

Databases Normalization. Algorithm to Determine the Closure of Attributes Set

Prof.dr. Nicolaie GIURGIȚEANU

Facultatea de Științe Economice, Universitatea din Craiova

Designing databases is one of the most laborious activities and so, one of the most boring parts in the life of an informatic system. Of course, the theorists have developed a system of axioms and rules in order to facilitate this activity. The main model of data representation, the relational model, has allowed the introduction into the databases theory, of some very efficient concepts. One of these concepts is that of functional dependency and multivalued dependency. These concepts are fundamental for the normalization of databases. In this paper, we present an algorithm, as a Prolog program, to determine the closure of a set of attributes, and also, another program that establishes whether a set of attributes is a primary key or not.

Keywords: database normalization, attributes set, system of axioms, system of rules.

Introducere

Proiectarea bazelor de date este una din activitățile cele mai plictisitoare din viața unui sistem informatic. De buna proiectare a bazelor de date depinde reușita ulterioară a realizării sistemului informatică. Deși este o activitate importantă cu care ar trebui să înceapă proiectarea oricărui sistem informatic, de multe ori ea este neglijată de marea majoritate a creatorilor de sisteme informatice și acest lucru are ca efect o slabă productivitate a muncii.

Conform bibliografiei actuale, cel puțin din punct de vedere teoretic, proiectarea bazelor de date se face în patru pași și anume: *specificații de cerere, proiectare conceptuală, proiectare logică și proiectare fizică.*

Specificațiile de cerere constau în obținerea tuturor cerințelor de la utilizatori, iar *proiectarea conceptuală* constă din traducerea acestor cerințe într-un model conceptual (de exemplu modelul Entități-Relații alias E-R). Rezultatele acestui pas se concretizează prin realizarea unei *scheme conceptuale*. Acest rezultat este preluat în etapa de *proiectare logică* și tradus într-un model de reprezentare a datelor (de exemplu modelul relațional) suportat de sistemul de gestiune a bazelor de date în care vom lucra. Rezultatul etapei de proiectare logică a bazelor de date se concretizează într-o schemă logică, ce va fi tradusă în etapa de *proiectare fizică*, într-o *schemă*

fizică adecvată configurației sistemului pe care se va exploata baza de date.

În lucrarea de față ne vom ocupa de proiectarea logică a bazelor de date relaționale și vom prezenta doi algoritmi, ambii în Visual Prolog, pentru determinarea închiderii unei mulțimi de atribute și, respectiv pentru stabilirea faptului că o mulțime de atribute este, sau nu, cheie pentru o schemă relațională dată.

Modelul relațional este unul din cele mai studiate modele și, în același timp, cel mai răspândit model. De fapt, vechile modele, ierarhic și rețea au fost de mult abandonate, poate noile modele, precum cel orientat obiect să constituie un adversar redutabil pentru cel relațional.

Una din cele mai importante ramuri ale teoriei bazelor relaționale este cea care se ocupă cu realizarea schemei relaționale. La rândul ei această ramură se bazează pe teoria dependențelor funcționale și restricțiilor impuse schemelor relaționale. Neformal, dependența funcțională se definește ca fiind proprietatea valorilor unor atribute de a permite reconstruirea, fără echivoc, a valorilor altor atribute.

Armstrong demonstrează că o familie de dependențe funcționale (sau echivalent, scheme relaționale) poate fi descrisă de închiderea unei mulțimi de atribute.

Scopul principal al lucrării noastre este de a

prezenta o modalitate de construire a închiderii unei mulțimi de atribute și de stabilire a unei chei primare pentru o schemă relațională dată și pentru o mulțime de dependențe funcționale, de asemenea, dată.

În a doua secțiune a lucrării vom prezenta conceptele, legate de dependențele funcționale, necesare pentru atingerea scopului lucrării. În partea a treia vom defini închiderea unei mulțimi de atribute, iar în partea a patra vom prezenta algoritmul Visual Prolog, pentru determinare închiderii și, respectiv al stabilirii unei chei primare.

Definiții de bază

În această secțiune vom prezenta, foarte pe scurt, conceptele principale ale teoriei dependențelor funcționale, necesare atingerii scopului propus. Alte concepte pot fi găsite în [01, 03, 04].

Fie U o mulțime de atribute (de exemplu: nume, adresa, studii etc. - universul discursului). Elementele lui U vor fi notate cu a, b, c, \dots, x, y, z sau, dacă este necesară o ordonare pe U , cu a_1, a_2, \dots, a_n . În unele cazuri, în locul literelor mici, se preferă literele mari. În teoria bazelor de date universul discursului se mai numește și schemă și atunci când vorbim de modelul relațional vom spune schemă relațională iar universul discursului va fi notat cu R .

