

Multidimensional data analysis using OLAP Technology (1)

Asist. Gianina RIZESCU

Catedra de Contabilitate și Informatică Economică, Universitatea "Dunărea de Jos" Galați

In this paper we present the main steps in creating an application using OLAP technology. The main goals of this application are: the quality analysis of the teaching process through the results obtained by the students of a university and the structure and dynamic analysis of the students in a university ("Dunărea de Jos" of Galați was taken as example). Thus, after a brief introduction, we will present de general architecture of the application followed by the description of the analysis, design, and populating stages of the data warehouse and also by the multidimensional data analysis using decision cubes.

Keywords: OLAP, OLTP, data warehouse, data mart, decision cub.

Introducere

Ca exemplificare practică a etapelor parcurse în realizarea oricărei aplicații utilizând tehnologia OLAP, s-a proiectat o aplicație informatică pentru universitatea „Dunărea de Jos” care își propune să realizeze două funcții majore:

- ✓ Analiza calității procesului de învățământ prin intermediul rezultatelor obținute de studenții universității
- ✓ Analiza structurii și dinamicii studenților din universitate

Necesitatea dezvoltării unei astfel de aplicații la nivel de universitate a apărut din următoarele considerente:

- ✓ Necesitatea obținerii unei imagini globale, la nivel de universitate, asupra structurii și dinamicii studenților. Această necesitate a apărut datorită concurenței acerbe care s-a declanșat între universități în ultimii ani, a creșterii numărului universităților particulare pe de o parte, iar pe de altă parte scăderii de la an la an a numărului de studenți. Acest lucru face necesară obținerea de informații suplimentare, integrate la nivel de universitate pentru elaborarea unor politici de marketing în atragerea unui număr cât mai mare de candidați.
- ✓ Necesitatea obținerii unei imagini globale și integrate asupra calității procesului de învățământ pentru a identifica și înlătura factorii care influențează în mod negativ acest proces pe de o parte, iar pe de altă parte pentru a găsi mijloacele necesare creșterii calității acestui proces

- ✓ Fiecare facultate utilizează propriul sistem de evidența a studenților și a rezultatelor acestora. Apare astfel necesitatea creării unei structuri de date integrate și standardizate la nivel de universitate pentru a se facilita analiza comparativă a datelor și pentru a se înlătura informațiile inconsistente și deseori contradictorii

- ✓ Oferirea unui instrument flexibil de analiză a datelor factorilor de decizie care să le ofere acestora independență față de sistemele tranzacționale a căror rapoarte și liste nu satisfac varietatea necesităților de informare.

Arhitectura generală a aplicației

Arhitectura generală a aplicației este prezentată în figura 1.

Referitor la această arhitectură generală se impun următoarele observații:

- ✓ Depozitul de date este de fapt un datamart (magazie, piață de date) datorită orientării sale pe un singur subiect. În continuare însă vom folosi noțiunea de depozit de date pentru ușurință exprimării și datorită faptului că un datamart este tot un depozit de date de mai mici dimensiuni și orientat pe subiect.

- ✓ Pentru crearea modulului de populare a depozitului de date s-a luat în considerare doar baza de date a unei singure facultăți deoarece conceperea unui sistem de extragere și de populare a depozitului de date cu date din sistemele tranzacționale ale tuturor facultăților implica alte aspecte care nu au fost tratate și o complexitatea greu de acoperit la momentul respectiv.

