

CAPITOLE SPECIALE DE INFORMATICĂ

## Curs 4: Calculul scorurilor în un sistem complet de extragere a informațiilor

18 octombrie 2018

Reamintim că în cursul precedent am prezentat modelul de spațiu vectorial în care fiecare document  $d$  din o colecție  $D$  de  $N$  documente care are vocabularul de termeni  $V = \{t_1, \dots, t_M\}$ , este reprezentat ca un vector

$$\vec{V}(d) = (\text{wf}_{t_1,d}, \dots, \text{wf}_{t_M,d})$$

unde  $\text{wf}_{t_i,d}$  este greutatea termenului  $t$  în documentul  $d$ . De obicei,  $\text{wf}_{t,d}$  se definește în unul din următoarele 2 feluri:

1.  $\text{wf}_{t,d} := \text{tf}_{t,d}$ , adică numărul de apariții ale termenului  $t$  în documentul  $d$ , sau
2.  $\text{wf}_{t,d} := \text{tf-idf}_{t,d} = \text{idf}_t \cdot \text{tf}_{t,d}$  unde  $\text{idf}_t = \log \frac{N}{\text{df}_t}$

Și cererile se reprezintă ca vectori: o cerere  $q = t_{i_1} \ t_{i_2} \ \dots \ t_{i_n}$  este descrisă de vectorul

$$\vec{V}(q) = (\text{w}_{t_1,q}, \dots, \text{w}_{t_M,q})$$

unde  $w_{t_j,q} = 1$  dacă  $j \in \{i_1, \dots, i_n\}$  și  $w_{t_j,q} = 0$  în caz contrar.

Valorile normalizate ale acestor vectori sunt vectorii de lungime 1 definiți astfel:

$$\vec{v}(d) := \frac{\vec{V}(d)}{|\vec{V}(d)|} \quad \text{și} \quad \vec{v}(q) := \frac{\vec{V}(q)}{|\vec{V}(q)|} = \frac{|\vec{V}(q)|}{\sqrt{n}}$$

Scorul (sau gradul) de similaritate cosinusoidală dintre  $d$  și  $q$  este cosinus unghiului dintre  $\vec{v}(q)$  și  $\vec{v}(d)$ , adică produsul scalar  $\vec{v}(q) \cdot \vec{v}(d)$ .

Sistemele complete de extragere a informațiilor facilitează determinarea rapidă a primelor  $K$  documente care se potrivesc cel mai bine cu o cerere de căutare  $q$ , adică acele documente  $d$  pentru care scorul  $\vec{v}(q) \cdot \vec{v}(d)$  este cât mai mare. Pentru performanță ridicată (adică timp de răspuns scurt, fără a compromite semnificativ calitatea răspunsului), se elimină niște calcule inutile și se folosesc metode euristice care nu garantează calculul exact al celor mai bune  $K$  potriviri, dar dau  $K$  răspunsuri apropiate de cele mai bune  $K$ .

În acest curs sunt prezentate (1) tehnici și metode euristice de calcul al documentelor cu cele mai bune scoruri de potrivire cu o cerere dată, (2) o arhitectură pentru un motor de căutare web, bazat pe generalizări ale indecșilor și metodelor de calcul al scorului de potrivire, și (3) extensii ale modului de interogare pentru modelul de reprezentare a documentelor în spațiu vectorial.

## Optimizare: evitarea normalizării vectorului $\vec{V}(q)$

Primele  $K$  documente care se potrivesc cel mai bine cu o cerere  $q$  sunt documentele  $d_1, \dots, d_K \in D$  astfel încât

- $\vec{v}(q) \cdot \vec{v}(d_i) \geq \vec{v}(q) \cdot \vec{v}(d_{i+1})$  pentru  $1 \leq i < K$ , și
- $\vec{v}(q) \cdot \vec{v}(d_K) \geq \vec{v}(q) \cdot \vec{v}(d)$  pentru toți  $d \in D \setminus \{d_1, \dots, d_K\}$ .

