

## Three of Normalization Glitches

Prof.dr. Marin FOTACHE

Catedra de Informatică Economică, Universitatea „Alexandru Ioan Cuza”, Iași

*Among the database designers there is a strong feeling that normalization haven't fulfilled one of it's main objectives - being an appropriate methodology for guiding any database designer in achieving a good schema, according to the application requirements. One of the reasons could be too much formalism. But this paper is focused on a different area, trying to propose solutions for three types of problems which occurs quite often in practice: the apparent transitivity of some FDs, canonical FDs, and some semantic relationships between attributes which cannot be expressed as FDs.*

**Keywords:** normalization, functional dependencies (FD), apparently transitive FD, 3<sup>rd</sup> normal form, Boyce-Codd normal form, non-canonic FD.

### Introducere

Toți autorii importanți atrag atenția asupra faptului că, deși riguros formalizată, normalizarea nu este un panaceu. Există numeroase reguli semantice neincorporabile în normalizare. Pe de altă parte, unele cazuri reclamă o examinare mai atentă a DF sau chiar o anumită doză de artificialitate. Prezenta lucrare își propune să identifice și să propună soluții trei tipuri situaționale ce pot afecta calitatea unei scheme de baze de date obținută în urma normalizării: DF aparent tranzitive, DF necanonice și relații semantice dintre atribute ce nu pot fi modelate sub forma DF.

### Dependențe funcționale aparent tranzitive

Problema a fost studiată tangențial și nu prea zgomotos. Dintre cele câteva lucrări în care este analizată, merită amintită ce a lui Paolo Atzeni și Stott Parker (Atzeni & Parker 82). Să luăm un exemplu foarte simplu: R {Profesor, Birou, Catedră, Telefon} care conține date despre “domiciliul” în facultate al fiecărui profesor. Un profesor are un singur birou (Birou identifică fără ambiguitate sala respectivă): (1) Profesor  $\longrightarrow$  Birou. Un profesor face parte dintr-o singură catedră: (2) Profesor  $\longrightarrow$  Catedră. Într-un birou pot “locui” mai mulți profesori: Birou  $\longleftarrow$  Profesor. Fiecare birou este arondat unei singure catedre (Economie politică, Contabilitate etc.): (3) Birou  $\longrightarrow$  Catedră. Fiecare birou are un număr de telefon unic: (4) Birou  $\longrightarrow$  Telefon.

Din (3) ar reieși că (2) este DF tranzitivă, iar relația ar trebui descompusă numai în R1 {Birou, Catedră, Telefon} și R2 {Profesor, Birou}. Apare însă o situație ce aruncă în aer algoritmul: un profesor are un birou alocat altei catedre din care face parte ! Situația este reală. Catedrele pot “împrumuta”, temporar sau definitiv, birouri sau câte un loc-două din anumite birouri. Pentru profesorii care ocupă un birou alocat altei catedre, schema bazei alcătuită din R1 și R2 furnizează eronat răspunsul la întrebarea: Din ce catedră face parte fiecare profesor ?

O soluție ar fi să se delimiteze, semantic, catedrele care dețin birouri și catedrele la care sunt afiliați profesorii: R {Profesor, CatedrăProfesor, Birou, CatedrăBirou, Telefon}:

Acum dependențele ar fi:

- (1) Profesor  $\longrightarrow$  Birou
- (2) Profesor  $\longrightarrow$  CatedrăProfesor
- (3) Birou  $\longrightarrow$  CatedrăBirou
- (4) Birou  $\longrightarrow$  Telefon

Noua schemă ar fi: R1 {Birou, CatedrăBirou, Telefon} și R2 {Profesor, CatedrăProfesor, Birou}.

Se poate obiecta că ambele atribute, CatedrăBirou și CatedrăProfesor, se referă la aceeași entitate, deci seamănă a redundanță. Cu toate acestea, soluția elimină problema pierderii de informații, deci poate fi însușită la elaborarea schemei bazei de date.

