

## CAPITOLE SPECIALE DE INFORMATICĂ

# Curs 2: Determinarea vocabularului de termeni și a listelor de postări

4 octombrie 2018

Reamintim că liste de indecsi inversați se construiesc în 4 pași:

1. Se colectează o mulțime  $D$  de documente care se vor indexa.
2. Se tokenizează textul din documentele colecției.
3. Tokenurile se preprocesează lingvistic.
4. Se indexează documentele în care apare fiecare termen.

Scopul acestui curs este să indice

- cum se determină unitățile de măsură a documentelor, și secvențele de caractere pentru conținutul lor.
- aspectele importante ale tokenizării și preprocesării lingvistice.
- implementarea listelor de postări și a listelor extinse de postări.
- structuri de date utile pentru interogări ce includ fraze și operatori de proximitate.

## 1 Determinarea documentelor și decodificarea secvenței de caractere

De obicei, documentele în format electronic sunt secvențe de octeți stocate în un fișier sau pe un browser web. Se folosesc mai multe metode de codificare/decodificare a secvențelor de caractere în secvențe de caractere: ASCII, Unicode UTF-8, UTF-16, etc. Preprocesarea acestor documente presupune:

1. Determinarea codificării corecte, urmată de decodificare. Determinarea codificării corecte este o problemă de clasificare în machine learning, care poate fi rezolvată în mai multe feluri: cu metode euristice, pe baza unei selecții a utilizatorului, sau pe baza metadatelor asociate cu documentul respectiv.

- Uneori, secvența de caractere produsă prin decodificare mai trebuie preprocesată. De exemplu, documentele XML conțin entități care se înlocuiesc cu caractere speciale. De exemplu, & se înlocuiește cu &.

In general, produsele software comerciale trebuie să știe să recunoască și să decodifice un număr mare de tipuri de documente și codificări. De obicei, decodificarea în un sir de caractere se face cu biblioteci software achiziționate cu licență în acest scop. Aceste biblioteci știu și cum să decodifice texte scrise în arabă sau japoneză, care nu sunt simple siruri de caractere scrise de la stânga la dreapta.

- Să se decidă ce unitate de măsură a documentelor (engl. *document unit*) să folosească. Adeseori, un document nu corespunde unui fișier: un fișier zip poate arhiva mai multe documente; mesajele din mail-boxul unui utilizator sunt toate stocate în un singur fișier; unele programe (precum `latex2html`) descompun un document în mai multe fișiere separate.

Prin urmare, se pune problema determinării *granularității*: cât de mare să fie un document; cum definim și cum determinăm o unitate de măsură pentru documente.

## 2 Determinarea vocabularului de termeni

Secvența de caractere produsă prin decodificare se desparte în tokenuri. De obicei, tokenurile se sunt separate prin semne de punctuație și secvențe de caractere speciale (eng. *whitespace*).

De exemplu, tokenizarea textului “Friends, Romans, Countrymen, lend me your ears.” produce secvența de tokenuri

Friends   Romans   Countrymen   lend   me   your   ears

Apoi, se rulează un proces de normalizare care extrage o secvență de **termeni** din lista de tokenuri. Termenii vor fi elementele incluse în dicționarul sistemului de IR.

Tokenizarea este o problemă dificilă: cum definim și cum determinăm tokenurile? De exemplu, cum tokenizarea frazei

“Mr. O’Neill thinks that the boys’ stories about Chile’s capital aren’t amusing.”

se poate face corect doar după ce se detectează limbajul în care este scris textul. O tehnică foarte eficientă de **detectie a limbajului** este cea bazată pe clasificatori ce utilizează subsecvențe scurte de caractere.

Alte probleme complexe rezolvate de algoritmii avansați de tokenizare sunt:

- interpretarea corectă a liniștei de subliniere. De exemplu *co-education* și *Hewlett-Packard* sunt considerate un singur token.
- recunoașterea grupurilor de cuvinte care constituie un singur token. De exemplu *Los Angeles* sau *New York University*.

- Segmentarea unor cuvinte compuse lungi. De exemplu, cuvântul german “Lebensversicherungsgesellschaftsangestellter” înseamnă “angajat al unei companii de asigurare pe viață” și poate fi descompus în mai multe tokenuri.
- Eliminarea unor cuvinte de legătură folosite frecvent. Cuvinte precum “și” sau articole precum “al” se omit din lista de tokenuri.
- ...