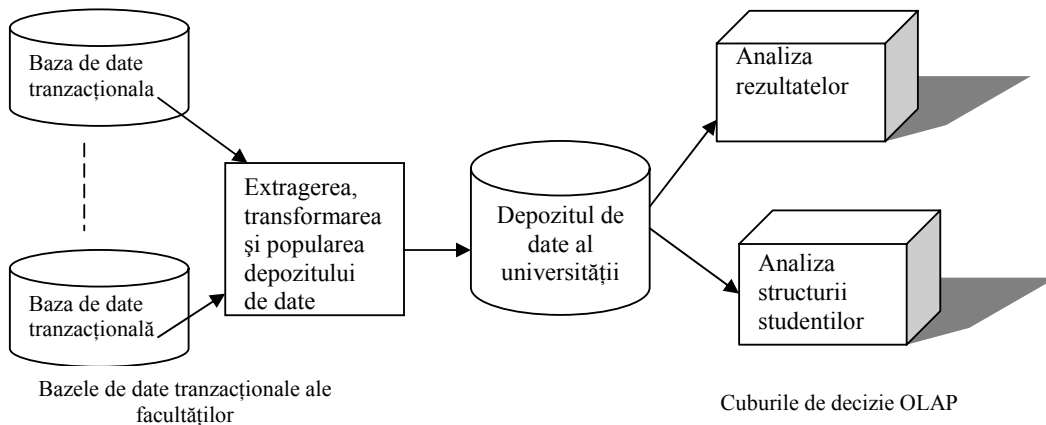


Fig. 1. Arhitectura generală a sistemului

Pentru realizarea aplicației s-a utilizat tehnologia oferită de Microsoft, SQL Server 7.0 din mai multe considerente:

- ✓ Oferă o soluție integrată și completă de dezvoltare a sistemelor de decizie bazate pe OLAP
- ✓ Oferă o soluție mai ieftină deoarece nu se plătește licența suplimentară pentru serverul OLAP acesta fiind furnizat odată cu Microsoft SQL Server 7.0
- ✓ Face posibilă implementarea depozitelor de date de către firmele mici și mijlocii care nu-și permit costuri suplimentare pentru alte sisteme de gestiune a bazelor de date
- ✓ SQL Server este mult mai ușor de utilizat decât majoritatea sistemelor de gestiune a bazelor de date de pe piață
- ✓ Oferă instrumente vizuale și intuitive de dezvoltare a aplicațiilor care permit crearea mult mai rapidă a lor.

În procesul de dezvoltare a depozitului de date au fost parcurse etapele care alcătuiesc ciclul de viață a oricărui depozit de date și anume: Analiză, Proiectare, Populare.¹

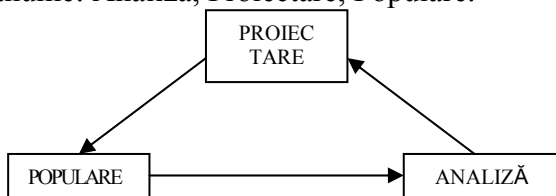


Fig. 2. Ciclul de viață al depozitelor de date

I. Analiza în vederea creării depozitului de date

Punctul de plecare în implementarea oricărei baze de date este analiza, în timpul căreia

trebuie să se răspundă la o serie de întrebări. Care sunt entitățile implicate? Care sunt relațiile între entități? Ce atribute trebuie asignate fiecărei entități? Ce funcții trebuie îndeplinite și cum afectează aceste funcții entitățile și atributele identificate?

Proiectarea logică sau faza de analiză a unui proiect în mod tipic are două puncte de plecare, care în general se suprapun:

- ✓ Procesele sau funcțiile. Ce funcții va trebui să îndeplinească viitorul sistem?
- ✓ Datele. Ce date sunt necesare pentru a sprijini realizarea funcțiilor afacerii?

Toate acestea se determină plecând bineînțeles de la cerințele utilizatorilor. Aceste lucruri sunt valabile atât pentru sistemele OLTP cât și pentru sistemele OLAP. Diferența constă în categoria de utilizatori pe care o are în vedere: dacă în sistemele tranzacționale utilizatorii țintă sunt cei executivi (în principal), în sistemele analitice ținta o reprezintă utilizatorii factori de decizie.

