

## Cursul 14

### Algoritmul PageRank

#### 14.1 Construcția lanțului Markov pe graful WEB. Algoritmul PageRank

Succesul extraordinar și dominația motorului Google se datorează în principal algoritmului PageRank, care exploatează structura linkurilor din WWW pentru a determina un indice de popularitate al fiecărei pagini, independent de interogarea formulată de utilizator.

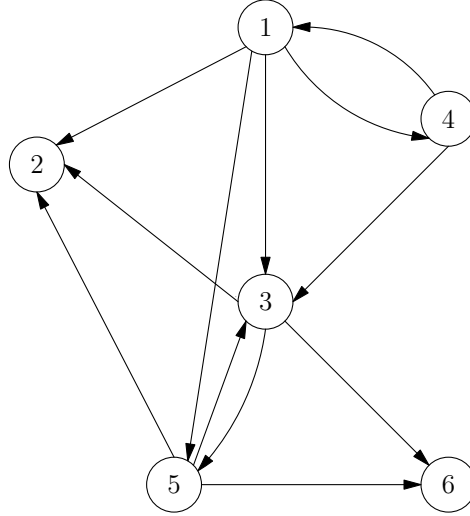
Documentele de pe WEB (paginile WEB) sunt identificate de aplicațiile software ale motorului, numite roboți sau *crawlere*. Documentele sunt apoi indexate. Modulul de indexare extrage cuvintele cheie, constituind așa numitul sac de cuvinte. Un alt modul, numit *query module* (modulul de interogare), convertește cererea formulată de utilizator, în limbaj natural, într-un vector cerere, cu care consultă indexul de conținut și extrage paginile relevante cererii. Modulul de ierarhizare ordonează descrescător aceste pagini în funcție de coeficienții de popularitate. PageRank-ul este un vector ale cărui coordonate sunt coeficienții de popularitate ai paginilor WEB identificate de crawler. Acest vector este distribuția de echilibru a unui lanț Markov definit pe graful WEB.

Să definim mai întâi lanțul Markov ce stă la baza algoritmului **PageRank**. Considerăm  $W = \{1, 2, \dots, m\}$  mulțimea tuturor paginilor WEB,  $H = (h_{ij})$  matricea de conectivitate a lui  $W$  sau matricea hyperlink:

$$h_{ij} = \begin{cases} 1, & \text{dacă există link în pagina } i \text{ către pagina } j, \\ 0, & \text{dacă nu există link în pagina } i \text{ către pagina } j. \end{cases}$$

Se spune că  $H$  este o matrice rară, căci are foarte multe zerouri (în medie, 3-10 elemente sunt nenule pe o linie). Suma elementelor de pe linia  $i$  a matricei  $H$  indică numărul de out-linkuri, adică numărul de linkuri din pagina  $i$  către alte pagini sau ea însăși. Notăm această sumă cu  $r_i = \sum_{j=1}^m h_{ij}$ .  $r_i$  se numește ordinul ieșirilor din pagina  $i$ . Suma elementelor de pe coloana  $i$  a matricei hyperlink indică numărul de in-linkuri ale paginii  $i$ , adică numărul de linkuri către pagina  $i$ .

Larry Page și Serghei Brin au definit un mers aleator pe graful WEB considerând că un surfer ajuns în pagina  $i$  alege cu aceeași probabilitate oricare din paginile către care



**Fig.14.1:** Graf orientat ilustrând linkurile între 6 pagini WEB.

aceasta are linkuri, prin urmare probabilitatea de a trece din pagina  $i$  în pagina  $j$  este:

$$p_{ij} = \begin{cases} \frac{1}{r_i}, & \text{dacă există link în pagina } i \text{ către pagina } j, \\ 0, & \text{dacă nu există link în pagina } i \text{ către pagina } j. \end{cases}$$

De exemplu, dacă

$$H = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

atunci ordinul de ieșire din pagina 5 este  $r_5 = 3$ , deci probabilitatea de a trece din pagina 5 în oricare din paginile  $\{1, 2, \dots, 6\}$  este  $p_{5j} = h_{5j}/3$ . Altfel spus, cu aceeași probabilitate de  $1/3$  un surfer poate trece din pagina 5 în pagina 2, 3 sau 6.

Vom exemplifica construcția propusă de L. Page și S. Brin prin modelul simplu de rețea izolată de pagini WEB (rețea intranet) din Fig. 14.1.

