

Data mining - de la sisteme de baze de date la sisteme cu baze de cunostinte

Prof.dr. Constanta BODEA, lect. Cristina IONITA,
asist. Anton CRETU, prep. Petrisor OPREA
Catedra de Informatica Economica, A.S.E. Bucuresti

Una dintre cele mai importante surse de achizitionare a cunostintelor o constituie, în prezent, bazele de date ale întreprinderii. Utilizarea instrumentelor de achizitionare si utilizare a cunostintelor în sistemele de baze de date a devenit o necesitate pentru toate organizatiile care folosesc un volum mare de date, colectate din diferite surse. Aceste instrumente au la baza o serie de metode, de la metodele statistice de interogare si raportare clasice pâna la metodele de învatare automata. Articolul prezinta câteva din aceste metode, punând în evidenta caracteristicile lor principale precum si avantajele si dezavantajele în utilizare. Pentru doua dintre aceste metode, analiza multidimensionala si calculul neuronal sunt prezentate rezultatele obtinute în cadrul unui studiu caz pentru domeniul vânzarilor.

Cuvinte cheie: achizitionarea cunostintelor, utilizarea cunostintelor, data mining, sisteme de baze de date multidimensionale, calcul neuronal.

1. Metode de achizitionare si utilizare a cunostintelor în sistemele de baze de date

Infrastructura informationala a organizatiilor este în continua crestere si tinde sa atinga un grad de maturitate caracterizat de un anumit nivel al calitatii datelor si al facilitatilor de acces la date. Pe masura ce infrastructura informationala se maturizeaza, creste necesarul de informatie relevanta, ceea ce reclama gasirea unor solutii de crestere a valorii datelor [3]. În forma în care sunt culese, datele au un nivel scazut de relevanta. Pentru a le utiliza, datele trebuie pregatite si prelucrate prin selectii, proiectii, reducerea dimensiunii, extragerea de pattern-uri si modele. Prin aceste operatii obtinem *cunostinte*, care sunt transmise managerilor pentru a fi folosite în procesul de luare a deciziilor sau sunt utilizate chiar de sistem pentru a ajunge la anumite concluzii relevante pentru manageri. Sistemele de baze de date (BD) tind sa devina sisteme de achizitionare si utilizare a cunostintelor (AUC).

Cei mai importanti pasi pentru realizarea AUC în domeniul BD sunt:

- a) regasirea si raportarea datelor relevante din BD;
- b) prelucrarea datelor (selectie, proiectie, reducerea dimensiunii);

- c) transformarea datelor (totalizare, normalizare);

- d) extragerea de pattern-uri si modele;

- e) utilizarea pattern-urilor si modelelor pre-determinate pentru a ajunge la cunostinte relevante (selectarea modelului, testarea si validarea modelului, dezvoltarea modelului). Pasii a)-d) sunt referiti, în general, drept etape ale procesului de achizitionare a cunostintelor, în timp ce pasul e) reprezinta procesul de utilizare a cunostintelor. Pasul d) este denumit si procesul de extragere a cunostintelor din date (*data mining*).

Cele mai raspândite metode pentru realizarea AUC în sistemele de BD sunt urmatoarele: metodele statistice si de raportare (metodele OLAP); metodele neuronale; metodele de învatare automata; metodele baze pe multimi aproximative (rough sets).

În figura 1 se prezinta trendul utilizarii acestor metode în sistemele de BD.

2. Stabilirea unei metrice pentru analiza comparativa a metodelor AUC

Evaluarea avantajelor utilizarii diferitelor metode de AUC poate fi realizata prin raportarea la cerintele domeniului de aplicabilitate si la resursele solicitate.

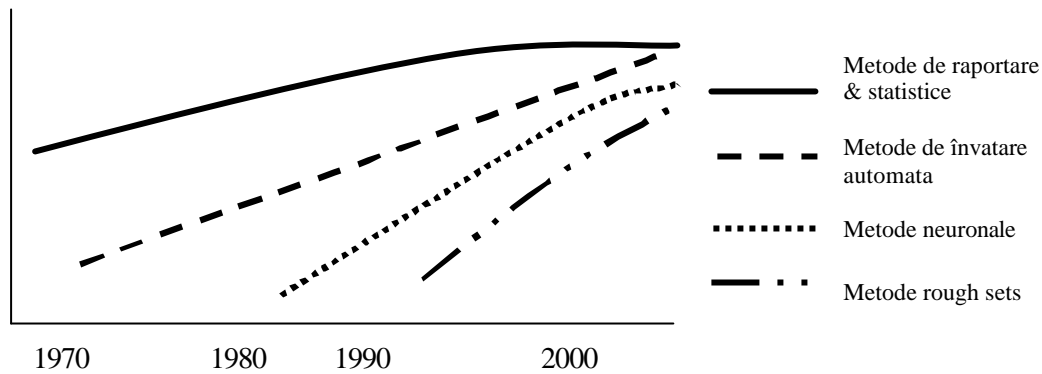


Fig. 1. Trendul utilizării metodelor AUC

În tabelul 1 sunt prezentate dimensiunile analizei comparative între principalele metode de AUC, dimensiuni grupate în patru clase: (1) calitatea rezultatelor și necesarul

