

## DATA MINING – TEHNOLOGII DEDICATE EXTRAGERII CUNOSTINTELOR

Obiective:

- însusirea tehnologiei Data Mining de extragere a cunostintelor din colectiile de date existente;
- însusirea unor tehnici Data Mining pentru obtinerea unor solutii în cadrul problemelor decizionale.

Concepte cheie: Data Mining; tehnici Data Mining.

Existenta unor volume imense de date a pus problema reorientării utilizării lor de la un proces de exploatare retrospectiv către unul prospectiv. Data Mining poate avea mai multe definitii, însă toate converg în esență către miezul problemei, și anume că acest concept reprezintă un proces de extragere de informații noi din colectiile de date existente. Termenul de dată are semnificatia de descriere a unui eveniment bine determinat care se produce în lumea reală și este perfect verificabil.

Prin tehnologia Data Mining se prelucrează date care referă perioade anterioare (date istorice), care sunt examinate și sunt deja cunoscute, pe baza lor constituindu-se un model. Acest model va putea fi aplicat situațiilor noi de același tip cu cele deja cunoscute. Informațiile care se pot obtine prin Data Mining sunt predictive sau descriptive. De exemplu directionarea acțiunilor de marketing pot constitui o problemă tipică predictivă. Detectarea fraudelor produse cu carduri bancare reprezintă o problemă tipică de aplicație descriptivă.

Dezvoltarea tehnicilor de Data Mining se explică prin acumularea de volume pe care acestea le-au derulat de-a lungul anilor. De asemenea, concurența tot mai acerbă precum și creșterea exigentelor pieței au determinat firmele să ia tot mai mult în considerare potențialul urias pe care îl oferă arhivele de date. Alături de arhivele de date memorate pe suporturi informatice mai există încă doi factori care au dus la necesitatea Data Mining: existenta și perfecționarea algoritmilor și a produselor program dedicate precum și creșterea capacității de memorare și prelucrare a calculatoarelor care permit tratarea corelativă a volumelor mari de date.

Este de remarcat că depozitele de date pot fi surse pentru Data Mining, iar rezultatele obtinute pot completa câmpurile înregistrărilor din depozitele de date, care apoi pot fi valorificate prin proiecțiile multidimensionale specifice OLAP.

Potențialul oferit de Data Mining se încorporează în procesele comerciale ale firmelor, iar căutarea informațiilor nu devine un scop în sine ci este utilă doar dacă este transformată ca acțiune. Astfel firmele pot alege să reacționeze sau nu la situațiile diverse create de realitate (diminuarea numărului de clienți, scăderea vânzărilor, pierderea unor piețe de desfacere etc.). Pasul următor după această alegere este exploatarea propriu-zisă a datelor utilizând diverși algoritmi.

De multe ori, acțiunea de Data Mining poate fi un eșec și nu o reușită, fiind posibil ca măsurile luate să nu fie adecvate informațiilor obtinute.

Toate elementele considerate anterior conduc spre ideea de ciclu în utilizarea Data Mining în cursul căruia sunt patru etape:

- definirea oportunităților comerciale și a datelor
- obtinerea de informații din colectiile de date existente prin tehnici Data Mining;
- adoptarea deciziilor și acțiunilor în urma informațiilor rezultate;

- cuantificarea cât mai corectă a rezultatelor concrete pentru a identifica și alte căi de exploatare a datelor.

### **Căutarea cunostintelor și verificarea ipotezelor**

Tehnicile de Data Mining se pot aplica atât ascendent, cât și descendent. Pentru abordarea descendentă se iau în considerare ipotezele formulate în prealabil prin alte mijloace. Abordarea ascendentă urmărește extragerea de cunostinte sau informații noi din date disponibile, această căutare putând fi dirijată sau nendirijată.

Căutarea dirijată presupune că se ia în considerare un atribut sau un câmp, ale cărui valori se explică prin celelalte câmpuri. Căutarea nendirijată identifică relațiile sau structurile din datele examinate fără a asigura prioritate unui câmp sau a altuia. Ceea ce se exploatează prin Data Mining sunt colecții de date constituite pentru alte scopuri (exemplu tranzacții derulate pe o perioadă de timp). Deseori la acest tip de date se adaugă și cele provenite din alte surse cum statistici oficiale care privesc evoluția în ansamblu a economiei, date privind concurența sau măsuri legislative. De aceea se folosește tot mai des noțiunea de informație ascunsă în sensul că este aproape imposibilă detectarea corelațiilor sau raporturile pe care datele le încorporează în mod intrinsec.