Definiția 2.1. O aplicație *dom* asociază fiecărui atribut $a \in R$ o mulțime de valori, numită domeniu de valori, care este notată cu $dom(a)$.

Definiția 2.2. O relație r peste R este o submulțime a produsului cartesian $\prod_{a \in R} dom(a)$.

Notația 2.1. Dacă relația r are drept schemă relațională pe R , atunci, pentru a preciza acest lucru, vom folosi o notație de forma $r(R)$.

Definiția 2.3. Numim tuplu, al unei relații r , un element al produsului Cartesien $\prod_{a \in R} dom(a)$.

Deci o relație poate fi privită ca o mulțime de tupluri, adică $r = \{t_1(b), t_2(b), \dots, t_n(b)\}$, unde: $t_i(b) \in \prod_{a \in R} dom(a)$, sau cu alte cuvinte $t_i(b) = (v_1, v_2, \dots, v_n)$, cu $b = (a_1, a_2, \dots, a_n)$ și $t_i(a_j) = v_j \in dom(a_j)$.

Definiția 2.4. O dependență funcțională, sau pe scurt o FD, este o expresie de forma $X \rightarrow Y$, cu $X, Y \subseteq R$. Spunem că această FD are loc pe r dacă și numai dacă oricare ar fi tuplurile t_1 și t_2 din r , ori de câte ori vom avea $t_1(X) = t_2(X)$, vom avea și $t_1(Y) = t_2(Y)$.

Cu alte cuvinte, o FD este o restricție, naturală sau impusă, unei scheme relaționale, care permite stabilirea valorilor unei mulțimi de atribute, plecând de la valorile altei mulțimi de atribute. În acest fel avem un criteriu de verificare (validare) a relației.

Observația 2.1. Să observăm că dacă $X \rightarrow R$, atunci, valorile diferite ale lui X , vor identifica în mod unic, valori diferite ale lui $R - X$, sau cu alte cuvinte X va identifica fără ambiguitate tuplurile lui r . Aceasta înseamnă că X va acționa ca o cheie pe R .

Observația 2.2. Dacă între mulțimile de atribute X și Y există o dependență funcțională de forma $X \rightarrow Y$, atunci vom spune că X determină funcțional pe Y sau că Y este determinat funcțional de X .

Observația 2.3. Dacă $X \cup Y = R$, atunci dependențele funcționale pot fi folosite pentru eliminarea unuia din cele două tupluri care se repetă.

O altă aplicabilitate a acestei restricții o vom vedea după prezentarea câtorva proprietăți ale dependențelor funcționale.

Proprietăți ale dependențelor funcționale

Să observăm că dacă X, Y și Z sunt trei mulțimi de atribute cu proprietatea că X determină funcțional pe Y și Y determină funcțional pe Z , atunci X determină funcțional și pe Z .

Demonstrarea acestei afirmații este banală și nu o vom reproduce. Am folosit această afirmație pentru a vedea că dependențele funcționale au o serie de proprietăți. Printre proprietățile cele mai importante ale dependențelor funcționale enumerăm:

A1. Dacă $Y \subseteq X$, atunci $X \rightarrow Y$ (axioma reflexivității).

A2. Dacă $X \rightarrow Y$ și $Y \rightarrow Z$, atunci $X \rightarrow Z$ (axioma tranzitivității).

A3. Dacă $X \rightarrow Y$ și V este o mulțime oarecare de atribute atunci $XV \rightarrow YV$ (axioma de dilatare).

A4. (axioma reuniunii) dacă $X \rightarrow Y$ și $X \rightarrow Z$, atunci $X \rightarrow YZ$.

A5. (axioma descompunerii) Dacă $X \rightarrow YZ$ și $Y \cap Z = \Phi$, atunci $X \rightarrow Y$ și $X \rightarrow Z$.

A6. (axioma pseudo-tranzitivității) Dacă $X \rightarrow Y$ și $WY \rightarrow Z$, atunci $XW \rightarrow Z$.

Observația 3.1. O mulțime de dependențe funcționale poate genera o altă mulțime de dependențe funcționale. În acest fel se poate spune că dependențele funcționale implică logic alte dependențe funcționale.

Definiția 3.1. Mulțimea dependențelor funcționale, logic implicate de mulțimea de dependențe funcționale F , o vom numi închiderea lui F și o vom nota cu F^+ .

