

II. FRAGMENTAREA RELAȚIILOR

1. Fragmentarea orizontală primară

Se consideră relația COUNTRIES_ALL.

Fie p_1 și p_2 două predicate:

p_1 : region_name= 'Europe'

p_2 : region_name= 'Americas'

Observație:

Se presupune că:

- aceste predicate sunt cel mai des utilizate în aplicații;
- probabilitățile de acces ale fiecărei aplicații la orice tuplu din fragmentele orizontale ce sunt definite pe baza acestor predicate sunt egale;
- există cel puțin o aplicație care să acceseze în mod diferit fragmentele determinate de aceste predicate.

Exercițiu:

- a) să se fragmenteze orizontal relația COUNTRIES_ALL pe baza celor două predicate;

Observație: Fragmentul COUNTRIES_AM va fi stocat pe stația 1 (*server1*), iar fragmentul COUNTRIES_EU va fi stocat pe stația 2 (*server2*).

- conectați-vă la *bd_am* și creați tabelul *countries_am_**** folosind informațiile din tabelul *countries_all* pentru care numele regiunii este *Americas*.

- Încercați comanda:

```
CREATE TABLE countries_eu_***@bd_eu AS
SELECT * FROM countries_all
WHERE region_name='Europe';
```

Ce observați?

- conectați-vă la *bd_eu* și creați tabelul *countries_eu_**** folosind informațiile din tabelul distant *countries_all* (numele legăturii este *bd_am*) pentru care numele regiunii este *Europe*.

- b) să se verifice că fragmentarea realizată este corectă;

- $Pr = \{p_1, p_2\}$ este mulțimea de predicate simple, deoarece predicatele simple p_1 și p_2 sunt cele mai des utilizate predicate în aplicații.
- Pr este completă și minimală.
 - Deoarece există probabilități de acces egale ale fiecărei aplicații la orice tuplu care aparține oricărui fragment orizontal ce este definit pe baza acestei mulțimi, rezultă că mulțimea Pr este completă.
 - Mulțimea Pr este minimală, deoarece toate predicatele mulțimii Pr sunt relevante (un predicat simplu determină divizarea unui fragment f în fragmentele f_i și f_j , dacă există cel puțin o aplicație care să acceseze

fragmentele f_i și f_j în mod diferit; acest predicat simplu va fi relevant pentru fragmentare).

- Se determină mulțimea M de predicate compuse pe baza mulțimii de predicate simple Pr .

$Pr = \{p_1, p_2\}$, unde

p_1 : region_name = 'Europe'

p_2 : region_name = 'Americas'

Prin interogarea tabelului COUNTRIES_ALL se deduce că

region_name \in {'Europe', 'Americas'}

Se construiește mulțimea de implicații I :

i_1 : (region_name= 'Europe') \Rightarrow \neg (region_name= 'Americas')

i_2 : \neg (region_name= 'Europe') \Rightarrow (region_name= 'Americas')

Se construiește mulțimea M care are patru predicate compuse:

m_1 : (region_name= 'Europe') \wedge (region_name= 'Americas')

m_2 : (region_name= 'Europe') \wedge \neg (region_name= 'Americas')

m_3 : \neg (region_name= 'Europe') \wedge (region_name= 'Americas')

m_4 : \neg (region_name= 'Europe') \wedge \neg (region_name= 'Americas')

Se observă că:

- predicatele m_1 și m_4 nu au sens în raport cu implicațiile mulțimii I și de aceea trebuie eliminate din mulțimea M ;
- $m_2 = p_1$ și $m_3 = p_2$.

Deci mulțimea de predicate care determină fragmentele este $M = Pr = \{p_1, p_2\}$

- Fragmentarea este corectă dacă sunt îndeplinite regulile de completitudine, reconstrucție și disjuncție.

Completitudinea: M este completă, deci proprietatea de completitudine este asigurată.

