

## CAPITOLE SPECIALE DE INFORMATICĂ

# Curs 1: Extragerea informațiilor. Modelul boolean și modelul boolean extins

27 septembrie 2018

**Extragerea informațiilor** (engl. *Information Retrieval, IR*) constă în găsirea și extragerea de material (de obicei, documente) nestructurat sau semi-structurat (de obicei, care conține text) care satisfacă o **necesitate de informare** din colecții mari (care, de obicei, sunt stocate pe calculatoare).

- **Datele nestructurate** sunt date fără o structură precisă care să poată fi procesată ușor cu un program pe calculator.
  - Exemple de date structurate: o bază de date relațională cu informații despre angajații unei companii.
  - Majoritatea datelor textuale au o structură lingvistică latentă: titlu, secțiuni, paragrafe, etc. Astfel de surse de informații se numesc *documente semistructurate*.
- Necesațile de informare sunt comunicate unui sistem de extragere a informațiilor prin intermediul unei **cereri** (engl. **query**) formulate de către utilizator.
- Un document este **relevant** pentru necesitatea de informare a unui utilizator dacă utilizatorul respectiv consideră că documentul respectiv conține informații valoroase pentru nevoia lui de informare.

Exemple tipice de probleme rezolvate de către sistemele de extragere a informațiilor:

- **Căutare semistrustructată**, de exemplu, găsirea documentelor al căror titlu conține “Java” și al căror conținut (corp) conține “threading.”
- **Clustering**, adică împărțirea unei colecții mari de documente în grupuri (clustere) care conțin documente cu conținut asemănător.
- **Clasificare**: dacă se dă (1) o mulțime de tematici ce corespund unor necesități de informare, sau (2) alte categorii de text (de exemplu, pentru cititori de anumite vârste), să se decidă la ce tematică sau categorie aparține fiecare document.

- Adesea se bazează pe algoritmi de învățare ce pornesc de la o mulțime de documente clasificate manual, și se speră că algoritmul este suficient de bun încât să învețe să clasifice corect documentele noi.
- **Căutare al-hoc** a documentelor din o colecție care sunt relevante pentru orice necesitate de informare comunicată sistemului printr-o cerere inițiată de către utilizator.

Documentele de extragere a informațiilor pot fi clasificate după mărimea colecției de documente pe care o analizează:

- Sistemele de **căutare web** analizează miliarde de documente stocate pe milioane de calculatoare. Un astfel de sistem este optimizat
  - să analizeze un volum cât mai mare de documente existente pe web (*crawling*)
  - să proceseze caracteristici specifice documentelor web: hiperlink-uri, tentative de manipulare a scorului de importanță, etc.
- Sistemele de **căutare pentru un mediu de afaceri, instituție sau un anumit domeniu**.
  - Colecțiile de documente analizate pot fi: documente interne ale unei corporații; o bază de date de patente; articole de cercetare despre biochimie; etc.
  - De obicei, aceste colecții sunt stocate în sisteme centralizate de fișiere, iar căutarea de informații în ele de face cu ajutorul câtorva calculatoare dedicate.
- Sisteme de **extragere a informațiilor personale** precum Spotlight (Mac OS X). Instant Search (Windows Vista), sau programe de email cu capacitate de clasificare automată a mesajelor, filtrare spam:
  - analizează o gamă largă de tipuri de documente stocate pe un calculator personal
  - sunt rapide, ocupă puțină memorie, și nu afectează semnificativ experiența de lucru a utilizatorului cu calculatorul

## 1 Modelul boolean de extragere a informațiilor

Se presupune dată o colecție de documente  $D$ . Are următoarele caracteristici:

1. Cererea de informare este o expresie booleană  $B$  care descrie ce termeni dorim să apară sau nu în documentele pe care le căutăm în colecția  $D$ . Se pot folosi conectorii logici AND, OR și NOT:

$$B ::= \langle \text{termen} \rangle \mid \text{NOT } B \mid B \text{ AND } B \mid B \text{ OR } B$$

2. Se consideră că un document este doar o mulțime de termeni.
3. Documentele din colecție sunt **indexate** înainte de a efectua cereri de informare. Indexarea se face astfel: se citesc toate documentele, se extrag toți termenii care apar în ele, și se indexează, adică se rețin toate documentele în care apare termenul respectiv.
  - Mulțimea tuturor termenilor din colecția  $D$  se numește **vocabular** sau **dicționar** și o notăm  $V$ .
  - Rezultatele indexării pot fi reținute în o **matrice de incidentă**  $V \times D$  în care elementul de la poziția  $(t, d)$  este 1 dacă termenul  $t$  apare în documentul  $d$ , și 0 în caz contrar.

Un exemplu de matrice de incidentă pentru termenii ce apar în câteva din operele lui Shakespeare este ilustrat mai jos:

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Anthony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0
...						

