

ASISTAREA INFORMATICĂ A DECIZIEI PRIN DEMERSUL IDENTIFICĂRII DE CUNOȘTINȚE ÎN BAZELE DE DATE - KDD (KNOWLEDGE DISCOVERY IN DATABASES)

PROF.UNIV.DR. IONESCU BOGDAN

Academia de Studii Economice Bucuresti, Adresa de corespondenta: Str 9 Mai nr. 5, bl. 36, sc.B, et. 2, ap. 25, sector 6, București, tel/fax 0212213446, mobil 0722 593695, e-mail bionescu@gmail.com

CONF.UNIV.DR. MIHAI FLORIN

Academia de Studii Economice Bucuresti, Adresa de corespondenta: str. Vintila Mihailescu 16, Bloc 70, ap 11, sector 6, Bucuresti, tel 0724255333, email fmihai@gmail.com

L'assistance informatique de la décision par la démarche de découverte des connaissances dans les bases des données.

La démarche KDD (découverte des connaissances contenues dans les bases des données) représente un domaine de la technologie d'information qui combine les approches des systèmes informatiques décisionnels, de reconnaître des modèles, d'analyse des données, de présentation de l'information, pour extraire automatiquement les relations, des concepts et des connaissances contenues dans les grandes bases des données. Le processus KDD met en évidence et interprète les schémas des connaissances existantes dans les données.

KDD este un domeniu emergent care combină tehnicile “mașinilor de învățare”, de recunoaștere a modelelor, de statistică, de vizualizare, pentru a extrage în mod automat inter-relațiile de concepte și module, conținute în marile baze de date. Sarcina de bază a tehnologiilor KDD este deci aceea de a extrage cunoștințe, plecând de la nivelurile cele mai de jos ale bazelor de date. KDD interpretează schemele de cunoștințe existente în bazele de date.

Figura următoare furnizează o idee despre locul acestor ansambluri față de o bază de date.

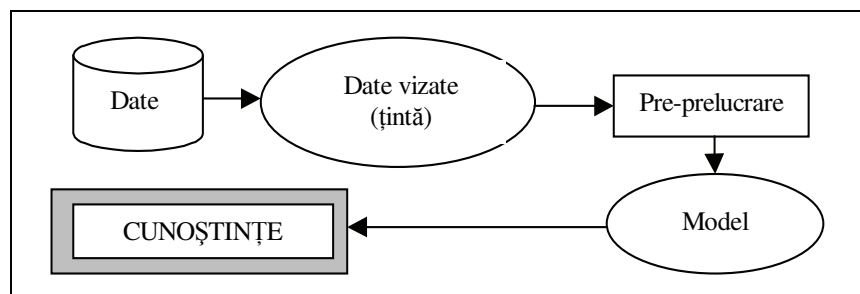


Figura 1. Legătura dintre date, modele, cunoștințe

Cunoștințele extrase sunt utilizate ca instrumente de asistare a deciziei, de exemplu, pentru operațiile de previziune și clasificare, pentru a rezuma conținutul unei baze de date sau pentru explicarea fenomenelor și tendințelor observate. Utilizarea sistemelor KDD permite degrevarea decidenților de problema analizei manuale a marilor cantități de date din

bazele de date. Aceste sisteme se dezvoltă și actualmente încep a fi implementate în întreprinderi, în mod particular, în domeniul finanțelor, al diagnosticului, al analizei piețelor, al detectării fraudelor, etc.

Facilitățile oferite de KDD devin extrem de importante în condițiile creșterii continue a masei de date, a informațiilor, a cunoștințelor disponibile. Astfel, capacitatea și rapiditatea procesului de introducere a datelor crește tot mai mult datorită noilor tehnologii, în timp ce capacitatea de analiză a acestor date și informații nu se dezvoltă în aceeași proporție. Deci, se produce în mod inevitabil o sub-optimizare a ansamblului de resurse informaționale. KDD își aduce o contribuție majoră la soluționarea acestei contradicții pentru optimizarea sistemului informațional.