### Dependențe tranzitive necanonice

Nici a o doua problemă pe care o vom ridica în continuare nu a beneficiat de prea multă atenție din partea celebriților normalizării, deși implicațiile sunt serioase. Mai toate lucrările recomandă folosirea exclusivă a depedențelor funcționale în formă canonică – cele în care destinația este alcătuită dintr-un singur atribut (de ex. [Saleh94], [Garcia-Molina s.a 02]). Lucrurile nu sunt grozav de complicate, identificarea depedențelor fiind o operațiune lesnicioasă. Din păcate, în unele situații este posibil ca tranzitivitatea să se stabilească prin depedențe ale căror destinații sunt compuse. La normalizarea relațiilor este recomandabil ca, în situația unor depedențe funcționale de genul:  $A \longrightarrow B$ ,  $A \longrightarrow C$ ,  $A \longrightarrow D$ ,  $A \longrightarrow E$ , depedențe ce derivă din rolul de cheie primară pe care îl

îndeplinește atributul A, să se verifice dacă nu cumva există fie o dependență funcțională simplă, de genul  $B \longrightarrow D$ , fie o dependență funcțională cu sursa compusă, de genul  $(B, C) \longrightarrow D$ . Dacă una din cele două afirmații de mai sus este adevărată, atunci dependența funcțională  $A \longrightarrow E$  nu este directă. De fapt, că problema poate fi pusă în legătură uneori cu forma normală Boyce-Codd.

Luăm în discuție tabela COMENZI, prezentată în figura 1, tabelă care gestionează comenzile pe care o firmă le-a trimis furnizorilor (în care se solicită materialele/produsele/serviciile de care are nevoie). Pe baza unei comenzi, un furnizor trimite produsele solicitate, întocmind o factură din care va fi remis un exemplar.

COMENZI					
NrComanda	NrFactura	CodFurnizor	DataFactura	Valoare	TVADeduct
1501	12344	12	20.06.2004	119000000	19000000
1502	11987	5	15.05.2004	67830000	10830000
1503	45963	2	14.05.2004	71400000	11400000
1504	12456	12	02.07.2004	67592000	10792000
1505	12344	5	27.06.2004	106691830	17034830
1505	12344	12	20.06.2004	101692640	16236640
1507	47890	2	29.05.2004	93886835	14990335

Fig.1. Relația COMENZI

Stabilim următoarele restricții:

- o comandă are un număr unic, stabilit de întreprinderea noastră, fiind întocmită pentru un singur furnizor;
- pentru o comandă, furnizorul va întocmi o singură factură;
- în schimb, o factură primită de la furnizor poate onora una, două sau mai multe comenzi.

Astfel, atributul NrComanda este cheia primară. Ca urmare, se poate scrie:

- (1) NrComanda  $\longrightarrow$  NrFactură
- (2) NrComanda  $\longrightarrow$  CodFurnizor
- (3) NrComanda  $\longrightarrow$  DataFactură
- (4) NrComanda  $\longrightarrow$  Valoare
- (5) NrComanda  $\longrightarrow$  TVADeduct

Facturile sunt întocmite de furnizori (care le numerotează independent între ei):

- (6) (NrFactură, CodFurnizor)  $\longrightarrow$  DataFactură
- (7) (NrFactură, CodFurnizor)  $\longrightarrow$  Valoare

(8) (NrFactură, CodFurnizor)  $\longrightarrow$  TVADeduct

Întrucât o factură primită poate onora mai multe comenzi:

(NrFactură, CodFurnizor)  $\not\rightarrow$  NrComanda  
În virtutea ultimelor trei depedențe funcționale, rezultă că dependențele (3), (4) și (5) sunt tranzitive,

NrComanda  $\longrightarrow$  (NrFactură, CodFurnizor)  
 $\longrightarrow$  DataFactură

NrComanda  $\longrightarrow$  (NrFactură, CodFurnizor)  
 $\longrightarrow$  Valoare

NrComanda  $\longrightarrow$  (NrFactură, CodFurnizor)  
 $\longrightarrow$  TVADeduct

după cum se observă în graful DF din figura 2.

