

Multidimensional Data Analysis Using OLAP Technology (2)

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 the 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.

În articolul anterior au fost prezentate etapele de analiză și proiectare a depozitului de date. În continuare vor fi prezentate popularea depozitului de date și analiza datelor utilizând cuburile de decizie.

III. Popularea depozitului de date

Popularea depozitului de date se face cu date preluate din sistemele tranzacționale care trebuie supuse unor procese de transformare pentru a se încadra în structura prefigurată a depozitului. Aceasta etapă va fi repetată periodic pentru a adăuga datele noi. SQL Server 7.0 oferă un instrument puternic pentru transformarea datelor, Data Transformation Services (DTS), care face această etapă mult mai ușoară. Încărcarea datelor în depozitul de date se face în câteva etape: extragerea datelor din sursele de date, transformarea datelor și popularea depozitului.

III.1 Planificarea transformărilor

Transformarea datelor implică câteva procese. Aceste procese trebuie proiectate pentru a asigura validitatea și integritatea, utilitatea și în anumite cazuri agregarea și calcularea unor noi valori bazate pe datele sursă. Aceste procese sunt clasificate și cunoscute ca: validarea, curățirea și transformarea datelor.

Validarea datelor asigură integritatea și validitatea. De exemplu, asigură faptul că orașul menționat se află în județul sau țara corespunzătoare. Validarea asigură de asemenea integritatea referențială a bazei de date în termeni de chei primare și străine.

Curățirea datelor analizează și integrează

datele provenite din surse diferite pentru ca în depozitul de date să ajungă o versiune unitară a acestora. Prin curățire se elimină discrepanțele sau conflictele între diferitele surse de date.

Transformarea datelor este procesul prin care datele extrase din sursele de date sunt convertite în formatul corespunzător încărcării lor în depozitul de date. Metadatele sunt în mod obișnuit utilizate pentru a memora caracteristicile mapărilor utilizate pentru a transforma datele sursă în Microsoft Repository. Metadatele definesc orice schimbare care este cerută înaintea încărcării datelor în depozitul de date. Transformarea datelor va înlătura anomaliile din sursele de date și va oferi date corespunzătoare pentru depozitul de date. Transformarea datelor poate interveni la nivel de înregistrare sau la nivel de câmp. Există trei modalități de bază utilizate în transformarea datelor: transformare structurală, transformare de conținut, transformare funcțională.

Transformarea structurală schimbă structura înregistrărilor sursă pentru a fi similară cu cea din baza de date destinație. Această modalitate transformă datele la nivel de înregistrare. Ceea ce înseamnă că întreaga înregistrare este transformată pentru ca structura acesteia să fie similară cu structura înregistrării din baza de date destinație.

Transformarea de conținut schimbă valorile înregistrărilor. Acest tip de transformare se aplică la nivel de câmp și schimbă valorile

câmpurilor utilizând algoritmi sau tabele pentru transformarea datelor (Lookup Table).

Transformarea funcțională creează noi valori pentru datele destinație pornind de la datele sursă. Această tehnică transformă datele la nivel de câmp. Aceste transformări pot fi folosite și în agregarea sau îmbogățirea datelor. Agregarea reprezintă calcularea unor valori derivate cum ar fi: totaluri sau medii ale valorilor câmpurilor din înregistrări diferite. Îmbogățirea combină două sau mai multe valori și creează unul sau mai multe noi atribute din una sau mai multe înregistrări care pot proveni din aceeași sursă sau surse diferite.

Microsoft SQL Server oferă câteva instrumente pentru toate operațiile prezentate mai sus. Aceste instrumente sunt: Data Transformation Services, Import/Export Wizard, programe sau script-uri ActiveX, DTS Lookup.

Transformarea datelor utilizând serviciul DTS implică în general planificarea și proiectarea transformărilor precum și crearea și executarea pachetului DTS ce va efectua aceste transformări în mod automat. DTS Designer permite definirea unor transformări complexe cu multe task-uri bazate pe constrângeri de precedență.

