



Volumul VI, Numărul 10 / 2004

ISSN 1454-9980

Procesul de data mining: descoperirea cunoștințelor ascunse

(pag. 19-24)

Bogdan ANASTASIEI

Volume VI, Issue 1 (10) / 2004

Cross-cultural
Management
Journal

PROCESUL DE *DATA MINING*: DESCOPERIREA CUNOȘTIȚELOR ASCUNSE

Bogdan ANASTASIEI

“Avem foarte multe date colectate în firmă, ce să facem cu ele acum?” Această problemă a devenit una obișnuită în multe organizații mari. Informația digitală este ieftin de obținut și relativ ieftin de stocat. Dar care este scopul stocării unei cantități de date atât de mari? În afara argumentului legat de modalitățile convenabile de stocare a datelor în format electronic, mai există și altul: firmele colectează date deoarece managerii “simt” că aceste date reprezintă un activ valoros, care ar putea fi folosit la un moment dat. În institutele de cercetare, datele se referă la observații asupra fenomenelor aflate sub studiu, culese cu atenție. În organizațiile economice, datele și informațiile se referă la piețe, concurenți, furnizori, distribuitori, clienți, precum și la informații interne legate de procesele de producție, organizarea muncii, depozitare etc.

Pentru multe firme din Occident, procesul de explorare și analiză a datelor nu este unul nou. Dar programele de *data mining* permit realizarea acestor analize mult mai rapid și, potențial, cu o eficiență sporită. În multe cazuri, tehnicile de explorare a datelor permit ca datele colectate pentru un anumit scop să poată fi folosite în multe alte scopuri.

Ca urmare a perfecționării modalităților de colectare și stocare a datelor, foarte multe companii mari s-au trezit deodată că stau pe un munte de date cărora ar trebui să le găsească o utilizare. De exemplu, companiile furnizoare de cărți de credite înregistrează permanent datele referitoare la tranzacțiile comerciale, supermarketurile înregistrau date privind cumpărăturile, utilizarea cupoanelor de reduceri ș.a.m.d. Apariția necesității de a valorifica toate aceste date era doar o chestiune de timp.

Prin *data mining* se înțelege procesul de extragere automată a unor informații cu caracter predictiv din marile baze de date. Ea poate prezice tendințele viitoare și identifica comportamente care le-au scăpat responsabililor din organizație.

Datele brute sunt rareori folositoare în forma în care se găsesc. Valoarea lor reală este dată de posibilitatea de a extrage din ele informații care să ajute la fundamentarea deciziilor sau la înțelegerea fenomenelor care guvernează viața organizației. În mod tradițional, analiza datelor era un proces strict manual. Unul sau mai mulți analiști deveneau foarte familiari cu datele și furnizau periodic rapoarte, cu ajutorul metodelor statistice. O asemenea abordare devine totuși dificil de pus în practică, pe măsură ce volumul de date crește până la dimensiuni uriașe. Cine ar putea “înțelege” o bază de date cu câteva milioane de obiecte, fiecare având câteva zeci de câmpuri? Pentru a complica situația, volumul de date crește într-un ritm atât de accelerat încât analiza manuală (chiar și atunci când este posibilă) nu mai poate ține pasul.

Comunitatea de cercetători și practicieni interesați de problema automatizării analizei datelor a crescut rapid. Două sunt denumirile sub care este cuprins acest proces de automatizare: “knowledge discovery in databases” (KDD) și “data mining”, pe care noi o vom mai numi aici și “explorarea datelor”. Primul workshop pe tema KDD a fost ținut în 1989; el s-a transformat ulterior într-o conferință internațională care reunește anual peste 500 de participanți.

Explorarea datelor ar putea fi de asemenea definită ca “procesul de căutarea a unui *pattern*

(model) în interiorul datelor”. Ea utilizează metode sofisticate de analiză statistică. Odată găsită, informația ascunsă în spatele datelor brute trebuie prezentată într-o formă accesibilă cu ajutorul rapoartelor, tabelelor, graficelor etc.

