



Baze de date

Curs 6 – Normalizarea bazelor de date

Sorina Preduț

sorina.predut@my.fmi.unibuc.ro

Universitatea din București



Proiectarea bazelor de date - recap.

- Metodele curente de proiectare a BD sunt divizate în 3 etape:
 1. crearea schemei conceptuale (curs3),
 2. crearea design-ului logic al bazei de date (curs3) și
 3. crearea design-ului fizic al bazei de date (cursul curent).



Proiectarea bazelor de date - recap.

1. **Crearea schemei conceptuale.** Este un design de nivel înalt care descrie datele și relațiile necesare pentru execuția operațiilor necesare, fiind independent de orice model de baze de date.
2. **Crearea design-ului logic al bazei de date.** La sfârșitul acestei etape vom avea un număr de tabele care vor permite stocarea și manipularea corectă a tuturor datelor necesare sistemului.
3. **Crearea design-ului fizic al bazei de date.** În această etapă designul logic este transformat într-o structură fizică eficientă.

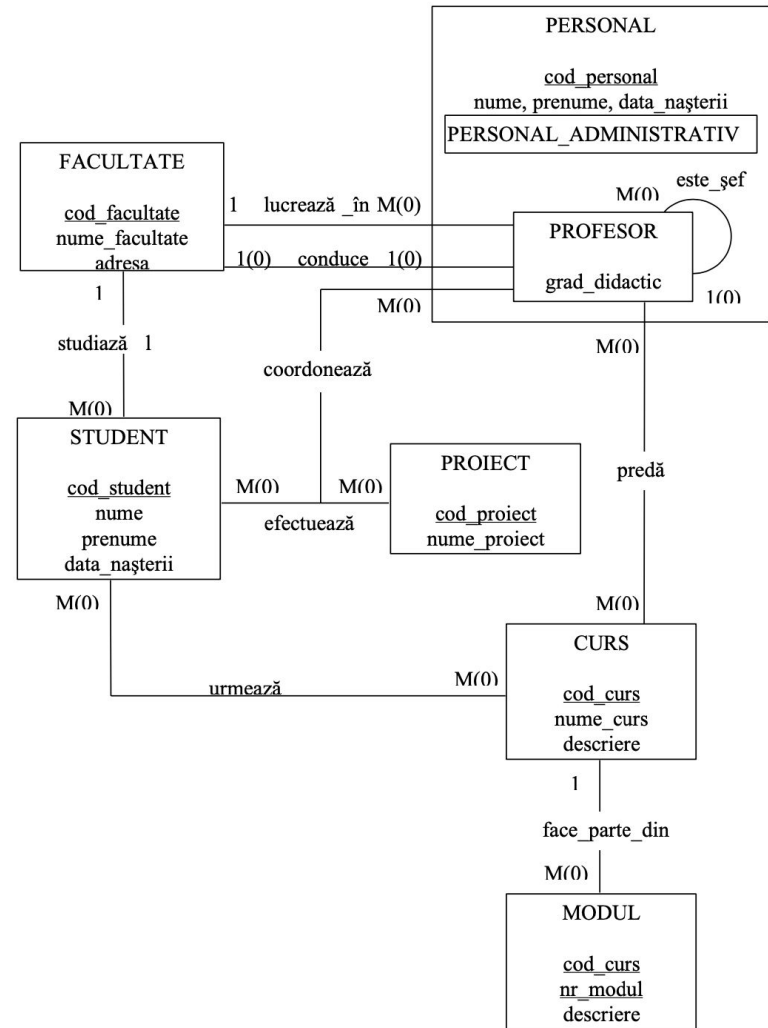


1. Schema conceptuală

- **Etapele obținerii modelului entitate-legătură:**
 - Identificarea entităților sistemului.
 - Identificarea relațiilor sistemului și stabilirea cardinalității acestora.
 - Identificarea atributelor entităților și relațiilor sistemului.
 - Stabilirea cheilor primare ale entităților.
 - Trasarea diagramei entitate-legătură.

Exemplu de ERD

Diagrama entitate-legătură a sistemului prezentat ca exemplu în cursul 3.



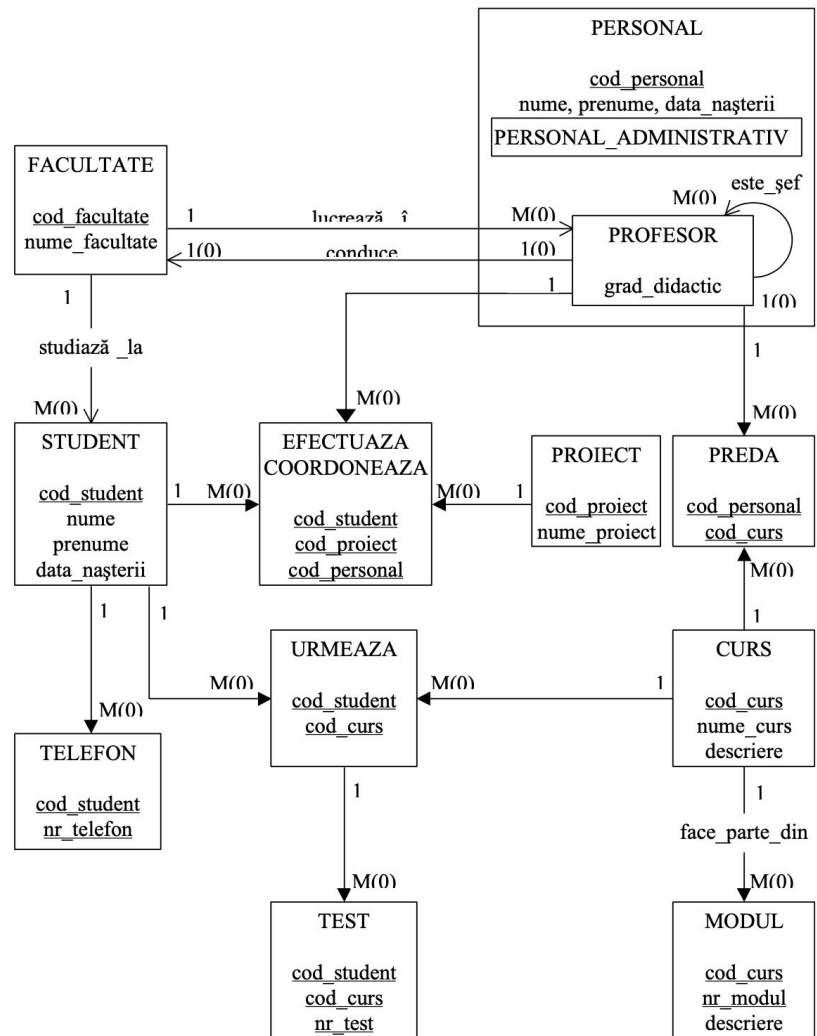


2. Design logic

- Se pornește de la schema conceptuală, mai precis de la modelul entitate-legătură, și se încearcă reprezentarea entităților și a legăturilor sub formă de tabele relaționale.

Exemplu de diagramă logică

Diagrama logică a BD pentru sistemul descris ca exemplu în cursul 3.





Exemplu de diagramă logică - cont.