Acum putem reveni la proprietatea axiomelor lui Armstrong de a fi complete și stabile. Completitudinea se traduce prin aceea că: cele trei axiome permit determinarea închiderii oricărei mulțimi de dependențe funcționale, iar stabilitatea se traduce prin aceea că: aplicarea lor nu duce la apariția de dependențe funcționale inexistente sau parazite.

Observația 3.2. Dacă se dă o mulțime de atribute X și o mulțime de dependențe funcționale F pe R , atunci dependențele logic implicate de F , care folosesc atribute din X , vor duce la determinarea unei mulțimi de atribute logic implicate de X .

Definiția 3.2. Mulțimea tuturor atributelor logic determinate de X , sub mulțimea de dependențe funcționale F , o vom numi închiderea lui X și o vom nota cu X_F^+ .

Definiția 3.3. Dacă X_F^+ este chiar R , atunci vom spune despre X că este o cheie (un identificator) pentru R .

Observația 3.3. Dacă X este o cheie pentru

Agenție	Avere	Oraș	Dosar	Client	Val împ
BCR Iași 1	9.000.000	Iași	17	Aioanei	2.000
BCR Craiova	8.700.000	Craiova	21	Georgescu	3.000
BCR Iași 2	7.500.000	Iași	23	Apetrei	2.500
BCR Cluj	9.600.000	Cluj-Napoca	45	Pătru	3.000
BCR Craiova	8.700.000	Craiova	37	Ionescu	2.500
BCR Iași 1	9.000.000	Iași	32	Sandu	3.500
BCR Oradea	7.500.000	Oradea	45	Manta	2.500

Fig. 1. Relația *împrumut*

Se constată că această relație suferă de următoarele anomalii.

1. **repeteția informației.** Se constată că pentru orice împrumut acordat de o agenție, trebuie să repetăm averea și orașul în care-și are sediul agenția respectivă. O redundanță de informație inacceptabilă.

R atunci orice tuplu al lui r va fi identificat fără ambiguități de o valoare a lui X .

Aplicații ale dependențelor funcționale

În secțiunea 2 am văzut una din aplicațiile dependențelor funcționale. În cele ce urmează vom arăta și alte aplicații ale acestora.

Pentru a vedea și alte aplicații ale dependențelor funcționale să considerăm schema relațională *Împrumut*=(*agenție, avere, oraș, dosar, client, valoare împrumut*), dintr-un sistem bancar. Semnificația atributelor acestei scheme relaționale este: *agenție* – sucursala unei bănci; *avere* - reprezintă capitalul sucursalei respective; *oraș* este orașul în care-și are sediul sucursala sau agenția respectivă; *client* este clientul care face un împrumut la bancă prin agenția respectivă; *dosar* este numărul dosarului cu care s-a acordat împrumutul, iar *valoare împrumut* este valoarea împrumutului acordat.

Să vedem ce avantaje ar avea o relație construită pe baza acestei scheme relaționale. În primul rând să observăm că dacă ne-ar interesa valoarea împrumutului acordat unui client vom rezolva direct această problemă fără alte prelucrări. De asemenea, dacă ne-ar interesa agenția care a acordat împrumutul unui client o putem rezolva fără alte prelucrări. Din punct de vedere al rapidității obținerii informației se pare că această schemă relațională este cea mai bună.

Să analizăm însă și alte aspecte legate de această schemă relațională. Pentru aceasta să considerăm relația *împrumut* prezentată în figura 1.

2. **anomalii de actualizare.** Repetarea informației nu are repercursiuni doar asupra dimensiunilor memoriei în care se păstrează informația ci are repercursiuni și asupra actualizărilor. De exemplu să presupunem că se mărește capitalul agenției BCR Iași 1. Toate tuplurile care se referă la această agenție tre-

buie să conțină același capital. Dacă dintr-un motiv oarecare nu toate tuplurile vor fi aduse la același capital, atunci am avea o bază de date care oferă capitaluri diferite ale aceleiași agenții. Acest lucru mai este cunoscut și sub denumirea de inconsistență a datelor. Acest lucru vrea să spună că din aceleași date nu obținem aceeași informație prin prelucrări și interpretări identice.

3. anomalii de inserare - inabilitatea reprezentării informației. Să presupunem că o cheie primară pentru această relație ar fi perechea (agenție, dosar). Orice nou tuplu inserat în relație trebuie să conțină o valoare pentru aceste atribute, deoarece integritatea existențială cere ca nici o cheie să nu fie parțial sau total NULL. Cu toate acestea dacă cineva vrea să insereze un nou client și un nou împrumut în baza de date, el nu va putea face acest lucru fără ca acest client să facă un împrumut de la o agenție oarecare, astfel încât atributele cheii primare să primească o valoare. Deci informațiile despre un nou client nu pot fi inserate decât dacă clientul face un împrumut. Acest tip de problemă este cunoscut ca anomalii de inserare.

