

## CAPITOLE SPECIALE DE INFORMATICĂ

# Curs 3: Scoruri de potrivire a documentelor. Calculul greutății termenilor. Modelul de spațiu vectorial

11 octombrie 2018

Modelul boolean de extragere a informațiilor are o limitare importantă: nu poate calcula un scor de importanță care să indice cât de mult răspunde un document unei cereri de informare. Este de dorit ca pentru fiecare pereche  $(d, q)$  în care

- $d$  este un document, și
- $q$  este o cerere de informare

să putem calcula un scor  $scor(d, q) \in \mathbb{R}$  care măsoară gradul de potrivire al lui  $d$  ca răspuns la cererea  $q$ : cu cât documentul  $d$  răspunde mai bine cererii  $q$ , cu atât  $scor(d, q)$  este mai mare.

De ce este util calculul unui scor de potrivire?

- În general, căutarea documentelor care se potrivesc cu o cerere  $q$  produce o listă uriașă de rezultate. Identificarea documentelor care răspund cel mai bine cererii  $q$  este dificilă.
- Scorurile de potrivire permit ordonarea și afișarea rezultatelor în ordine descrescătoare a importanței lor ca răspunsuri.

În acest curs sunt prezentate 3 tehnici folosite în implementarea sistemelor de extragere a informațiilor, care permit calculul unui scor de importanță:

1. Indexarea metadatelor caracteristice unui document. Acest proces produce indecsi parametrici și indecsi de zonă, care permit (1) căutari bazate pe metadate (de ex., titlu sau limbaj folosit); și (2) o metodă simplă de calcul al unui scor de importanță.
2. Calculul unei măsuri de importanță, numită **greutate**, pentru fiecare termen din dicționar.

- Reprezentarea documentelor ca vectori  $M$ -dimensionali de greutăți, unde  $V = \{t_1, \dots, t_M\}$  este dicționarul de termeni al colecției de documente. Fiecare document  $d$  este reprezentat de vectorul  $(g_1, g_2, \dots, g_M)$  unde  $g_i$  este scorul de importanță al lui  $t_i$  în documentul  $d$ . Această reprezentare ușurează calculul scorurilor de potrivire  $scor(d, q)$ .

## 1 Scoruri de potrivire bazate pe metadate

Majoritatea documentelor în format electronic au, pe lângă conținutul text, și **metadate**. Metadatele sunt forme specifice de date despre un document, și pot fi de 2 feluri:

**Câmpuri.** Fiecare câmp are un nume și o mulțime finită de valori posibile.

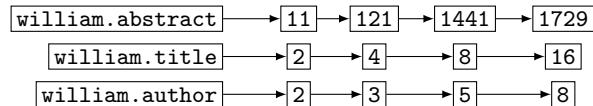
Exemple de nume de câmpuri de metadate de document sunt: data creării documentului; limba în care este scris documentul; etc.

Pentru fiecare câmp se construiește un **index parametric** care permite găsirea rapidă a tuturor documentelor pentru care câmpul respectiv are o anumită valoare.

**Zone.** O zonă are un nume și o valoare care este un text oarecare oricât de mare.

Exemple de nume de zone sunt: titlu; autor; rezumat (sau abstract); etc. Indexarea zonelor se poate face în 2 feluri:

- Se construiește un index inversat separat pentru fiecare zonă, de exemplu: un index pentru abstract-uri, altul pentru titluri, și altul pentru autori. De pildă, liste de indecsi inversați pentru termenul **william** în zonele din metadate ar putea arăta astfel:



- Se extind postările pentru termenii din zone cu specificări ale zonelor în care apar. De exemplu:

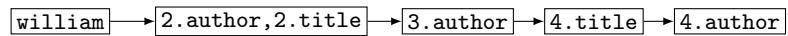


Figura 1 ilustrează interfața unui sistem de extragere a informațiilor cu indecsi parametри și de zone.

### 1.1 Scoruri bazate pe greutatea zonelor

Listele de postări pentru termenii din zonele de metadate permit o tehnică simplă de definire a unui scor de potrivire pentru zone în intervalul  $[0, 1]$ :

- Presupunem că există  $\ell$  zone de metadate distințe  $z_1, \dots, z_\ell$ . Fixăm o pondere de importanță  $g_i \in [0, 1]$  pentru fiecare zonă  $z_i$ . Suma ponderilor zonelor trebuie să fie 1, adică  $\sum_{i=1}^{\ell} g_i = 1$ .