Reconstrucția: Se observă că:

$COUNTRIES_ALL = COUNTRIES_AM \cup COUNTRIES_EU$

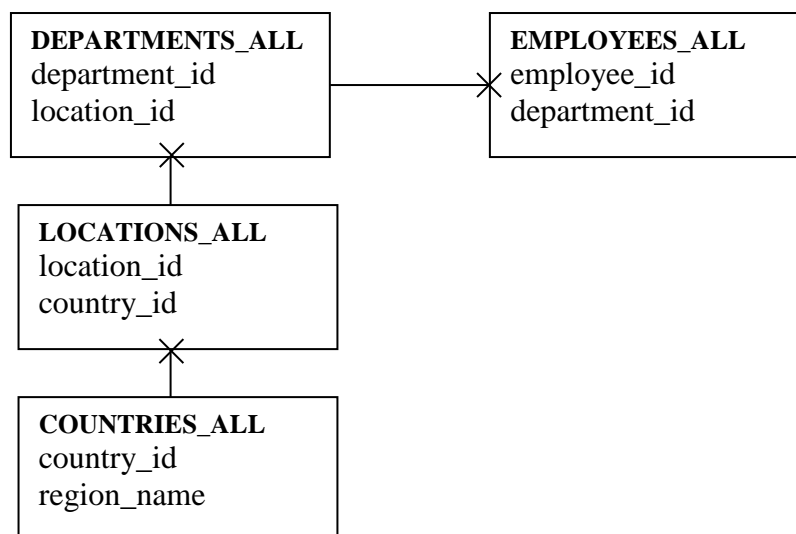
Disjuncția:

$COUNTRIES_AM \cap COUNTRIES_EU = \emptyset$

- să se afișeze relația globală (toate țările), utilizând cele două fragmente.
- să se verifice cele 3 reguli de corectitudine a fragmentării, folosind comenzi SQL.

2. Fragmentarea orizontală derivată

Se consideră următoarele relații:



Aceste relații vor fi fragmentate aceste relații, pornind de la predicatele:

p₁: region_name= 'Europe'

p₂: region_name= 'Americas'

Fragmentele vor fi stocate astfel:

- pe stația 1
LOCATIONS_AM, DEPARTMENTS_AM, EMPLOYEES_AM
- pe stația 2
LOCATIONS_EU, DEPARTMENTS_EU, EMPLOYEES_EU

Exercițiu:

- a) creați fragmentul *locations_am_**** utilizând următoarele comenzi.

```

CONNECT magi_bdd/parola@bd_am;
CREATE TABLE locations_am_*** AS
  SELECT a.*
  FROM   locations_all a, countries_am_*** b
  WHERE  a.country_id=b.country_id;

```

- b) creați fragmentele *departments_am_**** și *employees_am_****.

- c) creați fragmentul *locations_eu_**** utilizând următoarele comenzi.

```

CONNECT magi_bdd/parola@bd_eu;
CREATE TABLE locations_eu_*** AS
  SELECT a.*
  FROM   locations_all@bd_am a, countries_eu_*** b
  WHERE  a.country_id=b.country_id;

```

- d) creați fragmentele *departments_eu_**** și *employees_eu_****.

- e) verificarea corectitudinii fragmentării orizontale derivate;

- LOCATIONS_ALL

Completitudinea: Este asigurată, deoarece fiecărui tuplu din cele două fragmente ale relației LOCATIONS_ALL îi corespunde un tuplu din relația COUNTRIES_ALL.

Reconstrucția: Se observă că:

$$\text{LOCATIONS_ALL} = \text{LOCATIONS_AM} \cup \text{LOCATIONS_EU}$$

Disjuncția: Este asigurată, deoarece graful ce reprezintă diagrama relațiilor (COUNTRIES_ALL, LOCATIONS_ALL) și a legăturilor dintre acestea este un graf simplu (există o singură legătură ce pleacă sau vine pentru orice fragment).

- DEPARTMENTS_ALL similar cu LOCATIONS_ALL;
- EMPLOYEES_ALL

Completitudinea: similar cu LOCATIONS_ALL.

Reconstrucția: similar cu LOCATIONS_ALL.

Disjuncția: Graful ce reprezintă diagrama relațiilor (DEPARTMENTS_ALL, EMPLOYEES_ALL) și a legăturilor dintre acestea nu este un graf simplu. În aceste condiții se verifică valorile fiecărui tuplu.

