

## 2 Reguli de asociere și mulțimi frecvente de articole

Problema coșului de produse presupune că avem un mare număr de articole, e.g. “pâine”, “lapte”. Cumpărătorii pun în coșul lor de produse anumite submulțimi de articole iar noi vom afla ce articole sunt cumpărate împreună chiar dacă nu și de cine. Vânzătorii utilizează aceste informații pentru așezarea articolelor în rafturi și controlează modul în care un cumpărător tipic traversează magazinul.

În afara utilizării în marketing, același tip de problemă are următoarele utilizări:

1. Coș = documente; articole = cuvinte. Cuvintele aparând frecvent împreună în documente pot reprezenta fraze sau concepte legate între ele. Poate fi utilizat pentru colectarea de informații (intelligence).
2. Coș = propoziții, articole = documente. Două documente conținând aceleași propoziții pot reprezenta un plagiat sau “mirror sites” pe web.

### 2.1. Obiective pentru căutarea în coșuri de produse

1. *Regulile de asociere* sunt propoziții de forma  $\{X_1, X_2, \dots, X_n\} \Rightarrow Y$  însemnând că dacă găsim toate  $X_1, X_2, \dots, X_n$  în coș atunci sunt mari șanse să-l găsim și pe  $Y$ . Probabilitatea de a-l găsi pe  $Y$  pentru a accepta această regulă este numită *încrederea* respectivei reguli. În mod normal vom căuta doar reguli care au o încredere peste un anumit prag. Putem de asemenea cere ca încrederea să fie semnificativ mai mare decât în cazul în care articolele ar fi plasate aleator în coș. De exemplu putem găsi o regulă ca  $\{\text{lapte, unt}\} \Rightarrow \text{pâine}$  doar pentru că foarte multă lume cumpără pâine. Totuși exemplul bere/scutece arată că regula  $\{\text{scutece}\} \Rightarrow \{\text{bere}\}$  este verificată cu o încredere semnificativ mai mare decât a submulțimii de coșuri conținând bere.
2. *Ccauzalitate*. Ideal, am vrea să știm dacă într-o regulă de asociere prezența elementelor  $X_1, X_2, \dots, X_n$  efectiv “cauzează” (determină) cumpărarea lui  $Y$ . “Ccauzalitatea” este însă un concept echivoc. Cu toate acestea, pentru datele de tip coș de produse următorul test arată ce înseamnă cauzalitatea. Dacă scădem prețul scutecelor și creștem prețul berii putem ademeni cumpărătorii de scutece care au înclinația de a cumpăra bere din magazin, acoperind astfel pierderile din vânzarea scutecelor. Această strategie este valabilă deoarece “scutece determină bere”. Acțiunea inversă, micșorarea prețului la bere și mărirea prețului scutecelor nu va determina cumpărătorii de bere să cumpere scutece în număr mare și vom pierde bani.
3. *Mulțimi frecvente de articole (frequent itemsets)*. În multe situații (dar nu în toate) ne interesează doar regulile de asociere și cauzalitatea în ceea ce privește mulțimi de articole care apar frecvent în coșuri. De exemplu, nu putem conduce o bună strategie de marketing care implică produse pe care oricum nu le cumpără nimeni. Astfel, extragerea de cunoștințe din date pornește de la premiza că ne interesează doar mulțimile de articole cu un larg suport; i.e., ele apar împreună în multe coșuri de produse. Găsim apoi regulile de asociere sau cauzalitățile implicând doar articolele cu larg suport (i.e.  $\{X_1, X_2, \dots, X_n, Y\}$  trebuie să apară în cel puțin un anumit procent din coșuri, numit *prag de suport*)

### 2.2. Cadru pentru căutarea mulțimilor frecvente de articole

Utilizăm termenul *mulțimi frecvente de articole (frequent itemsets)* pentru “o mulțime de articole  $S$  care apare în cel puțin a “ $s$ ”-a parte din coșuri”, unde  $s$  este o constantă aleasă, de obicei 0.01 sau 1%.

Vom presupune că avem o cantitate de date care nu încap în memoria centrală a calculatorului. Fie sunt stocate într-o bază de date relațională (BDR), de exemplu o relație (tabelă) *Coșuri(IdCoș, articol)* fie ca un fișier de înregistrări de forma  $(\text{IdCoș}, \text{articol1}, \text{articol2}, \dots, \text{articol-n})$ . Când evaluăm timpul de rulare al algoritmilor:

- Numărăm trecerile prin date. Deoarece costul principal este dat adesea de timpul necesar citirii datelor de pe disc, numărul de citiri necesar pentru fiecare dată este adesea cea mai bună măsură a timpului de rulare al algoritmului.