de resurse, (2) cerințe de inginerie software, (3) restricții logistice, (4) aria de aplicabilitate și restricții de integrare [3].

Tabelul 1 - Analiza comparativă a principalelor metode de AUC

| Dimensiuni | Rezultate statistice și de raportare | Metode neuronale | Metode de învățare automată | Metode bazate pe mulțimi aproximative |
|-------------------------------------------------------------|--------------------------------------|------------------|-----------------------------|---------------------------------------|
| 1) calitatea rezultatelor și necesarul de resurse | | | | |
| acuratețe | medie | mare | medie | medie |
| explicabilitate | medie | mica | medie | medie |
| viteza / fiabilitate | mare | mare | mare | medie |
| toleranta la zgomote | mica | mare | medie | mica |
| toleranta la raritatea datelor | mica | mica | mica | medie |
| toleranta la complexitate | medie | mare | medie | mare |
| 2) cerințe de inginerie software | | | | |
| flexibilitate | mare | mare | mare | mare |
| scalabilitate | medie | medie | mare | medie |
| compactare | mica | mare | medie | mare |
| usurinta în utilizare | mare | medie | medie | medie |
| 3) restricții logistice | | | | |
| independența față de experți | medie | mare | medie | medie |
| usurinta în calcul | medie | medie | medie | mare |
| timp de dezvoltare | mare | mediu | mediu | mediu |
| 4) aria de aplicabilitate și restricții de integrare | | | | |
| pași procesului de AUC | a+b+c | c+d+e | c+d | c+d |
| variabilitatea domeniilor de utilizare | mare | mare | medie | medie |
| migrarea datelor | mare | medie | mica | medie |

Relevanța rezultatelor analizei comparative depinde de domeniul de aplicabilitate a metodelor, de instrumentele hardware și soft-

ware utilizate pentru implementarea metodelor.

Pentru a fundamenta analiza metodelor AUC în sistemele de BD am realizat un studiu de caz pentru un sistem tranzactional al vânzariilor. În acest sistem, datele despre vânzari sunt organizate într-o baza de date relationala (BDR) a carei structura este

redată în figura 2. Managerii au nevoie de o solutie informatica pentru descoperirea pattern-urilor de evolutie a vânzariilor si previzionarea vânzariilor pe diferite arii geografice, categorii de produse sau tipuri de clienti.