**Bibliographic Search**

| Search category   | Value  |
|---|--|
| <a href="#">Author</a>                                    | Example: Widom, J or Garcia-Molina<br><input type="text"/>   |
| <a href="#">Title</a>                                     | Also a part of the title possible<br><input type="text"/>  |
| <a href="#">Date of publication</a>                       | Example: 1997 or <1997 or >1997 limits the search to the documents appeared in, before and after 1997 respectively<br><input type="text"/> |
| <a href="#">Language</a>                                  | Language the document was written in<br>English <input type="button" value="▼"/>   |
| <a href="#">Project</a>                                   | ANY <input type="button" value="▼"/>   |
| <a href="#">Type</a>                                      | ANY <input type="button" value="▼"/>   |
| <a href="#">Subject group</a>                             | ANY <input type="button" value="▼"/>   |
| <a href="#">Sorted by</a>                                 | Date of publication <input type="button" value="▼"/>   |
| <input type="button" value="Start bibliographic search"/> |  |
| Find document via ID <input type="text"/>                 |  |

Figure 1: Căutare bazată pe indecsi parametrici și indecsi de zonă.

2. Fie  $s_i \in \{0, 1\}$  valoarea booleană de potrivire a zonei  $z_i$  a documentului  $d$  cu cererea  $q$ . Scorul (sau gradul) de potrivire dintre zonele lui  $d$  și  $q$  este dat de formula

$$scor(d, q) := \sum_{i=1}^{\ell} g_i \cdot s_i$$

De exemplu, dacă avem o colecție  $D$  de documente cu trei zone:  $z_1$  pentru **author**,  $z_2$  pentru **title** și  $z_3$  pentru **body**, și fixăm ponderile  $g_1 = 0.2$ ,  $g_2 = 0.3$  și  $g_3 = 0.5$  atunci:

- Scorul de potrivire  $scor(d, \text{shakespeare})$  pentru un document în care termenul **shakespeare** apare în zonele pentru **title** și pentru **body** este

$$0.2 \cdot s_1 + 0.3 \cdot s_2 + 0.5 \cdot s_3 = 0.2 \cdot 0 + 0.3 \cdot 1 + 0.5 \cdot 1 = 0.8$$

Această tehnică de calcul al unui scor de potrivire pentru zone se numește **extragere booleană gradată a informațiilor**.

Listele de indecsi inversați pentru zone de metadate pot fi folosite pentru calculul direct al scorurilor de potrivire. De exemplu, dacă  $q_1$  și  $q_2$  sunt două cereri, iar

- $p_1$  este lista de postări care indică documentele și zonele ce se potrivesc cu cererea  $q_1$  (metoda 2 de la pagina 2)
- $p_2$  este lista de postări care indică documentele și zonele ce se potrivesc cu cererea  $q_2$  (metoda 2 de la pagina 2)

atunci putem calcula direct un tablou  $scoruri[ ]$  care, pentru fiecare document  $d$  al căruia identificator  $docId$  apare în  $p_1$  și  $p_2$ , reține în  $scoruri[docId]$  scorul

de potrivire dintre  $d$  și cererea  $q_1 \text{ AND } q_2$ . Algoritmul de calcul este ilustrat în figura 2 în care presupunem că

1. colecția are  $N$  documente, și fiecare document are  $\ell$  zone cu ponderile  $g[1], \dots, g[\ell]$
2. metoda auxiliară  $\text{SCORDEZONĂ}(p_1, p_2)$  este apelată cu două postări  $p_1, p_2$  pentru același document  $d$ ;  $p_1$  indică zonele lui  $d$  care se potrivesc cu cererea  $q_1$ ;  $p_2$  indică zonele lui  $d$  care se potrivesc cu cererea  $q_2$ ; metoda întoarce valoarea  $\sum_{i=1}^{\ell} g_i \cdot s_{1,i} \cdot s_{2,i}$  unde  $s_{j,i}$  este 1 dacă zona  $i$  a lui  $d$  se potrivește cu cererea  $q_j$ , și 0 în caz contrar.

```

SCORURIDEZONĂ( $q_1, q_2$ )
1 float  $scoruri[N] = [0]$ 
2 constant  $g[\ell]$ 
3  $p_1 := postings(q_1)$ 
4  $p_2 := postings(q_2)$ 
5 //  $scoruri$  este un tablou cu o intrare inițializată pentru fiecare docID
6 //  $p_1, p_2$  sunt inițializați să se refere la începuturile listelor de postări
7 // Presupunem că  $g[ ]$  este inițializat cu ponderile corespunzătoare de zonă
8 while  $p_1 \neq Nil$  and  $p_2 \neq Nil$  do
9     if  $=\text{docID}(p_1) = \text{docID}(p_2)$  then
10         $scoruri(\text{docID}(p_1)) := \text{SCORDEZONĂ}(p_1, p_2, g)$ 
11         $p_1 := \text{next}(p_1)$ 
12         $p_2 := \text{next}(p_2)$ 
13    else if  $\text{docID}(p_1) < \text{docID}(p_2)$  then
14         $p_1 := \text{next}(p_1)$ 
15    else  $p_2 := \text{next}(p_2)$ 
16 return  $scoruri$ 

