

5. Căutarea pe web

Puncte importante :

- *Rangul paginii*, pentru descoperirea celor mai "importante" pagini de web, utilizat de Google.
- *Indecși și autorități*, o evaluare mai detaliată a importanței paginilor web utilizând o variantă a calculului de valori proprii utilizată pentru rangul paginii.

5.1. Rangul paginii

Intuitiv rezolvăm problema definiției "importanței" recursiv : o pagina este importantă dacă pagini importante conțin legături către ea.

Creăm o matrice stohastică a Internetului astfel :

1. Fiecare pagină i corespunde liniei i și coloanei i a matricii.
2. Dacă pagina j are n succesori (legături), atunci elementul i, j al matricii este $1/n$ dacă pagina i este unul dintre acești succesori ai paginii j și 0 altfel.

Intuiția care stă în spatele acestei matrici este :

- Să ne imaginăm că inițial fiecare pagină are o unitate de importanță. La fiecare pas fiecare pagină își împarte importanța între succesorii săi și primește noi fracțiuni de importanță de la predecesorii săi.
- Eventual, importanța fiecărei pagini atinge o limită care este componenta corespunzătoare ei din vectorul principal de valori proprii al matricii.
- Această importanță este de asemenea probabilitatea ca un navigator pe web, pornind de la o pagină aleatoare și urmând legături aleator alese din fiecare pagină, să ajungă la pagina în discuție după o lungă serie de legături.

Exemplul 5.1 : În 1839 Internetul consta din doar trei pagini : Netscape, Amazon și Microsoft. Legăturile între aceste trei pagini erau ca în Fig. 9.

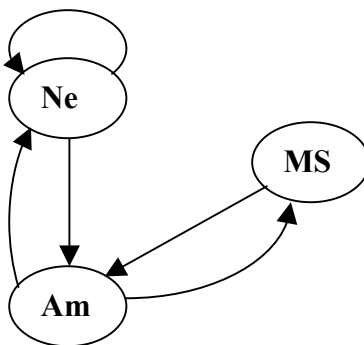


Figura 9: Internetul în 1839

Fie $[n, m, a]$ vectorul importanței pentru cele trei pagini : Netscape, Microsoft respectiv Amazon. Atunci ecuația care descrie valorile asimptotice ale acestor trei variabile este :

$$\begin{bmatrix} n \\ m \\ a \end{bmatrix} = \begin{bmatrix} 1/2 & 0 & 1/2 \\ 0 & 0 & 1/2 \\ 1/2 & 1 & 0 \end{bmatrix} \begin{bmatrix} n \\ m \\ a \end{bmatrix}$$

De exemplu, prima coloană a matricii reflectă faptul că Netscape își divide importanța între el însuși și Amazon. A doua coloană că Microsoft dă toată importanța sa către Amazon.

Putem rezolva ecuații ca aceasta începând cu aserțiunea că $n = m = 1$ și aplicând repetat matricea la estimarea curentă a acestor valori. Primele patru iterații dau următoarele estimări :

$$\begin{array}{rcl}
 n & = & 1 \quad 1 \quad 5/4 \quad 9/8 \quad 5/4 \\
 m & = & 1 \quad 1/2 \quad 3/4 \quad 1/2 \quad 11/16 \\
 a & = & 1 \quad 3/2 \quad 1 \quad 11/8 \quad 17/16
 \end{array}$$

La limită, soluția este $n = a = 6/5$; $m = 3/5$. Adică Netscape și Amazon au fiecare aceeași importanță și de două ori mai mare decât importanța Microsoft (ei, asta se întâmpla în 1839).

- De notat că nu putem să obținem niciodată valorile absolute ale lui n , m și a ci doar raportul lor, de vreme ce aserțiunea inițială că fiecare a pornit de la 1 a fost arbitrară.
- Deoarece matricea este stohastică (suma pe fiecare coloană este 1), procesul de *relaxare* de mai sus converge către vectorul principal de valori proprii.

5.2. Probleme cu grafuri reale ale Internetului

1. *Dead end*: o pagină care nu are succesori nu are către cine să-și trimită importanța. Eventual, toată importanța "se va scurge" din Internet
2. *Capcane*: un grup de una sau mai multe pagini care nu au legături către pagini din afara grupului vor acumula eventual toată importanța din Internet.