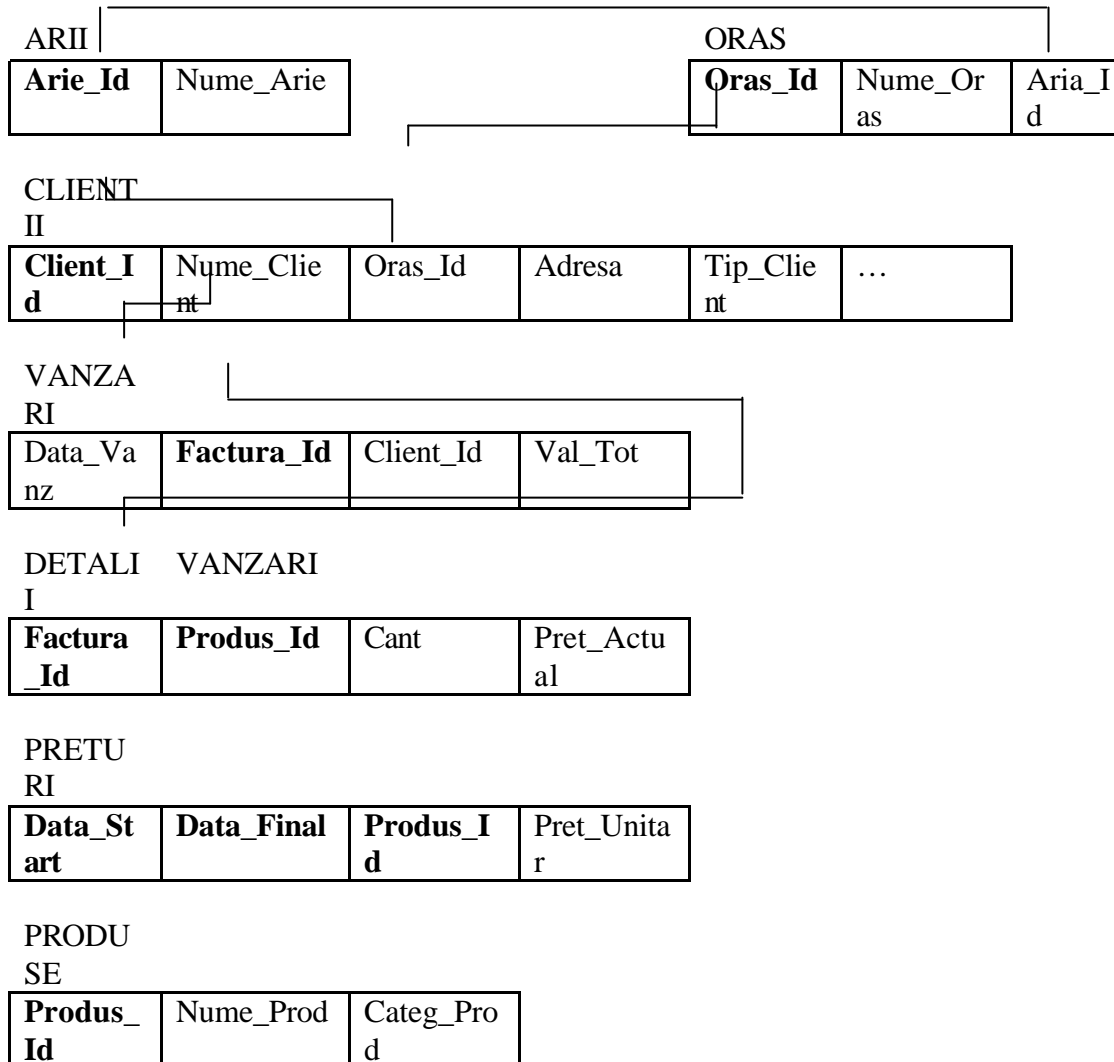


Fig. 2. Structura BDR pentru sistemul de vânzari

Pentru ilustrarea procesului de AUC din BDR pentru vânzari am ales doua metode: analiza multidimensionala si calculul neuronal. Urmatoarele sectiuni ale articolului prezinta rezultatele obtinute în analiza vânzariilor, utilizând ORACLE Express pentru analiza multidimensionala si toolbox-ul Neural Networks din MATLAB pentru calcul neuronal.

3. Analiza multidimensionala în sistemele de BD – Studiu de caz

Managerul de vânzari analizeaza activitatile economice în functie de factorii (“dimensionile”) care influenteaza evolutia activitatilor, cum ar fi: timpul, aria geografica, categoria de produs si tipul de client. Managerii nu gândesc în termenii relatiilor bidimensionale, plate, din BDR care presupun realizarea a numeroase jonctiuni pentru a obtine o imagine detaliata a vânzariilor care

sa surprinda toti factorii de influenta. Ei vizualizeaza datele despre vânzari sub forma de hipercub în care dimensiunile de analiza (timp, aria de vânzare, produsul, clientul) sunt muchiile hipercubului, iar datele despre vânzari (valoarea vânzariilor, cantitatea, pretul mediu) sunt celulele din interiorul hipercubului [9]. O solutie eficienta de stocare si prelucrare a hipercuburilor de date este oferita de bazele de date multidimensionale (BDM).

Pentru studiul de caz privind analiza vânzariilor s-a proiectat o BDM a carei structura este prezentata în figura 3. BDM are o tabela centrala de fapte care stocheaza informatiile privind valoarea vânzariilor, cantitatea si pretul mediu si mai multe tabele dimensiune pentru a stoca informatiile privind factorii care influenteaza evolutia vânzariilor. Unele din tabelele de dimensiune au o structura ierarhica (TIMP, PRODUS) iar altele sunt ierarhic legate între ele (CLIENT, ORAS, ARIA) pentru a putea realiza agregarea si

descompunerea datelor în raport cu aceste dimensiuni.

Pentru a ilustra tehnicile de analiza OLAP pe o BDM am ales instrumentele din suita ORACLE Express, care ofera facilitati de analize statistice si raportare avansate (analize pe serii de timp, regresii, simulari, prognoze).