➤ Tabelele asociate acestei diagrame sunt următoarele:

PERSONAL (cod_personal, nume, prenume, data_nastere, sex, stare_civila)

PERSONAL_ADMINISTRATIV (cod_personal, profesie, funcție)

PROFESOR (cod_personal, grad_didactic, titlu, sef, ore_predate, data_angajării, *cod_facultate*)

CURS (cod_curs, nume_curs, descriere, nr_ore)

PREDĂ (cod_personal, cod_curs)

MODUL (cod_curs, nr_modul, descriere)

FACULTATE (cod_facultate, nume_facultate, localitate, strada, nr, cod_postal *cod_decan*)



Exemplu de diagramă logică - cont.

STUDENT (cod_student, nume, prenume, data_nasterii, tara, localitate, strada, nr, cod_postal, studii_anterioare)

TELEFON (cod_student, nr_telefon, tip_telefon)

PROIECT (cod_proiect, nume_proiect, domeniu)

EFFECTUEAZA_COORDONEAZA (cod_student, cod_proiect, cod_personal)

URMEAZA (cod_student, cod_curs, nota_examen, nota_restantă, observatii)

TEST (cod_student, cod_curs, nr_test, nota_test, observatii)

- **Notă:** Atributele subliniate constituie CP a tabelului iar cele italice constituie chei străine.



2. Design fizic

- aka Normalizarea bazei de date.



Normalizarea bazei de date

- În general, pt. proiectare se creează mai întâi schema conceptuală care este transformată într-un design logic.
- În trecut, în locul acestei tehnici se utiliza o altă tehnică numită normalizare.
- **Normalizarea** constă din descompunerea unui tabel relațional în mai multe tabele care satisfac anumite reguli și stochează aceleași date ca și cel inițial.
- \exists 6 forme normale: prima formă normală, a doua formă normală, a treia formă normală, forma normală Boyce-Codd, a patra formă normală și a cincea formă normală.



Dependențe (nu dependențe!) funcționale

- Fie R un tabel relațional și X și Y 2 submulțimi de coloane ale lui R .
Vom spune că **X determină funcțional pe Y** sau **Y depinde funcțional de X** dacă \nexists 2 rânduri în tabelul R care să aibă aceleași valori pentru coloanele din X și să aibă valori diferite pentru cel puțin o coloană din Y .
Cu alte cuvinte, o valoare a lui X determină în mod unic o valoare a lui Y , adică \forall 2 rânduri din R care au aceeași valoare pentru X trebuie să ia aceeași valoare pentru Y .
Notăție: **$X \rightarrow Y$** . X se va numi **determinant**, iar Y **determinat**.
Spunem că dependența $X \rightarrow Y$ este **trivială** dacă toate elementele lui Y sunt și elemente ale lui X , $Y \subseteq X$.
- În curs simbolul \forall poate însemna oricare, orice sau pentru orice, după caz.

Exemplu

- VÂNZĂRI (cod_client, nume_client, nr_telefon, cod_comandă, data, cod_articol, nume_articol, cost_articol, cantitate)

VÂNZĂRI

cod_client	nume_client	nr_telefon	cod_comandă	data	cod_articol	nume_articol	cost_articol	cantitate
A1	Popescu	3215576	C1	12.05.99	P1	cămașă	100.000	2
A1	Popescu	3215576	C1	12.05.99	P3	tricou	50.000	1
A2	Ionescu	2325587	C2	13.05.99	P1	cămașă	100.000	3
A2	Ionescu	2325587	C2	13.05.99	P3	tricou	50.000	2
A2	Ionescu	2325587	C2	13.05.99	P2	pantaloni	200.000	1
A1	Popescu	3215576	C3	14.05.99	P3	tricou	50.000	3
A3	Georgescu	4555895	C4	14.05.99	P1	cămașă	100.000	1



Exemplu - cont.

- Dependențe funcționale:
 $\{\text{cod_articol}\} \rightarrow \{\text{nume_articol}, \text{cost_articol}\}$
 $\{\text{cod_comandă}\} \rightarrow \{\text{data}, \text{cod_client}, \text{nume_client}, \text{nr_telefon}\}$
 $\{\text{cod_client}\} \rightarrow \{\text{nume_client}, \text{nr_telefon}\}$
- **Notă:** dependența $\{\text{cod_comandă}\} \rightarrow \{\text{nume_client}, \text{nr_telefon}\}$ poate fi dedusă din dependențele $\{\text{cod_comandă}\} \rightarrow \{\text{cod_client}\}$ și $\{\text{cod_client}\} \rightarrow \{\text{nume_client}, \text{nr_telefon}\}$. Astfel de dependențe se numesc **dependențe tranzitive** și vor fi definite mai riguros ulterior.



Prima formă normală (1NF)

- Un **tabel relațional** este în **prima formă normală (1NF)** dacă fiecărei coloane îi corespunde o valoare indivizibilă (atomică), deci \forall valoare nu poate să fie o mulțime sau un tuplu de valori. În plus, nu pot să apară grupuri de attribute repetitive.
 - Tabelul VÂNZĂRI se află în prima formă normală.
-
- NF reprezintă abrevierea din lb. eng. pentru normal form



Exemple

- Pentru o coloană care conține date calendaristice, sub forma **zz-ll-aa** (doi digiți pentru zi, doi pentru lună și doi pentru an), se consideră că **valoarea** respectivă **nu se poate descompune** în ziua, luna și anul corespunzătoare acestei valori.
- Dacă vom considera că **adresa** este un atribut compus format din componentele țară, oraș, stradă, număr și cod, fiecare având semnificație proprie și putând fi folosite independent la interogarea bazei de date, atunci adresa va fi reprezentată în tabel prin 5 coloane în loc de una.



Exemple - cont.

- un tabel aflat în 1NF nu poate conține atribute sau grupuri de atribute repetitive.
De exemplu, în tabelul următor, dacă un student poate avea **mai multe numere de telefon**, atunci coloanele “telefon1”, “telefon2”, “telefon3” constituie un **grup de atribute repetitive**. Deci tabelul STUDENT nu este în prima formă normală, el conținând atât un atribut compus (“adresa”) cât și un grup de atribute repetitive.

Exemple - cont.

STUDENT								
cod_ student	nume	prenume	adresa	telefon1	telefon2	telefon3	materia	nota
101	Ionescu	Vasile	Romania, București, Str. Polizu 5, 7355	6245981	3215678	092659019	Drept	10
101	Ionescu	Vasile	Romania, București, Str. Polizu 5, 7355	6245981	3215678	092659019	Engleza	8
102	Pop	Costică	Romania, Cluj, Str. Unirii 7, 3551	3215469			Fizică	8

- Algoritmul 1NFA permite aducerea unei relații în 1NF prin eliminarea atributelor compuse și a celor repetitive.



Algoritmul 1NFA

1. Se înlocuiesc în tabel coloanele corespunzătoare atributelor compuse cu coloane ce conțin componentele elementare ale acestora.
2. Se plasează grupurile de attribute repetitive, fiecare în câte un nou tabel.
3. Se introduce în fiecare tabel nou creat la pasul 2 CP a tabelului din care a fost extras atributul respectiv. Prin urmare, în tabelul nou creat attributele introduse vor fi chei străine ce fac referință la tabelul din care au fost extrase.
4. Se stabilește CP a fiecărui nou tabel creat la pasul 2. Aceasta va fi creată din chei străine introdusă la pasul 3 plus una sau mai multe coloane adiționale.