### 3 Normalizarea tokenurilor

*Normalizarea* este procesul de împărțire a tokenurilor în clase de echivalență, adică mulțimi de tokenuri care au același înțeles sau înțelesuri foarte asemănătoare. De exemplu, tokenurile **mașină** și **automobil** sunt echivalente. Prin urmare, când căutam documente despre **mașină**, vrem să găsim și documente despre **automobil**.

O soluție la această problemă este să reprezentăm toți termenii din o clasă de echivalență cu un singur termen. Înlocuirea unui termen cu reprezentantul clasei sale de echivalență se numește **normalizare**.

Altă soluție este să reținem care termeni sunt sinonimi. De obicei, acest lucru se realizează în 2 feluri:

1. Se indexează toate tokenurile nenormalizate și se asociază o *listă de expansiune* pentru toți termenii admisi în o cerere de informare. De exemplu, **mașină** se va putea înlocui cu **mașină OR automobil**.
2. Expansiunea tokenurilor se face pe durata indexării. De exemplu, un document care conține **mașină** se indexează și pentru tokenul **automobil**.

Unele probleme de normalizare trebuie să țină cont de limbajul din care provin tokenurile:

- În engleză, accentele și diacriticile au de obicei un rol marginal ⇒ normalizarea se poate face eliminând diacriticile: “cliche” în loc de “cliché”, “naïve” în loc de “naïve.”

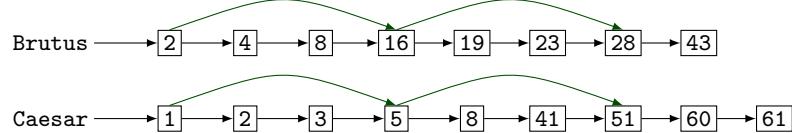
În alte limbi, precum spaniola, un accent poate schimba complet sensul unui cuvânt: “peña” înseamnă “faleză”, iar “pena” înseamnă “durere.”

...

### 4 Liste extinse de postări

În cursul 1 am vazut că intersecția a 2 liste de postări de lungimi  $m$  și  $n$  se poate face în timp  $O(m + n)$ . Dacă dorim timpi și mai scurți (sub-liniari) de calcul al intersecției, putem folosi liste cu pointeri de ocolire (engl. **skip lists**):

- ▶ pointerii de ocolire sunt “scurtături” efective care ne permit să evităm să procesăm porțiuni mari din liste de postări care nu apar în rezultatul intersecției. Figura de mai jos ilustrează 2 liste de postări cu pointeri de ocolire:



Apar 2 probleme noi: (1) unde se plasează pointerii de ocolire, și (2) cum se poate face intersecție rapidă cu pointeri de ocolire.

Algoritmul de intersecție cu pointeri de ocolire se poate defini în mai multe feluri. O posibilitate este ilustrată mai jos:

```

INTERSECTCUOCOLIRI( $p_1, p_2$ )
1  $raspuns := \langle \rangle$ 
2 while  $p_1 \neq Nil$  and  $p_2 \neq Nil$  do
3     if  $\text{docID}(p_1) = \text{docId}(p_2)$ 
4         adaugă  $\text{docId}(p_1)$  la  $raspuns$ 
5          $p_1 := \text{next}(p_1)$ 
6          $p_2 := \text{next}(p_2)$ 
7     else if  $\text{docID}(p_1) < \text{docId}(p_2)$ 
8         if  $\text{hasSkip}(p_1)$  and  $(\text{docID}(\text{skip}(p_1)) \leq \text{docId}(p_2))$ 
9             while  $\text{hasSkip}(p_1)$  and  $(\text{docID}(\text{skip}(p_1)) \leq \text{docId}(p_2))$  do
10             $p_1 := \text{skip}(p_1)$ 
11        else  $p_1 := \text{next}(p_1)$ 
12    else if  $\text{hasSkip}(p_2)$  and  $(\text{docID}(\text{skip}(p_2)) \leq \text{docId}(p_1))$ 
13        while  $\text{hasSkip}(p_2)$  and  $(\text{docID}(\text{skip}(p_2)) \leq \text{docId}(p_1))$  do
14             $p_2 := \text{skip}(p_2)$ 
15        else  $p_2 := \text{next}(p_2)$ 
16 return  $raspuns$ 
  
```

Observații:

- Pointerii de ocolire sunt folositorii pentru a face mai rapidă operația de intersecție; nu pot fi folosiți pentru operația de reunire.
- Mai mulți pointeri de ocolire implică
  - rute de ocolire mai scurte
  - mai multe teste de ocolire și mai mult spațiu ocupat de pointerii de ocolire.
- Ouristică simplă care funcționează bine în practică este ca pentru liste de postări de lungime  $P$  să alocăm  $\sqrt{P}$  pointeri de ocolire care să peste căte  $\sqrt{P}$  postări.