Expresia KDD poate admite două interpretări: “KD in Data” și “KD in Databases”. Pentru a surmonta dificultatea de a gestiona un imens volum de informații, există două abordări posibile:

- reimplementarea algoritmilor de învățare pentru a face operaționale bazele de date, în scopul manipulării unui volum mare de date;
- reducerea volumului de date și de atribute.

Acest ultim punct implică o prelucrare anterioară a datelor și arată importanța descoperirii de cunoștințe în date (KD in Data). În acest caz, pot interveni elemente statistice, precum și alte tehnici cum ar fi, avizul experților, existența și calitatea cunoștințelor în domeniu.

Pentru a avea maximum de beneficii din tehnologia KDD, o strategie mai completă este aceea de a utiliza cu preponderență descoperiri în date (KD in Data) și nu în bazele de date (in Databases). Dezvoltarea abordărilor specifice KDD necesită încă timp pentru o cercetare aprofundată. Astfel de abordări, vor constitui o provocare pentru statisticieni, cercetători și pentru informaticienii specialiști în domeniul bazelor de date. Procesul de pre-lucrare este un instrument pentru abordări KDD aplicative, în particular, reducând cantitatea de date și atributele aferente bazei de date.

Nu există o metodologie generală a tehnologiei KDD, aplicabilă tuturor cazurilor și domeniilor, dar există o anumită experiență ce poate permite definirea unui oarecare număr de etape generale, pentru a lăsa o marjă de manevră în funcție de condițiile în care se aplică procedura de descoperire de cunoștințe.

Aplicațiile de descoperire de cunoștințe sunt aplicații ce posedă următoarele *caracteristici*:

- Cunoștințele codate, ascunse în aplicații, sunt în sens larg extrase din datele întreprinderii. Aceste cunoștințe sunt în general validate și îmbogățite prin avizul dat de experții în domeniu;
- Rezultatul produs de către aplicații este re-lucrabil. Acest lucru este posibil numai dacă există un vocabular al întreprinderii în măsură să “înțeleagă” și să poată fi utilizat direct, fără o altă transformare pentru sarcinile de gestiune;
- Cunoștințele codate (ascunse) în aplicații, pot fi menținute și pot fi perfecționate de utilizatorii întreprinderii, de o așa manieră încât acestea pot fi modificate fără ca aplicația să fie recompilată.

Aplicațiile ce au ca obiect descoperirea de cunoștințe în date sau în bazele de date sunt concepute într-o primă abordare pentru utilizatorii informației întreprinderii. Aceste

aplicații trebuie să încurajeze utilizarea decizională a informației și, mai ales, trebuie să explicitizeze reprezentările cunoștințelor înconjurătoare.

Iată câteva probleme generale, ce trebuie luate în considerare în tehnologia KDD:

a) Construirea unui Data Warehouse

Se cunoaște faptul că în majoritatea întreprinderilor, datele informatice se găsesc dispersate în diferite sisteme operaționale. Aceste sisteme au fost concepute pentru a prelucra și actualiza date, dar această situație ce poartă amprenta eterogenității, nu este practic cea mai bună soluție pentru aplicațiile KDD. Pentru a dezvolta aplicațiile de “descoperire de cunoștințe”, o întreprindere trebuie să-și transpună datele existente din sistemele informatice operaționale într-un model de date unic, coerent și logic care să acopere ansamblul necesităților informaționale. Acest sistem consolidat constituie un veritabil Data Warehouse. A construi un Data Warehouse reprezintă o etapă de pre-prelucrare necesară dezvoltării aplicațiilor ce au ca obiect “descoperirea de cunoștințe”.

b) Descoperirea de cunoștințe

Această activitate începe cu înțelegerea sarcinilor de îndeplinit pentru a obține rezultatele așteptate și cu înțelegerea universului organizației în sensul statistic al termenului, plecând de la datele care “populează”, un Data Warehouse. Bazându-se pe această înțelegere, instrumentele de clasificare și vizualizare pot fi utile pentru a segmenta acest univers în subansambluri. Acest proces de segmentare este urmat de selectarea variabilelor independente pentru a construi și dezvolta un model, după care se pune problema alegerii metodologiilor de modelare. Pentru acest lucru, trebuie să se țină cont de anumiți factori cum ar fi: tipul datelor prelucrate și distribuția variabilelor independente. Acești factori trebuie să garanteze punerea în aplicare a acestor metodologii de modelare. La sfârșit, se dezvoltă un model, se testează și se evaluează funcționalitatea sa.

