

6. Extragerea de cunoștințe din web

Puncte importante :

1. *Numărarea dinamică a mulțimilor de articole*: Căutarea de mulțimi *interesante* de articole într-un spațiu mult prea mare pentru a se putea lua în considerare fiecare pereche de articole.
2. *"Cărți și autori"*: Intrigantul experiment al lui Sergey Brin de extragere de date relaționale din web.

6.1 Găsirea mulțimilor de articole neobișnuite

Problema este de a găsi mulțimi de cuvinte care apar "neobișnuit de des" împreună pe web, e.g. "New" și "York" sau {Ducea, de, York}.

- "Neobișnuit de des" poate fi definit în diverse moduri pentru a încorpora ideea că numărul de documente web conținând mulțimea de cuvinte este mult mai mare decât cel așteptat în cazul în care cuvintele ar fi fost alese la întâmplare, fiecare cuvânt cu probabilitatea sa de apariție într-un document.
- Un mod adecvat este *entropia per cuvânt din mulțime*. Formal, *interesul* unei mulțimi de cuvinte S este:

$$\frac{\log_2\left(\frac{\text{prob}(S)}{\prod_{w \in S} \text{prob}(w)}\right)}{|S|}$$

De notat că împărțim la dimensiunea lui S pentru a evita "efectul Bonferroni", în care sunt atât de multe mulțimi de o dimensiune dată încât unele, din motive probabilistice, par a fi corelate.

- Exemplu: Dacă *a*, *b* și *c* (cuvinte) apar fiecare în 1% din toate documentele și $S = \{a, b, c\}$ apar în 0.1% din documente, interesul lui S este $(\log_2(0.001/(0.01 \times 0.01 \times 0.01)))/3 = \log_2(1000)/3$ adică aproximativ 3.3.
- Problema tehnică: interesul nu este monoton sau "închis în jos" în modul de la produse cu larg suport. Asta înseamnă că putem avea o mulțime S cu o valoare mare a interesului și totuși unele sau chiar toate submulțimile sale stricte să nu fie interesante. Prin contrast, dacă S are suport larg, atunci toate submulțimile sale au cel puțin același suport.
- Problemă tehnică: Cu mai mult de 10^8 cuvinte diferite apărând pe web nu este posibil nici măcar să considerăm toate perechile de cuvinte.

6.2 Motorul DICE

DICE (dynamic itemset counting engine) vizitează repetat paginile web într-un mod de tip "round-robin". De fiecare dată numără aparițiile anumitor mulțimi de cuvinte și ale fiecărui cuvânt din aceste mulțimi. Numarul de mulțimi numărate este suficient de mic încât contorii lor încap în memoria centrală.

Din când în când, să spunem la fiecare 5000 de pagini, DICE își reconsideră mulțimile pentru care numără. Înlătură acele mulțimi care au cel mai mic interes și le înlocuiește cu alte mulțimi.

Alegerea noilor mulțimi se bazează pe *proprietatea muchiei grele (heavy edge property)* care este o observație justificată experimental că acele cuvinte care apar în mulțimi cu interes ridicat au probabilitatea mai mare să apară în alte mulțimi cu interes ridicat. Astfel, când selectează noi mulțimi pentru a începe numărarea, DICE este direcționat în favoarea cuvintelor care apar deja în mulțimi cu interes ridicat. Totuși, el nu se bazează exclusiv pe aceste cuvinte altfel nu ar putea niciodată să găsească mulțimi cu interes ridicat compuse din multele cuvinte pe care nu le-a considerat niciodată. Unele (dar nu toate) din construcțiile pe care le utilizează DICE pentru crearea noilor mulțimi sunt:

1. Două cuvinte aleatoare. Aceasta este singura regulă independentă de aserțiunea muchiei grele și ajută noi cuvinte să ajungă în mulțime.
2. Un cuvânt dintr-una din mulțimile interesante și un cuvânt aleator.

