

### 3. Căutare în date pentru suport scăzut, corelare puternică

Vom continua să considerăm modelul de date “coșuri de produse” și vom vizualiza datele ca o matrice booleană unde linii=coșuri și coloane=articole. Aserțiuni importante:

1. Matricea este foarte rară; aproape peste tot 0.
2. Numărul de coloane (articole) este suficient de mic pentru a putea stoca în memoria centrală ceva per coloană dar suficient de mare astfel încât nu putem stoca ceva per pereche de articole în memoria centrală (aceeași aserțiune pe care am făcut-o până acum privind regulile de asociere).
3. Numărul de linii este atât de mare încât nu putem stoca întreaga matrice în memorie chiar profitând de faptul ca e rară și comprimând-o (din nou aceeași aserțiune ca întotdeauna).
4. Nu suntem interesați de perechile sau mulțimile de coloane cu larg suport; în schimb dorim perechile de coloane puternic corelate.

#### 3.1. Aplicații

În timp de aplicațiile de marketing sunt interesate doar de produsele de larg consum (nu merită să se încerce promovarea obiectelor pe care oricum nu le cumpără nimeni), există un număr de aplicații care se potrivesc cu modelul de mai sus, de interes fiind în special problema perechilor de coloane/articole de consum restrâns dar puternic corelate:

1. Liniiile și coloanele sunt pagini de web;  $(r, c) = 1$  înseamnă că pagina corespunzătoare liniei  $r$  conține o legătură către pagina coloanei  $c$ . Coloanele similare pot fi pagini despre același domeniu.
2. La fel ca (1) dar pagina corespunzătoare coloanei  $c$  conține legături către pagina liniei  $r$ . Acum, coloane similare pot reprezenta copii multiple (mirror) ale unei pagini.
3. Linii = pagini web sau documente; coloane = cuvinte. Coloane similare reprezintă cuvinte care apar aproape mereu împreună, e.g. “frază”.
4. La fel ca (3) dar liniile sunt propoziții [iar coloanele sunt pagini web]. Coloane similare pot indica copii multiple ale unei unei pagini sau plagiat.

#### 3.2. Similaritate.

Să ne gândim la o coloană ca la multimea liniilor pentru care coloana conține 1. Atunci *similaritatea* a două coloane  $C_1$  și  $C_2$  este  $Sim(C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$ .

**Exemplul 3.1:**

$$\begin{array}{cc} 0 & 1 \\ 1 & 0 \\ 1 & 1 = 2/5 = 40\% \text{ similare} \\ 0 & 0 \\ 1 & 1 \\ 0 & 1 \end{array}$$

#### 3.3. Signaturi

Ideia principală: Se mapează (“dispersează”) fiecare coloană  $C$  într-o cantitate mică de date [*signatura*,  $Sig(C)$ ] astfel încât:

1.  $Sig(C)$  este suficient de mică pentru ca signaturile tuturor coloanelor să încapă în memoria centrală.
2. Coloanele  $C_1$  și  $C_2$  sunt puternic similare dacă și numai dacă  $Sig(C_1)$  și  $Sig(C_2)$  sunt puternic similare. (dar de notat că este nevoie să definim “similaritatea” pentru signaturi).

O idee care însă nu funcționează: Se iau aleator 100 de linii și șirul de 100 de biti ai coloanelor [pentru acele linii] este semnatura fiecărei coloane. Motivul [pentru care ideea nu funcționează] este că matricea este presupusă ca fiind rară deci multe coloane vor avea semnături formate doar din 0 chiar dacă ele nu sunt deloc similare.

Convenție utilă: dându-se două coloane  $C_1$  și  $C_2$ , ne vom referi la liniile lor ca fiind de patru tipuri –  $a$ ,  $b$ ,  $c$ ,  $d$  – în funcție de biții lor pe aceste coloane, după cum urmează:

| Tip | $C_1$ | $C_2$ |
|-----|-------|-------|
| a   | 1     | 1     |
| b   | 1     | 0     |
| c   | 0     | 1     |
| d   | 0     | 0     |

De asemenea vom utiliza  $a$  pentru “numărul de linii de tip  $a$ ”, ș.a.m.d.

- De notat că  $Sim(C_1, C_2) = a / (a+b+c)$ .
- Dar cum cele mai multe linii sunt de tip  $d$ , într-o selecție de, să spunem, 100 de linii alese aleator toate vor fi de tip  $d$ , deci similaritatea coloanelor doar pentru aceste 100 de linii nici nu este definită.

### 3.4. Dispersia de tip Min

Să ne imaginăm liniile permutate într-o ordine aleatoare. “Dispersăm” fiecare coloană  $C$  în  $b(C)$ , numărul primei linii în care coloana  $C$  are un 1.