Există un principiu cheie, numit *monotonicitate* (monotonicity) sau *principiul a-priori* care ne ajută să găsim mulțimile frecvente de articole:

- Dacă o mulțime de articole  $S$  este frecventă (i.e., apare cel puțin în a “ $s$ ”-a parte a coșurilor), atunci orice submulțime a lui  $S$  este de asemenea frecventă.

Pentru a găsi mulțimile frecvente de articole:

1. Procedăm nivel cu nivel, găsim întâi articolele frecvente (mulțimi de dimensiune 1), apoi perechile frecvente, tripletele frecvente, etc. În discuția noastră ne vom concentra pe găsirea perechilor frecvente deoarece:
  - a. Adeseori perechile sunt suficiente.
  - b. În multe mulțimi de date partea cea mai grea este găsirea perechilor; continuarea pe nivelele superioare necesită mai puțin timp decât găsirea perechilor frecvente.Algoritmii de acest tip utilizează o trecere per nivel.
2. Găsim toate *mulțimile frecvente de articole maximale* (i.e., mulțimile  $S$  a.i. nici o mulțime care include strict pe  $S$  nu este frecventă) într-o singură trecere sau în câteva treceri.

### 2.3. Algoritmii a-priori

Acest algoritm procedează nivel cu nivel.

1. Dându-se pragul de suport  $s$ , în prima trecere găsim articolele care apar în cel puțin a “ $s$ ”-a parte a coșurilor. Această mulțime este notată  $L_1$ , mulțimea articolelor frecvente.
2. Perechile de articole din  $L_1$  devin mulțimea  $C_2$  a *perechilor candidate* pentru a doua trecere. Sperăm că dimensiunea lui  $C_2$  nu este atât de mare pentru că altfel nu este suficient spațiu în memoria centrală pentru un contor numeric întreg al apariției fiecărei perechi. Perechile din  $C_2$  al căror contor ajunge sau depășește  $s$  formează mulțimea  $L_2$  a perechilor frecvente.
3. Tripletele candidate  $C_3$  sunt mulțimile  $\{A, B, C\}$  pentru care  $\{A, B\}$ ,  $\{A, C\}$  și  $\{B, C\}$  sunt în  $L_2$ . În a treia trecere sunt numărate aparițiile tripletelor din  $C_3$ ; cele al căror contor este cel puțin  $s$  formează mulțimea  $L_3$  a tripletelor frecvente.
4. Se poate merge oricât de departe se dorește (sau până mulțimile devin vide).  $L_i$  conține mulțimile frecvente de articole de dimensiune  $i$ ;  $C_{i+1}$  este mulțimea candidatelor de dimensiune  $i+1$  a.i. fiecare submulțime a lor de dimensiune  $i$  este inclusă în  $L_i$ .

### 2.4. De ce este util principiul a-priori

Să considerăm următoarea cerere SQL pe o relație (tabelă)  $Coșuri(IdCoș, articol)$  având  $10^8$  tupluri care conțin date despre  $10^7$  coșuri de câte 10 articole fiecare. Presupunem existența a 100.000 articole diferite (tipic pentru o rețea de magazine ca Wal-Mart de exemplu).

```
SELECT b1.articol, b2.articol, COUNT(*)
FROM Coșuri b1, Coșuri b2
WHERE b1.IdCoș = b2.IdCoș AND b1.articol < b2.articol
GROUP BY b1.articol, b2.articol
HAVING count(*) >= s;
```

Notă:  $s$  este pragul de suport iar al doilea termen al clauzei WHERE elimină perechile formate din același produs și apariția de două ori a aceleiași perechi.

În joinul  $Coșuri \bowtie Coșuri$  fiecare coș contribuie cu  $C_{10}^2 = 45$  de perechi astfel încât joinul are  $4,5 \times 10^8$  tupluri. A-priori “împinge clauza HAVING în jos pe arborele expresiei”, determinându-ne în primul rând să înlocuim  $Coșuri$  cu rezultatul cererii:

```
SELECT *
FROM Coșuri
GROUP BY articol
HAVING COUNT(*) >= s;
```

Dacă  $s = 0,01$  atunci cel mult 1000 de grupuri de articole pot trece de clauza HAVING. Motivul: în relația Coșuri sunt  $10^8$  linii conținând articole iar fiecare articol are nevoie de  $0,01 \times 10^7 = 10^5$  dintre acestea pentru a apărea în 1% din coșuri.