Notăm cu  $Q = (p_{ij})_{i,j=1,6}$ , matricea probabilităților de tranziție definite mai sus. Se observă din structura grafului de conectivitate că paginile 2 și 6 sunt pagini ce nu conțin linkuri către alte pagini. Acestea se numesc *dangling pages*. De exemplu, fișierele **pdf**, **ps** sau fișierele imagine sunt pagini dangling. Prin urmare, liniile 2 și 6 din matricea de tranziție au toate elementele nule și, astfel,  $Q$  nu este o matrice stohastică, deci nu poate fi interpretată ca matricea de tranziție a unui lanț Markov cu spațiul stărilor  $\{1, 2, 3, 4, 5, 6\}$ .

$$Q = \begin{array}{c|cccccc} & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 1 & 0 & 1/4 & 1/4 & 1/4 & 1/4 & 0 \\ 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 3 & 0 & 1/3 & 0 & 0 & 1/3 & 1/3 \\ 4 & 1/2 & 0 & 1/2 & 0 & 0 & 0 \\ 5 & 0 & 1/3 & 1/3 & 0 & 0 & 1/3 \\ 6 & 0 & 0 & 0 & 0 & 0 & 0 \end{array}$$

Pentru a remedia această problemă, Page și Brin au propus ca vector de probabilitate dintr-o pagină dangling  $i$  distribuția uniformă, considerând

$$p_{ij} = 1/m, \quad j = \overline{1, m}.$$

Adică, în mod artificial, se adaugă linkuri dintr-o pagină dangling către toate paginile WEB sau, echivalent, ajuns într-o pagină dangling, un navigator poate apoi alege cu o probabilitate uniformă orice pagină din WWW. Astfel, matricea stohastică obținută din matricea  $Q$  este:

$$\tilde{Q} = \begin{array}{c|cccccc} & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 1 & 0 & 1/4 & 1/4 & 1/4 & 1/4 & 0 \\ 2 & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} \\ 3 & 0 & 1/3 & 0 & 0 & 1/3 & 1/3 \\ 4 & 1/2 & 0 & 1/2 & 0 & 0 & 0 \\ 5 & 0 & 1/3 & 1/3 & 0 & 0 & 1/3 \\ 6 & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} \end{array}$$

Lanțul Markov definit de matricea stohastică  $\tilde{Q}$  nu este în general ireductibil (adică nu există drum de linkuri între orice două pagini sau, echivalent, graful WEB nu este tare conex) și pot exista traiectorii periodice, adică surferul, navigând conform matricei de tranziție  $\tilde{Q}$ , ar putea fi prins, ca într-o capcană, într-o mișcare aleatoare ciclică. Din acest motiv, dar și pentru că la un moment dat, și în realitate, orice surfer renunță să navigheze urmând linkurile din pagini, L. Page și S. Brin au introdus ipoteza că doar cu probabilitatea  $\alpha \in (0, 1)$  surferul navighează conform matricei  $\tilde{Q}$  și cu probabilitatea  $1 - \alpha$  ignoră linkurile și alege cu probabilitate uniformă oricare din paginile de pe WEB, introducând adresa URL în linia de comandă a browser-ului. Probabilitatea  $\alpha$  se numește *factor de damping*. În lucrarea inițială a fondatorilor Google,  $\alpha$  era menționat ca având valoarea 0.85. Cu această modificare matricea de tranziție este:

$$G = \alpha \tilde{Q} + (1 - \alpha) \underbrace{\begin{bmatrix} 1/m & 1/m & \dots & 1/m \\ 1/m & 1/m & \dots & 1/m \\ & & \vdots & \\ 1/m & 1/m & \dots & 1/m \end{bmatrix}}_E.$$

Matricea  $G$  se numește matricea Google, iar matricea  $E$ , de elemente identice  $1/m$ , se numește *matricea de teleportare*, deoarece surferul se teleportează din navigarea aleatoare urmând linkuri într-o "navigare artificială". Evident că și matricea  $E$  este o matrice stohastică, iar  $G$  fiind o combinație convexă de astfel de două matrice este o matrice stohastică (vezi proprietățile matricelor stohastice). Mai mult,  $G(i, j) > 0, \forall i, j = \overline{1, m}$ , deci matricea Google este ireductibilă și aperiodică.

Se presupune că matricea Google este cea mai "uriasă" matrice cu care se lucrează în vreo aplicație la ora actuală.

Lanțul Markov având:

- spațiul stărilor constituit din mulțimea paginilor WEB, de cardinal  $m$
- matricea de tranziție de tipul  $G$ , cu  $\alpha$  fixat
- distribuția inițială de probabilitate  $\pi_0$  (distribuția uniformă sau oricare alta)

este un lanț ireductibil și aperiodic, deci are o unică distribuție de echilibru  $\pi$ , numită vectorul PageRank.