I.1 Analiza cerințelor și a datelor necesare
Înainte începerii dezvoltării depozitului de date, domeniul afacerii trebuie să fi fost înțeles pe deplin. Scopul final al unui astfel de sistem este de a oferi utilizatorilor libertatea de a manipula datele liber, fără constrângeri externe. De aceea analiza este necesară pentru a ne asigura că în depozitul de date ajung toate datele necesare – la nivelul adecvat de detaliu – astfel încât acesta, depozitul de date, să fie capabil să răspundă la cât mai multe dintre interogările utilizatorilor. Deci, în faza de analiză trebuie identificate care sunt cele mai frecvente și cele mai importante interogări la care depozitul de date trebuie să răs-

¹ Dorin Zaharie, & co, „Sisteme informatice pentru asistarea deciziei”, Ed. DualTech, București 2001

pundă.

Astfel, pentru analiza calității procesului de învățământ s-a considerat că trebuie oferite răspunsuri pentru următoarele întrebări:

- ✓ Care sunt rezultatele studenților pe discipline sau grupuri de discipline în diferite perioade?
- ✓ Care sunt rezultatele comparative ale studenților pe specializări, secții, facultăți?
- ✓ Care sunt rezultatele comparative ale studenților pe orașe, județe, țări?
- ✓ Care sunt rezultatele comparative ale studenților pe sexe, stare civilă?
- ✓ Care sunt rezultatele comparative ale studenților în funcție de regimul de școlarizare: buget, taxă, zi, IDD(Învățământ Deschis la Distanță), IFR(Învățământ cu Frecvență Redusă)
- ✓ Rezultatele comparative ale studenților în orice combinație dintre criteriile definite mai sus.

Pentru analiza structurii și dinamicii studenților din universitate, s-a considerat necesar să se răspundă la următoarele întrebări:

- ✓ Care este structura și dinamica studenților pe orașe, județe, țări?
- ✓ Care este structura și dinamica studenților pe sexe, stare civilă?
- ✓ Care este structura și dinamica studenților pe specializări, secții, facultăți?
- ✓ Care este structura și dinamica studenților în funcție de regimul de școlarizare : buget, taxă, zi, IDD(Învățământ Deschis la Distanță), IFR(Învățământ cu Frecvență Redusă)?
- ✓ Structura și dinamica studenților pe orice combinație dintre criteriile de mai sus.

Identificarea datelor necesare nu este însă suficientă. În afară de această identificare se urmărește elaborarea unei structuri unice și integrate a datelor: tipul de date, modalitatea de exprimare a unor date, stabilirea modului de exprimare a timpului în analiza datelor știind că orice depozit de date are o dimensiune timp care se exprimă diferit în funcție de semnificația timpului pentru sistemul respectiv. Toate acestea se realizează prin utilizarea de convenții consistente în privința numelor, măsurătorilor, atributelor și semanticii.

II. Proiectarea depozitului de date

Structura depozitului de date are în vedere

identificarea precisă a datelor stocate și accesul rapid la ele. Pentru aceste deziderate, masa de informații care se va stoca în depozit trebuie astfel organizată încât să reflecte atât datele importante cât și contextul lor. Modelarea dimensională oferă suportul necesar pentru proiectarea structurii depozitului de date. Structura se implementează sub forma unei baze de date care să asigure atât stocarea unui volum imens de date cât și accesul rapid la ele.

Cel mai popular model pentru depozitele de date este modelul multidimensional. Acesta poate fi în formă de stea (star schema), de fulg de zăpadă (snow-flake schema) sau de constelație (constellation schema), în care se regăsesc datele cantitative (cantități valori) din tabelele de tranzacții agregate în principal pe unitatea de timp (zi, lună etc.) și apoi după alte criterii: (pe student, pe disciplină, pe specializare, pe facultăți etc.) și întotdeauna data calendaristică, primul criteriu de agregare. Aceste date cantitative centralizate sunt *măsurile ale activității* iar criteriile de agregare sunt denumite *dimensiuni*. Măsurile identificate prin dimensiuni sunt stocate într-o tabelă relațională denumită *tabela de fapte*. Codurile criteriilor de agregare sunt explicitate în tabele de tip nomenclator asociate tabelii de fapte (tabelele de dimensiuni), schema relațională căpătând forma de stea. Mai multe asemenea scheme de tip stea care folosesc aceleași nomenclatoare formează un model de tip *constelație* iar dacă dimensiunile se pot divide în subdimensiuni, atunci nomenclatoarele pot avea la rândul lor alte nomenclatoare formând o schemă ce are forma unui *fulg de zăpadă*.