III.2 Pachetul DTS pentru popularea automată a depozitului de date

Pachetul realizat pentru popularea automată a depozitului de date are următoarele caracteristici:

✓ Are ca sursă o singură bază de date și anume baza de date SQL Server, *Disertatie.mdf*

✓ Destinația o reprezintă depozitul de date care este tot o bază de date de tip SQL Server, *Diseratiemart.mdf*

✓ Extrage datele din baza de date sursă, le transformă și le încarcă în baza de date destinație

Pentru realizarea transformărilor, pachetul DTS creat are următoarea structură:

Conexiuni de date. Pentru realizarea transformărilor în pachet au fost necesare următoarele conexiuni:

✓ Conexiuni de date destinație: sunt cele trei conexiuni *disertație OLAP* care fac conexiunea cu depozitul de date

diseratiemart.mdf. Prima conexiune începând de la stânga a fost creată pentru realizarea transformărilor necesare populării tabelor de dimensiunii, a doua pentru popularea tablei de fapte *rezultate* iar a treia pentru popularea tablei de fapte *structură_studenți*.

✓ Conexiuni de date sursă realizează conectarea la baza de date sursă, *disertatie.mdf*. Acestea sunt:

- DW_tipexamen – creată pentru sprijinirea populării tablei de dimensiuni *tip_examen*

- DW_profesor - creată pentru sprijinirea populării tablei de dimensiuni profesor

- DW_timp - creată pentru sprijinirea populării tablei de dimensiuni *timp*

- DW_disciplină - creată pentru sprijinirea populării tablei de dimensiuni *disciplina*

- DW_student - creată pentru sprijinirea populării tablei de dimensiuni *student*

- DW_zonă - creată pentru sprijinirea populării tablei de dimensiuni *zona_geografică*

- DW_structura studentii - creată pentru sprijinirea populării tablei de dimensiuni *structura student*

- DW_timp_stud - creată pentru sprijinirea populării tablei de dimensiuni *timp_stud*

- DW_rezultate - creată pentru sprijinirea populării tablei de fapte *rezultate*

- DW_structură_studenți - creată pentru sprijinirea populării tablei de fapte *structură_studenți*

Taskuri. În cadrul pachetului a fost folosit un singur task *actualizarea tabelor* de tipul Execute SQL Task, care actualizează (șterge înregistrările) tablele de dimensiuni și tablele de fapte înaintea începerii procesului de transformare și încărcare a depozitului de date.

Fluxuri (Workflows). Au fost necesare următoarele fluxuri:

✓ Fluxuri pentru transformarea datelor (Data Transformation).

- Fluxul DW_tipexamen - *disertatieOALP* – extrage datele necesare, le transformă și le încarcă în tabela de dimensiuni *tip_examen*

- DW_profesor - extrage datele necesare, le transformă și le încarcă în tabela de dimensiuni *profesor*

- DW_timp - extrage datele necesare, le

transformă și le încarcă în tabela de dimensiuni *timp*

- DW_disciplină - extrage datele necesare, le transformă și le încarcă în tabela de dimensiuni *disciplina*
 - DW_student - extrage datele necesare, le transformă și le încarcă în tabela de dimensiuni *student*
 - DW_zonă - extrage datele necesare, le transformă și le încarcă în tabela de dimensiuni *zona_geografică*
 - DW_rezultate - extrage datele necesare, le transformă și le încarcă în tabela de fapte *rezultate*
 - DW_structură_studenți - extrage datele necesare, le transformă și le încarcă în tabela de fapte *structură_studenți*
- ✓ **Constrângeri de precedență** au fost necesare următoarele constrângeri:

- Constrângerile de precedență On Succes între task-ul *actualizarea tabelelor* și conexiunile de date sursă pentru ca procesul de încărcare a tabelelor de dimensiuni sa nu înceapă decât după actualizarea cu succes a acestora.
- Constrângerea de precedență On succes între conexiunea destinație *disertație OLAP* și conexiunea sursă *DW_rezultate* pentru ca procesul de populare a tabelii de fapte sa nu înceapă până când nu se încheie cu succes popularea tabelelor de dimensiuni.
- Constrângerea de precedență On Succes între conexiunea destinație *disertație OLAP* și conexiunea sursă *DW_structură_studenți* pentru ca popularea tabelii de fapte *structură_studenți* să nu înceapă până când popularea tabelii de fapte *rezultate* nu se încheie cu succes.