Statistica este pilonul central al întregii activități de explorare a bazelor de date. Atât în ceea ce privește validarea ipotezelor, cât și în ceea ce privește analiza exploratorie a datelor și realizarea predicțiilor, statistica are o importanță fundamentală.

Calitatea datelor este critică pentru obținerea unor rezultate consistente din analiza datelor. Primul pas în KDD și în *data mining* este construirea unui sistem de culegere a datelor într-o formă cât mai precisă și fidelă. Datele primare de calitate slabă pot compromite întregul proces.

Sistemul cu ajutorul căruia se realizează în fiecare zi activitățile din organizație precum preluarea comenzilor, înregistrările contabile, gestiunea stocurilor etc. poartă numele de “sistem operațional”, iar informațiile pe care acest sistem le creează poartă denumirea de “date operaționale”. Explorarea datelor este folosită pentru a da consistență și substanță acestor date operaționale, care sunt date brute, primare. Este un proces de “descoperire”, în urma căruia ies la iveală informații pe care datele brute nu le pot oferi în forma actuală. Nu sunt utilizate rapoarte cu o structură prestabilită, ci este permisă decidentului crearea de interogări a bazei de date pornind de la informațiile de care are nevoie în momentul respectiv.

Să luăm un exemplu de aplicare a procesului de *data mining* în vânzări. Utilizatorul (managerul de vânzări) poate începe prin a cere bazei de date informații cu privire la vânzările totale pe trimestrul curent și pe trimestrul anterior, precum și un calcul al diferenței și a procentajului de creștere sau descreștere a vânzărilor. El află, de exemplu, că vânzările au crescut, per total, cu 3,5% de la trimestru la trimestru. Mai departe, el dorește să afle cum stă evoluția vânzărilor pe fiecare din cele 42 de județe, și cere bazei de date un raport asemănător, numai că acum defalcat pe județe. Primind acest raport, el află că vânzările au crescut în majoritatea județelor, dar au scăzut în Bacău și Vrancea.

Pentru a afla mai mult, el cere situația vânzărilor pentru fiecare din cei 10 distribuitori din județul Bacău, unul din județele cu

probleme. Astfel, află că în acest județ, 8 din 10 distribuitori au înregistrat scăderi destul de drastice ale vânzărilor. Mai departe, el cere pentru fiecare din acești distribuitori, o situație a evoluției vânzărilor pe produse și observă că declinul cel mai accentuat se înregistrează la produsul X.

Urmărind situația vânzărilor pentru zona Vrancea, utilizatorul descoperă același lucru: vânzările la produsul X scad la majoritatea distribuitorilor din acest județ. Studiind vânzările produsului X și în județele unde se înregistrează creșteri de ansamblu ale cifrei de afaceri, managerul află că în cazul acestui produs, creșterile sunt foarte reduse (aproape de zero), iar în multe cazuri (cca. 30% din totalul județelor) s-au produs chiar scăderi ușoare.

Concluzia analizei: produsul X pare a fi un produs-problemă. Momentan, problema se manifestă mai acut în două zone, dar s-ar putea extinde în viitor. Așadar, deși per ansamblu vânzările firmei au crescut, există totuși o problemă importantă, care a fost identificată cu ajutorul procesului de *data mining*.

Iată câteva alte aplicații posibile ale procesului de *data mining* în organizațiile economice (și nu numai).