Rezultatele obținute sunt cu atât mai relevante cu cât ele se bazează pe un volum mare de date. Datele pot fi exploatate pentru a obține informații prin diverse tehnici cum sunt: rețele neuronale, arbori de decizie, algoritmi genetici, analiza grupurilor, raționamente bazate pe cazuri, analiza legăturilor. Aceste tehnici pot fi asociate cu tehnici statistice cum sunt regresiiile sau analiza factorială. Data Mining nu este capabilă, ca tehnică, să rezolve orice problemă de gestiune. De fapt ceea ce poate oferi se rezumă la câteva acțiuni cum sunt: clasificarea, estimarea, predicția, gruparea, analiza grupărilor, care folosite la locul potrivit pot deveni utile pentru o multime de probleme din domeniul decizional.

Destinația și caracteristicile acțiunilor oferite de Data Mining

Clasificarea are ca scop plasarea obiectelor prelucrate într-un grup limitat de clase predefinite. De exemplu, vânzarea unui produs nou se poate încadra într-una din următoarele categorii de risc: scăzut, mediu, ridicat. Obținute în mod clasificat vor fi reprezentate sub formă de înregistrări care la rândul lor sunt compuse din atribute sau câmpuri. Ca tehnici de Data Mining pentru clasificare sunt arborii de decizie și raționamentul bazat pe cazuri.

Estimarea va atribui o valoare unei variabile pe baza celorlalte date de intrare. Rezultatele obținute în urma estimării sunt valori continue. Pentru acest tip de prelucrări se pot utiliza rețelele neuronale.

Predicția poate clasa înregistrările luate în considerare în funcție de un anumit comportament sau o valoare viitoare estimată. De aceea se va recurge la o colecție de exemple care vizează date din trecut, în care valorile variabilei de previzionat sunt deja cunoscute. Cu ajutorul lor se va construi un model care va putea explica comportamentul observat. Aplicând acest model înregistrărilor care fac obiectul prelucrării, se va obține o predicție a comportamentului sau a valorilor acestora în viitor.

Gruparea poate duce la determinarea acelor obiecte care apar cel mai frecvent împreună. Un exemplu este „analiza cosului gospodăriei” în evaluările

statistice.

Analiza grupului urmărește o divizare a populației eterogene în grupuri mai omogene, care poartă numele de clustere.

În această tehnică nu se pleacă de la un set predeterminat de clase și nici din exemple din trecut. Segmentarea pe grupuri se face în funcție de similitudinile obiectelor.

### **Explorarea datelor – conținut și etape**

Programele care realizează implementarea algoritmilor pentru Data Mining nu sunt suficiente. Ele trebuie alimentate cu date care provin din diverse surse organizate pentru alte scopuri. De aceea este necesar un proces de curățare a acestora și de uniformizare pentru a fi explorate așa cum sunt ele furnizate de programe, conținutul lor trebuind a fi analizat de specialiști care vor identifica informațiile utile pe care acestea (rezultatele) le conțin. Având în vedere aceste particularități, tehnicile de Data Mining se pot utiliza numai în procese specifice complexe și de cele mai multe ori neliniare. Se pot astfel distinge etapele:

- definirea problemei;
- identificarea surselor de date;
- colectarea și selectarea datelor;
- pregătirea datelor;
- definirea și construirea modelului;
- evaluarea modelului;
- integrarea modelului.

Definirea problemei constă în sesizarea unei oportunități sau necesități de afaceri. De aceea se va delimita ceea ce urmează a fi rezolvat prin Data Mining, obiective urmărirea și rezultate scontate. Problema ce urmează a fi rezolvată prin Data Mining este o parte componentă a oportunității organizației, dar nu se identifică cu ea. De asemenea problema trebuie să primească o formă adecvată pentru a putea fi tratată cu această tehnică.