4. anomalii de ștergere – pierdere de informații importante. În anumite circumstanțe se poate pierde informație utilă prin efectuarea unei ștergeri. De exemplu, dacă dintr-un motiv oarecare clientul care face un împrumut la agenția BCR Oradea renunță la acesta, atunci prin ștergerea tuplului corespunzător se va pierde informația despre agenția BCR din Oradea.

Probleme ca cele de mai sus apar în primul rând datorită faptului că relația *Împrumut* conține atât informații despre agenții, cât și informații despre împrumuturile acordate clienților. Una din soluțiile folosite pentru eliminarea acestor probleme o constituie descompunerea relației și deci a schemei relaționale, în două sau mai multe sub-relații (subscheme). Descompunerea poate produce

beneficii importante mai ales în sistemele distribuite, informații distincte pot fi memorate în locuri distincte. Desigur că descompunerea relațiilor are și efecte negative asupra timpului de obținere a informației deoarece relațiile distincte trebuie legate, în vederea obținerii unui răspuns corect și complet. Din analiza relației de mai sus și schemei sale relaționale se constată că există următoarele restricții de tip dependențe funcționale:

agenție → *avere* ;

agenție → *oraș* ;

dosar → *val_imp* ;

dosar → *agenție*.

Toate aceste dependențe sunt naturale, intrinseci problemei respective. Pentru a evita marea majoritate a anomaliilor menționate mai sus am putea descompune schema relațională *Împrumut* în următoarele scheme relaționale:

Agenții=(*agenție, avere, oraș*) și

Credite=(*agenție, dosar, client, val_imp*).

Se poate observa că aceste noi scheme relaționale elimină o parte din anomalii și că ele se obțin pe baza dependențelor funcționale. În teoria bazelor de date se arată că există o serie de scheme relaționale cu anumite proprietăți standard, numite forme normale și aceste forme normale se obțin pe baza dependențelor funcționale.

În secțiunea următoare vom prezenta doi algoritmi descriși în Prolog, unul pentru determinarea închiderii unei mulțimi de atribute iar al doilea pentru determinarea unei chei de identificare.

Algoritmi Visual Prolog pentru lucru cu FD

Una din formele normale cele mai cunoscute și mai practice, din teoria bazelor de date, este forma normală Boyce-Codd (BCNF). Pentru descompunerea unei scheme relaționale non BCNF în scheme relaționale BCNF vom folosi un procedeu interactiv combinat cu algoritmul descris în Prolog pentru determinarea închiderii unei mulțimi de atribute și prezentat în figura 2.

```
domains
    lista=symbol*
database
    fd(lista, lista)
predicates
    incarcaBICD
    nondeterm insert(symbol, lista, lista)
    nondeterm concaten(lista, lista, lista)
    nondeterm closure(lista)
    nondeterm incl(lista, lista)
    nondeterm apartine(symbol, lista)
```

```

nondeterm aflista(lista)
nondeterm go
nondeterm append(lista,lista,lista)
nondeterm qs(lista,lista)
nondeterm partition(symbol,lista,lista,lista)
clauses
incarcaBICD:-
/* assertz(fd([agentie],[avere]),
assertz(fd([agentie],[oras]), */
assertz(fd([dosar],[val imp])),
assertz(fd([dosar],[agentie])),
assertz(fd([dosar],[client])).
closure(X):-
R=X, retract(fd(U,V)), incl(U,R),
concaten(V,R,RR),
qs(R,R1), qs(RR,R2), not(R1=R2),
consult("BICD"), closure(R2).
closure(X):-write("[",aflista(X),write("]"),nl.
aflista([]).
aflista([X|R]):-
write(X,","), aflista(R).
incl([],).
incl([X|R],Y):-
apartine(X,Y), incl(R,Y).
apartine(X,[X|_]).
apartine(X,[_|R]):-apartine(X,R).
concaten([],L,L):-!.
concaten([X|R],L2,L3):-
insert(X,L2,L4),concaten(R,L4,L3).
insert(X,[],[X]).
insert(X,L,L):-L=[X|_],!.
insert(X,[Y|R],[Y|L]):-insert(X,R,L).
append([],L2,L2):-!.
append([X|L1],L2,[X|L3]):-
append(L1,L2,L3).
qs([],[]).
qs([A],[A]).
qs([A,B|R],Ans):-
partition(A,[B|R],Left,Right),
qs(Left,Ls),qs(Right,Rs),
append(Ls,[A,L1],append(L1,Rs,Ans)).
partition(, [], [], []).
partition(A,[B|Rest],[B|L],R):-
A > B, partition(A,Rest,L,R).
partition(A,[B|Rest],L,[B|R]):-
A <= B, partition(A,Rest,L,R).
go:-
incarcaBICD,save("BICD"),write("Inchiderea lui [",aflista([dosar]),
write("] este "),closure([dosar])).
goal go.