II.1 Identificarea și definirea dimensiunilor

Dimensiunile determină modalitățile în care datele sunt organizate – cu alte cuvinte, cum datele pot și trebuie să fie filtrate, organizate, și grupate. Pentru determinarea dimensiunilor s-au avut în vedere următoarele:

- ✓ Caracteristicile dimensiunii;
- ✓ Caracteristicile tabelii de dimensiuni;

Caracteristicile dimensiunii.

Tabelele de dimensiuni trebuie proiectate ținând cont în permanență de utilizatorul final

și de necesitățile sale de analiză. Dimensiunile sunt constituite de tabele relaționale și ca orice tabelă relațională, sunt definite prin cheie primară și atribute. Acestea trebuie să conțină: atribute strâns corelate, denumiri clare, elemente de nume și adresă atomice, chei „surogat”. În loc de a se utiliza valorile cheilor din sistemul sursă este recomandat să se utilizeze aceste chei surogat, obținute prin generarea unor valori unice pentru identificarea entităților din tabelele de dimensiuni. Aceste chei sunt valori numerice, secvențiale care sunt gestionate de sistemul de gestiune a bazei de date. Nu sunt atribuite de către utilizator. Folosirea acestor chei oferă următoarele avantaje:

- ✓ Oferă independență față de sistemul sursă, înlăturând următoarele probleme:
 - modificarea aplicațiilor poate cauza și modificarea structurii cheilor
 - aplicațiile sursă pot reutiliza cheile
- ✓ Oferă o indexare mai eficientă

Caracteristicile tablei de dimensiuni.

Tabelele de dimensiuni au caracteristici comune, care le fac mult mai ușor de definit. La definirea unei table de dimensiuni există câteva reguli ce trebuie avute în vedere: Orice tabelă de dimensiuni are o cheie primară care este determinată unic pentru fiecare tabelă. Nu se utilizează chei compuse și nu se preiau cheile din aplicațiile sursă; Între o tabelă de dimensiuni și tabela de fapte există o relație unu la mai mulți. Unei înregistrări din tabela de dimensiuni îi corespunde mai multe înregistrări în tabela de fapte; Conține cel puțin un atribut în afară de cheia primară; Conține alte atribute care sunt utile pentru definirea diferitelor niveluri de agregare. Astfel de atribute sunt denumite *ierarhii* ale dimensiunii; Conține un număr limitat de înregistrări care crește încet în timp.

Având în vedere cele două obiective majore ale depozitului de date: analiza rezultatelor școlare ale studenților și analiza structurii studenților din universitate, și respectând pe cât posibil observațiile de mai sus, au fost identificate și create următoarele tabele de dimensiuni:

Student: conține toate atributele legate de student care ar putea interveni la un moment

dat în analiză: nr_matricol, numele și prenumele studentului, sex, stare civilă, vârsta, anul admiterii, media de admitere.

Timp: în acest caz, dimensiunea timp are o structură mai particulară, vând în vedere ca o analiză a rezultatelor școlare ale studenților au sens doar pe sesiune și pe an calendaristic, acestea sunt și atributele dimensiunii timp. Combinația dintre anul calendaristic și sesiune poate duce și la obținerea rezultatelor pe semestru.

Tip_examen: conține ca atribute examenul, tipul examenului și data acestuia, pentru a putea analiza rezultatele pe tipuri de examene și chiar pe fiecare examen în parte

Disciplina: conține atribute referitoare la discipline: denumire disciplină, categorie disciplină.

Profesor: conține toate atributele referitoare la profesor: numele, sex, funcție, vârstă, vechime. Introducând această dimensiune, analiza capătă o altă nuanță și anume, se poate analiza, pe lângă rezultatele studenților, performanțele didactice ale profesorilor.