- în marketing, la realizarea profilului consumatorilor. Caracteristicile clienților buni (celor care cumpără de la firmă, regulat și în cantități mari) sunt utilizate ca variabile de predicție. Ele îi vor ajuta pe marketeri să-i “țintească” pe noii clienți. Cu ajutorul explorării putem identifica, în baza de date cu clienți, o serie de modele care pot fi aplicate ca criterii de selectare pentru o bază de date cu clienți potențiali. Acei clienți potențiali, care se aseamănă cel mai mult cu “clientul ideal”, vor constitui cei mai buni candidați pentru o campanie de marketing direct, de exemplu. Este de așteptat ca beneficiile oferite de produsul firmei, să fie interesante și pentru ei. Acest fapt permite eficientizarea campaniilor de marketing, prin direcționarea lor către acei potențiali clienți care au cea mai mare probabilitate de cumpărare;
- explorarea datelor îi poate ajuta pe detașiști să înțeleagă cum arată “coșul de produse” al unui cumpărător individual (ce tip de produse și ce mărci sunt achiziționate simultan). Ei pot afla astfel ce produse trebuie expuse pe rafturi și cum trebuie

expuse. Explorarea poate de asemenea ajuta la măsurarea eficienței campaniilor promoționale ale magazinului;

- o altă utilizare obișnuită a explorării datelor o întâlnim în gestiunea relațiilor cu clienții. Determinând caracteristicile clienților pentru care este foarte mare probabilitatea de a ne părăsi în favoarea unui client, compania poate pune în practică acțiuni de reținere a clienților, știut fiind că reținerea clienților vechi costă mult mai puțin decât atragerea de noi clienți;
- detectarea fraudelor este un alt câmp de aplicație pentru *data mining*, de care sunt interesate în special companiile de telecomunicații, organizațiile care emit cărți de credit, bursele, agențiile guvernamentale. Suma totală pierdută din cauza fraudelor este enormă. Cu ajutorul explorării datelor, organizațiile pot identifica tranzacțiile potențial frauduloase și stopa pierderile până nu este prea târziu;
- organizațiile financiare utilizează explorarea datelor în scopul determinării caracteristicilor pieței și a diferitelor sectoare de activitate, prezicând astfel comportamentul viitor al titlurilor financiare ale diferitelor companii;
- o altă aplicație este cea din domeniul medical: explorarea datelor poate prezice eficiența intervențiilor chirurgicale, analizelor medicale, administrării unor medicații, serviciilor medicale ș.a.m.d.

Bazele de date utilizate pentru acest proces de explorare trebuie să fie deosebit de flexibile, deoarece utilizatorul poate face explorarea într-o manieră ad-hoc, fără a urma standarde prestabilite. În exemplul nostru anterior, managerul și-a structurat interogarea pe direcția județ-distribuitori-produs, dar altcineva ar putea cere foarte bine un raport pe direcția produs-distribuitori-județ sau produs-județ-distribuitori. Baza de date trebuie să cuprindă foarte multe referințe încrucișate, pentru a putea asocia cu ușurință câmpurile și extrage rapid rapoarte rezumative. O altă cerință importantă este aceea ca datele să fie cât mia “condensate” cu puțință, astfel încât timpul necesar obținerii răspunsului la o interogare să fie cât mai redus posibil (câteva sutimi de secundă). Dacă timpul necesar obținerii unui raport este prea mare (de ordinul secundelor), utilitatea procesului este pusă sub semnul întrebării.

Caracteristicile cheie ale unui bun proces de *data mining* sunt:

- ușurința utilizării. Procesul trebuie să fie ușor de realizat și să nu solicite o instruire prea îndelungată a utilizatorilor. Este important ca utilizatorul să se poată concentra asupra propriilor sale necesități și probleme, asupra modului de formulare a interogărilor, și nu asupra modului de utilizare a programului. Cu cât este mai dificilă rularea programului de *data mining*, cu atât este mai sigur că nu va fi folosit chiar de către persoanele care ar avea cea mai mare nevoie de el;
- accesibilitatea. Valoarea reală a explorării datelor se obține oferind acest instrument celor care au nevoie de el, care ar avea de aflat cele mai multe lucruri în urma analizei datelor. Managerii de vânzări, de exemplu, îl vor folosi pentru a evalua eficiența agenților de vânzări și a distribuitorilor, pe fiecare tip de client în parte. Managerii de achiziții îl vor folosi pentru a determina evoluția stocurilor pe fiecare tip de material. Pentru contabili, ar putea fi un instrument util în construirea bugetelor și a situațiilor financiare, și așa mai departe;
- răspuns rapid. Este foarte important ca analiza și furnizarea informației să fie realizată în cel mai scurt timp posibil (de ordinul zecimilor sau sutimilor de secundă). Numai în acest fel capacitățile procesului de *data mining* vor fi folosite la întregul lor potențial. Decidenții vor fi mai puțin înclinați spre a utiliza programul atunci când știu că trebuie să aștepte secunde bune pentru a obține rezultatul unei interogări, deoarece acest lucru ar însemna pentru ei o pierdere de timp importantă;
- informație la zi (actualizată). Folosirea unor date vechi de săptămâni sau luni reduce considerabil eficiența și utilitatea procesului. Cu cât datele sunt mai actuale, cu atât mai repede poate fi produs un plan de măsuri pentru rezolvarea problemelor sau exploatarea oportunităților identificate prin procesul de explorare a datelor. În exemplul nostru, combaterea declinului înregistrat la vânzările produsului X poate fi realizată cel mai eficient, dacă problema a fost descoperită chiar în momentul apariției ei și cât mai devreme după acest moment.

Beneficiile procesului de explorare a datelor sunt următoarele:

- costul mai redus al procesării datelor. Atunci când decidenții obțin datele necesare în urma unor analize făcute de ei înșiși, dispare necesitatea de a crea programe informatice specializate pentru obținerea și raportarea acestor date. Astfel, programatorii își pot concentra atenția asupra altor probleme legate de operarea sistemelor din firmă;
- reducerea cantității de hârtie ocazionate de tipărirea rapoartelor și situațiilor. Marea majoritate a programelor informatice oferă managerilor sumele totale ce reflectă volumul activității din domeniul lor (de exemplu, vânzările totale). În cazul procesului de explorare a datelor, managerul poate începe prin a solicita aceste sume totale pe ecranul computerului, apoi poate cere defalcarea lor în funcție de informația pe care dorește să o primească. Odată ce deține această informație, el va cere eventual tipărirea unui raport, foarte probabil mult mai util decât cel ce prezintă doar totalurile (justificând așadar mai bine costul hârtiei pe care a fost tipărit);
- reducerea încărcării sistemului operațional. Procesul de *data mining* utilizează de regulă baze de date care au fost realizate și condensate fie offline, fie în afara orelor de program din firmă. Întrucât procesele de sortare, selectare și raportare sunt mai ușoare pe aceste baze de date, încărcarea sistemului informatic al organizației va fi substanțial redusă.

Exemplele din acest articol sunt exemple simple, care totuși ne arată cât de util este a oferi decidenților posibilitatea de a privi datele din diverse unghiuri de vedere. Pe măsură ce întreprinderile își dau seama că procesul de *data mining* nu este unul extrem de complex și de costisitor, tot mai multe își vor include achiziționarea sistemelor de acest fel în planurile lor de dezvoltare viitoare.

Există trei clase de activități de *data mining*:

- descoperirea, care este procesul de căutare în bazele de date a unor modele ascunse, fără ipoteze și idei preconceptuate asupra naturii acestora. Cu alte cuvinte, inițiativa de a căuta aceste modele îi aparține

programului, înainte ca utilizatorul să formuleze interogările;