Identificarea surselor de date constă în stabilirea structurii generale a datelor necesare pentru rezolvarea problemei, precum și regulile de constituire a acestora și localizarea lor. Fiecare sursă de date va fi examinată pentru o familiarizare cu conținutul său și pentru identificarea incoerențelor sau a problemelor de definire.

Colectarea și selecția datelor este etapa în care se face extragerea și depunerea într-o bază comună a datelor care urmează a fi utilizate ulterior. Această etapă ocupă un timp mare, cam 80% din timpul total, iar existența depozitelor de date constituie un real avantaj.

În funcție de limitele echipamentelor de calcul folosite, de produsele program aplicate colecțiilor de date și nu în ultimul rând de bugetul disponibil se poate prelucra întregul fond de date disponibil sau un esantion. Dacă opțiunea aleasă este dirijată spre lucrul cu esantionare, atunci trebuie respectate toate regulile și cerințele de selecție a acestora.

Pregătirea datelor. Datele sunt de obicei stocate în colecții de date care au fost construite pentru alte scopuri. De aceea firesc este să existe o fază preliminară de pregătire înainte de extragere prin Data Mining. Transformările la care sunt supuse datele pentru Data Mining se referă la: valori extreme, valori lipsă, valori de tip text, tabele.

Tratarea valorilor extreme se poate face prin încadrarea între anumite limite cuprinse între medie și un număr de abatere standard prin excludere sau limitare sau prin izolarea vârfurilor.

În cazul valorilor lipsă se pot elimina câmpurile cu valori nule din înregistrări, sau se pot completa câmpurile cu date de valori medii, deoarece existența lor poate duce la o funcționare incorectă a algoritmilor de Data Mining.

Valorile de tip text ridică probleme întrucât separarea prin spații a cuvintelor duc la apariția de valori diferite. Din acest motiv este indicată eliminarea lor, dar dacă prelucrarea lor nu poate fi eliminată, soluția cea mai pertinentă este de codificare prin tabele de corespondență, în care să se evidențieze toate sirurile valide de caractere.

Rezumarea se aplică atunci când datele sunt considerate a reprezenta detalii neesențiale pentru rezolvarea problemei, sau când numărul de exemple este insuficient.

Codificarea incoerentă apare în momentul în care obiecte identice sunt reprezentate diferit în unele din sursele utilizate. Incompatibilitățile arhitecturale informatice se referă la diferențele existente între modul de reprezentare internă a valorilor datorat creării lor cu sisteme din generații diverse.

Definirea și construirea modelului este etapa care se apropie cel mai mult de noțiunea de Data Mining și se referă la crearea modelului informatic care va efectua exploatarea. Etapa de definire și construire a modelului este însoțită de faza de instruire sau învățare, depinzând de tehnicile de Data Mining utilizate.

Indiferent de aceste tehnici toate au de parcurs două etape: învățarea și testarea. Învățarea presupune existența unui set suficient de reprezentativ de exemple complete de la care se porneste pentru a identifica relațiile de legătură între valorile câmpurilor sau atributelor. Se consideră ca fiind încheiat procesul de învățare, în momentul în care rezultatele obținute prin model se apropie suficient de mult de soluțiile continuate de datele după care s-a învățat. Nu întotdeauna rezultatele sunt cele scontate și atunci modelul va fi supus testării cu date diferite de cele folosite pentru învățare, dar care apar în aceeași colecție. În această etapă sunt formulate alte două obiective, și anume: obținerea de date preclasate și distribuirea acestora în seturi de învățare, testare sau evaluare.

Evaluarea modelului are ca scop de a determina corect valorile în care modelul are capacitatea de a determina corect valorile pentru cazurile noi. Modelul va fi astfel aplicat asupra ultimei părți din datele preclasate care sunt dedicate evaluării. Procentul de eroare ce se stabilește acum va fi considerat că va fi acceptat și pentru datele noi.

Performanțele unui model se vor aprecia cu „matricea de confuzie” care are rolul de a compara situația reală cu cea pe care modelul o furnizează. Integrarea modelului este etapa în care se finalizează procesul, prin încorporarea modelului în SIAD ca element de bază, sau prin includerea sa într-un proces decizional general din organizație.