Exemple - cont.

STUDENT

cod_ student	nume	prenume	adresa	telefon1	telefon2	telefon3	materia	nota
101	Ionescu	Vasile	Romania, București, Str. Polizu 5, 7355	6245981	3215678	092659019	Drept	10
101	Ionescu	Vasile	Romania, București, Str. Polizu 5, 7355	6245981			Engleza	8
102	Pop	Costică	Romania, Cluj, Str. Unirii 7, 3551	3215469			Fizică	8

TELEFON

cod_ student	telefon
101	6245981
101	3215678
101	092659019
102	3215469



Observații

- 1NF este o cerință minimală a tuturor sistemelor relaționale. Sistemele de baze de date care nu respectă nici măcar această formă nu pot fi numite relaționale.
- Tabelele aflate în 1NF permit o referire simplă a datelor prin indicarea numelui tabelului, a coloanei și a cheii rândului din care face parte informația respectivă.
- Operatorii pentru aceste tabele sunt mai simpli și permit definirea unor tehnici de proiectare și utilizare a bazelor de date.



A doua formă normală (2NF)

- Fie R un tabel relațional și fie X și Y două submulțimi de coloane ale lui R .
O **dependență funcțională** $X \rightarrow Y$ se numește **totală** dacă
 \forall subset de coloane Z al lui X , $Z \subseteq X$, dacă $Z \rightarrow Y$ atunci $Z = X$.
Cu alte cuvinte, \nexists nici un subset Z al lui X , $Z \neq X$, pentru care $Z \rightarrow Y$.
O dependență funcțională care nu este totală se numește **parțială**.



A doua formă normală (2NF) - cont.

- Un tabel relațional R este în a doua formă normală (2NF) dacă și numai dacă:
 - R este în 1NF.
 - \forall coloană care depinde parțial de o cheie a lui R este inclusă în acea cheie.
- Cu alte cuvinte, a doua formă normală **nu permite dependențe funcționale parțiale față de cheile tabelului**, cu excepția dependențelor triviale, de incluziune.

Deci tabelul VÂNZĂRI nu este în 2NF, coloanele “nume_articol” și “cost_articol” depinzând parțial de CP a tabelului.



Regulă de descompunere

- Pentru a obține tabele relaționale în 2NF, tabelul inițial se descompune fără pierdere de informație astfel:

Fie $R(K1, K2, X, Y)$ un tabel relațional unde $K1$, $K2$, X și Y sunt submulțimi de coloane ale lui R a.î. $K1 \cup K2$ este o cheie a lui R , iar $K1 \rightarrow X$ este o dependență funcțională totală. Dacă $X \subset K1$ atunci tabelul este deja în 2NF, altfel tabelul R poate fi descompus prin proiecție în $R1(K1, K2, Y)$ - având cheia $K1 \cup K2$ - și $R2(K1, X)$ - având cheia $K1$.



Regulă de descompunere - cont.

- Această descompunere conservă nu numai datele, ci și dependențele funcționale, atât determinantul cât și determinatul dependenței eliminate regăsindu-se în tabelul nou creat.
- Folosind această regulă, **algoritmul 2NFA permite** aducerea în 2NF a unui tabel relațional aflat în 1NF prin **eliminarea dependențelor funcționale parțiale**.



Algoritmul 2NFA

1. \forall coloană X care depinde funcțional parțial de o cheie K , $K \rightarrow X$, și care nu este inclusă în K , se determină $K1 \subset K$ un subset al lui K , astfel încât **dependența** $K1 \rightarrow X$ este **totală** și se creează un nou tabel **$R1(K1, X)$** , adică un tabel format din determinantul ($K1$) și determinatul (X) acestei relații.
2. Dacă în tabelul $R \exists$ mai multe dependențe totale ca mai sus cu același determinant, atunci pentru acestea se creează un singur tabel format din determinant - luat o singură dată - și din determinații dependențelor considerate.



Algoritmul 2NFA - cont.

3. Se elimină din tabelul inițial R toate coloanele, X , care formează determinatul dependenței considerate.
4. Se determină CP a fiecărui tabel nou creat, $R1$. Aceasta va fi $K1$, determinantul dependenței considerate.
5. Dacă noile tabele create conțin alte dependențe parțiale, atunci se merge la pasul 1, altfel algoritmul se termină.



Exemplu

- Pentru tabelul VÂNZĂRI CP este {cod_comandă, cod_articol}.
Atributul “nume_articol” și “cost_articol” depind funcțional parțial de această CP și depind funcțional total de atributul “cod_articol” conținut în CP.
Același lucru se întâmplă și cu attributele “data”, “cod_client”, “nume_client” și “nr_telefon” care depind funcțional total numai de atributul “cod_comandă” conținut în CP.

Prin urmare vom avea următoarele **dependențe totale**:

{cod_articol} → {nume_articol, cost_articol}

{cod_comandă} → {data, cod_client, nume_client, nr_telefon}



Exemplu - cont.

- Tabelul VÂNZĂRI (cod_client, nume_client, nr_telefon, cod_comandă, data, cod_articol, nume_articol, cost_articol, cantitate) va fi descompus în

VÂNZĂRI_1 (cod_client, nume_client, nr_telefon, cod_comandă, data, cod_articol, cantitate)
și

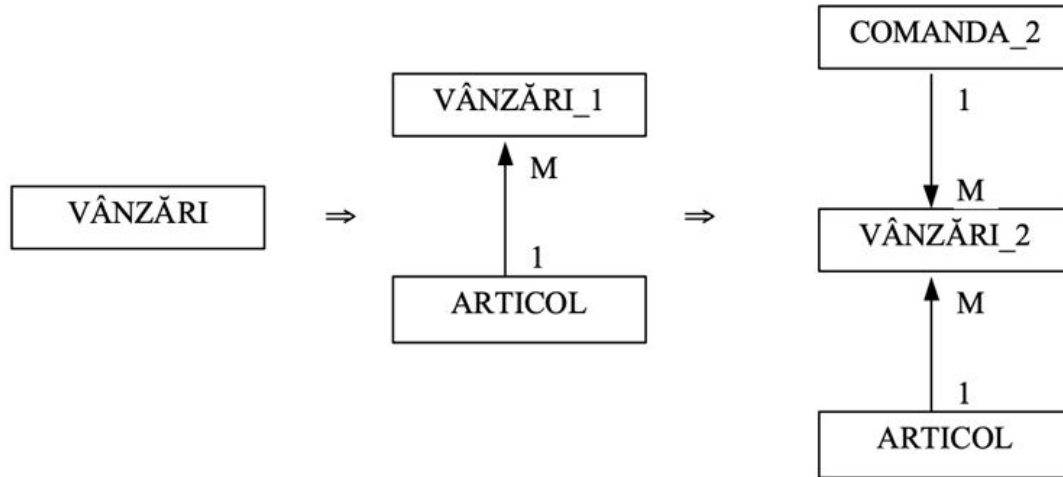
ARTICOL (cod_articol, nume_articol, cost_articol).