```

Figure 2: Calculul scorurilor ponderate de zonă pentru cererea  $q_1 \text{ AND } q_2$ .

## 1.2 Învățarea scorurilor

Determinarea ponderilor de zonă  $g_1, \dots, g_\ell$  se face adesea cu un algoritm de învățare care funcționează astfel:

- Primește un set  $T$  de documente clasificate deja de către niște operatori umani. Setul  $T$  se numește set de **exemple de antrenare**.
- Determină ponderile  $g_1, \dots, g_\ell$  care minimizează o măsură de eroare a scorurilor calculate pentru exemplele de antrenare.
- Ponderile calculate se vor folosi pentru calculul scorurilor unor documente noi.

### Exemplu ilustrat

Fie  $D$  o colecție de documente cu două zone:  $T$  pentru titlu și  $C$  pentru conținutul propriu-zis (corpul) documentului. În acest caz, avem de fixat o valoare  $g \in [0, 1]$  cu care calculăm scorul de potrivire

$$scor(d, q) = g \cdot s_T(d, q) + (1 - g) \cdot s_C(d, q).$$

În formula de mai sus avem  $s_T(d, q) = 1$  dacă zona  $T$  a lui  $d$  se potrivește cu cererea  $q$ , și  $s_T(d, q) = 0$  în caz contrar. Deasemenea,  $s_C(d, q) = 1$  dacă zona  $C$  a lui  $d$  se potrivește cu cererea  $q$ , și  $s_C(d, q) = 0$  în caz contrar.

În acest caz, o mulțime de exemple de antrenare este o mulțime finită de tripleți  $\Phi_j = (d_j, q_j, r(d_j, q_j))$  formate din un document  $d_j$ , o cerere  $q_j$  și o decizie de relevanță:  $r(d_j, q_j) = 1$  dacă documentul  $d_j$  se consideră relevant pentru cererea  $q$ , și  $r(d_j, q_j) = 0$  în caz contrar.

Se consideră că **eroarea de calcul al scorului** pentru exemplul de antrenare  $\Phi_j$  este  $\epsilon(g, \Phi_j) := (r(d_j, q_j) - scor(d_j, q_j))^2$ , iar **eroarea totală de calcul al scorurilor de antrenare** este

$$\sum_{j=1}^M \epsilon(g, \Phi_j) \quad \text{dacă } T = \{\Phi_1, \Phi_2, \dots, \Phi_M\}.$$

Se dorește să se aleagă valoarea  $g \in [0, 1]$  pentru care eroarea totală ia valoarea minimă. Fie

- $n_{ijr}$  numărul exemplelor de antrenare  $(d_j, q_j, r(d_j, q_j))$  cu  $s_T(d_j, q_j) = i$  și  $s_C(d_j, q_j) = j$  și  $r(d_j, q_j) = 1$ .
- $n_{ijn}$  numărul exemplelor de antrenare  $(d_j, q_j, r(d_j, q_j))$  cu  $s_T(d_j, q_j) = i$  și  $s_C(d_j, q_j) = j$  și  $r(d_j, q_j) = 0$ .

Rezultă că eroarea totală este o expresie polinomială de gradul 2 în  $g$ :

$$\sum_{j=1}^M \epsilon(g, \Phi_j) = (n_{01r} + n_{10n}) \cdot g^2 + (n_{10r} + n_{01n})(1 - g)^2 + n_{00r} + n_{11n}$$

Deci, avem de rezolvat o problemă de minimizare pătratică. Valoarea lui  $g$  pentru care derivata acestei expresii este 0 este

$$\frac{n_{00r} + n_{11n}}{n_{00r} + n_{00n} + n_{11r} + n_{11n}} \in [0, 1]$$

Rezultă că aceasta este valoarea optimă a lui  $g$  pe care o calculează algoritmul de învățare.