Rândul din matrice corespunzător unui termen  $t$  este un sir de  $n$  biți, unde  $n$  este numărul de documente din  $d$ . De exemplu, sirul de biți pentru Brutus este 110100, pentru Caesar este 110111, și pentru Calpurnia este 010000.

4. Permite căutarea ad-hoc folosind expresii booleene cu termeni ca să comunice necesitățile de informare ale utilizatorilor.

Necesitatea de informare pentru documentele care conțin termenii **Brutus** și **Caesar** dar nu conțin **Calpurnia**, poate fi comunicată cu expresia booleană

**Brutus AND Caesar AND NOT Calpurnia**

iar răspunsul se obține prin operații booleene pe vectori de biți:

```
110100 AND 110111 AND NOT 010000 =
110100 AND 110111 AND 101111 = 100100
```

Documentele returnate ca răspuns la cererea utilizatorului sunt primul și al patrulea, adică Anthony and Cleopatra și Hamlet, fiindcă rezultatul 100100 are primul și al patrulea bit cu valoarea 1.

### 1.1 Statistici de estimare a eficacității unui sistem de IR

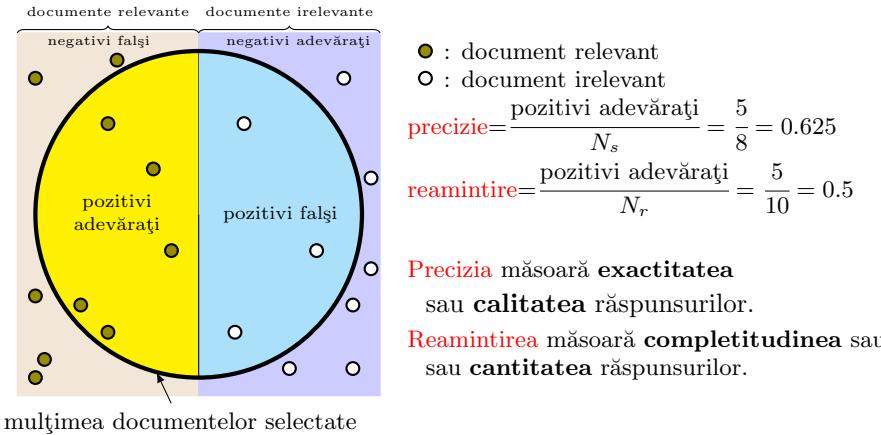
**Eficacitatea** unui sistem de extragere a informațiilor se estimează folosind următoarele măsurători statistice:

- **Precizia:** ce procent din rezultatele returnate de sistem sunt documente relevante pentru necesitatea de informare a utilizatorului?
  - **Reamintirea** (engl. **recall**): ce procent din documentele relevante din colecție sunt printre rezultatele returnate de sistem ca răspuns la necesitatea de informare a utilizatorului (adică pozitivi adevărați)?

Exemplu: Fie  $S$  este un sistem de extragere a informațiilor din o colecție  $D$  de  $N = 19$  documente. Presupunem că un utilizator comunică cererea  $Q$  sistemului  $S$  și că:

- ▶  $D$  conține  $N_r = 10$  documente relevante pentru  $Q$
  - ▶ Sistemul  $S$  returnează  $N_s = 8$  documente ca răspunsuri la cerere.  $N_s$  este numărul documentelor selectate de către  $S$  din  $D$  ca răspunsuri.
  - ▶ 5 din cele 8 răspunsuri sunt documente relevante pentru  $Q$  (sau *pozitivi adevărați*). Celelalte 3 răspunsuri sunt documente irelevante pentru  $Q$  (sau *pozitivi falsi*).

Situatia din acest exemplu este ilustrata in diagrama de mai jos.



## 1.2 Liste de indecșii inversați

Sunt o alternativă de reprezentare a rezultatului indexării termenilor care, de obicei, ocupă mult mai puțin spațiu de memorie decât matricea de incidentă.

De exemplu, dacă  $D$  este o colecție de  $N \approx 10^6$  documente alcătuite fiecare din  $\approx 10^3$  termeni, iar fiecare termen ocupă în medie 6 octeți, atunci toată colecția  $D$  are o mărime de  $\approx 6$  GB. Putem estima că  $D$  conține  $N \approx 5 \cdot 10^5$  termeni distincți.

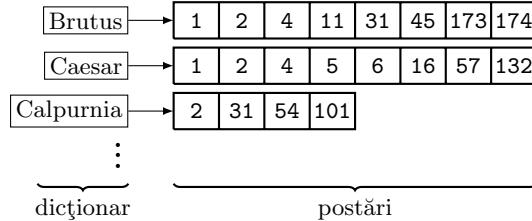
**veste rea:** matricea de incidență pentru colecția  $D$  de documente din exemplu ocupă  $5 \cdot 10^5 \times 10^6 \approx 0.5 \cdot 10^{12}$  biți: nu începe în memoria RAM a unui calculator obișnuit

**veste bună:** matricea conține  $\approx 10^6$  elemente care sunt 1. Rezultă că  $\approx 99.8\%$  din elementele matricii sunt 0  $\Rightarrow$  este mai economic să reținem în liste pozițiile la care apare 1 în matrice.