Zona geografică: conține ca atribute orașul, județul, regiunea, țara ceea ce va permite analiza rezultatelor și a structurii studenților după zona lor de proveniență.

Structură_organizatorică: conține atributele legate de apartenența unui student la o anumită grupă. Utilizarea acestei dimensiuni face posibilă urmărirea studenților în dinamica trecerii lor de la o grupă la alta, de la un an de studiu la altul. Cuprinde următoarele atribute: facultate, secție, grupă, subgrupă, anul de studiu.

Timp_stud: conține ca atribute, anul calendaristic și anul universitar necesare pentru analiza numărului și structurii studenților.

Pentru toate dimensiunile s-au folosit ca și chei primare cheile surogat pentru a asigura independența acestora de baza de date tranzacțională și cu ajutorul cărora se va asigura și legătura dimensiunilor cu tabela de fapte.

II.2 Definirea tablei de fapte și a măsurilor

Înainte de definirea tablei de fapte trebuie stabilit nivelul de detaliu până la care se dorește analizarea datelor. Atributele tablei de fapte trebuie alese cu grijă deoarece acestea

stau la baza tuturor deciziilor viitoare. Granularitatea tabelului de fapte va determina configurația atributelor acestuia.

Nivelul de sumarizare ales pentru tabela de fapte va avea un profund efect asupra depozitului de date. Cel mai mic nivel, posibil de granularitate este acela de a stoca faptele la nivelul atomic tranzacțional. Acesta oferă același nivel de detaliu ca și sistemul sursă. Stocarea faptelor la acest nivel de detaliu prezintă avantajul major că permite analiza până la cele mai mici niveluri de detaliere, că poate răspunde la interogări neașteptate și că depozitul de date este mult mai flexibil la introducerea de noi date. Cel mai mare dezavantaj îl reprezintă faptul că adesea este aproape imposibil ca depozitul de date să poată stoca un asemenea volum de date, și deoarece cele mai multe interogări se fac la un anumit nivel de sumarizare, datele detaliate nu vor face decât să încetinească timpul de răspuns la aceste interogări. Granularitatea tabelului de fapte va determina profunzimea analizelor care se pot face asupra depozitului de date. O tabelă de fapte în general, are o serie de caracteristici, de care de asemenea am ținut cont în realizarea ei: Cheie primară a tabelului de fapte se compune din cheile primare ale tabelurilor de dimensiuni; Măsurile sunt reprezentate prin coloane numerice care reflectă nivelul de detaliu al datelor stocate în tabela de fapte; Cheile primare ale tabelurilor de dimensiuni sunt declarate chei străine în tabela de fapte pentru relațiile unu la mai mulți între tabelurile de dimensiuni și tabela de fapte; Multe înregistrări din tabela de fapte nu vor avea valori pentru toate cheile străine asociate dimensiunilor; Conține un număr foarte mare de înregistrări, de ordinul sutelor de mii și chiar milioane. Nivelul de detaliu stabilit pentru tabela de fapte determină în mod direct și dimensiunea acesteia și implicit timpul necesar de răspuns la o interogare.

Având în vedere cele două obiective majore ale depozitului de date: analiza calității procesului de învățământ și analiza structurii studenților din universitate, respectând pe cât posibil observațiile de mai sus, au fost identificate două tabele de fapte:

Tabela **Rezultate** stochează toate rezultatele

obținute de studenți. Măsura acestor rezultate este dată de notă, credite, note peste 5, note peste 8 care reprezintă măsurile tabelii de fapte. Aceste măsuri vor sta la baza calculării indicatorilor de analiză a calității procesului de învățământ.

Tabela **Structură Studenți** stochează faptele legate de structura studenților ținând cont de dinamica acestora pe parcursul anilor. Măsura tabelii este dată de numărul de studenți, măsură ce va fi utilizată pentru calcularea indicatorilor necesari pentru analiza structurii studenților.