- Deși 99% dintre articole sunt eliminate de algoritmul a-priori nu trebuie să concluzionăm că relația Coșuri care rezultă are doar  $10^6$  tuple. În fapt, *toate* tuplele pot fi pentru produse cu larg suport. Totuși, în situațiile reale, micșorarea relației Coșuri este substanțială și dimensiunea joiului scade cu pătratul acestei micșorări.

## 2.5. Îmbunătățiri ale algoritmului a-priori

Sunt de două tipuri:

1. Micșorarea dimensiunii mulțimilor candidat  $C_i$  pentru  $i \geq 2$ . Această opțiune este importantă chiar și pentru găsirea perechilor frecvente deoarece numărul de elemente trebuie să fie suficient de mic pentru ca un contor de apariții pentru fiecare să fie ținut în memoria centrală.
2. Contopirea încercărilor de găsire a  $L_1, L_2, L_3, \dots$  în doar una sau două treceri în loc de o trecere per nivel.

## 2.6. Algoritmul PCY

Park, Chen și Yu au propus, utilizând o tabelă de dispersie, să determine la prima trecere (când este calculat  $L_1$ ) că multe perechi nu sunt perechi frecvente posibile. Profită de faptul că memoria centrală este uzual *mult* mai mare decât numărul de articole. În timpul celor două faze pentru găsirea lui  $L_2$  memoria centrală este ocupată ca în Fig. 1.

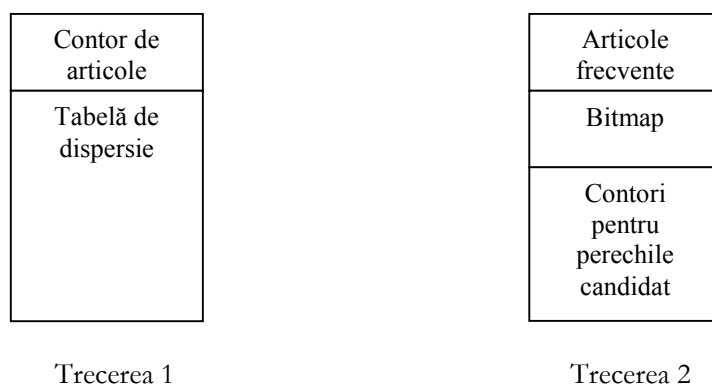


Figura 1: Două treceri ale algoritmului PCY

Presupunem că datele sunt stocate într-un fișier de înregistrări constând dintr-un identificator IdCoș și o listă cu articolele sale.

1. Trecerea 1:
    - a. Se număra aparițiile fiecărui articol.
    - b. Pentru fiecare coș constând din articolele  $\{i_1, \dots, i_k\}$ , se aplica funcția de dispersie fiecărei perechi asociind-o unei intrări a tabelii de dispersie și se incrementează contorul acesteia cu 1.
    - c. La sfârșitul trecerii, se determină  $L_1$ , articolele cu contorul cel puțin  $s$ .
    - d. De asemenea la sfârșitul trecerii se determină acele intrări ale tabelii de dispersie cu contor cel puțin  $s$ .
      - Punct cheie: o pereche  $(i, j)$  nu poate fi frecventă decât dacă este dispersată într-o intrare frecventă astfel încât perechile care sunt dispersate în alte intrări nu vor fi candidate în  $C_2$ .
- Se înlocuiește tabela de dispersie cu un bitmap având un bit per buchet: 1 dacă buchetul a fost frecvent, 0 altfel.