- Probabilitatea ca  $b(C_1) = b(C_2)$  este  $a / (a+b+c)$  deoarece valoarea de dispersie este aceeași dacă prima linie cu un 1 în vreuna din coloane este de tip  $a$  și este diferită dacă prima astfel de linie este de tip  $b$  sau  $c$ . De notat că această probabilitate este aceeași cu  $Sim(C_1, C_2)$ .
- Dacă repetăm experimentul cu o nouă permutare a liniilor de un număr mare de ori, să zicem 100, obținem o semnătură constând din 100 de numere de linii pentru fiecare coloană. “Similaritatea” acestor liste (fracțiune a pozițiilor în care ele sunt egale) va fi foarte apropiată de similaritatea coloanelor.
- Observație importantă: nu trebuie să permutăm fizic liniile, ceea ce ar duce la multe treceri prin întreaga cantitate de date. În schimb citim liniile într-o ordine oarecare și dispersăm fiecare linie [numărul acesteia] utilizând (să zicem) 100 de funcții de dispersie diferite. Pentru fiecare coloană memorăm cea mai mică valoare a funcției de dispersie a unei linii în care coloana are un 1, independent pentru fiecare dintre cele 100 de funcții de dispersie. După parcurgerea tuturor liniilor vom avea pentru fiecare coloană primele linii în care coloana are 1 dacă liniile ar fi fost permutate în ordinea dată de fiecare dintre cele 100 de funcții de dispersie.

### 3.5. Dispersia senzitivă la localizare

Problema: avem semnăturile fiecărei coloane în memoria centrală iar semnături similare înseamnă cu mare probabilitate coloane similare, dar pot fi totuși atât de multe coloane încât a face ceva care este proportional cu pătratul numărului de coloane, chiar și în memoria centrală, este prohibitiv. Dispersia senzitivă la localizare (DSL, LSH în engleză) este o tehnică destinată a fi utilizată în memoria centrală pentru a aproxima mulțimea de perechi de coloane similare cu o complexitate mult mai mică decât cea pătratică.

Scopul: în timp proporțional cu numărul de coloane să se elimine cea mai mare parte a perechilor de coloane din mulțimea posibilelor perechi similare.

1. Considerăm semnatura ca fiind o coloană de întregi.
2. Partizionăm liniile signaturilor în *benzi*, să spunem  $l$  benzi de câte  $r$  linii fiecare.
3. Dispersăm coloanele din fiecare bandă în intrări ale unei table de dispersie. O pereche de coloane este o pereche-candidat dacă ambele sunt dispersate în aceeași intrare în vreo bandă.

4. După identificarea candidaților se verifică fiecare pereche candidat  $(C_i, C_j)$  examinând pentru similaritate  $Sig(C_i)$  și  $Sig(C_j)$ .

**Exemplul 3.2:** Pentru a vedea efectul DSL să considerăm date cu 100.000 de coloane și semnături constând din 100 de întregi fiecare. Signaturile ocupă 40Mb de memorie, nu atât de mult la standardele actuale. Să presupunem că vrem perechile care sunt 80% similare. Vom examina signaturile în loc de coloane, deci în mod real vom identifica coloanele ale caror *signaturi* sunt 80% similare – deci nu chiar același lucru.

- Dacă două coloane sunt 80% similare atunci probabilitatea ca ele să fie identice în una dintre benzile de 5 întregi este  $(0,8)^5 = 0,328$ . Probabilitatea ca ele să nu fie identice în *nici una* dintre cele 20 de benzi este  $(1-0,328)^{20} = 0,00035$ . Astfel toate mai puțin aproximativ 1/3000 dintre perechile cu semnături 80% similare vor fi identificate ca și candidate.
- Acum, să presupunem că două coloane sunt doar 40% similare. Atunci probabilitatea ca ele să fie identice într-o bandă este  $(0,4)^5 = 0,01$  iar probabilitatea ca ele să fie identice în cel puțin una dintre cele 20 de benzi nu este mai mare ca 0,2. Astfel, putem ignora cel puțin 4/5 dintre perechi care nu vor deveni candidate dacă 40% este similaritatea tipică a coloanelor.
- În fapt, cele mai multe perechi vor fi cu mult mai puțin decât 40% similare astfel încât realmente eliminăm o parte importantă a perechilor de coloane care nu sunt similare.

### 3.6. Dispersia $k$ -Min

Dispersia Min ne cere să dispersăm fiecare număr de linie de  $k$  ori dacă vrem semnături de  $k$  întregi. În loc de asta, în cazul *dispersiei  $k$ -min* dispersăm fiecare linie o singură dată și, pentru fiecare coloană luăm ca semnătură numerele primelor  $k$  linii în care acea coloană are un 1.