Se observă că

$$\begin{aligned} \vec{v}(q) \cdot \vec{v}(d) \geq \vec{v}(q) \cdot \vec{v}(d') &\Leftrightarrow \frac{\vec{V}(q)}{|\vec{V}(q)|} \cdot \vec{v}(d) \geq \frac{\vec{V}(q)}{|\vec{V}(q)|} \cdot \vec{v}(d') \\ &\Leftrightarrow \vec{V}(q) \cdot \vec{v}(d) \geq \vec{V}(q) \cdot \vec{v}(d') \end{aligned}$$

Deci, primele  $K$  documente care se potrivesc cel mai bine cu o cerere  $q$  sunt documentele  $d_1, \dots, d_K \in D$  pentru care  $\vec{V}(q) \cdot \vec{v}(d)$  iau valorile cele mai mari. Deoarece componentele nenele ale lui  $\vec{V}(q)$  sunt 1, rezultă că algoritmul de calcul al celor mai bune potriviri din cursul 3 poate fi optimizat astfel:

```

SCORCOSINUSOIDALRAPID(q)
1 float scoruri[N] = [0]
2 for fiecare  $d$  din colecția  $D$  do
3   inițializează lungime[d]
4 for fiecare termen  $t$  din cererea  $q$  do
5   obține lista de postări a lui  $t$ 
6   for fiecare pereche  $(d, tf_{t,d})$  din lista de postări a lui  $t$  do
7     scoruri[d] += wf_{t,d}
8 for fiecare document  $d \in D$  do
9   scoruri[d] := scoruri[d]/lungime[d]
10 return primele  $K$  componente din scoruri[ ]

```

Figure 1: Un algoritm de calcul mai rapid al scorurilor în modelul vectorial.

Extragerea primelor  $K$  componente din  $scoruri[ ]$  se poate face în 2 pași: (1) se sortează documentele din  $D$  în ordine descrescătoare a scorurilor (în timp  $O(N \log N)$ ) și se extrag primele  $K$  elemente din lista sortată de documente (în timp constant, fiindcă  $K$  este o valoare constantă). O alternativă mai eficientă este să se plaseze documentele  $d$  cu  $scoruri[d] > 0$  în o structură *heap*. Dacă în colecție sunt  $J$  documente  $d$  cu  $scoruri[d] > 0$ , construcția structurii heap se face folosind  $2 \cdot J$  pași de comparare, iar apoi extragerea primelor  $K$  elemente cu cele mai mari scoruri se poate face folosind  $\log J$  pași de comparare.

# 1 Calculul inexact al primelor $K$ potriviri

Idee: există tehnici cu timp de calcul mult mai scurt care, pentru o cerere  $q$ , extrag  $K$  documente din o colecție, pentru care probabilitatea este mare să fie printre primele  $K$  documente cu scor de potrivire maxim.

D.p.d.v. al utilizatorului, aceste tehnici de potrivire inexactă sunt adesea un lucru bun:

- timpul de răspuns este mult mai mic,
- măsura de similaritate cosinusoidală este doar o aproximare a gradului de relevanță perceput de utilizator.

Vom prezenta metode euristicice care au probabilitate mare să returneze  $K$  documente cu scoruri cosinusoidale apropriate de cele ale primelor  $K$  potriviri cosinusoidale. Aceste metode operează în 2 pași:

1. Se determină o mulțime de candidați  $A$  care are mai mult de  $K$  elemente, adică  $K < |A| \ll N$ . Nu este necesar ca  $A$  să conțină primele  $K$  potriviri pentru cererea  $q$ , dar sunt sănse mari ca  $A$  să conțină multe documente cu scoruri apropriate de primele  $K$  scoruri de potrivire.
2. Se returnează  $K$  documente din  $A$  care au cele mai mari scoruri de potrivire cu  $q$ .

## a) Eliminarea indecșilor

Această metodă euristică se aplică pentru cereri  $q = t_1 \dots t_n$  cu cel puțin 2 termeni de interogare (adică  $n > 1$ ).

- Idee de bază: se consideră doar documente  $d \in D$  care conțin cel puțin un termen din  $q$ .
- Alte euristică posibile:
  - Se consideră doar documente  $d \in D$  care conțin termeni din  $t$  cu  $\text{idf}_t$  mai mare decât o valoare prestabilită. În general, termenii  $t$  cu  $\text{idf}_t$  mic au liste de postări lungi  $\Rightarrow$  mulțimea de documente pentru care se calculează scoruri devine mult mai mică.
  - Se consideră doar documente  $d \in D$  care conțin mulți termeni din  $q$ .