În procesul de descoperire a cunostintelor deosebit de eficiente sunt facilitatile oferite utilizatorilor de a sectiona hipercuburile de date si respectiv de a agrega datele din hipercuburi pe o singura dimensiune sau pe mai multe, operatii cunoscute sub numele de *drill down* si respectiv de *roll-up*. Astfel, managerul poate începe analiza cu vizualizarea agregata a datelor. Daca observa un fenomen si doreste sa descopere factorii care au influentat producerea acestuia, poate începe un proces de sectionare succesiva a datelor pe o singura dimensiune sau corelat pe mai multe dimensiuni.

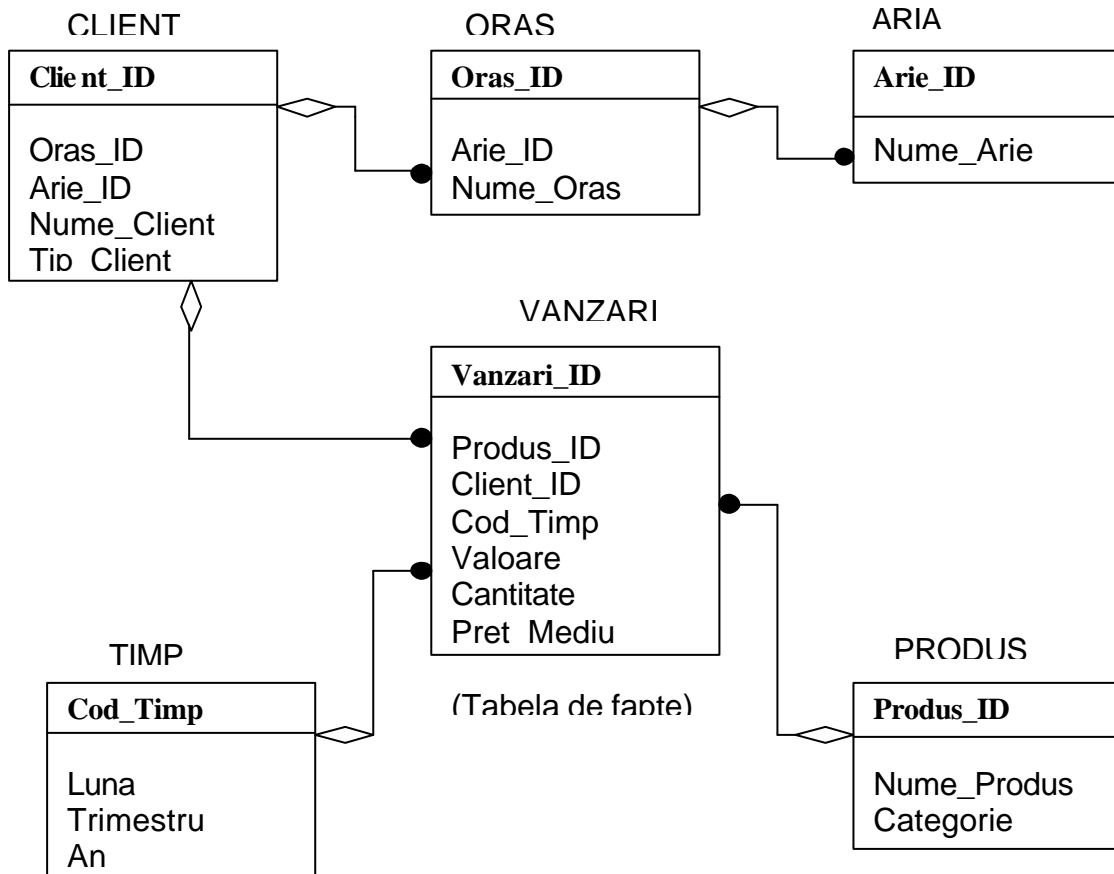


Fig. 3. Structura BDM pentru vânzari

De exemplu, analizând valoarea agregata a vânzarilor pe fiecare trimestru din ultimii cinci ani (1994 -1998) din figura 4 managerul poate sesiza vizual ca, desi tendinta generala este de crestere a vânzarilor, în trimestrul 1 si trimestrul 4 vânzarile sunt mai scazute. Sectionând datele privind

valoarea trimestriala a vânzarilor pe categorii de produse (figura 5) managerul va sesiza ca produsele din categoria 3 prezinta o sezonabilitate a vânzarilor, înregistrând o scadere a vânzarilor mai mica în trimestrul 4 si mai accentuata în trimestrul 1.