### **Rationamentul bazat pe cazuri**

Prin această tehnică se caută o rezolvare a problemelor apărute prin analogie cu experiența acumulată. Această metodă se poate aplica pentru clasificări și pentru predicții. Cazurile pe care este bazat rationamentul sunt memorate ca înregistrări

compuse din setul de atribute care descriu fiecare caz. Un caz nou este prezentat tot ca o înregistrare, numai că în câmpurile în care valoarea trebuie determinată sunt vide. Pentru a determina aceste valori se caută înregistrările cu care înregistrarea „caz nou” se aseamănă și conținutul acestora se consideră a fi răspunsul. Prin urmare se poate afirma că există două funcții fundamentale de prelucrare:

- a) măsurarea distanței dintre membrii fiecărui cuplu de înregistrări, pentru a afla vecinele cele mai apropiate;
- b) combinarea rezultatelor obținute de la „vecine” în răspunsul propus pentru cazul curent.

Măsurarea distanței dintre câmpuri. Se numește distanță expresia modulului în care se evaluează similitudinea. Distanța are ca proprietăți: poate fi definită și se prezintă ca un număr real; distanța de la un element la el însuși este totdeauna nulă; sensul de măsurare este fără semnificație în maniera că distanța de la elementul A la elementul B este egală cu distanța de la B la A și nu există un punct C intermediar lui A și B prin a cărei parcurgere să se scurteze drumul de la A la B.

Ca moduri de calcul pentru distanța câmpurilor numerice se enumeră:

- diferența între valori absolute  $|A-B|$ ;
- pătratul diferenței  $(A-B)^2$ ;
- diferența între valori absolute normalizată  $|A-B|$  (diferența maximă). Ultima variantă produce rezultate cu valori cuprinse între 0 și 1. Măsurarea distanței între înregistrări. Când apare necesitatea de a considera simultan mai multe câmpuri ale înregistrării, se calculează distanța pentru fiecare câmp în parte, iar rezultatul se combină într-o valoare mică care reprezintă distanța înregistrării respective.

Se vor enumera câteva procedee de combinare a distanței câmpurilor: însumarea, însumarea normalizată (suma distanțelor/suma maximă), distanța euclidiană (rădăcina pătrată din suma pătratelor distanțelor). Distanța euclidiană evidențiază cel mai bine înregistrările pentru care toate câmpurile sunt vecine. Combinarea rezultatelor presupune aflarea celor mai apropiați vecini, iar soluția problemei se obține prin combinarea răspunsurilor obținute de la aceștia.

Fiecare vecin poate avea diverse variante de răspuns, dar se vor lua în calcul doar cei care sunt mai apropiați. Rezultatul ce obține majoritatea va fi atribuit cazului curent. Cerința minimă este ca numărul votanților să fie impar, pentru a evita situațiile de nedeterminare.

Metodele care se bazează pe vot dau rezultate satisfăcătoare în situațiile în care răspunsurile așteptate sunt de tip enumerativ. O altă soluție posibilă este interpolarea valorilor înregistrărilor vecine care însă introduce o aplatizare a rezultatelor care se înscriu între cele două limite folosite în calcul. De asemenea, se poate constata că rezultate bune se obțin prin metode de regresie statistică aplicate asupra valorilor date de vecinii cei mai apropiați. Se obține ecuația unei drepte sau a unei curbe care permite calcularea mai precisă a valorilor aferente cazului curent.

Se poate concluziona că raționamentul bazat pe cazuri este o tehnică de Data Mining suficient de bună și care se poate aplica unui mare număr de probleme, caz în care conduce la soluții acceptabile. Toate acestea sunt valabile dacă volumul de date pe care se bazează este bine ales și concludent. Ca avantaj pentru această metodă se pot enumera:

- aplicarea unui mare număr de tipuri de date, pe structuri de date complexe, iar câmpurile tip text sunt mai bine tratate decât în alte tehnici;
- luarea în considerare a oricât de multor câmpuri;
- rezultatele obținute sunt explicite;
- elementele de noutate care apar în procesul de învățare sunt ușor de înglobat și de folosit în raționamente.

Ca orice metodă prezintă și unele dezavantaje dintre care se pot menționa: volumul mare de memorie și resursă timp de prelucrare relativ mare, și de asemenea, timpul de prelucrare mare pentru aplicarea funcțiilor de distanță asupra tuturor înregistrărilor și câmpurilor necesare pentru obținerea rezultatelor.