c) Proiectarea unei “ieșiri reprelucrabile”

Această activitate este una de post-prelucrare, ce vine să răspundă la următoarea întrebare: “care vor fi rezultatele procesului de descoperire de cunoștințe?”. Printre rezultatele reprelucrabile, ar trebui să se găsească “cadre” de informație (utile aplicațiilor de descoperire de cunoștințe pentru analize ulterioare), motoare ale bazelor de date active (utile pentru generarea rapoartelor asupra operațiunilor cheie ale întreprinderii) și agenți (utili pentru implementarea normală a operațiunilor comerciale).

d) Actualizarea procesului de descoperire de cunoștințe

Pentru a realiza menținerea în exploatare a tehnologiei KDD, precum și a criteriului de perfecționare a aplicațiilor ce au la bază descoperirea de cunoștințe, se remarcă aplicațiile de analiză denumite generic “post-hoc”. În cazul acestor aplicații cunoștințele sunt reprezentate sub o formă cauzală abstractă, utilizatorul interacționând cu aplicația prin procedee de modificare a datelor. Astfel, prin procedee de actualizare, modelul este recalculat, pentru că acesta se bazează pe noi date.

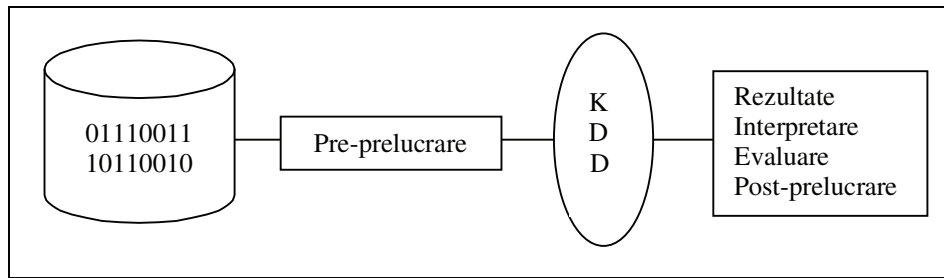


Figura 2. Procese de pre-prelucrare și post-prelucrare

În general, se poate spune că trebuie ținut cont de următoarele elemente atunci când se proiectează aplicații KDD:

- discuții ample cu utilizatorul pentru a se descoperi natura și finalitatea aplicației. Etapa cea mai importantă în acest proces o constituie selectarea unei variabile dependente care furnizează într-adevăr pentru întreprindere informația prelucrabilă. În paralel, se regăsește activitatea de achiziție a tuturor variabilelor (din coloanele bazelor de date) pentru a consolida modelul. Câteva tipuri de probleme pot avea nevoie de numai 50 de variabile, în timp ce altele au nevoie de 500 sau 1000 de variabile;
- instrumentul KDD trebuie să fie capabil să prelucreze un mare volum de date, adăugând noi date la variabilele existente (problemă de consolidare). Câteva aplicații KDD pot genera rezultate utile, numai exploatând zeci de mii de înregistrări, în timp ce majoritatea aplicațiilor exploatează curent sute de mii până la milioane de înregistrări.
- instrumentul KDD trebuie să fie ușor de utilizat. Aceasta înseamnă că interfața trebuie să furnizeze view-urile relațiilor dintre date, calitatea previziunii și a modelelor, într-un mod direct și intuitiv. Ideal ar fi ca utilizatorii să nu aibă nevoie de cursuri de formare în acest domeniu, pentru a învăța relațiile între tehnologiile informaticii decizionale sau pentru a învăța cum funcționează interfețele. Trebuie reflectat asupra faptului că cei mai mulți beneficiari ai KDD sunt decidenți și că aceștia nu au timp de investit pentru a afla cum funcționează noua tehnologie a informaticii decizionale;
- instrumentul KDD trebuie să manipuleze variabile numerice și categoriale, chiar dacă datele sunt răspândite în diverse surse sau chiar lipsesc, iar relațiile dintre date sunt adesea non-lineare.