Pentru a vedea de ce similaritatea acestor semnături este aproape aceeași cu similaritatea coloanelor din care derivă să examinăm Fig. 4. Această figură reprezintă signaturile  $Sig_1$  și  $Sig_2$  pentru coloanele  $C_1$  și respectiv  $C_2$  în cazul în care liniile au fost permutate în ordinea valorilor funcției de dispersie și liniile de tip  $d$  (în care nici o coloană nu are 1) sunt omise. Astfel, vedem doar liniile de tip  $a, b$  și  $c$  și indicăm că o linie este în semnatura printr-un 1.

|             |      |      |          |
|-------------|------|------|----------|
|             | Sig1 | Sig2 |          |
| 100 de<br>1 | 1    | 1    | 100 de 1 |
|             | 1    | 0    |          |
|             | 1    | 1    |          |
|             | 0    | 1    |          |
|             | .    | .    |          |
|             | .    | .    |          |
|             | .    | .    |          |
|             | 1    | 1    |          |
|             | 1    | 0    |          |
|             | 1    | 1    |          |

Figura 4: Exemplu de semnături pentru două coloane în cazul dispersiei  $k$ -Min

Să presupunem că  $c \geq b$  astfel încât situația tipică (pentru  $k = 100$ ) este cea din Fig. 4: cele 100 de linii pentru prima coloană includ unele linii care nu sunt printre cele 100 ale celei de-a doua coloane. Atunci o estimare a similarității lui  $Sig_1$  și  $Sig_2$  poate fi calculată astfel:

$$| Sig_1 \cap Sig_2 | = \frac{100a}{a+c}$$

deoarece, în medie, primele 100 de linii din  $C_2$  care sunt și în  $C_1$  este  $a/(a + c)$ . De asemenea:

$$| Sig_1 \cup Sig_2 | = 100 + \frac{100c}{a+c}$$

Motivul este că toate cele 100 de linii din  $Sig_1$  sunt în reuniune. În plus, liniile din  $Sig_2$  care nu sunt în  $Sig_1$  sunt de asemenea în reuniune, iar acestea sunt în medie în număr de  $100c/(a+c)$ . Astfel, similaritatea lui  $Sig_1$  cu  $Sig_2$  este:

$$| Sig_1 \cap Sig_2 | / | Sig_1 \cup Sig_2 | = \frac{100a}{a+c} / (100 + \frac{100c}{a+c}) = \frac{a}{a+2c}$$

Observăm că dacă  $c$  este apropiat ca valoare de  $b$  atunci similaritatea signaturilor este apropiată de similaritatea coloanelor. În fapt, dacă două coloane sunt foarte similare, atunci  $b$  și  $c$  sunt ambele mici comparate cu  $a$  și similaritățile signaturilor și coloanelor *trebuie* să fie apropiate.

### 3.7. Amplificarea lui 1 (DSL Hamming)

În cazul în care coloanele nu sunt rare ci au aprox. 50% de 1 nu avem nevoie de dispersie Min; o colecție aleatoare de linii servește ca semnătură. *DSL Hamming* construiește o serie de matrici, fiecare având jumătate din numărul de linii ale precedentei, aplicând operatorul SAU (OR) la câte două linii succesive din matricea precedentă, ca în Fig. 5.

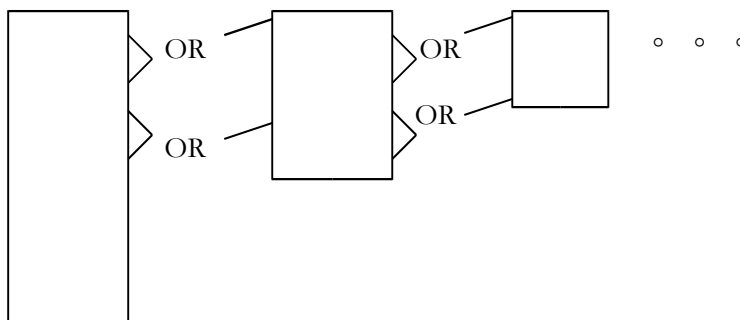


Figura 5: Construcția unei serii de matrici exponențial mai mici și mai dense

- Nu exista mai mult de  $\log n$  matrici, unde  $n$  este numărul de linii. Numărul total de linii în toate matricile este  $2n$  și pot fi calculate toate într-o singură trecere prin matricea originală, stocandu-le pe cele mari pe disc.
- În fiecare matrice, se aleg ca perechi candidat acele coloane care:
  1. Au o densitate de 1 să zicem între 20% și 80%
  2. E posibil să fie similare bazat pe testul DSL
- Observam ca plaja de densitate 20% - 80% ne garantează că doua coloane care sunt cel puțin 50% similare vor fi considerate împreună în cel puțin o matrice, în afara cazului nefericit în care densitatea lor relativă se schimbă din cauza operației OR care combină doi de 1 într-unul singur.
- O a doua trecere prin datele originale confirmă care dintre candidate sunt întradevăr similare.
- Aceasta metodă exploateaza o idee care poate fi folositoare și în alte părți: coloanele similare au un număr similar de 1, deci nu are rost compararea coloanelor al caror număr de 1 este foarte diferit