- f) verificați cele 3 reguli de corectitudine a fragmentării orizontale derivate utilizând comenzi SQL;
- g) folosind fragmentele create anterior, să se listeze codul angajatului, numele departamentului, orașul, țara și regiunea.

3. Fragmentarea verticală

Exercițiu:

a) Să se aplice algoritmi de fragmentare verticală pentru relația JOBS_ALL (job_id, job_title, min_salary, max_salary), știind că:

- cele mai des utilizate cereri sunt următoarele:
 - q_1 : SELECT min_salary FROM jobs_all WHERE max_salary = 10000;
 - q_2 : SELECT min_salary FROM jobs_all WHERE min_salary > 10000;
 - q_3 : SELECT max_salary FROM jobs_all WHERE min_salary > 10000;
 - q_4 : SELECT job_title FROM jobs_all WHERE max_salary = 40000;
- $nr_acc_l(q_k)$ reprezintă numărul de accesări ale atributelor (A_i, A_j) pentru fiecare execuție a cererii q_k pe stația S_l ; (presupunem că $nr_acc_l(q_k) = 1 \quad \forall k \in \overline{1,4}, \forall l \in \overline{1,2}$.)
- $fr_acc_l(q_k)$ reprezintă frecvența de acces a cererii q_k pe stația S_l .

Stația 1	Stația 2
$fr_acc_1(q_1) = 20$	$fr_acc_2(q_1) = 0$
$fr_acc_1(q_2) = 0$	$fr_acc_2(q_2) = 10$
$fr_acc_1(q_3) = 20$	$fr_acc_2(q_3) = 0$
$fr_acc_1(q_4) = 10$	$fr_acc_2(q_4) = 5$

Soluție:

Fie $A_0 = \text{job_id}$, $A_1 = \text{job_title}$, $A_2 = \text{min_salary}$, $A_3 = \text{max_salary}$.

$$\text{ref}(q_i, A_j) = \begin{cases} 1, & \text{dacă } A_j \text{ este referit de } q_i \\ 0, & \text{altfel} \end{cases}$$

- Valorile folosirii atributelor sunt definite de matricea VA, $VA(i, j) = \text{ref}(q_i, A_j)$.

VA	A ₁	A ₂	A ₃
q ₁	0	1	1
q ₂	0	1	0
q ₃	0	1	1
q ₄	1	0	1

- Matricea afinității atributelor:

$$af(A_i, A_j) = \sum_{k \in K} \sum_{\forall S_l} nr_acc_l(q_k) f_acc_l(q_k),$$

unde $K = \{k \mid \text{ref}(q_k, A_i) = 1 \wedge \text{ref}(q_k, A_j) = 1\}$;

$$af(A_1, A_2) = 0$$

$$af(A_1, A_3) = fr_acc_1(q_4) + fr_acc_2(q_4) = 10 + 5 = 15$$

$$\begin{aligned} af(A_2, A_3) &= fr_acc_1(q_1) + fr_acc_2(q_1) + fr_acc_1(q_3) + fr_acc_2(q_3) \\ &= 20 + 0 + 10 + 0 = 40 \end{aligned}$$

$$af(A_1, A_1) = fr_acc_1(q_4) + fr_acc_2(q_4) = 10 + 5 = 15$$

$$\begin{aligned} af(A_2, A_2) &= fr_acc_1(q_1) + fr_acc_2(q_1) + fr_acc_1(q_2) + fr_acc_2(q_2) + fr_acc_1(q_3) \\ &\quad + fr_acc_2(q_3) = 50 \end{aligned}$$

$$\begin{aligned} af(A_3, A_3) &= fr_acc_1(q_1) + fr_acc_2(q_1) + fr_acc_1(q_3) + fr_acc_2(q_3) + fr_acc_1(q_4) \\ &\quad + fr_acc_2(q_4) = 55 \end{aligned}$$

AA	A ₁	A ₂	A ₃
A ₁	15	0	15
A ₂	0	50	40
A ₃	15	40	55

Matricea afinității atributelor este utilizată pentru realizarea fragmentării verticale. Acest proces implică două etape:

- gruparea atributelor cu afinități mari;
- divizarea relației în raport cu grupurile de attribute formate.