O particularitate a sistemului constă în faptul că același student, va face parte în ani universitari diferiți din alte grupe și ani de studii, deci pentru a obține o situație reală a componentei structurale a studenților și a rezultatelor acestora, ei trebuie urmăriți în dinamica trecerii lor de la un an de studiu la altul. Acest lucru se realizează cu tabelele de dimensiuni *Student* și *Structură organizatorică* care permit încărcarea în tabelele de fapte a situației reale la un moment dat.

II.3 Schema conceptuală a depozitului de date

Modelele cele mai utilizate în faza de concepție a unui depozit de date sunt modelele dimensionale care regrupează datele din tabelele relaționale în scheme de tip: stea, fulg de zăpadă și constelație.

Modelul stea este cel mai comun model de date, în care depozitul de date conține un tabel central voluminos (tabelul de fapte) și un set de tabele însoțitor (tabelele dimensiune) pentru fiecare dimensiune. Tabelul de fapte cuprinde cea mai mare parte a datelor fără redundanțe. Graficul asociat seamănă cu o stea în care tabelele dimensiune sunt afișate radial în jurul tabelului de fapte central.

Modelul fulg de zăpadă este o variantă a modelului stea, unde o parte din tabelele dimensiune sunt normalizate, astfel datele sunt împărțite în tabele suplimentare. Rezultă o schemă reprezentată într-un graf similar unui fulg de zăpadă. Diferența majoră între modelul „fulg de zăpadă” și modelul stea este că tabelele dimensiune din modelul „fulg de zăpadă” pot fi păstrate în forma normalizată ceea ce determină o redundanță redusă. Ase-

menea tabele sunt ușor de întreținut și se economisește spațiu de stocare, deoarece un tabel dimensiune mare poate deveni enorm când structura dimensională este inclusă în coloane. Totuși această economie de spațiu este neglijabilă în comparație cu volumul foarte mare de date din tabelul de fapte. Mai mult, structura „fulg de zăpadă” poate reduce eficacitatea „browsing-ului” când mai multe „join-uri” trebuie executate la o interogare. De aceea, schema fulg de zăpadă este mai puțin răspândită.

Modelul constelație. Aplicațiile sofisticate pot solicita tabele multiple de fapte care partajează tabelele dimensiune. Acest gen de schemă poate fi văzută ca o colecție de stele

și, de aici, denumirea de schemă galaxie sau constelație de fapte (fact constellation). Pentru depozitele de date, schema constelație de fapte este în mod curent utilizată.

Depozitul de date pentru aplicația realizată are două tabele de fapte: *rezultate* și *structură studenți* care au ca dimensiuni comune (*zonă geografică*, *student* și *structură organizatorică*) de aceea schema depozitului de date va fi de tip constelație așa cum se poate observa în figura 4. Schema depozitului de date a fost obținută în urma unui proces de denormalizare a bazei de date tranzacționale care este tot o bază de date relațională, SQL Server. (figura 3)