3. Două cuvinte din două mulțimi interesante diferite.
4. Reuniunea a două mulțimi interesante a căror intersecție are dimensiunea 2 sau mai mult.
5. $\{a, b, c\}$ dacă toate mulțimile $\{a, b\}$, $\{a, c\}$ și $\{b, c\}$ sunt găsite ca fiind interesante.

Bineînțeles, în general sunt mult prea multe opțiuni de a aplica cele de mai sus în toate modurile posibile astfel încât se utilizează o selecție aleatoare a opțiunilor dând o anumită șansă fiecăreia dintre ele.

6.3 Cărți și autori

Ideea principală este de a căuta pe web fapte de un anumit tip, de genul celor care ar putea forma o relație de genul *Cărți(titlu, autor)*. Procesarea este sugerată de Fig. 13.

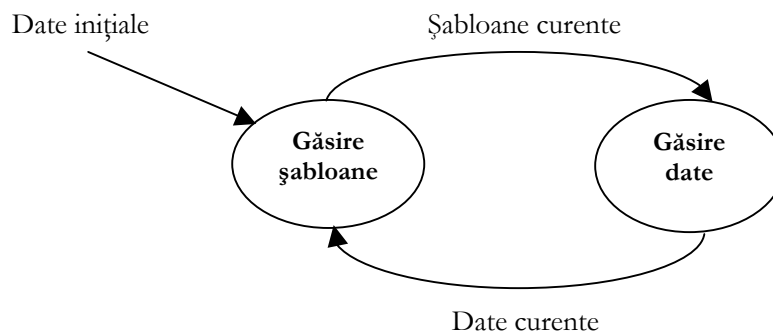


Figura 13: Extragerea relațiilor din web

1. Se pornește de la un eșantion al tuplelor care se doresc găsite. În exemplul discutat în lucrarea lui Brin au fost folosite cinci exemple de titluri de cărți și autori ai acestora.
2. Fiind dată o mulțime de exemple cunoscute, se caută pagini unde apar aceste date pe web. Dacă se găsește un șablon care identifică un număr de tupluri cunoscute și este suficient de specific încât e puțin probabil să identifice prea mult, atunci se acceptă acest șablon.
3. Fiind dată o mulțime de șabloane acceptate, se caută date care satisfac aceste șabloane și se adaugă la mulțimea datelor cunoscute.
4. Se repetă pașii (2) și (3) de un număr de ori. În exemplul citat au fost utilizate patru ciclări care au dus la 15,000 de tuple; aprox. 95% au fost perechi adevărate titlu-autor.

6.4 Cum arată un șablon?

Noțiunea constă în cinci elemente:

1. *Ordinea*; i.e. dacă titlul apare în text înaintea autorului sau vice-versa. Într-un caz general, în care tuplele au mai mult de 2 componente, ordinea va fi dată de permutarea componentelor.
2. *Prefixul adresei web (URL)*.
3. *Prefixul* textului, care apare înaintea primului dintre titlu și autor
4. *Mijlocul*: text care apare între cele două elemente de date.
5. *Suffixul* textului care urmează după al doilea dintre cele două elemente de date. Atât prefixul cât și suffixul au fost limitate la 10 caractere.

Exemplul 6.1 : Un șablon posibil poate consta din următoarele:

1. Ordinea: titlul și apoi autorul.
2. Prefixul URL: `www.stanford.edu/class`
3. Prefixul, mijlocul și sufixul de forma următoare:

`<I>titlu</I> de autor<P>`

Aici prefixul este `<I>`, mijlocul este `</I> de` (inclusiv spațiul după "de") și sufixul este `<P>`. Titlul este orice apare între prefix și mijloc; autorul este orice apare între mijloc și sufix.