Algoritmul de grupare al atributelor determină grupurile de attribute similare din punct de vedere al afinității.

Matricea legăturilor de afinitate se generează în trei pași:

- inițializarea – se plasează și se fixează o coloană arbitrară a matricei AA în matricea LA ;
- iterația – se alege fiecare coloană, neutilizată încă, din AA (se va alege cea coloană care aduce cea mai mare contribuție la măsura afinității globale) și se plasează în LA pe una dintre pozițiile disponibile; acest pas se repetă până când au fost alese toate coloanele matricei AA ;
- ordonarea liniilor – se permută liniile matricei LA astfel încât să fie păstrată simetria.

Se consideră că $LA(0, j) = LA(i, 0) = LA(n + 1, j) = LA(i, n + 1) = 0$.

Contribuția globală a afinității atributelor obținută prin plasarea atributului A_k între attributele A_i și A_j va fi dată prin formula:

$$cont(A_i, A_k, A_j) = 2leg(A_i, A_k) + 2leg(A_k, A_j) - 2leg(A_i, A_j),$$

$$\text{unde } leg(A_x, A_y) = \sum_{i=1}^n af(A_i, A_x)af(A_i, A_y).$$

CA	1 (A ₁)	2 (A ₂)	3
1	15	0	
2	0	50	
3	15	40	

$index \leftarrow 3$

$$\begin{aligned} cont(A_0, A_3, A_1) &= 2*leg(A_0, A_3) + 2*leg(A_3, A_1) - 2*leg(A_0, A_1) = \\ &= 2* \sum_{i=1}^4 af(A_i, A_0)af(A_i, A_3) + 2* \sum_{i=1}^4 af(A_i, A_3)af(A_i, A_1) - \\ &\quad - 2* \sum_{i=1}^4 af(A_i, A_0)af(A_i, A_1) = 2100 \end{aligned}$$

$$\begin{aligned} cont(A_1, A_3, A_2) &= 2*leg(A_1, A_3) + 2*leg(A_3, A_2) - 2*leg(A_1, A_2) = \\ &= 2* \sum_{i=1}^4 af(A_i, A_1)af(A_i, A_3) + 2* \sum_{i=1}^4 af(A_i, A_3)af(A_i, A_2) - \\ &\quad - 2* \sum_{i=1}^4 aff(A_i, A_1)aff(A_i, A_2) = 9300 \end{aligned}$$

$$\begin{aligned} cont(A_2, A_3, A_0) &= 2*leg(A_2, A_3) + 2*leg(A_3, A_0) - 2*leg(A_2, A_0) = \\ &= 2* \sum_{i=1}^4 af(A_i, A_2)af(A_i, A_3) + 2* \sum_{i=1}^4 af(A_i, A_3)af(A_i, A_0) - \\ &\quad - 2* \sum_{i=1}^4 af(A_i, A_2)af(A_i, A_0) = 8400 \end{aligned}$$

$loc = 2 \quad (A_1, A_3, A_2)$

LA	1 (A ₁)	2 (A ₃)	3 (A ₂)
1	15	15	0
2	0	40	50
3	15	55	40

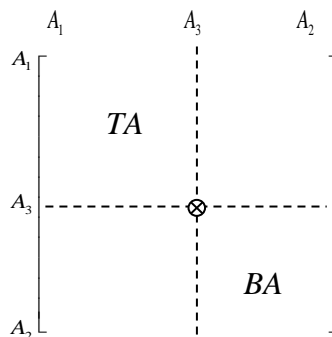
Se permută liniile, corespunzător permutării realizate la coloane (astfel matricea rămâne simetrică).

LA	1 (A ₁)	2 (A ₃)	3 (A ₂)
1 (A ₁)	15	15	0
2 (A ₃)	15	55	40
3 (A ₂)	0	40	50

Algoritmul de divizare a unei relații presupune găsirea mulțimilor de atribute care sunt accesate independent sau sunt accesate des de mai multe mulțimi distincte de aplicații.