## 2. Trecerea 2:

- a. Memoria centrală conține o listă cu toate articolele frecvente, i.e.  $L_1$ .
- b. Tot memoria centrală conține un bitmap reprezentând rezultatele dispersiei din prima trecere.
  - Punct cheie: intrările tabeli de dispersie utilizează 16 sau 32 de biți pentru un contor dar sunt comprimate la un singur bit. Astfel, chiar dacă tabela de dispersie ocupă aproape întreaga memorie centrală la prima trecere, bitmapul său nu ocupă mai mult de 1/16 din memoria centrală la trecerea 2.
- c. În final, memoria centrală conține de asemenea o tabelă cu toate perechile candidat și contorii asociați lor. O pereche  $(i, j)$  poate fi candidată în  $C_2$  doar dacă *toate* condițiile următoare sunt adevărate:
  - i.  $i$  este în  $L_1$ .
  - ii.  $j$  este în  $L_1$ .
  - iii.  $(i, j)$  este dispersată într-o intrare frecventăUltima condiție diferențiază PCY de a-priori clasic și reduce necesarul de memorie în trecerea 2.
- d. În timpul trecerii 2 luăm în considerare fiecare coș și fiecare pereche de articole din el, efectuând testul de mai sus. Dacă o pereche îndeplinește toate cele trei condiții, se incrementează contorul acesteia din memorie sau se crează unul dacă acesta nu exista încă.
  - Când este mai performant PCY decât a-priori? Când sunt prea multe perechi de articole din  $L_1$  pentru a încăpea într-o tabelă de perechi candidate și de contori asociați în memoria centrală iar numărul de intrări frecvente ale tabeli de dispersie din algoritmul PCY este suficient de mic pentru a reduce dimensiunea lui  $C_2$  suficient pentru a încăpea în memoria centrală (chiar și fără 1/16 din ea consumată de bitmap).
  - Când o mare parte a intrărilor tabeli de dispersie nu vor fi frecvente în PCY? Când sunt puține perechi frecvente iar cea mai mare parte a perechilor sunt atât de puțin frecvente încât chiar dacă contorii tuturor perechilor care sunt dispersate într-o aceeași intrare sunt adunați nu sunt mari șanse să se obțină o valoare egală sau mai mare ca  $s$ .

## 2.7. Extensia “Iceberg” a PCY

1. *Tabele de dispersie multiple*: se împarte memoria între două sau mai multe tabele de dispersie, ca în Fig. 2. La trecerea 2 se reține în memorie câte un bitmap pentru fiecare dintre acestea; de notat că spațiul necesar pentru aceste bitmapuri este exact același cu cel necesar în bitmapul unic din PCY deoarece numărul total de intrări reprezentate este același. Pentru a fi candidată la  $C_2$  o pereche:
  - a. Constă din articole din  $L_1$ , și
  - b. Este dispersată într-o intrare frecventă în *fiecare* tabelă de dispersie.

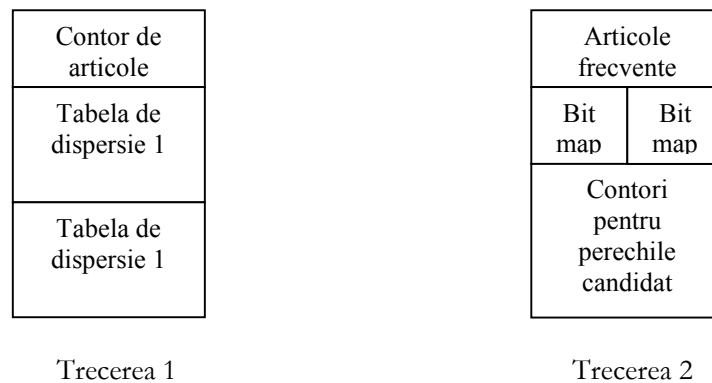


Figura 2: Utilizarea memoriei pentru tabele de dispersie multiple

2. Tabele de dispersie iterate *Multistage*: În locul verificării candidatelor în trecerea 2 se crează o altă tabelă de dispersie (altă funcție de dispersie) în trecerea 2 dar sunt dispersate doar acele perechi care îndeplinesc condițiile pentru PCY; i.e., ambele articole sunt din  $L_1$  și sunt dispersate într-un buchet frecvent în prima trecere. În a treia trecere, păstrăm câte un bitmap pentru fiecare tabela de dispersie și tratăm o pereche ca o candidată din  $C_2$  doar dacă:
- Ambele articole sunt în  $L_1$ .
  - Perechea a fost dispersată într-o intrare frecventă în prima trecere.
  - Perechea a fost dispersată de asemenea într-o intrare frecventă în trecerea 2.

Figura 3 arată utilizarea memoriei. Aceasta schemă poate fi extinsă la mai multe treceri, dar există o limită deoarece eventual memoria devina plină de bitmapuri și nu mai putem contoriza perechile candidat.