PageRank-ul,  $\pi$ , este limita șirului  $(\pi_n)$ , definit prin  $\pi_n^T = \pi_0^T G^n$ . Limita este aceeași indiferent de distribuția inițială de probabilitate  $\pi_0$ , adică indiferent cu ce probabilitate surferul alege pagina din care începe navigarea.  $\pi(j)$  se numește PageRank-ul paginii  $j$  și reprezintă șansa asimptotică pe care o are pagina  $j$  de a fi vizitată de navigatorul aleator sau proporția din timpul de navigare pe care surferul ar petrece-o vizitând pagina  $j$ . Deci,  $\pi(j)$  este un indice de popularitate al paginii.

Când un utilizator introduce cuvinte cheie în bara de căutare, motorul Google caută paginile ce conțin cuvintele cheie și le afișează în ordinea descrescătoare a PageRank-ului lor.

Remarcăm că PageRank-ul unei pagini este independent de interogarea formulată de utilizator. Ea depinde doar de structura grafului WEB și se poate calcula offline. PageRank-ul se calculează la intervale regulate de timp. Până în 2008 se calcula lunar, dar acum se actualizează la intervale mai scurte de timp.

Vectorul PageRank se calculează numeric, folosind așa numita metodă a puterii, adică se calculează recursiv. Pornind de la  $\pi_0$  și  $G$ , se determină distribuțiile (sau PageRank-ul la  $n$  pași de navigare)  $\pi_n^T = \pi_{n-1}^T G$ . Se consideră că metoda a atins stadiul de convergență (adică s-a ajuns la echilibru) într-o etapă  $n$  în care  $\|\pi_n - \pi_{n-1}\| < \varepsilon$ , unde  $\varepsilon$  este un număr pozitiv foarte mic, prescris.

Pseudocodul algoritmului de calcul al PageRank-ului este:

```

1: function PageRank(G, m);
2:    $\pi = [1/m, 1/m, \dots, 1/m]$ ; //Distributia initiala de probabilitate
3:    $\text{eps} = 10^{-7}$ ;
4:   do
5:      $\pi' = \pi$ ;
6:      $\pi = \pi' * G$ ;
7:   while ( $\|\pi - \pi'\| \geq \text{eps}$ );
8:   return  $\pi$ ;
9: end function
```

S-a demonstrat că viteza de convergență a metodei puterii este aceeași cu rata de convergență a lui  $\alpha^n$ , unde  $\alpha$  este factorul de damping.

**Implicații asupra PageRank-ului.** Din punctul de vedere al vitezei de convergență ar fi preferabil un factor  $\alpha$  cât mai apropiat de zero. În acest caz, ținând seama că matricea Google este  $G = \alpha\tilde{Q} + (1 - \alpha)E$ , ar rezulta că se acordă o pondere redusă,  $\alpha$ , navigării conform linkurilor din graful WEB (cu modificarea pentru pagini dangling) și o pondere mai mare navigării artificiale, conform matricei de teleportare  $E$ . Cu alte cuvinte, în acest caz PageRank-ul asociat nu ar reflecta popularitatea reală a paginilor WEB. De aceea o valoare rezonabilă, așa cum a fost ea aleasă inițial de Larry Page și Serghei Brin,  $\alpha = 0.85$ , conduce la rezultate mai apropiate de realitate și la o viteză de convergență suficient de bună (un reprezentant Google a declarat că metoda puterii converge după 100-200 de iterații). Dacă vreți să aflați PageRank-ul unor pagini WEB intrați aici:

[http://www.prchecker.info/check\\_page\\_rank.php](http://www.prchecker.info/check_page_rank.php)

## 14.2 Pagerank-ul personalizat

Pentru o ierarhizare personalizată a paginilor WEB, matricea  $E$  se calculează luând în considerare vectorul personalizat  $w$ , care este un vector probabilist  $w = [a_1, a_2, \dots, a_m]^T$  ale cărui coordonate reprezintă probabilitatea ca surferul, ce iese din navigarea conform linkurilor, să aleagă pagina  $1, 2, \dots, m$  din WEB. Cu alte cuvinte, el nu alege o pagină în mod uniform, ci are anumite preferințe, identificate de motor în decursul timpului. Astfel, matricea de teleportare va fi  $E = ew^T$ , unde  $e = [1, 1, \dots, 1]^T$ , iar matricea Google corespunzătoare este

$$G = \alpha\tilde{Q} + (1 - \alpha)ew^T.$$

Distribuția de echilibru corespunzătoare este PageRank-ul personalizat.