O altă extensie utilă a listelor simple de postări sunt **listele de indecsi poziționali**. O listă de indecsi poziționali pentru un termen  $t$  reține mai multe informații:

- frecvența lui  $t$  în colecția de documente, adică în câte documente apare  $t$ .
- pentru fiecare document cu identificatorul  $docID$ , o postare de forma

$$docId, n : \langle pos_1, pos_2, \dots \rangle$$

unde

- $pos_1, pos_2, \dots$  sunt pozițiile (sortate crescător) unde apare termenul  $t$  în documentul cu identificatorul  $docId$ .
- $n$  este numărul de aparții al lui  $t$  în documentul cu identificatorul  $docId$ .

De exemplu, listele de indecsi poziționali ai termenilor **to** și **be** în o colecție de documente ar putea arăta astfel:

**to,993427:**

$\langle 1,6 : \langle 7, 18, 33, 72, 86, 231 \rangle;$   
 $2,5 : \langle 1, 17, 74, 222, 255 \rangle;$   
 $4,5 : \langle 8, 16, 190, 429, 433 \rangle;$   
 $\dots \rangle$

**be,178239:**

$\langle 1,2 : \langle 17, 25 \rangle;$   
 $4,5 : \langle 17, 191, 291, 430, 434 \rangle;$   
 $5,3 : \langle 14, 19, 101 \rangle;$   
 $\dots \rangle$

Listele de indecsi poziționali permit găsirea documentelor care conțin fraze, adică secvențe de termeni.

De exemplu, găsirea documentelor care conțin “**to be**” poate decurge astfel: se caută pozițiile în lista de indecsi poziționali ai lui **be** care sunt imediat după o poziție în lista de indecsi poziționali ai lui **to** în același document. De pildă, pentru listele ilustrate mai sus, rezultatul căutării va fi

**to:**  $\langle \dots, 4 : \langle \dots, 429, 433 \rangle, \dots \rangle$   
**be:**  $\langle \dots, 4 : \langle \dots, 430, 434 \rangle, \dots \rangle$

Un algoritm mai general, care găsește toate perechile de poziții din 2 liste de indecsi poziționali care sunt la cel mult  $k$  cuvinte distanță una de cealaltă, este ilustrat în figura 1.

```

POSITIONALINTERSECT( $p_1, p_2, k$ )
1  $raspuns := \langle \rangle$ 
2 while  $p_1 \neq Nil$  and  $p_1 \neq Nil$  do
3   if docID( $p_1$ ) = docId( $p_2$ )
4      $l := \langle \rangle$ 
5      $pp_1 := pozitii(p_1)$ 
6      $pp_2 := pozitii(p_2)$ 
7     while  $pp_1 \neq Nil$  do
8       while  $pp_2 \neq Nil$  do
9         if  $|pos(pp_1) - pos(pp_2)| \leq k$ 
10          adaugă  $pos(pp_2)$  la  $l$ 
11        else if  $pos(pp_2) > pos(pp_1)$ 
12          break
13         $pp_2 := next(pp_2)$ 
14      while  $l \neq \langle \rangle$  and  $|l[0] - pos(pp_1)| > k$  do
15        şterge  $l[0]$  din  $l$ 
16      for each  $ps \in l$  do
17        adaugă  $\langle docId(p_1), pos(pp_1), ps \rangle$  la  $raspuns$ 
18         $pp_1 := next(pp_1)$ 
19         $p_1 := next(p_1)$ 
20         $p_2 := next(p_2)$ 
21      else if docID( $p_1$ ) < docID( $p_2$ )
22         $p_1 := next(p_1)$ 
23      else  $p_2 := next(p_2)$ 
24 return  $raspuns$ 

```

Figure 1: Algoritm de intersecție de proximitate dintre 2 liste extinse de postări.

## Bibliografie

1. Capitolul 2 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>