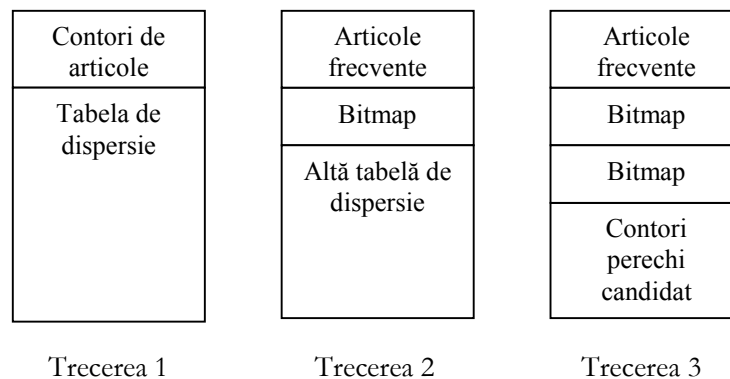


Figura 2: Utilizarea memoriei pentru tabele de dispersie Multistage

- Când sunt utile tabelele de dispersie multiple? Când cele mai multe dintre intrările tabelii de dispersie de la prima trecere a PCY au contori cu mult sub pragul  $s$ . Atunci putem dubla contorii intrărilor și totuși cele mai multe dintre ele vor fi sub prag.
- Când sunt utile tabelele de dispersie Multistage? Când numărul intrărilor frecvente din prima trecere este mare, dar nu toate (e.g. 50%). Atunci, a doua dispersie cu unele dintre perechile ignorate poate reduce numărul de intrări frecvente semnificativ.

## 2.8. Toate mulțimile frecvente de articole în două treceri

Metodele de mai sus sunt cele mai bune când dorim perechile frecvente, cazul cel mai comun. Dacă dorim toate mulțimile frecvente maximale de articole, inclusiv mulțimi mari, sunt necesari prea mulți pași. Există mai multe abordări pentru obținerea tuturor mulțimilor frecvente de articole în două treceri sau mai puțin. Fiecare se bazează într-un fel sau altul pe faptul că datele sunt aleatoare.

1. *Abordarea simplă*: Se ia un eșantion de date de dimensiunea memoriei centrale. Se rulează un algoritm nivel cu nivel în memoria centrală (deci nu sunt costuri de I/O) și se speră că eșantionul ne va conduce la adevăratele mulțimi frecvente.
  - De notat că pragul  $s$  trebuie scalat prin micșorare; e.g. dacă eșantionul este de 1% din date, se utilizează  $s/100$  ca prag de suport.
  - Se poate face o trecere completă prin date pentru a verifica dacă mulțimile frecvente de articole ale eșantionului sunt cu adevărat frecvente, dar vor fi pierdute mulțimile de articole care sunt frecvente în ansamblul datelor dar nu și în eșantion.
  - Pentru a minimiza falsele negative se poate scădea puțin pragul pentru eșantion găsimdu-se mai multe candidate pentru trecerea prin ansamblul datelor. Risc: prea multe candidate pentru a încăpea în memoria centrală.

2. *SON95* (Savasere, Omnecinski și Navathe în VLDB 1995; referit de Toivonen). Se citesc submulțimi ale datelor în memoria centrală aplicându-se abordarea simplă pentru descoperirea mulțimilor candidat. Fiecare coș este parte a uneia dintre aceste submulțimi. În a doua trecere o mulțime este candidată dacă a fost identificată ca și candidată în una sau mai multe submulțimi ale datelor.
  - Punct cheie: O mulțime nu poate fi frecventă în ansamblul datelor dacă nu este frecventă în cel puțin o submulțime a acestora.
3. *Algoritmul lui Toivonen:*
  - a. Se ia un eșantion care încapă în memoria centrală. Se rulează abordarea simplă pe aceste date dar cu un prag micșorat astfel încât să fie improbabilă pierderea vreunei adevărate mulțimi frecvente de articole (e.g. dacă eșantionul este de 1% din date se folosește  $s/125$  ca prag de suport).
  - b. Se adaugă candidatelor din eșantion *marginea negativă*: acele mulțimi de articole  $S$  astfel încât  $S$  nu este identificată ca frecventă în eșantion dar *orice* submulțime strictă maximală a lui  $S$  este identificată astfel. De exemplu, dacă  $ABCD$  nu este frecventă în eșantion dar  $ABC$ ,  $ABD$ ,  $ACD$  și  $BCD$  sunt frecvente în eșantion, atunci  $ABCD$  este în marginea negativă.
  - c. Se face o trecere prin date, numărând toate mulțimile frecvente de articole și marginea negativă. Dacă nici o mulțime din marginea negativă nu este frecventă în ansamblul datelor, atunci mulțimile frecvente de articole sunt exact acele candidate care sunt deasupra pragului.
  - d. Din păcate, dacă există o mulțime din marginea negativă care devine frecventă, atunci nu știm dacă vreuna dintre mulțimile care o conțin nu este de asemenea frecventă, astfel încât întregul proces trebuie să fie repetat (sau acceptăm că există și nu ne interesează câteva fals negative).