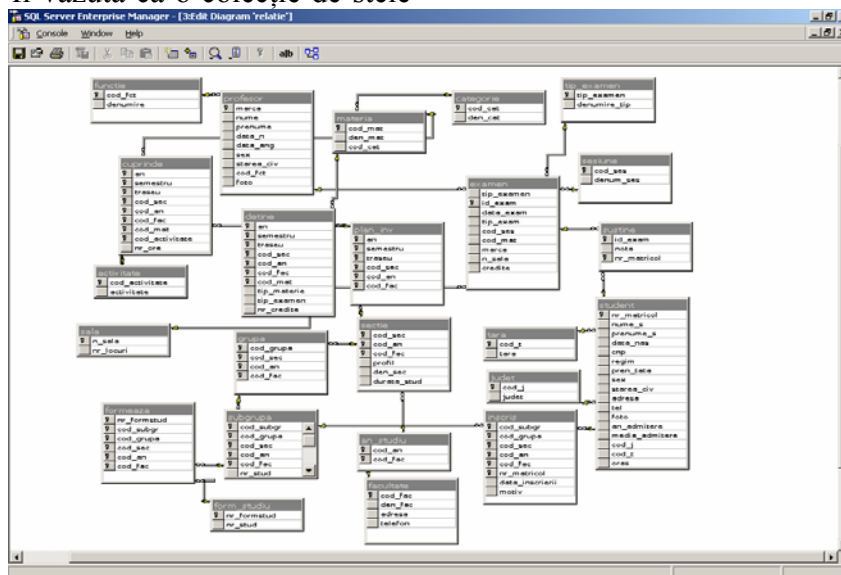


Fig. 3. Schema bazei de date tranzacționale

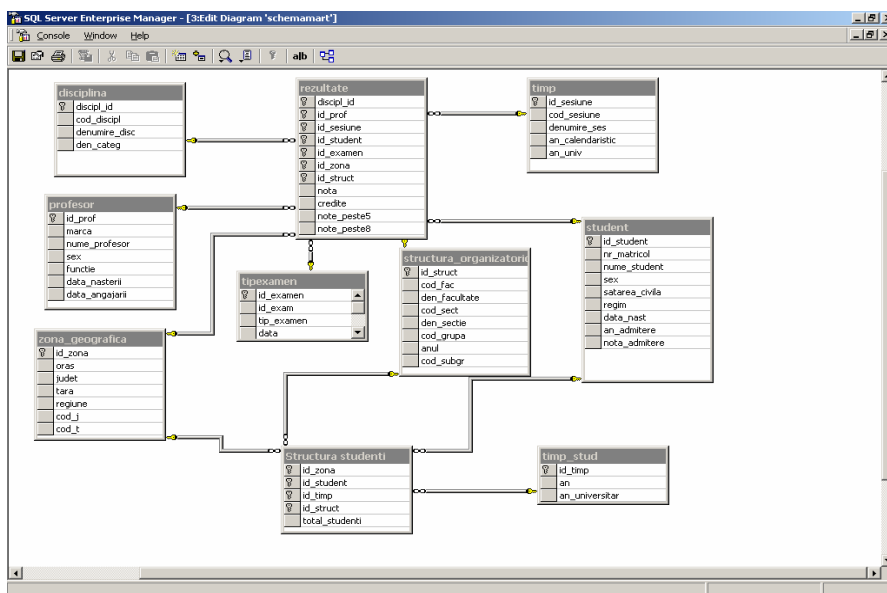


Fig. 4. Schema conceptuală a depozitului de date

În această primă parte am prezentat etapele de analiză și proiectare a depozitului de date, urmând ca în cea de-a doua parte a articolului să prezentăm popularea depozitului de date și analiza multidimensională a datelor utilizând cuburile de decizie.

Bibliografie

1. Albescu F., Bojan I. „Management information systems and decision support systems” Ed. DualTech, București 2001
2. Airinei D, „Sisteme informatice de asistare a deciziei”, curs on-line
3. Airinei D. „Depozite de date”, curs on-line
4. Bain T. &co, „Professiona SQLServer 2000 Datawarehousing with Analysis Services”, Wrox Press, Londra 2000
5. Connolly T., „Baze de date”, Teora, București 2001
6. Tanrikorur Tuturu, “Enterprise DSS architecture: a hybrid approach”, DM Review, February 1998
7. William McKnight, “Way a data warehouse?”, McKnight Associates, Inc, April 2002
8. Zaharie D & co, „Sisteme informatice pentru asistarea deciziei”, Ed. DualTech, București 2001
9. ***”Microsoft SQL Server 7.0 Data Warehousing Training Kit”, MS Press
10. ***Microsoft SQL Server 7.0 Data warehousing strategy
11. ***Microsoft SQL Server 7.0 OLAP Services
12. ***MS Press - MCSE Training Kit-SQL 7.0 Data Warehousing
13. <http://www.billinmon.com>
14. <http://www.datawarehouse.com>
15. <http://www.datawarehousingonline.com>
16. <http://www.dmreview.com>
17. <http://www.dw-institute.com>
18. <http://www.olapconcil.com>
19. <http://www.intelligententerprise.com>
20. <http://www.olapreports.co>