Intuiția pentru care acest șablon poate fi bun este că există probabil o mulțime de liste cu referințe bibliografice în paginile cursurilor de la Stanford. □

Pentru a se focaliza pe șabloanele care pot fi corecte, Brin a utilizat o serie de constrângeri asupra șabloanelor, după cum urmează:

- Definim *specificitatea* unui șablon ca fiind produsul lungimilor prefixului, mijlocului, sufixului și prefixului URL. În mare, specificitatea măsoară cât de posibil este să găsim date care corespund șablonului; cu cât specificitatea este mai mare, cu atât ne așteptăm la mai puține apariții ale acestuia în date.
- Apoi șablonul trebuie să îndeplinească două condiții pentru a fi acceptat:
 1. Trebuie să fie cel puțin 2 elemente de date cunoscute care apar conform aceluși șablon.
 2. Produsul specificității șablonului cu numărul de apariții de date conform acestuia trebuie să depășească un anumit prag T (nespecificat).

6.5 Apariții ale datelor

O apariție a unui tuplu este asociată cu un șablon după care acestea apar; i.e., același titlu și autor pot să apară după diferite șabloane. Astfel, o apariție a datelor constă în:

1. Un anumit titlu și autor.
2. Adresa Internet completa (URL) și nu doar prefixul ca în cazul șablonului.
3. Ordinea, prefixul, mijlocul și sufixul șablonului după care au apărut titlul și autorul respectiv.

6.6 Găsirea aparițiilor pornind de la datele cunoscute

Dacă avem câteva perechi autor-titlu cunoscute, primul pas în găsirea de noi șabloane este căutarea pe web pentru a vedea unde apar aceste titluri și autori. Să presupunem că există un index al web-ului astfel încât dându-se un cuvânt putem găsi toate (legăturile către) paginile conținând acel cuvânt. Metoda utilizată este esențial a-priori:

1. Găsim (legături către) toate paginile conținând oricare dintre autorii cunoscuți. Cum numele autorilor constau în general din 2 cuvinte, se utilizează indexul pentru fiecare dintre prenume și nume și se verifică dacă aparițiile sunt consecutive în document
2. Găsim (legături către) toate paginile conținând oricare dintre titlurile cunoscute. Se pornește prin găsirea paginilor conținând toate cuvintele titlului și apoi verificând că aceste cuvinte apar în ordine în pagină.
3. Se intersectează mulțimile de pagini care au un autor și un titlu în ele. Doar pentru aceste pagini este nevoie să se facă operația de căutare pentru găsirea șabloanelor în care se găsește o pereche cunoscută autor-titlu. Pentru prefix și sufix se iau cele 10 caractere alăturate acestora sau mai puține dacă nu există 10.

6.7 Construcția șabloanelor din aparițiile de date

1. Se grupează aparițiile de date după ordinea și mijlocul lor. De exemplu, un grup din acest "group-by" poate corespunde ordinii "titlu-apoi-autor" și mijlocului " de ".
2. Pentru fiecare grup se găsește cel mai lung prefix, sufix și prefix URL comun.
3. Dacă testul de specificitate pentru acest șablon este îndeplinit, se acceptă șablonul.
4. Dacă testul de specificitate *m* este îndeplinit, se încearcă spargerea grupului în doua prin extinderea lungimii prefixului URL cu un caracter și apoi se repetă pasul (2). Dacă este imposibil să spargem grupul (pentru că există doar un URL) atunci am eșuat în a produce un nou șablon din acel grup.

Exemplul 6.6 : Să presupunem că grupul conține trei URL-uri:

```
www.stanford.edu/class/cs345/index.html  
www.stanford.edu/class/cs145/index.html  
www.stanford.edu/class/cs340/readings.html
```

Prefixul comun este `www.stanford.edu/class/cs`. Dacă trebuie să spargem grupul, atunci următorul caracter, 3 sau 1, sparge grupul în două, cu acele apariții ale datelor din prima pagină (pot fi multe astfel de apariții) mergând într-un prim grup și aparițiile din celelalte două pagini în celălalt.

6.8 Găsirea aparițiilor pornind de la șabloane

1. Se găsesc toate URL-urile care se potrivesc cu prefixul URL al cel puțin unui șablon.
2. Pentru fiecare astfel de pagină se parcurge textul folosind o expresie regulată construită din prefixul, mijlocul și sufixul șablonului.
3. Se extrage din fiecare potrivire titlul și autorul, după ordinea specificată în șablon.