## 2 Calculul frecvenței și al greutății termenilor

În modelul boolean de extragere a informațiilor contează doar dacă un termen apare sau nu într-un document. O variantă îmbunătățită ar fi să considerăm că

un document sau zonă este cu atât mai relevant(ă) pentru o cerere exprimată de un termen  $t$  cu cât  $t$  apare mai frecvent în documentul sau zona respectivă.

Prin urmare, vom calcula **frecvența unui termen în un document**:

$\text{tf}_{t,d}$  := numărul de apariții al termenului  $t$  în documentul  $d$ .

**Observație:** calculul frecvențelor termenilor ne permite să interpretăm fiecare document ca pe o colecție de termeni, în care numărul de apariții al termenilor contează, dar ordinea în care apar termenii nu contează. De exemplu, documentul cu conținutul “Ion este mai înalt decât George” este identic în această interpretare cu documentul cu conținutul “George este mai înalt decât Ion.”

Nu toți termenii sunt la fel de importanți pentru a descrie conținutul specific unui document. De exemplu, este firesc ca termenul **auto** să apară foarte des într-o colecție de documente referitoare la industria auto, și să nu transmită informații noi. Pentru a reduce relevanța termenilor care apar frecvent în colecții și care nu au relevanță mare pentru descrierea conținutului documentului, se calculează valorile următoare:

- **frecvența de colecție a unui termen:**

$\text{cf}_t$  := numărul total de apariții al lui  $t$  în colecția de documente  $D$ .

- **frecvența de document a unui termen:**

$\text{df}_t$  := numărul total de documente din colecția  $D$  în care apare  $t$ .

În general, aceste valori se comportă diferit. De exemplu, frecvențele de colecție și de document ale termenilor **try** și **insurance** au valorile ilustrate mai jos în o colecție de documente de știri ale agenției Reuters:

| termen    | cf    | df   |
|-----------|-------|------|
| try       | 10422 | 8760 |
| insurance | 10440 | 3997 |

Pentru a reduce relevanța termenilor care apar în multe documente, se definește factorul numit **frecvență inversă de document** a lui  $t$ :

$$\text{idf}_t := \log \frac{N}{\text{df}_t}$$

unde  $N$  este numărul total de documente în colecția  $D$ . Se observă că  $\text{idf}_t$  ia valori mari pentru termeni ce apar în puține documente din colecție, și ia valori mici pentru termeni ce apar în multe documente din colecție. Figura 3 ilustrează acest comportament.

În final, se definește **greutatea tf-idf** a unui termen  $t$  în un document  $d$ :

$$\text{tf-idf}_{t,d} := \text{tf}_{t,d} \cdot \text{idf}_t.$$

Se observă că  $\text{tf-idf}_{t,d}$  atribuie o greutate

- mare când  $t$  apare des în un număr mic de documente,
- redusă când  $t$  apare rar în  $d$ , sau apare în multe documente,
- foarte redusă când  $t$  apare în majoritatea documentelor din colecție.

| termen    | $\text{df}_t$ | $\text{idf}_t$ |
|-----------|---------------|----------------|
| car       | 18165         | 1.65           |
| auto      | 6723          | 2.08           |
| insurance | 19241         | 1.62           |
| best      | 25235         | 1.5            |

Figure 3: Exemplu de valori df și idf pentru termeni din o colecție de 806791 documente a agenției de știri Reuters.

### 3 Modelul de spațiu vectorial pentru documente

Calculul scorurilor de potrivire al documentelor cu cereri de informare se simplifică dacă reprezentăm fiecare document  $d$  ca pe un vector de dimensiune  $M$

$$\vec{V}(d) := (\text{tf-idf}_{t_1,d}, \dots, \text{tf-idf}_{t_M,d})$$

unde  $V = \{t_1, \dots, t_M\}$  este dicționarul de termeni ai colecției de documente. Documentele devin vectori într-un spațiu  $M$ -dimensional care are o axă pentru fiecare termen, iar componenta  $i$  a vectorului are valoarea  $\text{tf-idf}_{t_i,d}$  care indică relevanța lui  $t_i$  în conținutul lui  $d$ .