Tabelul VÂNZĂRI_1 poate fi descompus în

VÂNZĂRI_2 (cod_comandă, cod_articol, cantitate) și

COMANDA_2 (cod_comandă, data, cod_client, nume_client, nr_telefon).

Exemplu - cont.



- Un tabel care are CP formată dintr-un singur atribut este automat în 2NF. Prin urmare, algoritmul 2NFA nu se poate aplica decât în cazul în care CP a unui tabel este o cheie compusă.



A treia formă normală (3NF)

- Deși tabelul COMANDA_2 este în 2NF, se observă că încă mai există redundanță în date - tuplul (A1, Popescu, 3215576) apare de două ori.
Intuitiv, aceasta se explică prin faptul că attributele “nume_client” și “nr_telefon” depind indirect de CP a tabelului, dependența făcându-se prin intermediul atributului “cod_client”. Aceste dependențe indirecte vor fi îndepărtate în 3NF.



A treia formă normală (3NF) - cont.

- Un tabel relațional **R** este în a treia formă normală (3NF) dacă și numai dacă:
 - R este în 2NF.
 - \forall coloană **A** neconținută în nici o cheie a lui R, dacă \exists un set de coloane **X** a.î. $X \rightarrow A$, atunci fie **X** conține o cheie a lui R, fie **A** este inclusă în **X**.
- A doua condiție din definiție **interzice dependențele funcționale totale față de alte coloane în afara celor care constituie chei ale tabelului.**
Prin urmare, un tabel este în 3NF dacă \forall coloană care nu este conținută într-o cheie depinde de cheie, de întreaga cheie și numai de cheie.



A treia formă normală (3NF) - cont.

- În mod evident tabelul COMANDA_2 nu este în 3NF existând dependențele
 $\{\text{cod_client}\} \rightarrow \{\text{nume_client}\}$ și
 $\{\text{cod_client}\} \rightarrow \{\text{nr_telefon}\}.$

Pe de altă parte, tabelele VÂNZĂRI_2 și ARTICOL sunt în 3NF.



A treia formă normală (3NF) - cont.

- Cea de-a doua condiție din definiția de mai înainte se formulează folosind noțiunea de **dependență tranzitivă**.
Fie R un tabel relațional, X o submulțime de coloane a lui R și A o coloană a lui R .
Spunem că **A este dependentă tranzitiv de X** dacă \exists o submulțime de coloane Y care nu include A și nu determină funcțional pe X astfel încât $X \rightarrow Y$ și $Y \rightarrow A$.
Dacă în această definiție se dorește să se evidențieze și Y atunci se spune că A depinde funcțional de X prin intermediul lui Y și se scrie $X \rightarrow Y \rightarrow A$.



A treia formă normală (3NF) - cont.

- De exemplu, în tabelul COMANDA_2 coloanele “nume_client” și “nr_telefon” depind tranzitiv de CP “cod_comandă” prin intermediul coloanei “cod_client”.
Folosind această definiție, condiția ca un tabel să fie în 3NF se poate reformula astfel:
- Un **tabel relațional R** este în a treia formă normală (3NF) dacă și numai dacă:
 - R este în 2NF.
 - \forall coloană neconținută în nici o cheie a lui R nu este dependentă tranzitiv de nici o cheie a lui R.



Regulă de descompunere

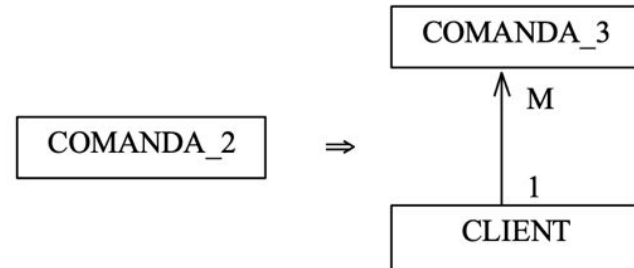
- Pentru a obține tabele relaționale în 3NF, tabelul inițial se descompune fără pierdere de informație după următoarele reguli.

Fie $R(K, X, Y, Z)$ un tabel relațional unde K este o cheie a lui R , iar X, Y și Z sunt submulțimi de coloane ale lui R .

- Dacă există **dependența tranzitivă** $K \rightarrow X \rightarrow Y$, atunci R se poate descompune în $R1(K, X, Z)$ - având cheia K - și $R2(X, Y)$ - având cheia X .
- Dependența tranzitivă poate fi mai complexă. Fie $K1 \subset K$ o parte a cheii K astfel încât există dependența tranzitivă $K \rightarrow K1 \cup X \rightarrow Y$. În acest caz, R poate fi descompus în $R1(K, X, Z)$ - având cheia K - și $R2(K1, X, Y)$ - având cheia $K1 \cup X$.

Regulă de descompunere - cont.

- Descompunerile corespunzătoare regulilor de mai înainte conservă nu numai datele, ci și dependențele funcționale, determinantul și determinatul dependențelor eliminate regăsindu-se în tabelele nou create.
- Un exemplu de aplicare a primei reguli este descompunerea tabelului COMANDA_2 în COMANDA_3 (cod_comandă, data, cod_client) și CLIENT (cod_client, nume_client, nr_telefon)



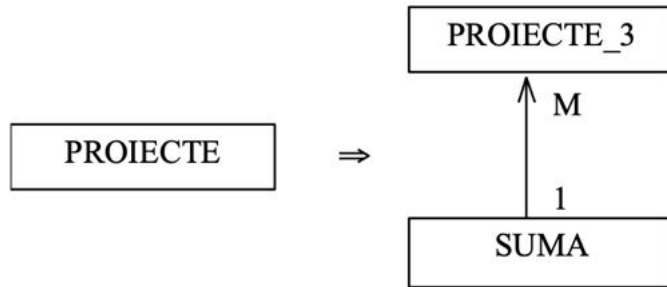


Exemplu

- Pentru a exemplifica cea de-a doua regulă, considerăm PROIECTE (cod_angajat, cod_proiect, rol_în_proiect, suma_obținută), care stochează date privind repartizarea pe proiecte a angajaților unei firme. Presupunem că suma obținută de un angajat depinde de proiectul respectiv și de rolul angajatului în acel proiect, deci avem dependența $\{\text{cod_proiect}, \text{rol_în_proiect}\} \rightarrow \{\text{suma_obținută}\}$.

Exemplu - cont.

- Aplicând regula a doua, tabelul PROIECTE se descompune în:
PROIECTE_3 (cod_angajat, cod_proiect, rol_în_proiect) și
SUMA (cod_proiect, rol_în_proiect, suma_obținută).



- Folosind această regulă, algoritmul 3NFA permite aducerea în 3NF a unui tabel relațional aflat în 2NF prin **eliminarea dependențelor funcționale tranzitive**.