Exemplul 5.2 : Să presupunem că Microsoft încearcă să profite că este un monopol înlăturând toate legăturile din situl său. Noul Internet este ca în Fig. 10 iar matricea descriind tranzițiile este :

$$\begin{bmatrix} n \\ m \\ a \end{bmatrix} = \begin{bmatrix} 1/2 & 0 & 1/2 \\ 0 & 0 & 1/2 \\ 1/2 & 0 & 0 \end{bmatrix} \begin{bmatrix} n \\ m \\ a \end{bmatrix}$$

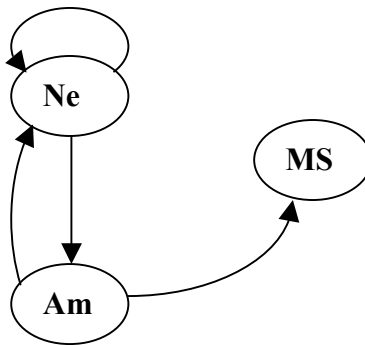


Figura 10: Microsoft devine dead end.

Primii patru pași ai soluției iterative sunt :

$$\begin{array}{rcl}
 n & = & 1 \quad 1 \quad 3/4 \quad 5/8 \quad 1/2 \\
 m & = & 1 \quad 1/2 \quad 1/4 \quad 1/4 \quad 3/16 \\
 a & = & 1 \quad 1/2 \quad 1/2 \quad 3/8 \quad 5/16
 \end{array}$$

Eventual, fiecare dintre n , m și a devin 0; i.e. toată importanța se scurge afară.

Exemplul 5.3 : Supărat de decizie, Microsoft decide să nu folosească decât legături către el însuși de acum încolo. Acum, Microsoft a devenit o capcană. Noul Internet este în Fig. 11, iar ecuația de rezolvat este:

$$\begin{bmatrix} n \\ m \\ a \end{bmatrix} = \begin{bmatrix} 1/2 & 0 & 1/2 \\ 0 & 1 & 1/2 \\ 1/2 & 0 & 0 \end{bmatrix} \begin{bmatrix} n \\ m \\ a \end{bmatrix}$$

Primii patru pași ai soluției sunt:

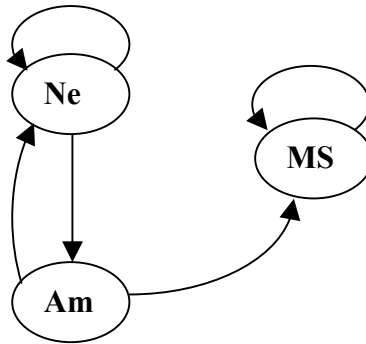


Figura 11: Microsoft devine o capcană.

$$\begin{array}{rcl}
 n & = & 1 \quad 1 \quad 3/4 \quad 5/8 \quad 1/2 \\
 m & = & 1 \quad 3/2 \quad 7/4 \quad 2 \quad 35/16 \\
 a & = & 1 \quad 1/2 \quad 1/2 \quad 3/8 \quad 5/16
 \end{array}$$

Acum m converge la 3 iar $n = a = 0$.

5.3. Soluția Google pentru dead end și capcane

În loc de a aplica matricea direct, "taxăm" fiecare pagină cu o fracțiune din importanța sa curentă și distribuim importanța taxată în mod egal tuturor paginilor.

Exemplul 5.4 : Dacă folosim o taxă de 20% ecuația din exemplul 5.3 devine:

$$\begin{bmatrix} n \\ m \\ a \end{bmatrix} = 0.8 \begin{bmatrix} 1/2 & 0 & 1/2 \\ 0 & 1 & 1/2 \\ 1/2 & 0 & 0 \end{bmatrix} \begin{bmatrix} n \\ m \\ a \end{bmatrix} + \begin{bmatrix} 0.2 \\ 0.2 \\ 0.2 \end{bmatrix}$$

Soluția acestei ecuații este $n = 7/11$; $m = 21/11$; $a = 5/11$.

- De notat că suma celor trei valori nu este 3 dar obținem o distribuție mult mai rezonabilă a importanței decât în Exemplul 5.3.

5.4. Procedee anti-spam la Google

"Spamming" este în acest context încercarea multor situri web de a părea că sunt despre un subiect care atrage navigatorii fără ca într-adevăr să fie despre acel subiect.

- Google, ca și alte motoare de căutare, încearcă să potrivească cuvintele din cererile de căutare cu cuvinte din pagini web. Cu toate acestea, Google, spre deosebire de alte motoare de căutare, tinde să creadă ceea ce spun alții în textul legăturilor despre o pagină web făcând mai greu pentru aceasta să pară ca fiind despre ceva ce nu este.
- Utilizarea rangului paginii pentru a măsura importanța în locul unei măsuri mult mai naive ca "numărul de legături către acea pagină" protejează de asemenea împotriva spamului. Măsura naivă poate fi înșelată de un spammer care creează 1000 de pagini care se referă între ele în timp ce rangul paginii recunoaște că nici una dintre acestea nu au importanță reală.