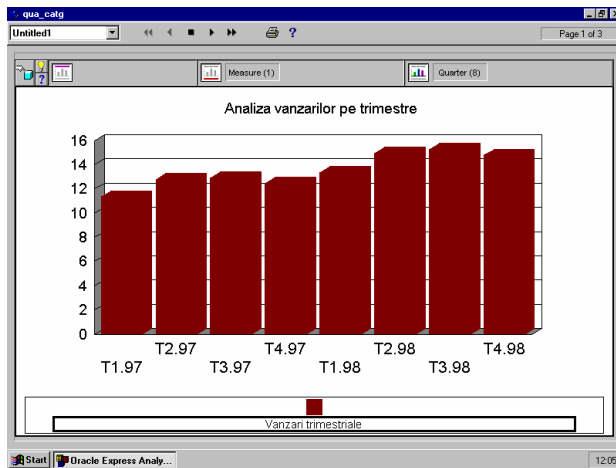


Fig. 4. Analiza vânzarilor pe trimestre

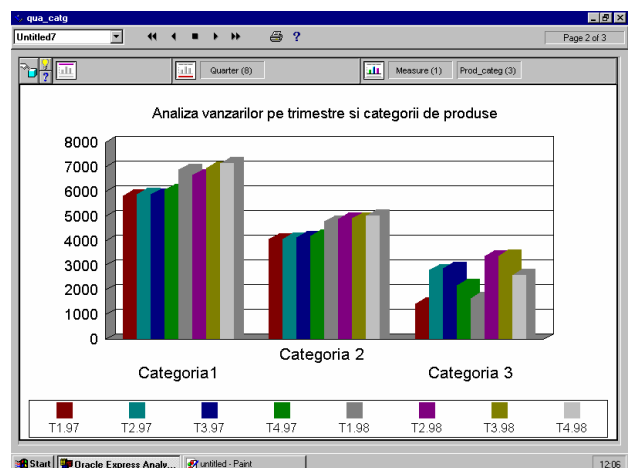


Fig. 5. Analiza vânzarilor pe trimestre si pe categorii de produse

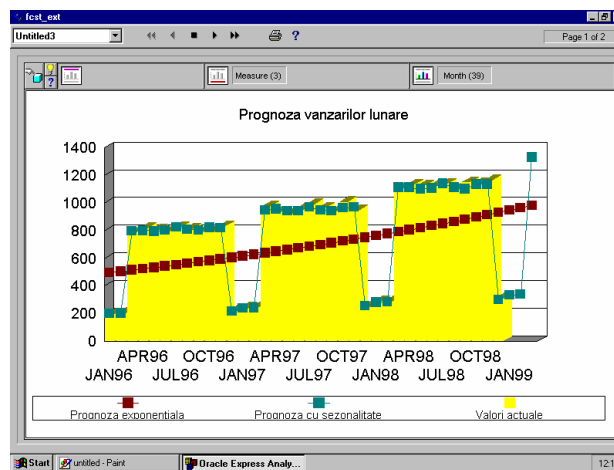


Fig. 6. Rezultatele procesului de prognoza pentru produsele din categoria 3

Proгноza datelor în Express presupune mai întâi identificarea vizuala a trendului de evolutie si apoi alegerea metodei de prognoza: liniara, exponentiala sau cu sezonabilitate. Pentru prognoza vânzarilor pe primele trei luni din anul 1999, pe baza datelor privind vânzarile lunare din perioada 1996-1998 am folosit doua metode de prognoza:

exponentiala si cu sezonabilitate. Rezultatele celor doua metode de prognoza sunt prezentate comparativ în graficul din figura 6 numai pentru produsele din categoria 3. Se observa ca prognoza cu sezonabilitate de tip Holt-Winters aproximeaza datele cu o exactitate mai mare. Rezultatele valorice ale prognozei sunt prezentate în tabelul 3.

4. Calcul neuronal în sisteme de baze de date – studiu de caz

Scopul propus este de a determina funcțiile care aproximează trendul vânzătorilor pe cele 3 categorii analizate de produse, în vederea prognozarilor vânzătorilor lunare.

Pentru realizarea acestui scop am definit câte o arhitectura de rețea neuronală feed-forward multinivel pentru fiecare din cele 3 categorii de produse, folosind algoritmul faster-backpropagation pentru antrenarea acestor rețele. Fiecare rețea este constituită din:

- un nivel de intrare, cu 5 unități; 3 unități de intrare semnifică vânzătorii pentru categoria analizată de produse în ultimele 3 luni, iar 2 unități de intrare semnifică vânzătorii pentru categoria analizată de produse în aceeași lună a ultimilor 2 ani;
- un nivel de ieșire, cu un element, semnificând nivelul vânzătorilor pentru categoria analizată de produse în luna curentă;
- un nivel intermediar (ascuns), conținând 5 elemente de procesare. Numărul neuronilor de pe nivelul ascuns a rezultat în urma mai multor sesiuni de antrenare, luând în considerare timpul necesar antrenării, precum și o aproximare cât mai exactă a trendului.