Aceste liste se numesc **liste de indecsi inversați** sau **liste de postări**.

- pentru fiecare termen  $t \in V$  reținem lista de identificatori de documente (docID) în care apare  $t$ .
- identificatorii de documente sunt reținuți în ordine crescătoare în liste.

De acum încolo vom presupune că fiecare document are un număr serial unic care este o valoare întreagă pozitivă. Acest număr se numește **identificator de document**. De exemplu, listele de indecsi inversați pentru termenii din operele lui Shakespeare pot arăta astfel:



Listele de indecsi inversați se construiesc în 4 pași:

1. Se colectează documentele care urmează să fie indexate  $\Rightarrow$  o colecție de  $N$  documente  $D$  cu identificatorii  $d_1, \dots, d_N$ .
2. Se transformă texte documentelor în liste de token-uri. Acest pas se numește **tokenizare**.

De exemplu, tokenizarea textului

“Friends, Romans, countrymen. So let it be with Caesar.”

este [Friends] [Romans] [countrymen] [So] [let] [it] [be] [with] [Caesar]

3. Se face o preprocesare lingvistică a token-urilor  $\Rightarrow$  o listă de termeni normalizați numiți **termeni**. Termenii sunt elementele care se indexează. De exemplu, lista de termeni pentru textul ilustrat mai sus este

[friend] [roman] [countryman] [so] ...

4. Indexarea propriu-zisă. După pasul 3 se vor cunoaște

- mulțimea  $V$  a tuturor termenilor care apar în textele din colecția  $D$ . Această mulțime de termeni  $V$  se numește **dicționar** și se sortează crescător.
- lista de identificatori de documente în care apare fiecare termen  $t$ . Aceste liste se sortează crescător.
- valori statistice despre textele documentelor. De exemplu, **frecvența de document** a fiecărui termen  $t$ , care reprezintă numărul de documente în care apare termenul  $t$ .

### 1.3 Procesarea cererilor booleene

Procesarea unei cereri booleene  $B_1 \text{ AND } B_2$  produce o listă de postări care satisfac cererile  $B_1$  și  $B_2$ . Această listă se determină în felul următor:

1. Se determină listele de postări  $p_1$  pentru  $B_1$ , și  $p_2$  pentru  $B_2$ .
2. Se intersectează listele  $L_1$  și  $L_2$ . Intersecția se poate face rapid, în timp  $O(m + n)$  unde  $m$  este lungimea lui  $p_1$  și  $n$  este lungimea lui  $p_2$ , cu algoritmul următor:

```

INTERSECT( $p_1, p_2$ )
1 raspuns := ⟨⟩
2 while  $p_1 \neq \text{Nil}$  and  $p_2 \neq \text{Nil}$  do
3     if docID( $p_1$ ) = docID( $p_2$ )
4         adaugă docID( $p_1$ ) la raspuns
5          $p_1 := \text{next}(p_1)$ 
6          $p_2 := \text{next}(p_2)$ 
7     else if docID( $p_1$ ) < docID( $p_2$ )
8          $p_1 := \text{next}(p_1)$ 
9     else  $p_2 := \text{next}(p_2)$ 
10 return raspuns
```

Procesările celorlalte tipuri de cerei booleene

$$\text{NOT } B \quad \text{și} \quad B_1 \text{ OR } B_2$$

se pot defini în mod asemănător.

## 2 Extensiile ale modelului boolean

Modelul boolean de extragere a informațiilor este prea limitat fiindcă (1) fiecare document este interpretat doar ca o mulțime de termeni distinți, și (2) nu poate calcula un scor de importanță al documentelor, care să indice care documente sunt răspunsuri mai importante la o cerere de căutare.

Căteva din aceste limitări au fost eliminate de un **model boolean extins** de extragere a informațiilor. Acest model are un **operator de proximitate** cu

care putem preciza cât de aproape trebuie să fie 2 termeni în documentele care ne interesează.

Westlaw (<http://www.westlaw.com/>) este un sistem comercial de căutare într-o colecție de zeci de terabytes de documente despre legi și legislație. Westlaw are un model boolean extins de extragere a informațiilor. Exemple interesante de cereri de căutare mai sofisticate cu Westlaw sunt:

- "trade secret" /s disclos! /s prevent /s employe!  
căutarea de informații despre legi de prevenire a obținerii de secrete comerciale de către angajați care au lucrat anterior la o companie competitoare.
- disab! /p access! /s work-site work-place (employment /3 place)  
Cerințe pentru persoane cu handicap care își caută un loc de muncă.