite a programului.

Pe langa aceasta implementare, s-a constatat impactul major pe care ponderea initiala a valorii $A(u)$ (respectiv $H(u)$) o are asupra rezultatului final, autoritatilor indreptandu-se spre valori unitare in cazul echivalarii initiale a autoritatilor si hub-urilor cu 1.

5. Bibliografie

- [1] R. Lempel and S. Moran. 2001. SALSA: the stochastic approach for link-structure analysis. *ACM Trans. Inf. Syst.* 19, 2 (April 2001), 131–160. <https://doi.org/10.1145/382979.383041>
- [2] Allan Borodin, Gareth O. Roberts, Jeffrey S. Rosenthal, and Panayiotis Tsaparas. 2005. Link analysis ranking: algorithms, theory, and experiments. *ACM Trans. Internet Technol.* 5, 1 (February 2005), 231–297. <https://doi.org/10.1145/1052934.1052942>
- [3] Najork, Marc. (2007). Comparing the effectiveness of hits and salsa. 157-164. 10.1145/1321440.1321465
- [4] Marc A. Najork. 2007. Comparing the effectiveness of hits and salsa. In *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management (CIKM '07)*. Association for Computing Machinery, New York, NY, USA, 157–164. <https://doi.org/10.1145/1321440.1321465>
- [5] Farahat, Ayman, Thomas LoFaro, Joel C. Miller, Gregory Rae, and Lesley A. Ward. “Authority Rankings from HITS, PageRank, and SALSA: Existence, Uniqueness, and Effect of Initialization.” *SIAM Journal on Scientific Computing* 27, no. 4 (2006): 1181–1201. doi:10.1137/S1064827502412875
- [6] Mellah, Mohamed & Amine, Abdelmalek & Hamou, Reda & Kumar, A.V.. (2014). Link Analysis for Communities Detection on Facebook. *International Journal of Data Mining And Emerging Technologies*. 4. 10.5958/2249-3220.2014.00017.2.
- [7] <http://pubs.sciepub.com/ajss/3/2/3/index.html>, accesat: 22.05.2022, 14:00
- [8] https://docs.oracle.com/cd/E56133_01/2.4.0/reference/algorithms/salsa.html, accesat: 22.05.2022, 14:00
- [9] <https://snap.stanford.edu/data/index.html#web>, accesat: 23.05.2022, 14:00
- [10] <https://slideplayer.com/slide/6060463/>, accesat: 22.05.2022, 14:00