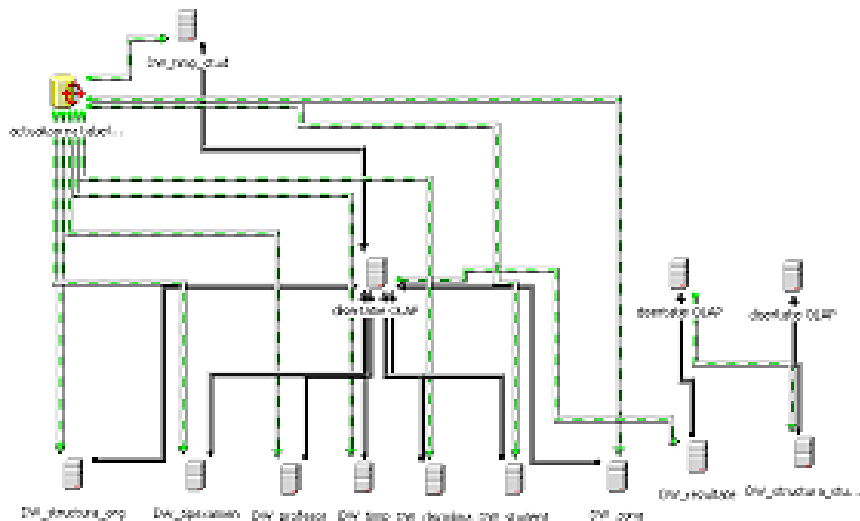


Fig.1. Pachetul DTS pentru popularea depozitului de date

IV Analiza datelor utilizând cuburile de decizie

IV.1 Proiectarea cuburilor: dimensiuni, măsuri, fapte

Proiectarea unui cub presupune identificarea și proiectarea pe rând a dimensiunilor, faptelor și a măsurilor.

Dimensiunile sunt elemente de bază în definirea unui cub. Cuburile de date sunt construite din dimensiuni. Ele sunt cele care oferă posibilitățile de vizualizare a datelor. O dimensiune poate conține un singur criteriu de vizualizare a datelor, de exemplu după stu-

dent, dar în mod obișnuit o dimensiune definește mai multe criterii organizate ierarhic, după care se pot vizualiza datele. Deci dimensiunile sunt cele care permit analiza „în adâncime” (drill-down) a datelor plecând de la nivel scăzut de detaliere și înaintând progresiv spre niveluri din ce în ce mai ridicate de detaliere.

Există două tipuri de dimensiuni care trebuie analizate atunci când se definesc acestea, și anume: dimensiuni private (Private Dimensions) și dimensiuni publice (Shared Dimensions). Dimensiunile Publice sunt

create independent de orice cub iar dimensiunile private sunt create odată cu cubul de decizie și sunt salvate în librăria cubului de decizie respectiv

Se optează pentru dimensiunile private în cazul în care dimensiunea respectivă implementează o logică de vizualizare a datelor valabilă doar în contextul cubului respectiv. În toate celelalte cazuri este bine să se definească dimensiuni publice care să asigure viziuni unitare și consistente asupra datelor.

Măsurile sunt valori numerice care dau măsura analizelor efectuate asupra datelor. Obligatoriu orice tabelă de fapte din depozitul de date trebuie să conțină măcar o măsură. Asupra măsurilor din tabela de fapte se pot aplica următoarele funcții de agregare: sum (efectuează suma elementelor), min (calculează valoarea minimă), max (calculează valoarea maximă), count (calculează numărul de elemente), avg (care calculează media aritmetică).

Aceste măsuri se pot defini atât pentru vizualizarea lor în cub cât și pentru obținerea unor membri calculați.

Membri calculați pot fi definiți pe baza datelor stocate deja în cub. Aceasta înseamnă că valorile membrilor calculați sunt calculate la cerere ele nefiind salvate fizic în cubul de decizie.

Agregarea datelor. Agregările sunt date precalculate care permit scăderea timpului de răspuns la interogări. Prin memorarea valorilor precalculate, server-ul nu trebuie să acceseze toate datele sursă necesare efectuării calculului respectiv. Utilizarea agregărilor este vitală pentru scăderea timpului de răspuns al sistemelor OLAP. În sistemele OLAP cuburile sunt cele care memorează datele agregate. Dimensiunile sunt coordonatele după care se interoghează cubul de decizie. Agregările sunt memorate, stocate la intersecția dintre dimensiuni. Fiecare intersecție, denumită celulă stochează o valoare agregată.