Pentru a realiza simularea evoluției vânzătorilor pentru categoriile analizate de produse am utilizat toolbox-ul Neural Networks din MATLAB. Datele necesare au fost exportate din baza de date în fișiere text, de unde au fost importate în MATLAB.

În faza de antrenare am utilizat pentru fiecare rețea 36 seturi de antrenare, care conțin ca outputuri dorite vânzătorii pe 36 luni (din ianuarie 1996 până în decembrie 1998) plus cele 5 inputuri cu semnificația menționată anterior. Funcția de activare a neuronilor este de tip logistica sigmoidă, atât pentru neuronii de pe nivelul ascuns, cât și pentru neuronii de pe nivelul de ieșire. Pentru a realiza o aproximare cât mai bună și mai rapidă a trendului, datele din seturile de antrenare (atât inputurile cât și outputurile) au fost normalizate astfel:

$$X_{n_{ij}} = (X_{ij} - X_{\min_j}) / (X_{\max_j} - X_{\min_j}), \quad i=1, Q; \quad j=1, R+S$$

unde: Q este numărul de seturi de antrenare; R este numărul de inputuri; S este numărul de outputuri; X_{ij} reprezintă datele preluate din seturile de antrenare; X_{\max_j} , X_{\min_j} reprezintă marginile superioare și inferioare (determinate pe fiecare input și output); $X_{n_{ij}}$ reprezintă datele normalizate.

Pe parcursul antrenării se urmărește minimizarea erorii globale, calculată ca suma a erorilor (diferențelor) între outputurile dorite și outputurile obținute din propagarea inputurilor în rețea. Cele 3 rețele au fost antrenate până când eroarea a scăzut sub o limită specificată. Numărul epocilor de instruire și eroarea globală pentru fiecare din cele 3 rețele sunt prezentate în tabelul 2.

Tabelul 2.

| Retea/Categorii de produs | Numar epoci | Eroare globala |
|---------------------------|-------------|----------------|
| Categoria 1 | 160000 | 0.0067 |
| Categoria 2 | 160000 | 0.0064 |
| Categoria 3 | 180000 | 0.0057 |

Evoluția erorii pentru cea de a 3-a categorie de produse este prezentată în figura 7, iar funcția care aproximează trendul vânzătorilor pentru categoria a 3-a de produse (rezultată în urma antrenării rețelei) este prezentată în figura 8.

În faza de simulare, pentru fiecare categorie de produse am utilizat ca seturi de test inputurile din seturile de antrenare, pentru a verifica dacă aproximarea este corectă. După aceasta am realizat câteva prognoze ale vânzătorilor lunare, începând cu ianuarie 1999. Câteva din rezultatele simularilor sunt prezentate în tabelul 3.

5. Concluziile analizei comparative

Modelele multidimensionale de organizare a datelor oferă un suport eficient pentru realizarea prelucrărilor analitice de tip OLAP și pentru fundamentarea deciziilor. Sistemele OLAP ce implementează BDM sunt orientate către decidenți și permit acestora

sa realizeze cu usurinta sectionari si grupari ale datelor supuse analizei si sa vizualizeze datele din diferite perspective pentru a descoperi factorii (dimensiunile) care influenteaza un anumit comportament. Sistemele OLAP ofera facilitati de analiza statistica avansate cum ar fi analize pe serii de timp,

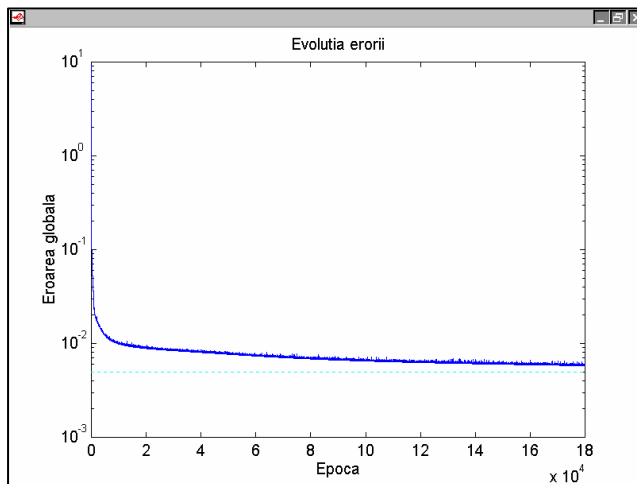


Fig. 7. Evolutia erorii la antrenarea rețelei neuronale pentru categoria 3

Calculul neuronal prezinta o serie de avantaje, precum: identificarea unor functii de aproximare foarte complexe, acuratete sporita a rezultatelor, o mare varietate de input-uri si output-uri, independenta fata de experti.