Economia de spațiu și diminuarea redundanței sunt mai degrabă simbolice, ba chiar lucrurile se complică prin introducerea unei restricții referențiale și, normal, prin apariția de două ori a atributelor de legătură NrFactură și CodFurnizor. Există, totuși, un merit indiscu-

tabil al noilor relații. În relația COMENZI nu se poate prelua nici o factură fără a-i cunoaște comanda (eterna problemă a restricției de entitate – nici un atribut component al cheii nu poate avea nule). Practica economico-financiară de la noi, dar și din alte părți, cunoaște suficiente cazuri în care firmele primesc facturi direct, fără a înainta o comandă prealabilă. Utilizând proaspetele tabele, FACTURI\_PRIMITE și COMENZI\_FACTURI, o factură fără comandă va apărea numai în prima din cele două relații. Deși, la prima vedere, am adăugat, și nu eliminat (așa cum, generos, sună unul dintre obiectivele normalizării) redundanță, noua structură este mai “sănătoasă” din punct de vedere relațional.

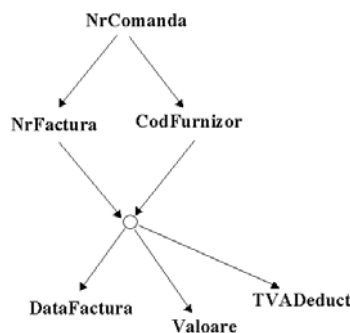


Fig.2. DF tranzitive necanonice

Pe baza tranzitivității, tabela COMENZI se poate sparge în două relații FACTURI\_PRIMITE și COMENZI\_FACTURI, ca în figura 3.

FACTURI PRIMITE					COMENZI FACTURI		
NrFactura	CodFurnizor	DataFactura	Valoare	TVADeduct	NrComanda	NrFactura	CodFurnizor
12344	12	20.06.2004	119000000	19000000	1501	12344	12
11987	5	15.05.2004	67830000	10830000	1502	11987	5
45963	2	14.05.2004	71400000	11400000	1503	45963	2
12456	12	02.07.2004	67592000	10792000	1504	12456	12
12344	5	27.06.2004	106691830	17034830	1505	12344	5
47890	2	29.05.2004	93886835	14990335	1505	12344	12
					1507	47890	2

Fig.3. Eliminarea DF tranzitive din tabela COMENZI

Practic, avantajul celor două relații din figura 3 este cu atât mai evident cu cât sunt numărul facturilor ce onorează două sau mai multe comenzi este mai mare.

Structura finală îndeplinește și cerințele formei normale Boyce-Codd ca fiecare sursă de DF să fie cheie candidată. Putem institui și o regulă: atunci când descompunerea se face urmând filiera clasică:

- stabilirea cheii primare a relației inițiale (universale), apoi
- eliminarea, dintre dependențele funcționale ce decurg din calitatea de cheie primară a relației inițiale, a celor parțiale, apoi
- eliminarea, din toate relațiile obținute în pasul anterior, a tuturor dependențelor tranzitive

**PROFI DISCIPLINE**

IdProf	NumeProf	Catedră	Cod Disc	DenDisc	LbStrProf	LbStrDisc
111	Airinei Dinu	Informatică Economică	AI1101	Introducere în informatica economică	FR	FR
112	Fotache Marin	Informatică Economică	AI1101	Introducere în informatica economică	EN	EN
112	Fotache Marin	Informatică Economică	AI1101	Introducere în informatica economică	FR	FR

Fig.4. Cursuri predate în limbi străine

tive (în care surse sunt cheile primare ale respectivelor relații), trebuie verificat dacă în relația inițială mai existau dependențe cu sursa compusă ce nu apar drept cheie primară sau candidat în nici una dintre relațiile finale obținute.

**Relații semantice imposibil de preluat în dependențe funcționale**

Un alt aspect delicat al normalizării ține de dificultatea preluării în schemă, sub formă de dependențe funcționale, a unor relații semantice. Pentru o primă ilustrare să luăm relația PROFII\_DISCIPLINE din figura 4 ce conține date legate de profesori și cursurile predate într-o facultate.

Un profesor este identificat prin IdProf și face parte dintr-o singură catedră:

- (1) IdProf  $\longrightarrow$  NumeProf
- (2) IdProf  $\longrightarrow$  Catedră.