În mod normal agregările îmbunătățesc performanțele unui cub cu toate acestea la definirea lor trebuie avuți în vedere următorii factori:

✓ Spațiul pe disc. De la un anumit punct

spațiul pe disc ocupat suplimentar de către agregări nu se justifică în raport cu performanțele obținute. De aceea se determina și se creează mai întâi cele mai importante agregări și apoi în timp se adaugă noi agregări care devin necesare sau se pot elimina cele care sunt utilizate mai rar

Timpul necesar pentru procesarea cubului. Cu cât mai multe agregări sunt adăugate cubului cu atât timpul necesar procesării acestuia va fi mai mare

„Explozia datelor” (Data Explosion). Explozia volumului de date este o problemă atunci când se definesc prea multe agregări. Dacă toate agregările posibile ar fi calculate, spațiul necesar stocării cubului ar fi enorm. OLAP Services dispune de un algoritm avansat pentru determinarea optimului între performanță și spațiu necesar pentru stocare.

Ținând cont de considerațiile prezentate anterior, pentru analiza calității procesului de învățământ s-a creat cubul de analiză rezultate care are următoarele caracteristici:

Dimensiuni: pentru proiectarea cubului de decizie au fost identificate și definite următoarele dimensiuni:

Profesor cu următoarea ierarhie: grad didactic, sex, nume profesor;

Disciplină cu ierarhia: denumire categorie, denumire disciplină

Structură organizatorică: facultate, secție, an, grupa, subgrupa

Examen: tip examen, examen

Timp: an, sesiune

Regim: regim

Student: sex, stare civilă, student

Zona geografică: țara, regiune, județ, oraș

Ca **măsuri** au fost definite următoarele:

✓ *Nota, note peste 5, note peste 8, credite* care sunt preluate din depozitul de date

✓ *nr note și suma* utilizate pentru obținerea membrului calculat *media*

✓ *minim, maxim* utilizate pentru obținerea notelor minimă, respectiv maximă

✓ *note5 note8* utilizate pentru obținerea membrilor calculați *promovabilitate*, respectiv *rezultate peste 8*.

Ca schemă conceptuală a depozitului de date

a fost ales modelul de tip stea așa după cum se poate vedea în figura 2.

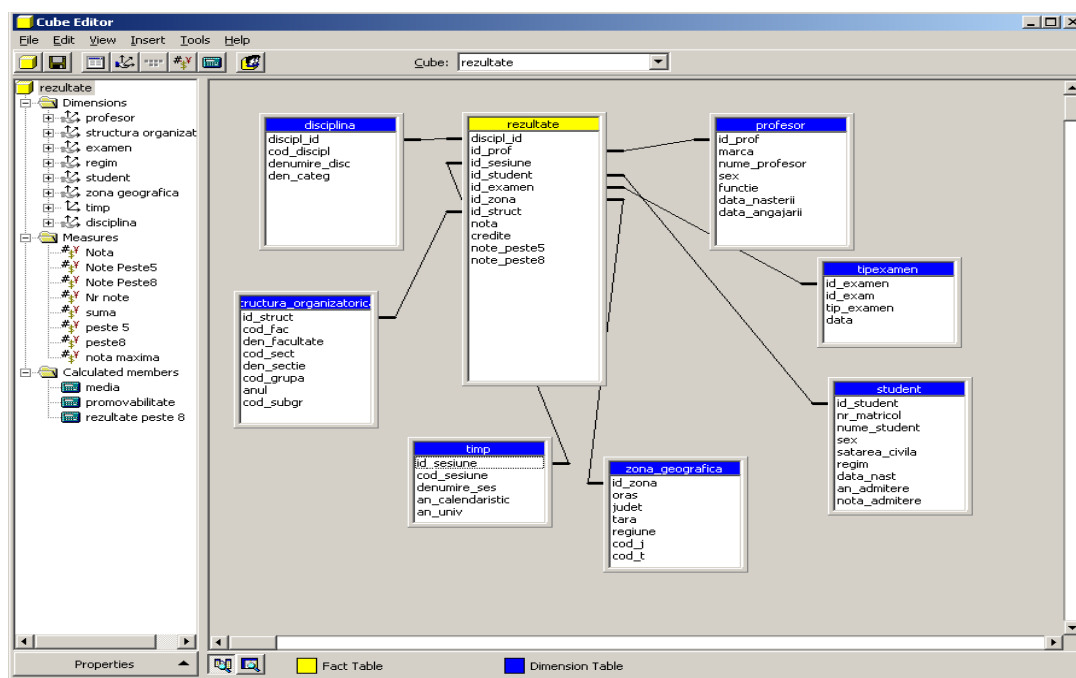


Fig 2. Schema conceptuală a cubului de date rezultate

Pentru analiza numărului și structurii studenților universității s-a creat cubul de analiză *structură studenți* care este definit de următoarele elemente:

Dimensiuni: pentru proiectarea cubului de decizie au fost identificate și definite următoarele dimensiuni:

✓ *Structură organizatorică cu următoarea ierarhie: facultate, secție, an, grupa, subgrupa*

✓ *Regim: regim*

✓ *Zona geografic cu ierarhia: țara, regiune, județ, oraș*

✓ *An: anul calendaristic și anul universitar*

✓ *Stare civilă: starea civilă*

✓ *Sex: sex*

Ca **măsuri** au fost definite următoarele:

✓ *Total studenți* calculează numărul total de studenți

Membri calculați:

✓ *Pondere fete* care calculează ponderea fetelor în numărul total de studenți

✓ *Pondere băieți* care calculează ponderea băieților în numărul total de studenți

✓ *Pondere buget* care calculează ponderea studenților la buget în numărul total de studenți

✓ *Pondere taxă* care calculează ponderea studenților cu taxă în numărul total de studenți

Ca schemă conceptuală a depozitului de date a fost ales de asemenea modelul de tip stea așa după cum se poate vedea în figura 3.

IV.2 Alegerea unei metode de stocare: MOLAP, ROLAP sau HOLAP

Datele din cuburile de analiză pot fi stocate în trei modalități: multidimensional, relațional sau hibrid. Modalitatea de stocarea trebuie aleasă cu grijă deoarece fiecare dintre acestea prezintă atât avantaje cât și dezavantaje.

MOLAP (Multidimensional OLAP) stochează atât datele cât și agregările în structuri multidimensionale, numite cuburi. Aceste structuri sunt deci stocate în afara depozitului de date. Datele stocate în cuburile OLAP au următoarele caracteristici:

✓ Toate măsurile sunt memorate în aceeași înregistrare ceea ce reduce timpul de acces la date

✓ Câmpurile cu valori nule nu sunt stocate, se elimină astfel problema „împrăștierii” (sparsity) datelor

- ✓ Row segments pot avea dimensiuni de 64 KB ce îmbunătățesc performanțele cubului
- ✓ Folosește indecși speciali, de tip bitmap

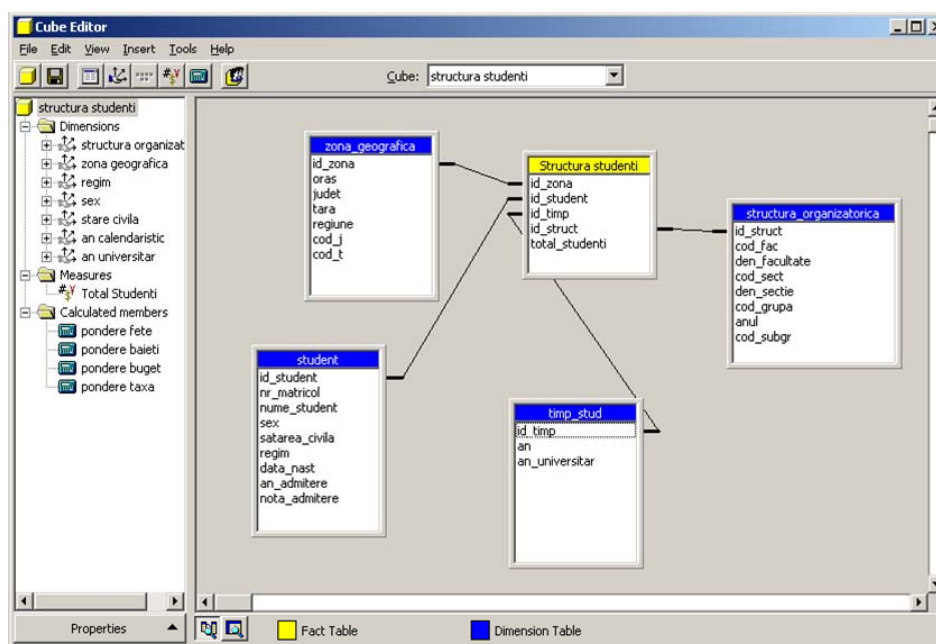


Fig.3. Schema conceptuală a cubului de date *statistici studenți*

Stocarea cuburilor în structuri MOLAP este cea mai potrivită pentru cele care sunt interogate frecvent și care necesită un timp de răspuns redus

ROLAP (Relational OLAP) stochează datele și agregările în tabele relaționale. Aceste tabele sunt memorate în aceeași bază de date în care sunt memorate și tabelele depozitului de date. Deoarece datele sunt stocate în depozitul de date, cuburile ROLAP nu necesită un alt spațiu de stocare. Având în vedere această modalitate de stocare, se pot utiliza și comenzi standard Transact-SQL în afara instrumentului OLAP Services, pentru interogarea cubului deși, numele tabelor și a coloanelor sunt greu de utilizat fiind generate automat de către sistem.