#### 3.1 Similaritatea documentelor

Modul standard de calcul al similarității a 2 documente este dat de **similaritatea cosinusoidală** a vectorilor  $\vec{V}(d_1)$  și  $\vec{V}(d_2)$ :

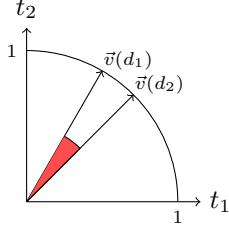
$$\text{sim}(d_1, d_2) := \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)| \cdot |\vec{V}(d_2)|}.$$

Reamintim că, dacă  $\vec{x} = (x_1, \dots, x_M)$  și  $\vec{y} = (y_1, \dots, y_M)$  sunt doi vectori de aceeași dimensiune  $M$ , atunci

- **produsul lor scalar** este  $\vec{x} \cdot \vec{y} := \sum_{i=1}^M x_i y_i$
- **lungimea** vectorului  $\vec{x}$  este  $|\vec{x}| := \sqrt{\sum_{i=1}^M x_i^2}$
- $\frac{\vec{x} \cdot \vec{y}}{|\vec{x}| \cdot |\vec{y}|}$  este valoarea cosinusului unghiului  $\theta$  dintre vectorii  $\vec{x}$  și  $\vec{y}$ .

De exemplu, dacă  $\vec{v}(d_1) = \vec{V}(d_1)/|\vec{V}(d_1)|$  și  $\vec{v}(d_2) = \vec{V}(d_2)/|\vec{V}(d_2)|$  atunci

- $\vec{v}(d_1)$  și  $\vec{v}(d_2)$  au lungimea 1, iar cosinusul unghiului dintre ei este produsul scalar  $\vec{v}(d_1) \cdot \vec{v}(d_2)$ .



- Cosinusul unghiului roșu dintre  $\vec{v}(d_1)$  și  $\vec{v}(d_2)$  reprezintă similaritatea dintre  $d_1$  și  $d_2$ . Se observă că

$$-1 \leq \text{sim}(d_1, d_2) \leq 1.$$

Dacă  $D = \{d_i \mid 1 \leq i \leq N\}$  este o colecție de documente și vrem să găsim documentele din  $D$  cele mai similare cu un document dat  $d$ , putem proceda astfel: calculăm toate produsele scalare  $\vec{v}(d_i) \cdot \vec{v}(d)$ , unde

$$\vec{v}(d_i) = \frac{\vec{V}(d_i)}{|\vec{V}(d_i)|} \quad \text{și} \quad \vec{v}(d) = \frac{\vec{V}(d)}{|\vec{V}(d)|}$$

și se returnează documentele  $d_i$  pentru care produsele scalare iau valorile cele mai mari.

O colecție cu  $N$  documente și  $M$  termeni poate fi reprezentată ca o matrice  $M \times N$  în care rândurile reprezintă  $M$  termeni iar coloanele reprezintă documentele din colecție.

### 3.2 Reprezentarea vectorială a cererilor

Considerăm expresii booleene de forma

$$q = t_{i_1} \ t_{i_2} \ \dots \ t_{i_n}$$

ca cereri de informare pentru documente care conțin termenii  $t_{i_1}, \dots, t_{i_n}$ . Pre-supunem că extragerea informațiilor se face din o colecție de  $N$  documente cu un vocabular de  $M$  termeni. În acest caz, cererea  $q$  se reprezintă cu vectorul

$$\vec{v}(q) := (v_1, \dots, v_M) \quad \text{unde } v_i = \begin{cases} \frac{1}{\sqrt{n}} & \text{dacă } i \in \{i_1, \dots, i_n\}, \\ 0 & \text{în caz contrar.} \end{cases}$$

Scorul de potrivire al cererii  $q$  cu un document  $d$  este

$$\vec{v}(q) \cdot \vec{v}(d) \quad \text{unde } \vec{v}(d) = \frac{\vec{V}(d)}{|\vec{V}(d)|}.$$

### Exemplu ilustrat

Presupunem că  $D$  este o colecție fictivă de  $N = 10^6$  documente în care apar doar termenii `auto`, `best`, `car` și `insurance`. Deasemenea, considerăm că cererea  $q$  este

`best car insurance`

și că avem valorile statistice indicate în tabelul de mai jos:

| termen    | cerere |       |     |           | document |    |           | produs |
|-----------|--------|-------|-----|-----------|----------|----|-----------|--------|
|           | tf     | df    | idf | $w_{t,q}$ | tf       | wf | $w_{t,d}$ |        |
| auto      | 0      | 5000  | 2.3 | 0         | 1        | 1  | 0.41      | 0      |
| best      | 1      | 50000 | 1.3 | 1.3       | 0        | 0  | 0         | 0      |
| car       | 1      | 10000 | 2.0 | 2.0       | 1        | 1  | 0.41      | 0.82   |
| insurance | 1      | 1000  | 3.0 | 3.0       | 2        | 2  | 0.82      | 2.46   |