Pentru ca tehnologia descoperirii de cunoștințe să fie capabilă să extragă informația din marile baze de date, trebuie ca aceasta să interacționeze îndeosebi cu sistemul de gestiune al bazelor de date, decât a obliga utilizatorul să-și extragă datele din aceste SGBD-uri și a le menține în afara bazelor de date. Această interacțiune va facilita gestionarea datelor. Dacă mai mulți utilizatori doresc să consulte în același moment bazele de date, aceștia riscă să nu aibă acces la informațiile dorite. Din păcate, multe dintre sistemele de gestiune a bazelor de date nu sunt încă pregătite să rezolve această problemă.

Printre altele, în materie de KDD poate fi relevată următoarea contradicție: pe de o parte, se asistă la creșterea în popularitate a descoperirii de cunoștințe în bazele de date, iar pe de altă parte, există totuși puține aplicații care au fost încununare de succes. Există mai multe rațiuni care explică acest fapt:

- succesul depinde, în primul rând, de *alegerea unei aplicații adecvate*. O astfel de aplicație trebuie să necesite descoperirea de rezultate care se știu extrase din date și ale căror prelucrări sunt cunoscute;
- succesul depinde în egală măsură de *conlucrarea aplicațiilor informatice existente și operaționale în întreprindere și metodele proprii descoperirii de cunoștințe*. În acest sens, există mai multe metode și abordări, dar pentru moment nici una nu se detașează net;
- *tranziția de la sistemul de explorare KDD la aplicația obișnuită*, face față aceluiași probleme ca de altfel toate transferurile de expertiză, de la cercetare la aplicațiile curente;
- publicațiile referitoare la KDD există, dar ele se adresează unui mic număr de inițiați.

Utilizarea aplicațiilor KDD este eficientă în situațiile în care volumul datelor necesare pentru a extrage o soluție este mic sau poate fi redus utilizând diferite tehnici statistice. În general, performanța algoritmilor ce reies din aceste aplicații asupra ansamblului de date introduse este non-lineară. Pentru a se putea servi de aceste aplicații, utilizatorii trebuie să stăpânească caracteristicile algoritmilor utilizați sau să aibă cunoștințe temeinice asupra tehnicilor statistice. Aceasta reprezintă o limită importantă de care trebuie să se țină seama și care trebuie rezolvată printr-o tehnologie ce permite aplicații mai ușor de utilizat și care beneficiază de interfețe mult mai conviviale.

Noile tehnologii de asistare a deciziei nu vin să înlocuiască exhaustiv sistemele clasice de asistare a deciziei, ci vin să le completeze, pentru ca decizia (în speță cea financiar-contabilă) – să poată considera și alte elemente “ascunse în date”. Aceste tehnologii ale informaticii decizionale, pot conlucra atât cu sistemele informatice clasice existente în întreprindere – pe care le exploatează, dar mai ales cu tehnologiile inteligenței artificiale care la rândul lor utilizează expertizele furnizate de noile abordări ale informaticii decizionale.

BIBLIOGRAFIE:

1. B. IONESCU, Sisteme de asistare a deciziei financiar contabile, Ed. InfoMega, Bucuresti, 2002
2. F. MIHAI, Sisteme informationale financiar – contabile, Ed. InfoMega, Bucuresti, 2004
3. KLEIN M., METHLIE L.B., Knowledge-based Decision Support Systems with applications in business, John Wiley&Sons, New York, 1999
4. PALLER Allan, LASKA R., The EIS book: Information System for top managers, Journal of Information Systems Management, Vol. 7, No.4, 1998
5. POWER D.J., Justifying a Data Warehouse project, <http://dss.cba.uni.edu/papers/dsscontext/index.html>