Comparând aceste doua metode de AUC în raport cu clasa de cerinte "calitatea rezultatelor si necesarul de resurse" prezentata generic în tabelul 1 am ajuns la urmatoarele concluzii:

1. *Acuratetea rezultatelor.* În analiza multidimensionala acuratetea depinde de abilitatea utilizatorilor de a selecta si utiliza corect modelele predefinite de analiza statistica, pe baza evaluarii vizuale a caracteristicilor de comportament ale faptelor. În multe din cazuri nu exista un model pre-

regresii, simulari, prognoze care au la baza modele predefinite. Instrumentele de tip OLAP precum ORACLE Express sunt rapide, usor de utilizat si gestioneaza volume mari de date de ordinul megabytes si gigabytes într-un timp relativ scurt.

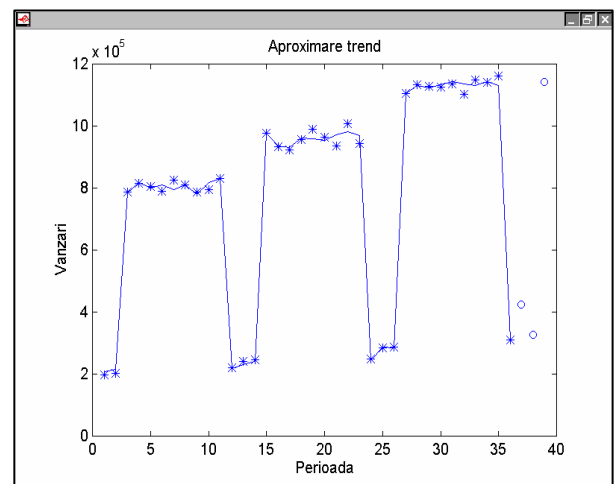


Fig. 8. Aproximarea trendului cu ajutorul rețelei neuronale pentru categoria 3

definit corespunzator, iar utilizatorii trebuie sa aleaga unul inadecvat. În studiul de caz prezentat, utilizarea unei prognoze exponentiale pentru valoarea vânzariilor genereaza o eroare inacceptabila pentru produsele din categoria 3, în timp ce prognoza cu sezonalitate asigura o eroare acceptabila (tabelul 3). În calculul neuronal, acuratetea rezultatelor depinde de datele de intrare si de calitatea proiectarii rețelei neuronale. O retea neuronală bine proiectata poate identifica si aproxima orice caracteristica a datelor de intrare, astfel încât acuratetea va fi foarte mare. În tabelul 3 se prezinta rezultatele valorice, care confirma acuratetea mare pentru calculul neuronal în domeniul vânzariilor fata de analiza multidimensionala.

Tabelul 3: Analiza comparativa a rezultatelor valorice ale prognozei

| Categoria de produs | Luna | Valoarea actuala | Express | | Matlab Valori prognozate |
|----------------------------|-------|------------------|--------------------------|--------------------------------|-----------------------------|
| | | | Prognoza exponentiala | Prognoza cu sezonalitate | |
| 1 | IAN98 | 229657 6 | 2133153 | 2279460.38 | 2296419 |
| | FEB98 | 227994 5 | 2160568.75 | 2332578.49 | 2282065 |
| | MAR98 | 232228 3 | 2188336.86 | 2281832.24 | 2332941 |
| | ... | | ... | ... | ... |
| | IAN99 | - | 2486423.10 | 2619696.72 | 2562881 |
| | FEB99 | - | 2518389.17 | 2655538.02 | 2514373 |
| | MAR99 | - | 2550745.94 | 2654901.27 | 2599749 |
| Procentul de eroare | | | 2.84 | 1.00 | 0.76 |
| 2 | IAN98 | 160395 6 | 1503090.74 | 1636561.04 | 1604627 |
| | FEB98 | 155445 0 | 1521920.48 | 1566818.46 | 1548870 |
| | MAR98 | 164008 2 | 1540986.10 | 1595496.05 | 1641169 |
| | ... | | ... | ... | ... |
| | IAN99 | - | 1745284.92 | 1940963.41 | 1826476 |
| | FEB99 | - | 1767148.71 | 1873342.82 | 1710930 |
| | MAR99 | - | 1789286.40 | 1915427.07 | 1753852 |
| Procentul de eroare | | | 3.07 | 1.25 | 0.69 |
| 3 | IAN98 | 284998 | 764641 | 282066.16 | 286905 |
| | FEB98 | 286231 | 778461 | 285667.17 | 285212 |
| | MAR98 | 110590 7 | 792530.93 | 1113696.34 | 1104408 |
| | ... | | ... | ... | ... |
| | IAN99 | - | 948002.36 | 335888.51 | 424112 |
| | FEB99 | - | 965136.50 | 340898.56 | 325655 |
| | MAR99 | - | 982580.31 | 1332685.64 | 1142169 |
| Procentul de eroare | | | 65.59 | 1.64 | 1.57 |