## b) Liste de campioni

Listele de campioni se calculează pentru termeni din dicționar și depind de alegerea unui număr întreg pozitiv  $r > 0$ : lista de campioni a lui  $t \in V$  constă din  $r$  documente  $d$  care au cele mai mari valori pentru  $\text{wf}_{t,d}$ .

**Idee:** construcția lui  $A$  bazată pe liste de campioni se face astfel:  $A$  este reuniunea listelor de campioni pentru termenii din cererea  $q$ . Calculul scorurilor de potrivire se face doar pentru documente din  $A$ .

Valoarea lui  $r$  se fixează înainte de construcția indecsilor, și poate varia în funcție de termen: de pildă,  $r$  poate fi mai mare pentru termeni care apar mai rar.

### c) Scoruri și ordini statice de calitate

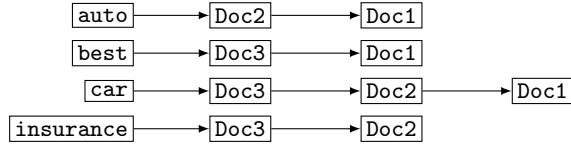
Numeiroase motoare de extragere a informațiilor calculează pentru fiecare document  $d$  o măsură  $g(d)$  a calității documentului, care nu depinde de vreo cerere. Astfel de măsuze se numesc scoruri statice de calitate. De obicei,  $0 \leq g(d) \leq 1$ . Valoarea lui  $g(d)$  poate fi determinată în mai multe feluri, ce de exemplu din numărul de comentarii favorabile primite de la vizitatori pe web.

Scorul net de potrivire a unui document  $d$  cu o cerere  $q$  este o combinație între  $g(d)$  și  $\text{wf}_{t,d}$ . De exemplu, poate fi

$$\text{scor-net}(q, d) = g(d) + \frac{\vec{V}(q) \cdot \vec{V}(d)}{|\vec{V}(q)| \cdot |\vec{V}(d)|}.$$

Pentru procesarea eficientă a listelor de postări, de exemplu pentru intersecția a două liste de postări, este necesar ca postările tuturor termenilor să fie ordonate la fel, de exemplu:

- în ordine crescătoare a identificatorilor de document, sau
- în ordine descrescătoare a valorilor statice  $g(d)$ . Exemplul de mai jos este pentru o colecție de trei documente Doc1, Doc2 și Doc3 cu scorurile statice de calitate  $g(1) = 0.25$ ,  $g(2) = 0.5$ ,  $g(3) = 1$ :



Valorile statice  $\{g(d) \mid d \in D\}$  ne permit să definim liste globale de campioni pentru o valoare fixată parametrului  $r$ : la fiecare termen  $t$  asociem ca postări postări  $r$  documente care au cele mai mari valori pentru  $g(d) + \text{tf-idf}_{t,d}$ . Listele de postări ale tuturor termenilor trebuie ordonate după o ordine comună; de exemplu în ordine crescătoare a identificatorului de document, sau în ordine descrescătoare a scorului static de calitate  $g(d)$ .

### d) Ordonare bazată pe impact

Algoritmi de calcul al rangului bazați pe gradarea după termeni nu necesită ca postările termenilor să fie ordonate după o ordine comună.

**Idee de bază:** ordonarea bazată pe impact ordinează documentele din lista de postări a unui termen  $t$  în ordine descrescătoare a valorilor lui  $\text{tf}_{t,d}$ . Prin urmare ordonarea documentelor poate varia de la o listă de postări a unui termen la lista de postări a altui termen.

Acest mod de ordonare a documentelor permite reducerea semnificativă a mulțimii de documente pentru care calculăm scoruri de potrivire:

1. Putem fixa criterii de oprire a traversării unei liste de postări pentru un termen  $t$ , de exemplu (1) după ce s-a traversat un număr de  $r$  postări; sau (2) după ce  $\text{tf}_{t,d}$  a devenit mai mică decât o anumită valoare.
2. Termenii care apar în cererea  $q$  se traversează în ordine descrescătoare a  $\text{idf}_t$ . Această euristică se bazează pe observația că termenii  $t$  din  $q$  care contribuie cel mai mult la scorurile de potrivire a documentelor sunt cei pentru care  $\text{idf}_t$  are valoare mare.  
 $\Rightarrow$  putem fixa criterii de încetare a acumulării de valori pentru scorurile documentelor, pentru termeni din  $q$  care contribuie cu valori mici.