Un curs este identificat prin cod și arondat unei singure catedre:

- (3) CodDisc  $\longrightarrow$  DenDisc
- (4) CodDisc  $\longrightarrow$  Catedră.

În schimb:

- un profesor predă mai multe discipline: IdProf  $\text{---}/\rightarrow$  CodDisc

- o disciplină poate fi predată de mai mulți profesori: CodDisc  $\text{---}/\rightarrow$  IdProf

- un profesor poate predă în mai multe limbi străine: IdProf  $\text{---}/\rightarrow$  LbStrProf

- un curs poate fi predat în mai multe limbi străine: CodDisc  $\text{---}/\rightarrow$  LbStrDisc

- un profesor poate predă același curs în mai multe limbi străine: (IdProf, CodDisc)  $\text{---}/\rightarrow$  LbStrDisc

- un profesor poate predă în aceeași limbă străină două sau mai multe cursuri:

(IdProf, LbStrDisc)  $\text{---}/\rightarrow$  CodDisc și

(IdProf, LbStrProf)  $\text{---}/\rightarrow$  CodDisc

Practic, graful DF se prezintă ca în figura 5.

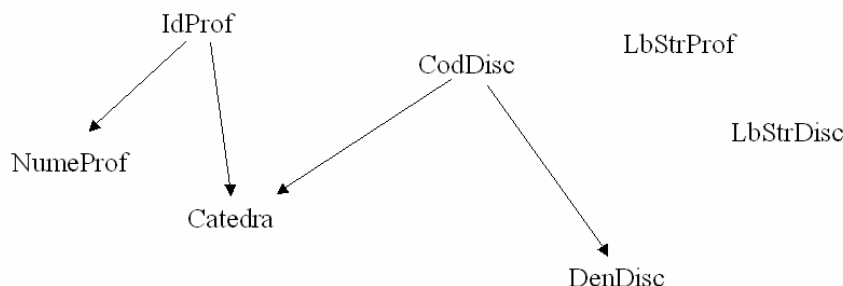


Fig.5. Graful DF pentru relația PROFIL\_DISCIPLINE

Nemulțumirea legată de informația reprezentată prin graf ține de faptul că un profesor știe anumite limbi străine, iar o disciplină este predată în câteva limbi, nicidecum în toate. Cu alte cuvinte, schema de mai sus ridică probleme la stocarea informațiile: care sunt limbile străine cunoscute de fiecare profesor și care sunt limbile străine în care se predă fiecare disciplină.

Atributele LbStrProf și LbStrDisc rămân suspendate în graf. Ce-i de făcut ? O idee, nu întotdeauna rezonabilă, este a le plasa într-o relație împreună cu celelate "rădăcini" ale grafului, obținându-se, astfel, schema:

PROFI {IdProf, NumeProf, Catedra}, DISCIPLINE {CodDisc, DenDisc, Catedra} și PROFIL\_DISC\_LBSTR { IdProf, CodDisc, LbStrProf, LbStrDisc}

Predicatul care ar exprima oricare tuplu al ultimei relații ar fi: un profesor (IdProf) care cunoaște o limbă străină (LbStrProf) predă o disciplină (CodDisc) într-una din limbile străine pe care le vorbește (LbStrDisc). Firește, apare întrebarea: Ce relație trebuie să existe între valorile LbStrProf și LbStrDisc

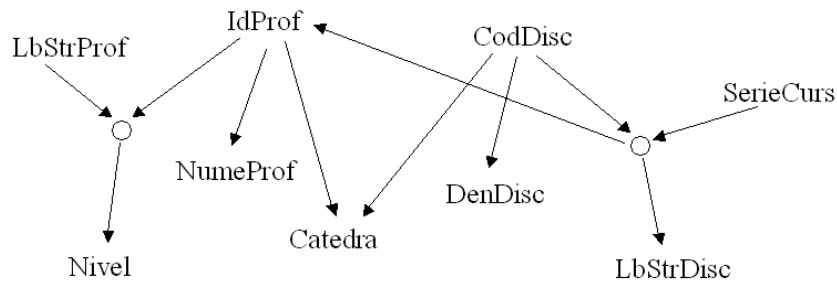
într-un tuplu, știind că disciplina este predată și în alte limbi decât LbStrProf, iar profesorul știe și alte limbi străine decât LbStrDisc ?