2. *Explicabilitatea proceselor de prelucrare necesare pentru obtinerea rezultatelor.* În analiza multidimensională majoritatea prelucrarilor tin doar de competența utilizatorilor, astfel încât explicabilitatea este evidentă. În calculul neuronal explicabilitatea este scăzută din cauza specificității schemei de reprezentare a cunosțințelor și a gradului mare de integrare a pașilor de prelucrare. Pentru a crește gradul

de explicabilitate în calculul neuronal, am extras regulile din rețelele neuronale, după faza de învățare, aplicând câteva metode uzuale ([1], [6]), dar rezultatele au fost neconcludente.

3. *Viteza de raspuns.* În analiza multidimensională, selectarea modelului de previziune și calcularea valorilor previzionate s-au realizat rapid. În abordarea neuronală, faza de antrenare a rețelei a fost de durată

mai mare (tabelul 2), dar apoi valorile previzionate au fost rapid calculate.

4. *Toleranta la zgomote.* Metodele de previziune predefinite în ORACLE Express

au tratat foarte bine aparitia zgomotului în date (tabelul 4). În faza de antrenare, o retea neuronală pentru prognoza nu poate identifica zgomotul în datele reale.

Tabelul 4 - Tratarea zgomotului în prognoza valorilor

| Luna | Valoarea actuala | Prognoza exponentiala | Prognoza cu sezonalitate |
|----------------------------|------------------|-----------------------|--------------------------|
| NOE96 | 1,708,483 | 1,913,989.51 | 1,712,180.14 |
| DEC96 | 1,759,802 | 1,933,564.86 | 1,747,569.26 |
| IAN96 (zgomot) | 11,831,74 6 | 1,953,340.42 | 4,435,909.49 |
| FEB96 | 1,969,470 | 1,973,318.23 | 4,868,775.84 |
| MAR96 | 1,975,974 | 1,993,500.36 | 2,419,024.36 |
| APR96 | 2,008,928 | 2,013,888.90 | 2,020,737.69 |
| Procentul de eroare | | 8.06% | 16.69% |

5. *Toleranta la complexitate.* Modelele de analiza implementate în ORACLE Express pot manipula datele cu trend si sezonalitate. Pentru dinamici mai complexe de evolutie aceste metode sunt improprii. Abordarea neuronală poate trata orice tip de dinamica. În concluzie, rezultatele obtinute în studiul de caz confirma aprecierile generale facute în analiza comparativa a principalelor metode de AUC.

Bibliografie

- [1] Bodea C. – *Inteligenta artificiala si sisteme expert*, Editura Infocrec, Bucuresti, 1998.
- [2] Brand E, Gerritsen R. – *Data Mining and Knowledge Discovery*, in: DBMS, Data Mining Solutions Supplement, November, 1998, <http://www.dbmsmag.com/9807m01.html>.
- [3] Dhar V., Stein R. – *Seven Methods for Transforming Corporate Data into Business Intelligence*, Prentice Hall, 1997.

[4] Fayyad U.M., Piatetsky-Shapiro, Smyth P., Uthurusamy R. – *Advances in Knowledge Discovery and Data Mining*, AAAI Series, 1998.

[5] Information Discovery Inc. – *Perspective on Data Mining*, <http://www.datamining.com/datamine/dm-ka.htm>

[6] Kasabov N.(ed) – *Brain-like Computing and Intelligent Information Systems*, Springer-Verlag, 1998.

[7] Munakata T. – *Fundamentals of the New Artificial Intelligence*, Springer-Verlag, 1998

[8] *** - *MATLAB. High Performance Numeric Computation and Visualisation Software*, Reference Guide, The MathWorks Inc., Natick, Massachusetts, 1992-1998.

[9] *** - *ORACLE Data Warehousing*, SAMS Publishing, 1997

[10] *** - *Express Analyzer*, Oracle Publishing, 1998