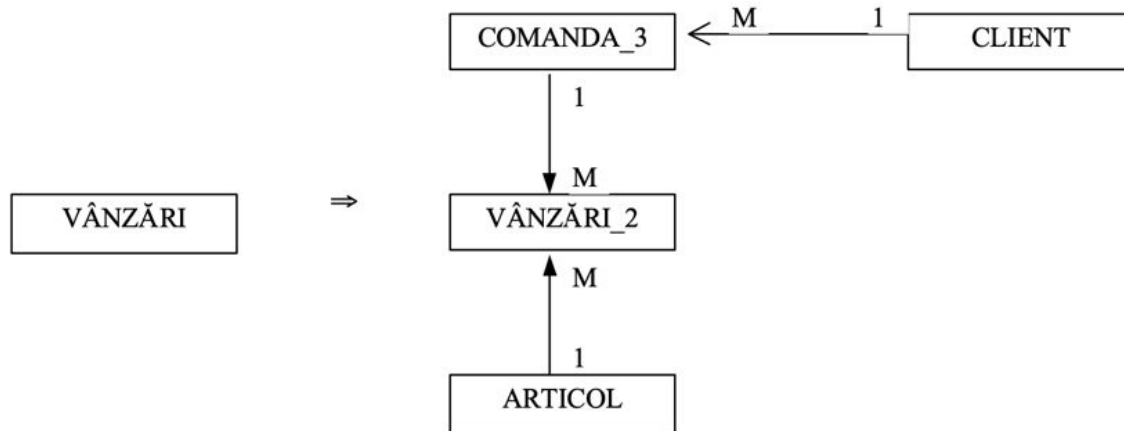


Algoritmul 3NFA

1. Pentru fiecare dependență funcțională tranzitivă $K \rightarrow X \rightarrow Y$ se transferă coloanele din X și Y într-o nouă relație.
2. Se determină CP a fiecărei noi relații create la pasul 1, aceasta fiind formată din coloanele din X .
3. Se elimină din relația principală coloanele din Y .
4. Dacă tabelele rezultate conțin alte dependențe tranzitive, atunci se merge la pasul 1, altfel algoritmul se termină.

Exemplu

- Aplicând algoritmi 2NFA și 3NFA, tabelul VÂNZĂRI a fost descompus în tabelele VÂNZĂRI_2, ARTICOL, COMANDA_3 și CLIENT, care sunt toate în 3NF. În aceste tabele nu mai există nici redundanță în date și nici anomalii de actualizare.





Forma normală Boyce-Codd (BCNF)

- Nu toate tabelele aflate în 3NF sunt lipsite de redundanță în date și anomalii de actualizare. Pentru a ilustra această situație să considerăm următorul exemplu.

O companie de transporturi efectuează curse, în care poate folosi unul sau mai mulți șoferi - de exemplu prima jumătate a cursei conduce un șofer, cea de-a doua alt șofer - și mai multe dintre autobuzele aflate în dotare, cu condiția ca într-o cursă un șofer să conducă un singur autobuz.

Pe de altă parte însă, un autobuz este repartizat unui șofer și deci nu poate fi condus decât de acesta.



Forma normală Boyce-Codd (BCNF) - cont.

- Această situație poate fi modelată printr-un tabel TRANSPORTURI (cod_cursă, cod_șofer, cod_autobuz, loc_plecare, loc_sosire).
- În acest tabel avem dependențele:
 $\{cod_cursă, cod_șofer\} \rightarrow \{cod_autobuz, loc_plecare, loc_sosire\}$
 $\{cod_autobuz\} \rightarrow \{cod_șofer\}$,
iar cheile tabelului sunt $\{cod_cursă, cod_șofer\}$ și $\{cod_cursă, cod_autobuz\}$.
- Tabelul TRANSPORTURI este în 3NF și totuși în acest tabel există redundanță în date datorită dependenței $\{cod_autobuz\} \rightarrow \{cod_șofer\}$.



Forma normală Boyce-Codd (BCNF) - cont.

- Forma Boyce-Codd elimină acest tip de redundanțe.
Intuitiv, un tabel R este în BCNF dacă fiecare determinant al unei dependențe funcționale este cheie candidată a lui R.
Tabelul TRANSPORTURI de mai înainte nu este în BCNF, “cod_autobuz” nefiind o cheie a lui R.
- Un **tabel relațional este în forma normală Boyce-Codd (BCNF)** $\Leftrightarrow \forall$ dependență funcțională totală $X \rightarrow A$, unde X este un subset de coloane iar A o coloană neconținută în X, X este o cheie a lui R.



Forma normală Boyce-Codd (BCNF) - cont.

- \forall tabel care are cel mult 2 coloane este în BCNF.
- \forall tabel relațional se poate descompune fără pierdere de informație în tabele aflate în BCNF, dar nu același lucru se poate spune despre descompunerea cu păstrarea dependențelor funcționale, după cum vom vedea în continuare.



Algoritmul BCNFA

1. \forall dependență **non-cheie** $X \rightarrow Y$, unde X și Y sunt subseturi de coloane ale lui R , se creează 2 tabele. Una dintre ele va fi formată din coloanele $\{X, Y\}$, iar cealaltă va fi formată din toate coloanele inițiale, mai puțin coloanele Y .
2. Dacă tabelele rezultate conțin alte dependențe non-cheie, atunci se merge la pasul 1, altfel algoritmul se termină.

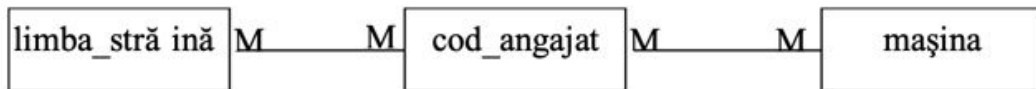


Exemplu

- Aplicând acest algoritm, tabelul TRANSPORTURI se va descompune în tabelele TRANSPORTURI_BC (cod_cursă, cod_autobuz, loc_plecare, loc_sosire) și AUTOBUZ (cod_autobuz, cod_șofer).
- Descompunerea s-a făcut fără pierdere de informație, dar a fost pierdută dependența funcțională $\{\text{cod_cursă}, \text{cod_șofer}\} \rightarrow \{\text{cod_autobuz}\}$.

A patra formă normală (4NF)

- Dacă BCNF elimină redundanțele datorate dependențelor funcționale, **4NF determină redundanțele datorate dependențelor multivaloare**.
- Pentru a ilustra acest tip de redundanțe, să considerăm tabelul ANGAJAȚI (cod_angajat, limba_străină, mașina) aflat în BCNF.
- Un angajat poate cunoaște mai multe limbi străine și poate avea mai multe mașini, dar \nexists nici o legătură între limba străină și mașină. Cu alte cuvinte, redundanța datelor din tabelul ANGAJAȚI este cauzată de existența a două relații N:M independente. 4NF va înlătura aceste relații N:M independente.





A patra formă normală (4NF) - cont.

- Fie R un tabel relațional, X și Y două submulțimi de coloane ale lui R și $Z = R - X - Y$ mulțimea coloanelor din R care nu sunt nici în X nici în Y .
Spunem că \exists o **dependență multivaloare** Y de X sau că X **determină multivaloare** pe Y , și notăm $X \twoheadrightarrow Y$, dacă, \forall valoare a coloanelor lui X , sunt asociate valori pentru coloanele din Y care nu sunt corelate în nici un fel cu valorile coloanelor lui Z .

A patra formă normală (4NF) - cont.