Se consideră matricea legăturilor de afinitate LA .

$$Q = \{q_1, q_2, q_3, q_4\}$$



TQ – aplicații ce accesează doar atributele din TA ;

BQ – aplicații ce accesează doar atributele din BA ;

OQ – aplicații ce accesează atât atribute din TA cât și atribute din BA

$$CQ = \sum_{q_i \in Q} \sum_{\forall S_j} fr_acc_j(q_i), \text{ costul total al accesării oricărui atribut } A;$$

$$CTQ = \sum_{q_i \in TQ} \sum_{\forall S_j} fr_acc_j(q_i), \text{ costul accesării unui atribut din } TA \text{ de pe orice stație};$$

$$CBQ = \sum_{q_i \in BQ} \sum_{\forall S_j} fr_acc_j(q_i), \text{ costul accesării unui atribut din } BA \text{ de pe orice stație};$$

$$COQ = \sum_{q_i \in OQ} \sum_{\forall S_j} fr_acc_j(q_i), \text{ costul accesării unui atribut din } TA \text{ sau din } BA \text{ de pe orice stație}.$$

$n = 1$

$TQ_1 = \emptyset$ (aplicațiile care accesează doar atributul A_1)

$BQ_1 = \{q_1, q_2, q_3\}$ (aplicațiile care accesează doar atributele A_3 și A_2)

$OQ_1 = \{q_4\}$

$CTQ_1 = 0$

$$CBQ_1 = \sum_{q_i \in BQ_1} \sum_{\forall S_j} fr_acc_j(q_i) = fr_acc_1(q_1) + fr_acc_2(q_1) + fr_acc_1(q_2) + fr_acc_2(q_2) + fr_acc_1(q_3) + fr_acc_2(q_3) = 50$$

$$COQ_1 = \sum_{q_i \in OQ_1} \sum_{\forall S_j} fr_acc_j(q_i) = fr_acc_1(q_4) + fr_acc_2(q_4) = 15$$

$$z_1 = CTQ_1 * CBQ_1 - COQ_1^2 = -225$$

$n = 2$

$TQ_2 = \{q_4\}$ (aplicațiile care accesează doar atributele A_1 și A_3)

$BQ_2 = \{q_2\}$ (aplicațiile care accesează doar atributul A_2)

$OQ_2 = \{q_1, q_3\}$

$$CTQ_2 = \sum_{q_i \in TQ_2} \sum_{\forall S_j} fr_acc_j(q_i) = fr_acc_1(q_4) + fr_acc_2(q_4) = 15$$

$$CBQ_2 = \sum_{q_i \in BQ_2} \sum_{\forall S_j} fr_acc_j(q_i) = fr_acc_1(q_2) + fr_acc_2(q_2) = 10$$

$$COQ_2 = \sum_{q_i \in OQ_2} \sum_{\forall S_j} fr_acc_j(q_i) = fr_acc_1(q_1) + fr_acc_2(q_1) + fr_acc_1(q_3) + fr_acc_2(q_3) = 40$$

$$z_2 = CTQ_2 * CBQ_2 - COQ_2^2 = -1450$$

Rezultă că optim este $z = z_1 = -225$, deci punctul de partiționare este 1.

A_0 este cheia primară și va fi conținută de ambele fragmente. Deci cele două fragmente sunt:

$\{A_0, A_1\}$ și $\{A_0, A_2, A_3\}$.

b) Să se creeze cele două fragmente obținute la punctul a).

- `jobs_am_***` pe `server1 (bd_am)`;
- `jobs_eu_***` pe `server2 (bd_eu)`.

c) Utilizând comenzi *SQL*, să se verifice corectitudinea fragmentării verticale realizate.

În urma fragmentărilor s-au obținut:

STAȚIA 1	STAȚIA 2
COUNTRIES_AM DEPARTMENTS_AM EMPLOYEES_AM JOBS_AM LOCATIONS_AM JOB_GRADES_ALL	COUNTRIES_EU DEPARTMENTS_EU EMPLOYEES_EU JOBS_EU LOCATIONS_EU