### e) Reducere bazată pe grupuri

În reducerea bazată pe grupuri (engl. *cluster pruning*) se efectuează un pas de preprocesare care produce grupuri de documente, iar atunci când se primește o cerere de informare, se iau în considerare doar documentele dintr-un număr mic de grupuri, pentru care se calculează scorurile de potrivire cu  $q$ .

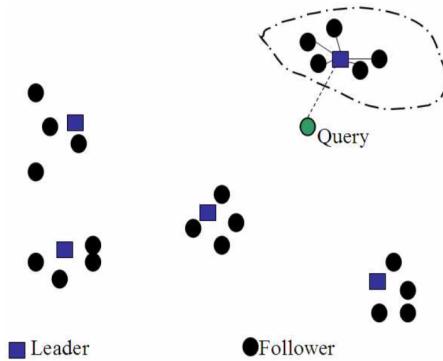
Preprocesarea unei colecții de  $N$  documente cu această metodă se face în 2 pași:

1. Se aleg la întâmplare  $\sqrt{N}$  documente din colecție. Documentele alese se numesc **lideri**.
2. Pentru fiecare document  $d$  din celelalte documente se calculează liderul cel mai apropiat de  $d$ .

Documentele care nu sunt lideri se numesc **următori**. Numărul estimat de următori ai fiecărui lider ales aleator este  $N/\sqrt{N} = \sqrt{N}$ .

Procesarea unei cereri de interogare  $q$  se face astfel:

1. Se detectează documentul lider  $L$  cel mai similar cu  $q$  (se calculează  $\sqrt{N}$  scoruri cosinusoidale)
2. Mulțimea de candidați  $A$  este mulțimea de următori ai lui  $L$ .



## 2 Componentele unui sistem complet de extragere a informațiilor

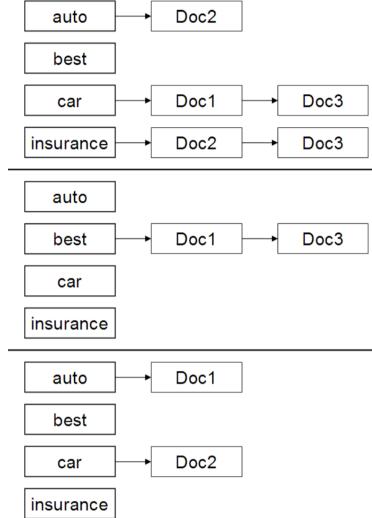
Un sistem rudimentar complet de extragere a informațiilor se poate dezvolta pornind de la modelul de spațiu vectorial pentru documente și de la ideile prezentate de calcul al scorurilor de potrivire cu o cerere de informare.

### 2.1 Indecși organizați pe niveluri

**Idee de bază:** Euristicile prezentate pot produce multimi de candidați  $A$  cu  $|A| < K$ . Indecși organizați pe niveluri rezolvă această problemă, după cum este ilustrat în exemplul de mai jos:

- Pe nivelul 1 reținem postări de documente  $d$  pentru termeni  $t$  cu  $\text{tf}_{t,d} > 20$ ,
- Pe nivelul 2 reținem postări de documente  $d$  pentru termeni  $t$  cu  $\text{tf}_{t,d} > 10$ ,
- Pe nivelul 3 reținem restul de postări de documente.

O astfel de situație este ilustrată în figura de mai jos.



Dacă eșuăm să obținem  $K$  rezultate de pe nivelul 1, continuăm cu procesarea documentelor de pe nivelul 2, și.a.m.d. Pe fiecare nivel, postările sunt ordonate după identificatorii de document.

### 2.2 Proximitatea termenilor din cereri

**Idee de bază:** sunt mai relevante documentele în care termenii din cerere au apariții apropiate. De exemplu, fie  $q$  o cerere de informare conține  $k$  termeni  $t_1, t_2, \dots, t_k$ , și  $w$  lățimea (adică, numărul de termeni) celui mai mic fragment

din  $d$  care conține toți termenii din  $q$ . Cu cât  $\omega$  este mai mic, cu atât este documentul mai relevant pentru cererea  $q$ .