- Cu alte cuvinte $X \twoheadrightarrow Y \Leftrightarrow \forall u \text{ și } v, 2 \text{ rânduri ale lui } R \text{ cu } u(X) = v(X), \exists s \text{ și } t, 2 \text{ rânduri ale lui } R \text{ a.î. } s(X) = u(X), s(Y) = u(Y), s(Z) = v(Z) \text{ și a.î. } t(X) = u(X), t(Y) = v(Y), t(Z) = u(Z), \text{ unde prin } u(X) \text{ am notat valoarea coloanelor } X \text{ corespunzătoare rândului } u, \text{ etc.}$
În mod evident, dacă $X \twoheadrightarrow Y$ atunci și $X \twoheadrightarrow Z$.
- Dependența multivaloare se mai numește și **multidependență**.

	X	Y	Z
u	x	y1	z1
v	x	y2	z2
s	x	y1	z2
t	x	y2	z1



A patra formă normală (4NF) - cont.

- În tabelul ANGAJAȚI avem dependențele multivaloare
 $\text{cod_angajat} \twoheadrightarrow \text{limba_străină}$ și
 $\text{cod_angajat} \twoheadrightarrow \text{mașina}$.
Pe de altă parte însă, nici “limba_străină” și nici “mașina” nu depinde funcțional de
“cod_angajat”.



A patra formă normală (4NF) - cont.

- Un tabel relațional **R** este în a patra formă normală (4NF) dacă și numai dacă:
 - R este în BCNF.
 - \forall dependență multivaloare $X \twoheadrightarrow Y$ este de fapt o dependență funcțională $X \rightarrow Y$.
- Condiția a doua arată că dacă \exists o dependență multivaloare $X \twoheadrightarrow Y$, atunci \forall coloană, A a lui R va depinde funcțional de coloanele din X, $X \rightarrow A$ - aceasta deoarece existența unei dependențe multivaloare $X \twoheadrightarrow Y$ implică și existența unei dependențe multivaloare $X \twoheadrightarrow (R - X - Y)$.
Ținând cont de faptul că R este în BCNF, înseamnă că \exists o cheie candidată a lui R inclusă în X.
- Un tabel relațional **R** este în a patra formă normală (4NF) \Leftrightarrow
 \forall dependență multivaloare $X \twoheadrightarrow Y \exists$ o cheie a lui R inclusă în X.



Regulă de descompunere

- Un tabel poate fi adus în 4NF prin descompunere fără pierdere de informație astfel:
Fie $R(X, Y, Z)$ un tabel relațional în care \exists o **dependență multivaloare** $X \twoheadrightarrow Y$ astfel încât X nu conține nici o cheie a lui R .
Atunci tabelul R poate fi descompus prin proiecție în două tabele $R1(X, Y)$ și $R2(X, Z)$.
- Aplicând această regulă tabelul ANGAJAȚI se descompune în tabelele ANGAJAȚI_4A (cod_angajat, limba_străină) și ANGAJAȚI_4B (cod_angajat, mașina).



Algoritmul 4NFA

1. Se identifică **dependențele multivaloare** $X \twoheadrightarrow Y$ pentru care X și Y nu conțin toate coloanele lui R și X nu conține nici o cheie a lui R .
Se poate presupune că X și Y sunt disjuncte datorită faptului că din $X \twoheadrightarrow Y$ rezultă $X \twoheadrightarrow (Y - X)$.
2. Se înlocuiește tabelul inițial R cu două tabele, primul format din coloanele $\{X, Y\}$, iar celălalt din toate coloanele inițiale, mai puțin coloanele Y .
3. Dacă tabelele rezultate conțin alte dependențe multivaloare, atunci se face transfer la pasul 1, altfel algoritmul se termină.



A cincea formă normală (5NF)

- A cincea formă se întâlnește destul de rar în practică, ea având mai mult valoare teoretică. Dacă în a patra formă normală sunt eliminate relațiile N:M independente, a cincea formă normală are ca scop **eliminarea relațiilor N:M dependente**.

Redundanțele datorate unor astfel de relații pot fi înlăturate prin descompunerea tabelului în 3 sau mai multe tabele.



Exemplu

- Considerăm tabelul LUCRĂTOR_ATELIER_PRODUS (cod_lucrător, cod_atelier, cod_produs), aflat în 4NF. Aparent acest tabel este ilustrarea unei relații de tip 3 care există între lucrător, atelier și produs. Dacă însă presupunem că între lucrător și atelier, lucrător și produs, atelier și produs există relații N:M, atunci în tabel pot exista redundanțe în date care pot fi înlăturate prin descompunerea tabelului LUCRĂTORI în 3 tabele:
LUCRĂTOR_ATELIER (cod_lucrător, cod_atelier),
LUCRĂTOR_PRODUS (cod_lucrător, cod_produs) și
ATELIER_PRODUS (cod_atelier, cod_produs).

LUCRĂTOR_ATELIER_PRODUS

Cod_lucrător	cod_atelier	cod_produs
L1	A1	P1
L2	A1	P2
L2	A1	P1
L2	A2	P1

LUCRĂTOR_ATELIER

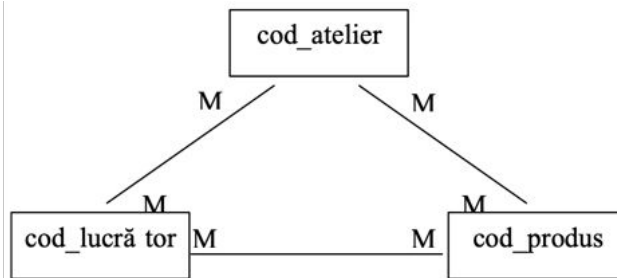
Cod_lucrător	cod_atelier
L1	A1
L2	A1
L2	A2

LUCRĂTOR_PRODUS

Cod_lucrător	cod_produs
L1	P1
L2	P2
L2	P1

ATELIER_PRODUS

Cod_atelier	cod_atelier
A1	P1
A1	P2
A2	P1

**R_12**

Cod_lucrător	cod_atelier	cod_produs
L1	A1	P1
L2	A1	P2
L2	A1	P1
L2	A2	P2
L2	A2	P1

R_13

Cod_lucrător	cod_atelier	cod_produs
L1	A1	P1
L1	A1	P2
L2	A1	P2
L2	A1	P1
L2	A2	P1

R_23

Cod_lucrător	cod_atelier	cod_produs
L1	A1	P1
L1	A2	P1
L2	A1	P2
L2	A1	P1
L2	A2	P1



Exemplu - cont.

- Tabelul LUCRĂTOR_ATELIER elimină redundanța (L2, A1), tabelul LUCRĂTOR_PRODUS elimină redundanța (L2, P1), în timp ce tabelul ATELIER_PRODUS elimină redundanța (A1, P1).
- Tabelul inițial LUCRĂTOR_ATELIER_PRODUS nu poate fi reconstituit din compunerea a doar două din tabelele componente, unde tabelul **R_12** rezultă din compunerea LUCRĂTOR_ATELIER și LUCRĂTOR_PRODUS, **R_13** rezultă din compunerea LUCRĂTOR_ATELIER și ATELIER_PRODUS, iar **R_23** rezultă din compunerea LUCRĂTOR_PRODUS și ATELIER_PRODUS.
- Tabelul inițial LUCRĂTOR_ATELIER_PRODUS poate fi obținut prin compunerea tuturor celor trei tabele componente, de exemplu el rezultă prin compunerea lui R_12 cu ATELIER_PRODUS.