În acest exemplu greutatea  $w_{t,q}$  a unui termen  $t$  în cererea  $q$  este valoarea idf (care este 0 pentru termenii care nu apar în cerere, ca de exemplu `auto`). În documente, se consideră că greutatea  $w_{t,d}$  a unui termen  $t$  în un document  $d$  este valoarea normalizată a lui  $tf_{t,d}$ . Deci

$$\begin{aligned}\vec{v}(q) &= (0, 1.3, 2.0, 3.0) \\ \vec{v}(d) &= (0.41, 0.041, 0.82)\end{aligned}$$

deci scorul de potrivire al lui  $q$  cu  $d$  este  $\vec{v}(q) \cdot \vec{v}(d) = 0 + 0 + 0.82 + 2.46 = 3.28$ .

### 3.3 Găsirea primelor $K$ documente care se potrivesc cel mai bine cu o cerere

De obicei, un sistem de extragere a informațiilor are o colecție de documente reprezentate ca vectori  $\{\vec{v}(d_i) \mid 1 \leq i \leq N\}$ , o cerere  $q$  reprezentată ca un vector  $\vec{v}(q)$ , și își cere să găsească primele  $K$  documente din colecție care se potrivesc cel mai bine cu cererea  $q$ . Pseudocodul din figura 4 este pentru algoritmul de bază care calculează scorurile de potrivire a documentelor din  $D$  cu cererea  $q$ . În acest pseudocod,  $lungime[d] = \sqrt{\sum_{i=1}^M wf_{t_i,d}}$  este lungimea euclidiană a vectorului  $\vec{v}(d) = (wf_{t_1,d}, \dots, wf_{t_M,d})$  iar în pasul 7,  $wf_{t,d}$  este de obicei fie  $tf_{t,d}$  sau  $tf\text{-}idf_{t,d}$ .

Sunt două posibilități de memorare a listei de postări pentru un termen  $t$ :

$$1. \quad t \rightarrow [d_1 \mid tf_{t,d_1} \mid wf_{t,d_1}] \rightarrow \dots \rightarrow [d_n \mid tf_{t,d_n} \mid wf_{t,d_n}] \rightarrow \dots$$

adică să se rețină valoarea lui  $wf_{t,d_i}$  (un număr `float` sau `double`) împreună cu postarea lui  $d_i$  pentru termenul  $t$ .

$$2. \quad t \rightarrow [N/df_t] \rightarrow [d_1 \mid tf_{t,d_1}] \rightarrow \dots \rightarrow [d_n \mid tf_{t,d_n}] \rightarrow \dots$$

```

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

```

Figure 4: Algoritmul de bază pentru calculul scorurilor de potrivire în spațiul vectorial de documente.

Prima intrare din lista de postări este specială: reține valoarea lui  $N/\text{df}_t$  (un **float** sau **double**). Această codificare a listei de postări ocupă mai puțin spațiu și permite (1) calcului lui  $\text{idf}_t := \log \frac{N}{\text{df}_t}$ , și (2) calculul valorilor  $\text{tf-idf}_{t,d_i} := \text{idf}_t \cdot \text{tf}_{t,d_i}$ .

Presupunând că  $q = t_{i_1} t_{i_2} \dots t_{i_n}$ , se observă că, la pasul 9, algoritmul memorizează în  $scoruri[d]$  valoarea

$$\frac{\sum_{k=1}^n \text{wf}_{t_{i_k},d} \cdot w_{t_{i_k},q}}{lungime[d]}$$

Valoarea lui  $scoruri[d]$  se calculează incremental: după ce termenul  $t_{i_j}$  este ales din cererea  $q$  în bucla **for** (pașii 4–7), fiecare element de tablou  $scoruri[d]$  va reține valoarea  $\sum_{k=1}^j \text{wf}_{t_{i_k},d} \cdot w_{t_{i_k},q}$ . Acest calcul incremental al scorurilor de potrivire se numește *gradare după termeni* (engl. *term-at-a-time scoring*) sau **acumulare**, iar elementele tabloului  $scoruri[ ]$  se numesc **acumulatori**.

## Bibliografie

1. Capitolul 6 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>