O soluție acceptabilă în multe cazuri constă în introducerea unor atribute suplimentare, tocmai pentru a "fixa" sub formă de depedente restricțiile respective. Spre exemplu, pentru a "formaliza" restricția un profesor știe anumite anumite limbi străine, se poate introduce un atribut, Nivel, pentru a ilustra gradul de comunicare al profesorului în limba străină respectivă: (IdProf, LbStrProf)  $\longrightarrow$  Nivel.

De asemenea, pentru a indica în ce limbi străine este disponibil un curs, se introduce atributul SerieCurs. La un obiect, o serie de curs se constituie pentru o singură limbă străină și case în responsabilitatea unui singur profesor:

(CodDisc, SerieCurs)  $\longrightarrow$  LbStrDisc,  
(CodDisc, SerieCurs)  $\longrightarrow$  IdProf.

Graful obținut nu mai conține atribute "suspendate" - vezi figura 6 - iar schema desprinsă din acesta este mult mai bună:



**Fig.6.** Noul graf DF pentru relația PROFI\_DISCIPLINE ameliorată

PROFI {IdProf, NumeProf, Catedra},  
DISCIPLINE {CodDisc, DenDisc, Catedra},  
PROFI\_LBSTR {IdProf, LbStrProf, Nivel}

și

SERII\_CURS\_LBSTR {CodDisc,  
SerieCurs, LbStrDisc, IdProf}.

Deși oarecum artificială, tehnica introducerii de atribute suplimentare, care nu sunt evidențiate în momentul construirii relației universale, este una salutară de multe ori, fiind folosită pe scară largă de practicieni. În plus, putem vorbi de un argument suplimentar în favoarea ideii de ciclicitate a normalizării, de reluare a o serie de pași până se obține o schemă considerată mulțumitoare.

### Concluzii

Prezenta lucrare încearcă să concilieze, pe calapodul a trei categorii situaționale, teoria și practica normalizării bazelor de date relaționale. Deși fără o încărcătură metodologică copleșitoare, cele trei probleme expuse sunt mai mult decât episodice pentru proiectanții de baze de date. Departate de a cădea în păcatul multor lucrări din domeniu, ce rezolvă problemele printr-o formulă matematică sau o teoremă și două leme, soluțiile prezentate sunt mai degrabă comune, simple, dar și aplicabile.

### Bibliografie

- [Atzeni & Parker 82] Atzeni, P., Parker, S. - *Assumptions in relational Database Theory*, Proc. of the 1<sup>st</sup> ACM SIGACT-SIGMOD symposium on Principles of database systems, Los Angeles, 1982
- [Date04] Date, C.J. - *An Introduction to Database Systems*, 8th edition, Pearson Addison-Wesley, Boston, 2004
- [Elmasri & Navathe 00] Elmasri, R., Navathe, S.R. - *Fundamentals of Database*

*Systems*, Addison-Wesley, Reading, Massachusetts, 2000

[Fotache01-1] Fotache, M. - *Câteva probleme ale formelor normalizate 2 și 3*, Net Report nr.104, mai, 2001

[Fotache01-2] Fotache, M. - *Cui îi mai este frică de formele normalizate Boyce-Codd, 4 și 5 ?*, Net Report nr.105, iunie, 2001

[Fotache04] Fotache, M. - *Do practitioners need database normalization ?*, Proc. of the Central and East European Conference in Business Information Systems, Ed. Risoprint, Cluj-Napoca, 2004

[Garcia-Molina s.a. 02] Garcia-Molina, H., Ullman, J., Widom, J. - *Database Systems. The complete Book*, Prentice Hall, Upper Saddle Riner, New Jersey, 2002

[Saleh94] Saleh, I. - *Les bases de données relationnelles. Conception et réalisation*, Hermes, Paris, 1994