Join-dependență

- Dependența funcțională și dependența multivaloare, și implicit regulile de descompunere pentru formele normale 1NF-4NF, permit **descompunerea prin proiecție a unui tabel relațional în 2 tabele relaționale**.

Totuși, regulile de descompunere asociate acestor forme normale nu dau toate descompunerile posibile prin proiecție a unui tabel relațional.

\exists tabele care nu pot fi descompuse în 2 tabele, dar pot fi descompuse în 3 sau mai multe tabele fără pierdere de informație.

Astfel de descompuneri, în 3 sau mai multe tabele, sunt tratate de 5NF.

Pentru a arăta că un tabel se poate descompune fără pierderi de informație a fost introdus conceptul de **join-dependență** sau **dependență la compunere**.



Join-dependență - cont.

- Fie R un tabel relațional și R_1, R_2, \dots, R_n o mulțime de tabele relaționale care nu sunt disjuncte - au coloane comune - a.î. U coloanelor din R_1, R_2, \dots, R_n este mulțimea coloanelor din R .
Se spune că R satisface **join-dependența** $\{R_1, R_2, \dots, R_n\}$ dacă R se descompune prin proiecție pe R_1, R_2, \dots, R_n fără pierdere de informație, adică tabelul inițial poate fi reconstruit prin compunerea naturală pe attribute comune ale tabelelor rezultate.
- În exemplul anterior, tabelul LUCRĂTOR_ATELIER_PRODUS satisface dependența de uniune $\{LUCRĂTOR_ATELIER, LUCRĂTOR_PRODUS, ATELIER_PRODUS\}$.



Join-dependență - cont.

- Join-dependența este o generalizare a dependenței multivaloare.
Mai precis, **dependența multivaloare corespunde join-dependenței cu două elemente.**
Într-adevăr, dacă în relația $R(X, Y, Z)$ avem multidependență $X \twoheadrightarrow Y$, atunci avem și join-dependență $\{(X \cup Y), (X \cup Z)\}$.
Invers, dacă avem join-dependență $\{R1, R2\}$, atunci avem și dependență multivaloare $(R1 \cap R2) \twoheadrightarrow (R1 - (R1 \cap R2))$.



A cincea formă normală (5NF) - cont.

- Un tabel relațional R este în a cincea formă normală (5NF) dacă și numai dacă \forall join-dependență $\{R_1, R_2, \dots, R_n\}$ este consecința cheilor candidate ale lui R , adică fiecare dintre R_1, R_2, \dots, R_n include o cheie candidată a lui R .



A cincea formă normală (5NF) - cont.

- \forall tabel relațional care este în 5NF este în 4NF deoarece \forall dependență multivaloare poate fi privită ca un caz particular de dependență la uniune.
Trecerea de la 4NF la 5NF constă în identificarea join-dependențelor cu mai mult de 3 elemente și descompunerea tabelului inițial prin proiecție pe aceste componente.
5NF are o importanță practică redusă, cazurile când apare în practică fiind extrem de rare.
- \forall tabel relațional poate fi descompus fără pierderi de informație într-o mulțime de tabele relaționale care sunt în 5NF.
Un tabel în 5NF nu conține anomalii ce pot fi eliminate luând proiecțiile pe diferite submulțimi de coloane ale acestuia.



Concluzii

- Normalizarea este procesul de transformare a structurilor de date și are ca scop **eliminarea redundanțelor și promovarea integrității datelor**.
Normalizarea este un pilon de bază al BD relaționale.
În general, un set de structuri de date nu sunt considerate relaționale decât dacă sunt complet normalizate.
- **Normalizarea** datelor este împărțită în **6** etape, numite **forme normale (NF)**.
Fiecare NF are asociat atât un criteriu cât și un proces.
Criteriul unei anumite NF este mai restrictiv decât al uneia inferioare, a. î.
∀ tabel relațional care este într-o anumită NF este și în NF inferioară.
Procesul asociat unei NF se referă la trecerea unor structuri de date din NF inferioară în forma normală curentă.



Concluzii - cont.

- $1NF \rightarrow 2NF$ elimină dependențele funcționale parțiale față de chei.
- $2NF \rightarrow 3NF$ elimină dependențele funcționale tranzitive față de chei.
- $3NF \rightarrow BCNF$ elimină dependențele funcționale pentru care determinantul nu este cheie.
- $BCNF \rightarrow 4NF$ elimină toate dependențele multivaloare care nu sunt și dependențe funcționale.
- $4NF \rightarrow 5NF$ elimină toate join-depenențele care nu sunt implicate de o cheie.



Concluzii - cont.

- În practică, cel mai adesea apar primele forme normale (1NF - 3NF), a. î. un tabel relațional aflat în 3NF este de obicei complet normalizat.

În special 5NF apare extrem de rar, având valoare practică foarte scăzută.



Denormalizare

- În principiu, denormalizarea este exact procesul invers normalizării.
- Este procesul de creștere a redundanței datelor, care are ca scop creșterea performanței sau simplificarea programelor de manipulare a datelor.

Totuși, trebuie reținut că o denormalizare corectă **nu înseamnă nicidecum a nu normaliza structurile de date inițiale.**

Din contră, denormalizarea **are loc după ce structurile bazei de date au fost complet normalizate**, și se face prin **selectarea strategică a acelor structuri unde denormalizarea aduce avantaje semnificative.**

Pe de altă parte, orice denormalizare trebuie însoțită de introducerea de măsuri suplimentare, care să asigure integritatea datelor chiar și în cazul unei BD care nu este complet nenormalizată.



Creșterea performanței

- Principalul obiectiv al denormalizării este creșterea performanței programelor de manipulare a datelor.
Una dintre cele mai întâlnite situații de acest gen este denormalizarea folosită pentru **operații sau calcule efectuate frecvent**.
- Să presupunem că pentru înregistrarea tranzacțiilor unui magazin care vinde articole la comandă se folosesc tabelele
VÂNZĂRI_2 (cod_comandă, cod_articol, cantitate),
ARTICOL (cod_articol, nume_articol, cost_articol),
COMANDA_3 (cod_comandă, data, cod_client),
CLIENT (cod_client, nume_client, nr_telefon)



Creșterea performanței - cont.

- Aceste tabele sunt în 5NF și conțin toate datele necesare pentru obținerea oricărui raport privind vânzările din magazin - de exemplu, în funcție de articol, dată, cantitate, client, etc. Totuși, să presupunem că majoritatea rapoartelor cerute de conducerea magazinului utilizează **cantitatea totală vândută într-o lună pentru fiecare articol**. Tabelele de mai înainte conțin toate datele necesare pentru obținerea acestei informații, de exemplu ea se poate obține printr-o simplă interogare SQL. Pe de altă parte însă, aceasta ar însemna ca **totalul pentru fiecare articol să fie recalculat de fiecare dată când este cerut**.



Creșterea performanței - cont.