Numeroase motoare de căutare (inclusiv Google) calculează măsuri de proximitate a termenilor din cereri  $q$ , care afectează scorul de potrivire al unui document cu  $q$ .

### 2.3 Construirea funcțiilor de analiză a cererilor și de calcul al scorului

**Stare de fapt:** interfețele majorității sistemelor de extragere a informațiilor au operatori complecși de exprimare a cererilor de informare.

- Acești operatori nu sunt gândiți să fie folosiți de utilizatori.
- De obicei, se presupune că un utilizator comunică o cerere de informare printr-un **text liber** (engl. *free text*), care este o secvență de termeni.
- Textul liber este analizat și transformat în una sau mai multe cereri de informare exprimate cu operatorii complecși al sistemului de extragere a informațiilor.

De exemplu, cererea prin text liber **rising interest rates** a  $K = 10$  potriviri, poate fi procesată astfel:

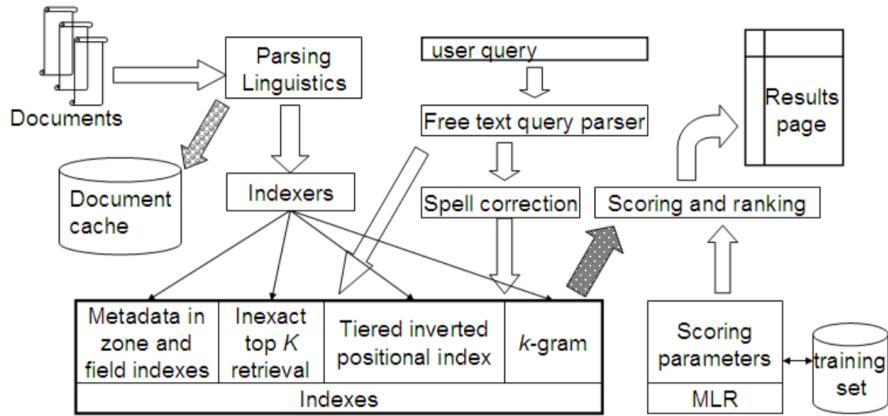
1. Mai întâi se caută documentele în care apare fraza **rising interest rates**, și li se calculează rangul folosind reprezentarea ca vectori.
2. Dacă sunt mai puțin de  $K$  documente care conțin această frază, se caută documentele care conțin frazele **rising interest** și **interest rates**. Din nou, se poate calcula rangul documentelor găsite folosind reprezentarea ca vectori.
3. Dacă nu am găsit încă  $K$  răspunsuri, se continuă cu cererile ce constau din termenii individuali ai cererii inițiale.

Fiecare din cei 3 pași produce o listă de documente cu scoruri aferente de potrivire. Fiecare scor calculat trebuie să combine contribuții ce provin din valoarea scorului în spațiu vectorial, scorul static de calitate, proximitatea în documente a termenilor din cerere, și alti factori potențiali — de exemplu, un document poate să apară în listele unor pași diferenți. Pentru a rezolva această problemă, trebuie definită o funcție de agregare a scorurilor de potrivire produse de pași diferenți:

- Unele sisteme oferă utilizatorului un toolkit cu care pot să-și configureze (1) modul de reducere a cererii din text liber în cereri concrete, și (2) funcția de agregare a scorurilor calculate de pași diferenți. Acest mod de configurare este adecvat pentru colecții de documente care nu se modifică frecvent: o dată configurat, este de așteptat ca sistemul de extragere a informațiilor să funcționeze satisfăcător.

- Pentru colecții de documente care se modifică frecvent, ca de exemplu documentele indexate de motoarele de căutare web, este imposibilă calibrarea manuală a funcțiilor de analiză a cererilor și de calcul al scorului. În acest caz, se folosesc tehnici de machine learning.

## 2.4 Arhitectura unui model complet de căutare



## Bibliografie

1. Capitolul 7: *Computing scores in a complete search system*  
din  
Christopher D. Manning, Prabhakar Raghavan, Hinrich Schütze: *An Introduction to Information Retrieval*. Ediție online (c) 2009 Cambridge UP.  
<http://nlp.stanford.edu/IR-book/pdf/irbookonlinereading.pdf>