- modelarea predictivă, care presupune utilizarea modelelor descoperite anterior în scopul prezicerii viitorului. Programul permite utilizatorului să introducă înregistrări care au unele câmpuri cu valoare necunoscută; valoarea acestora urmează a fi estimată de program tot cu ajutorul pattern-urilor identificate anterior;
- analiza extremelor, este procesul ce permite detectarea, pe baza modelelor identificate în prima etapă, a situațiilor anormale sau neobișnuite ce caracterizează obiectele din baza de date. Pentru a descoperi cazurile neobișnuite, va trebui să se răspundă mai întâi la întrebările: “Ce înseamnă normalul? Care este regula?”. Cazurile care înregistrează o deviere, cu o anumită mărime, de la această regulă, vor fi considerate cazuri neobișnuite. Analiza extremelor merge așadar mai departe decât descoperirea și modelarea predictivă, ajutându-ne nu doar să obținem așa numitele “cunoștințe obișnuite”, ci să identificăm și situațiile anormale.

Putem vorbi de trei tipuri de activități de explorare a datelor care se desfășoară într-o organizație: explorare episodică, explorare strategică și explorare continuă.

Explorarea episodică presupune efectuarea analizei datelor cu ocazia unui anumit eveniment: o campanie de marketing, de exemplu. Se pot analiza datele legate de această campanie pentru a se înțelege o serie de relații de cauzalitate (de exemplu, relația dintre cheltuielile pentru publicitate și volumul vânzărilor) și se poate folosi această analiză pentru a planifica și previziona rezultatele viitoarelor campanii de marketing. De regulă, analiștii sunt cei care realizează activități de *data mining* episodice.

Explorarea strategică presupune luarea în considerare a unui set de date mai extins, pentru a dobândi o înțelegere de ansamblu asupra unor parametri care caracterizează activitatea firmei (de exemplu, profitabilitatea). În acest caz, explorarea strategică poate răspunde la întrebări de genul: “Care sunt principalele noastre surse de profit?” sau “Care este relația dintre portofoliul nostru de produse și segmentele noastre de piață?”.

În cazul explorării continue, obiectivul este acela de a înțelege ce schimbări au avut loc în organizație și în mediul ei extern pe o anumită perioadă de timp și care sunt cei mai importanți factori care au influențat acele schimbări. Managerul s-ar putea întreba, de exemplu: “cum au evoluat vânzările în ultimul semestru, pe total și pe tipuri de produse?” sau “cum s-a modificat nivelul de satisfacție a consumatorilor noștri în ultimele luni și din ce motive?”.

Procesul de explorare a datelor are cinci componente principale, care sunt de fapt grupe de tehnici de analiză a datelor. Iată care sunt acestea:

- clasificarea (sau *clustering*-ul);
- asocierea;
- analiza secvențială;
- rețelele neuronale;
- arborii de decizie

În ceea ce urmează le vom discuta pe fiecare din ele pe scurt.

Clasificarea analizează un set de date și stabilește un ansamblu de reguli pe baza cărora vor fi grupate datele obținute în viitor. Software-ul de *data mining* identifică automat clasele (sau “ciorchinii”), studiind pattern-ul datelor existente. Odată ce au fost generate clasele se poate stabili, pe baza unor caracteristici precizate, cărei clase aparține un anume obiect din baza de date studiată. De exemplu, o clasă poate fi un segment de piață. Atunci când firma are un client nou ea poate stabili, cunoscând caracteristicile acestuia, în ce segment se încadrează.

O regulă de asociere este o regulă care implică existența unor relații între diferite seturi de obiecte din baza de date. Procesul de asociere duce la descoperirea unor asemenea reguli, la diferite niveluri de abstracțiune, între obiectele bazei de date. De exemplu, se poate afla că scăderea vânzărilor este totdeauna însoțită de o serie de simptome (evenimente din interiorul sau din afara organizației) și pe această bază se pot studia mai departe motivele care stau în spatele acestei asocieri.

Analiza secvențială duce la descoperirea acelor evenimente care se petrec întotdeauna într-o anumită secvență. Ea studiază datele care apar în tranzacții diferite (spre deosebire de asociație, care studiază relațiile dintre datele ce apar în cadrul aceleiași tranzacții). În urma

analizei secvențiale s-ar putea descoperi că majoritatea clienților, care în decursul unei săptămâni cumpără produsele A și B, săptămâna viitoare cumpără produsele C și D, de exemplu.