- Deoarece este puțin probabil ca acest total să se schimbe după încheierea lunii respective, repetarea acestora ar putea fi înlocuită cu un tabel suplimentar ARTICOL_LUNA (cod_articol, luna, cantitate_totală)
Datele conținute în acest nou tabel sunt **redundante**, ele putând fi obținute în orice moment din datele conținute în tabelele VÂNZĂRI_2 și COMANDA_3, dar prezintă avantajul că **folosirea lor este mai facilă și mai rapidă decât repetarea calculelor.**
- Pe de altă parte însă, folosirea acestui tabel suplimentar are și **dezavantaje**, putând duce la **pierderea integrității datelor.**
Modificarea datelor din tabelele VÂNZĂRI_2 și COMANDA_3 nu se reflectă automat în datele din noul tabel.



Creșterea performanței - cont.

- De exemplu, dacă după efectuarea calculelor și inserarea totalului în noul tabel se efectuează o vânzare suplimentară, această vânzare nu este reflectată în datele din noul tabel.

De aceea, în situațiile în care BD nu este total normalizată, trebuie luate măsuri suplimentare pentru păstrarea integrității datelor.



Creșterea performanței - cont.

- De exemplu, pentru ca orice modificare a tabelelor inițiale să fie reflectată automat în noul tabel, în Oracle se pot crea așa numitele **triggere** (declanșatoare) ale bazei de date bazate pe tabelele VÂNZĂRI_2 și COMANDA_3 care se declanșează după fiecare actualizare (INSERT, UPDATE, DELETE) a acestor tabele și care actualizează datele din tabelul ARTICOL_LUNA. Aceste triggere pe bază, la rândul lor, vor încetini operațiunile de actualizare din tabelele pe care sunt bazate.



Creșterea performanței - cont.

- Datorită problemelor legate de redundanță și integritate, orice denormalizare trebuie să fie gândită cu atenție pentru a fi siguri că avantajele acesteia în privința performanțelor și a simplității programelor depășesc efortul necesar pentru impunerea integrității.
- În exemplul anterior, se pot crea tabele suplimentare pentru totaluri pe lună sau pe zi. Dacă aceste totaluri sunt utilizate rar, este mai convenabil ca ele să fie calculate direct din tabele normalizate, fără a mai introduce alte date redundante.
- În concluzie, fiecare caz trebuie studiat cu atenție pentru a vedea dacă problemele inerente asociate denormalizării (redundanța, menținerea integrității) sunt compensate de creșterea în performanță pentru anumite situații.



Simplificarea codului

- Un alt motiv invocat pentru folosirea denormalizării este acela al simplificării codului pentru manipularea datelor.
Cu alte cuvinte, pentru un dezvoltator, un singur tabel poate părea uneori mai ușor de utilizat decât două sau mai multe.
De exemplu, să considerăm tabelul
`STOCURI(cod_depozit, cod_material, nume_material, cantitate)`
utilizat de o firmă pentru a înregistra cantitățile de materiale existente în fiecare din depozitele sale.
Evident, acest tabel nu este în 2NF și poate fi normalizat prin descompunerea în 2 tabele:
`STOCURI_2(cod_depozit, cod_material, cantitate)`
`MATERIAL(cod_material, nume_material)`



Simplificarea codului - cont.

- Totuși, dacă se dorește aflarea tuturor depozitelor în care există ciment, de exemplu, pentru dezvoltator este mai convenabil să folosească varianta nenormalizată, în acest caz interogarea SQL corespunzătoare fiind

```
SELECT cod_depozit, cantitate  
FROM stocuri  
WHERE nume_material = 'CIMENT';
```

- Evident, interogarea de mai sus este mai simplă decât varianta în care se folosesc cele două tabele normalizate

```
SELECT cod_depozit, cantitate  
FROM stocuri_2, material  
WHERE stocuri_2.cod_material = material.cod_material AND  
      nume_material = 'CIMENT';
```



Simplificarea codului - cont.

- Pe de altă parte însă, diferența de performanță (timp de execuție) dintre cele două variante de interogări este neglijabilă, de multe ori chiar interogarea pe structuri normalizate va fi mai rapidă decât cealaltă, astfel încât singurul avantaj al folosirii structurii nenormalizate este simplitatea.
- Dacă pentru un dezvoltator simplificarea invocată mai sus este nesemnificativă - orice dezvoltator trebuie să fie capabil să scrie o interogare pe mai multe tabele de tipul celei de mai sus - problema se pune cu mai mare stringență în cazul în care utilizatorii pot să-și scrie propriile interogări ale bazei de date.



Simplificarea codului - cont.

- În acest caz, pentru a păstra atât avantajele normalizării și pentru a permite în același timp simplificarea interogărilor, se poate crea o vedere bazată pe cele două tabele normalizate

```
CREATE VIEW stocuri AS
SELECT cod_depozit, stocuri_2.cod_material, nume_material, cantitate
FROM stocuri_2, material
WHERE stocuri_2.cod_material = material.cod_material;
```
- O astfel de vedere “ascunde” tabelele normalizate și permite utilizatorilor scrierea unor interogări mai simple.

Pe de altă parte însă, performanța programelor poate scădea mult în cazul vederilor complexe, în acest caz fiind uneori de preferat a se utiliza tabelele de bază pentru creșterea performanței.



Concluzii

- Designul logic al BD are un **impact** profund asupra **performanțelor și structurii sistemului**. Deși este posibil a compensa un design prost prin scrierea unui volum mare de cod, acest lucru nu este în nici un caz de dorit, mai ales că în majoritatea cazurilor el va duce inevitabil la scăderea performanțelor sistemului.
- În mod ideal, dezvoltarea oricărei aplicații trebuie să înceapă de la o structură complet normalizată, care să constituie o temelie solidă a acesteia.
- Denormalizarea, dacă ea este necesară, trebuie făcută numai după ce se obține o structură complet normalizată.
Ca regulă generală, denormalizarea este folosită rar și numai după o analiză atentă.



Temă

1. Se dă tabelul aflat în FN1:

INFO_STUDENT(cod_student, nume_student, data_nastere, cod_facultate, nume_facultate, cod_curs, nume_curs, data_inceput, nr_examen, nota)

Se presupune că:

- Un student poate studia la o singură facultate;
- Un student poate urma mai multe cursuri;
- Data_inceput reprezintă data la care studentul a început cursul;
- Între curs și facultate nu există legătură;
- În cadrul unui curs, studentul dă mai multe examene (numerotate prin intermediul coloanei nr_examen). Nota reprezintă nota luată de student la examenul cu numărul nr_examen din cadrul cursului cu codul cod_curs.



Temă - cont.

Să se determine cheia primară a tabelului și să se aducă în 2NF și 3NF, specificându-se cheile primare ale tabelelor rezultate.

Notă: nu se punctează decât rezolvările în care se parcurg toți pașii asociați normalizării (i.e. identificarea dependențelor parțiale și transitive față de cheie, descompunerea în tabele pentru eliminarea acestor dependențe, etc.).

2. Să se illustreze printr-un exemplu cum se transformă attributele repetitive (multivaloare) ale unei relații mulți-la-mulți (din modelul entitate-legătură) la crearea design-ului logic al unei baze de date relaționale.
3. Să se dea un exemplu de tabel relațional în care există o dependență multivaloare (multidependență) între attributele (coloanele) sale, care nu este dependență funcțională.



Bibliografie

F. Ipate, M. Popescu, Dezvoltarea aplicațiilor de baze de date în Oracle 8 și Oracle Forms 6, Editura ALL, 2000.