5.5. Indecși și autorități

Intuitiv, definim "index" și "autoritate" într-un mod mutual recursiv: un index conține legături către multe autorități iar o autoritate este referită de mulți indecși.

- Autoritățile pot fi pagini care oferă informații despre un subiect, e.g. pagina Quest despre proiectul IBM de data mining.
- Indecșii sunt pagini care nu furnizează informații ci spun unde se găsesc informații, e.g. pagina cursului CS345.
- Utilizează o formalizare matricială similară cu cea de la rangul paginii dar fără restricția stochastică. Numărăm fiecare legătură ca 1, indiferent de câți succesori sau predecesori are o pagină.
- Aplicarea repetată a matricii duce la divergență, dar putem introduce un factor de scalare pentru a ține valorile calculate pentru gradul de "autoritate" sau de "indexare" pentru fiecare pagină între limite finite.

Definim matricea A ale cărei linii și coloane corespund paginilor web având elementul $A_{ij} = 1$ dacă pagina i referă pagina j și 0 altfel.

- De notat că A^T , transpusa lui A , arată ca matricea utilizată pentru calculul rangului paginilor dar A^T are 1 acolo unde matricea pentru rang are fracții.

Fie a și h doi vectori iar componenta lor i corespunde gradului de autoritate respectiv indexare a paginii i . Fie λ și μ factorii de scalare corespunzători care vor fi calculați mai târziu. Atunci putem afirma că:

1. $h = \lambda A a$. Adică gradul de indexare al fiecărei pagini este suma gradelor de autoritate ale tuturor paginilor referite, scalată cu λ .
2. $a = \mu A^T h$. Adică gradul de autoritate al fiecărei pagini este suma gradelor de indexare ale tuturor paginilor care o referă, scalată cu μ .

Din (1) și (2) putem deduce folosind substituția, două ecuații care leagă vectorii a și h doar de ei înșiși:

$$a = \lambda \mu A^T A a; \quad h = \lambda \mu A A^T h$$

Ca urmare, putem calcula h și a prin relaxare, obținând vectorul principal de valori proprii al matricilor AA^T și respectiv $A^T A$

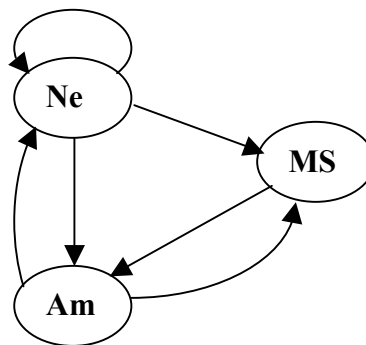


Figura 12: Internetul pentru exemplul 5.5

Exemplul 5.5 : Considerăm Internetul ca în Fig. 12. Matricile relevante sunt:

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \quad AA^T = \begin{bmatrix} 3 & 1 & 2 \\ 0 & 0 & 1 \\ 2 & 0 & 2 \end{bmatrix} \quad A^T A = \begin{bmatrix} 2 & 2 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 2 \end{bmatrix}$$

Dacă utilizăm $\lambda = \mu = 1$ și considerăm că vectorii $h = [h_n, h_m, h_a]$ și $a = [a_n, a_m, a_a]$ sunt inițial fiecare $[1, 1, 1]$, primele trei iterații ale ecuațiilor pentru a și h sunt:

$$\begin{aligned}
 \mathbf{a}_n &= 1 \quad 5 \quad 24 \quad 114 \\
 \mathbf{a}_m &= 1 \quad 5 \quad 24 \quad 114 \\
 \mathbf{a}_a &= 1 \quad 4 \quad 18 \quad 84
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{h}_n &= 1 \quad 6 \quad 28 \quad 132 \\
 \mathbf{h}_m &= 1 \quad 2 \quad 8 \quad 36 \\
 \mathbf{h}_a &= 1 \quad 4 \quad 20 \quad 96
 \end{aligned}$$

De exemplu, vectorul \mathbf{a} , scalat corespunzător, va converge către un vector în care $\mathbf{a}_n = \mathbf{a}_m$ și fiecare dintre aceste numere este mai mare ca \mathbf{a}_a în raportul $1 + \sqrt{3} : 2$ sau aproximativ 1.36.