Structura datelor ROLAP este formată din tabele și indecși. Dimensiunile sunt stocate în tabele iar măsurile în coloane. Fiecare nivel al unei dimensiuni este indexat.

Nu se poate utiliza această modalitate de stocare a datelor dacă sursa de date pentru cub este de tipul OLE DB Provider for OLAP Services (adică este multidimensională) deoarece ROLAP necesită o bază de date relațională ca sursă.

Stocarea datelor cubului într-o structură ROLAP se recomandă pentru cuburile a căror date nu sunt interogate frecvent. De exemplu dacă 80% dintre utilizatori solicită datele doar din ultimul an în interogările lor, datele istorice pot fi mutate într-o structură ROLAP pentru a micșora spațiul general de stocare.

HOLAP (Hybrid OLAP) este o combinație între MOLAP și ROLAP. Astfel, HOLAP stochează datele în tabele relaționale în aceeași bază de date ca și depozitul de date iar agregările în structuri multidimensionale în afara depozitului de date.

Stocarea datelor cubului într-o structură de tip HOLAP este recomandată în cazul interogărilor frecvente asupra datelor agregate ce au la bază un volum mare de date. De exemplu, rezultatele anuale ale studenților pot fi stocate într-o structură MOLAP iar rezultatele fiecărei sesiuni pot fi stocate într-o structură ROLAP.

După alegerea modalității de stocare și procesarea cubului, datele devin disponibile și pot fi interogate. Desigur că operația de procesare a cubului trebuie repetată periodic pentru actualizarea datelor din cub. Vizuali-

zarea datelor se face cu instrumentul Cube Browser care permite interogarea și vizualizarea rapidă și facilă a rezultatelor interogărilor într-un format standard.

Concluzii

Depozitele de date sunt concepute special pentru sprijinirea luării deciziilor. Ele au ca obiectiv regruparea datelor, agregarea și sintetizarea lor, organizarea și coordonarea datelor provenind din surse diferite, integrarea și stocarea acestora pentru a da decidenților o imagine adecvată care să permită regăsirea și analiza eficace a informațiilor necesare. Interogările obișnuite într-un depozit de date sunt mai complexe și mai variate decât cele din sistemele de gestiune a bazelor de date. Ele se aplică asupra unor volume foarte mari de date și presupun calcule complexe (analiza tendinței, medii, dispersii etc.) care necesită adesea agregări (group by). Tehnologia OLAP este cea care „se descurcă” cu toți acești factori critici (volumul mare de date, complexitatea calculelor și timpul de răspuns), transformând volumul imens de date stocate și gestionate în depozite în informații utile, just-in-time procesului de decizie

Trebuie să subliniem că această aplicație nu se dorește un sistem deoarece la realizarea ei nu s-au luat în calcul factorii care țin de implementare și problemele ce deriva de aici. A fost creat mai mult pentru a demonstra aplicabilitatea unui astfel de instrument de analiza a datelor (OLAP) în domeniu. Limita este dată în primul rând de faptul că nu a fost dezvoltată partea de interfață cu utilizatorul iar modulul pentru popularea depozitului de date a fost creat luând în considerare doar datele din baza de date tranzacțională a unei singure facultăți (Facultatea de Științe Economice și Administrative). Stadiul actual de dezvoltare al sistemelor tranzacționale aferente facultăților este insular (nu au nici un fel de legătură unele cu altele) și primitiv. Sunt facultăți la care aplicațiile de evidență ale studenților folosesc ca structuri de date, fișiere independente. În aceste condiții crea-

rea și chiar proiectarea unui sistem de populare cu date din toate sistemele tranzacționale ar fi atins o întindere și o complexitate apreciabile și ar fi depășit preocupările autorului. Aplicația se vrea un punct de plecare în dezvoltări ulterioare dacă o astfel de soluție de genul depozit de date se consideră utilă.

Bibliografie

1. Albescu F., Bojian I. „Management information systems and decizion 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”, MS Press
11. ****”Microsoft SQL Server 7.0 OLAP Services”, MS Press
12. ****”MCSE Training Kit-SQL 7.0 Data Warehousing”, MS Press
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.com>
21. <http://www.bijonline.com>