Pe măsură ce procesul de *data mining* devine mai obișnuit în organizații, rețelele neuronale și arborii de decizie se bucură de tot mai multă considerație. Deși rețelele neuronale sunt mai complexe în felul lor, utilizarea lor nu cere cunoștințe de statistică prea multe și prea avansate.

Rețelele neuronale utilizează un număr mare de parametri pentru a construi un model care preia și combină un set de inputuri în scopul prognozării comportamentului unei variabile cantitative sau categoriale. Valoarea fiecărui nod al rețelei se calculează ca medie ponderată a valorilor nodurilor anterioare. Procesul de construire a modelului implică aflarea acelor ponderi care produc cele mai precise prognoze, utilizând date reale pentru a “antrena” rețeaua. Cea mai comună metodă de “antrenare” este compararea datelor calculate în rețea cu valorile corecte cunoscute. După fiecare comparare, ponderile sunt ajustate și valorile sunt calculate din nou. După o perioadă îndelungată de asemenea comparări și ajustări, o rețea neuronală devine de regulă un foarte bun predictor.

Arborii de decizie reprezintă o serie de reguli a căror aplicare ne conduce către o anumită clasă sau valoare din baza de date. De exemplu, să presupunem că o bancă vrea să-i clasifice pe solicitanții de credite în două categorii: cei cu risc scăzut și cei cu risc ridicat. Un arbore de decizie construit în acest scop ar putea arăta în felul următor:

- are solicitantul un venit mai mare de 50 000 de dolari pe an? Dacă da, următoarea întrebare care se pune este: “Cât de mare este nivelul datoriei sale în prezent?”. Dacă persoana respectivă este foarte îndatorată, ea prezintă un risc mare la creditare; dacă este puțin îndatorată, atunci prezintă un risc redus;
- dacă solicitantul are un venit mai mic de 50 000 de dolari pe an, întrebarea următoare este: “Care este vechimea sa în muncă?”. Dacă este mai mare de 5 ani, se consideră că solicitantul prezintă un risc mic la creditare; dacă vechimea este mai mică de 5 ani, atunci este considerat un solicitant cu risc mare.

Arborii de decizie sunt folosiți deseori, deoarece au un grad de precizie rezonabil și, spre deosebire de rețelele neuronale, sunt mai ușor de înțeles și necesită un timp de construire mai redus. Ar mai trebui precizat că atât arborii de decizie, cât și rețelele neuronale pot fi utilizate pentru construirea unor modele de regresie.

Etapele unui proces de explorare a datelor (*data mining*) sunt în general următoarele:

- definirea problemei. Pentru o utilizare cât mai eficientă a tehnicilor de explorare trebuie mai întâi specificate foarte clar obiectivele studiului. De exemplu, în cazul unei campanii de *direct mailing*, obiective precum “creșterea numărului de răspunsuri”, “creșterea ratei răspunsurilor” sau “diminuarea costului per răspuns” sunt obiective diferite, care vor necesita metode de analiză diferite. Definirea clară a

obiectivelor ne va ajuta ulterior și la măsurarea eficienței procesului;

- selectarea tehnicii care va fi utilizată, pornind de la problema specifică de rezolvat;
- selecția și pregătirea datelor. Aceasta este etapa care consumă cel mai mult timp (între 50 și 85 la sută din timpul total alocat proiectului). După cum spuneam anterior, calitatea rezultatului final depinde foarte mult de acuratețea datelor de intrare;
- construirea modelului, ținând seama de tehnicile de analiză alese și de natura datelor pe care le avem la dispoziție;
- prezentarea rezultatelor finale, sub forma unor rapoarte conținând text, tabele și grafice;
- monitorizarea modelului, în vederea măsurării eficienței și utilității sale.