```

Fig. 2. Algoritm pentru determinarea închiderii unei mulțimi de atribute

Rulând programul prezentat în figura 2 vom obține închiderea mulțimii de atribute formată din atributul *dosar* și se constată că aceasta este chiar schema relațională *Credite*. Am adăugat, în acest program încă o dependență funcțională pe care o impunem noi și anume *dosar* → *client*. Modifi-

când programul de mai sus prin renunțarea la comentarii și calculând închiderea atributului *agenție* se constată că aceasta este chiar schema relațională *Agenții*. Algoritmul pentru determinarea unei chei de identificare este prezentat în figura 32.

```

domains
lista=symbol*
database
fd(lista,lista)
schema(lista)
predicates
incarcaBICD
nondeterm go
incarcaschema
nondeterm insert(symbol,lista,lista)
nondeterm concaten(lista,lista,lista)
nondeterm closure(lista)
nondeterm incl(lista,lista)
nondeterm apartine(symbol,lista)
nondeterm aflista(lista)
nondeterm append(lista,lista,lista)
nondeterm qs(lista,lista)
nondeterm partition(symbol,lista,lista,lista)
nondeterm verificaKey(lista)
clauses
incarcaschema:-
assert(schema([a,b,c,g,h,i])).
incarcaBICD:-

```

```

    assertz(fd([c,g],[h])),    assertz(fd([a],[b])),    assertz(fd([a],[c])),
    assertz(fd([b],[h])),    assertz(fd([c,g],[i])).
closure(X):-
  R=X,
  retract(fd(U,V)),          incl(U,R),          concaten(V,R,RR),
  qs(R,R1),                  qs(RR,R2),          not(R1=R2),
  consult("BICD"),          closure(R2).
closure(X):-verificaKey(X).
aflista([]).
aflista([X|R]):-
  write(X," "),
  aflista(R).
incl([],_).
incl([X|R],Y):-
  apartine(X,Y),
  incl(R,Y).
apartine(X,[X|_]).
apartine(X,[_ |R]):-apartine(X,R).
concaten([],L,L):-!.
concaten([X|R],L2,L3):-
  insert(X,L2,L4),concaten(R,L4,L3).
insert(X,[],[X]).
insert(X,L,L):-L=[X|_],!.
insert(X,[Y|R],[Y|L]):-insert(X,R,L).
verificaKey(L):-
  schema(L),                qs(L1,L2),          qs(L,L3),          L2=L3,
  write(" este identificator al schemei relationale n ["),
  aflista(L),write("]"),nl.
verificaKey(_):-
  schema(L),
  write(" nu este identificator al schemei relationale •n ["),
  aflista(L),write("]"),nl.
go:-
  incarcaschema,incarcaBICD,save("BICD"),write("["),aflista([dosar]),
  write("]"),closure([dosar]).
append([],L2,L2):-!.
append([X|L1],L2,[X|L3]):-
  append(L1,L2,L3).
qs([],[]).
qs([A],[A]).
qs([A,B|R],Ans):-
  partition(A,[B|R],Left,Right),
  qs(Left,Ls),qs(Right,Rs),
  append(Ls,[A],L1),append(L1,Rs,Ans).
partition(_,[],[],[]).
partition(A,[B|Rest],[B|L],R):-
  A > B,partition(A,Rest,L,R).
partition(A,[B|Rest],L,[B|R]):-
  A <= B,partition(A,Rest,L,R).
goal go.

```

Fig. 3. Algoritm pentru determinarea unei chei de identificare

Se poate constata, rulând acest algoritm, că *agenție* este o chei de identificare pentru schema relațională *Agenții*, iar *dosar* este chei de identificare pentru schema *Credite*.

Bibliografie

- Giurgițeanu N., *Proiectarea bazelor de date relaționale*. SITECH, Craiova 1997
- Giurgițeanu, N., Popa, S., *Using Prolog to verify the Functional and Multivalued Dependencies*, International Workshop IE&SI, Timișoara, 23-24 Mai 2003, 157-162
- Elmasri R., Navathe S., *Fundamentals of database Systems*. Benjamin/Cummings, 1994
- Codd E.F. A relational model for large shared data banks. Comm. ACM 13 (1970), 377-387.