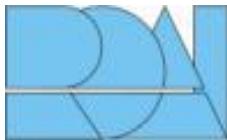


Mihai Istrate



**Centrul de Cercetare în Inteligență Aplicată – „Nicolae
Țăndăreanu”**

**Research Center for Applied Intelligence – „Nicolae
Țăndăreanu”**

www.rcai.eu

Colecția « Computer Science »

Coordonatori colecție:

Daniela Dănciulescu – Director RCAI – Universitatea din Craiova
Gabriel Stoian – Universitatea din Craiova

Comitetul științific:

Gheorghe Grigoraș, Universitatea Alexandru Ioan Cuza din Iași
Viorel Negru, Universitatea de Vest din Timișoara
Petrica Pop Sitar, Universitatea Tehnica din Cluj-Napoca
Mihaela Păun, Academia de Studii Economice din București
Cristian Kevorchian, Universitatea din București
Claudiu-Ionuț Popîrlan, Universitatea din Craiova

Inițiată în 2001, sub egida Centrului de Cercetare în Inteligență Aplicată – „Nicolae Țăndăreanu”, colecția „Computer Science” reunește contribuții valoroase – publicații științifice de înaltă ținută, studii, teze de doctorat, etc. – continuând astfel tradiția publicării în volum separat a seriei 100 (lucrările conferinței anuale AIDC – Artificial Intelligence and Digital Communications), a seriei 200 (Computer Science Fundamentals) și a seriei 300 (Research Reports in Artificial Intelligence).

Rămânând fidelă misiunii sale de dezvoltare și promovare a cunoașterii în domeniul științei calculatoarelor și a tehnologiilor digitale, colecția „Computer Science” reprezintă o sursă reală de informare și își propune lărgirea spectrului publicistic prin dezvoltarea de noi serii tematice.

Propunerile pentru publicare se vor adresa comitetului științific la adresa: office@rcai.eu

Mihai Istrate

**METODE DATA MINING APLICATE
ÎN E-COMMERCE ȘI INTERNET**

- MONOGRAFIE -



**Editura UNIVERSITARIA
Craiova, 2021**

Referenți științifici:

Prof.univ.dr. Daniela Dănciulescu
Lect.univ.dr. Claudiu-Ionuț Popirlan

Copyright © 2021 Editura Universitară
Toate drepturile sunt rezervate Editurii Universitaria

**Descrierea CIP a Bibliotecii Naționale a României
ISTRATE, MIHAI**

Metode data mining aplicate în e-commerce și internet / Mihai Istrate. - Craiova : Universitară, 2021
Conține bibliografie
ISBN 978-606-14-1705-6
004

© 2021 by Editura Universitară

Această carte este protejată prin copyright. Reproducerea integrală sau parțială, multiplicarea prin orice mijloace și sub orice formă, cum ar fi xeroxarea, scanarea, transpunerea în format electronic sau audio, punerea la dispoziția publică, inclusiv prin internet sau prin rețelele de calculatoare, stocarea permanentă sau temporară pe dispozitive sau sisteme cu posibilitatea recuperării informațiilor, cu scop comercial sau gratuit, precum și alte fapte similare săvârșite fără permisiunea scrisă a deținătorului copyrightului reprezintă o încălcare a legislației cu privire la protecția proprietății intelectuale și se pedepsesc penal și/sau civil în conformitate cu legile în vigoare.

LISTA DE FIGURI

1.1 Elementele specifice problemelor întâlnite la Data mining	24
1.2 Data mining - utilizări	25
1.3 Tendențe Data Mining	37
1.4 Data Mining automatizează	38
1.5 Data Mining prin Web	40
1.6 Pașii parcurși în Data mining	40
2.1 Obiectivele utilizării tehnologiilor Internet	44
2.2 Principalele paradigme de trecere la derularea pe suport informatic a afacerilor	45
3.1 Web structure mining	56
3.2 Semantic Web mining	63
3.3 Conceptul și operațiile Web mining semantic	63
3.4 Utilaje pentru agenți pentru Web mining semantic	64
3.5 Legătura strânsă între Data Mining și XML/ baza de date RDF	64
3.6 Legături libere între Data Miner și XML / baze de date RDF	65
3.7 Extragerea structurii și apoi utilajelor	65
3.8 Ontologii și Web Data Mining	66
3.9 Revizuirea agenților și utilajelor Web	67
3.10 Utilajele Web și Web-ul semantic	68
3.11 Utilajele web semantic și interoperabilitate	68
3.12 Utilaje pentru interoperabilitate semantică a bazelor de date XML	69
3.13 Servicii și utilaje Web	69
3.14 Evoluția utilajelor Web-ului semantic	70
3.15 Structura utilajelor Web	76
3.16 Structura unui sistem de Web usage mining	81
3.17 Cadru web usage mining	87

4.1	Un exemplu de graf hyperlink	98
4.2	Lanț Markov periodic, cu $k = 3$	100
5.1	Arhitectura generală a programului FastStats	108
5.2	Pachetul software SurfStats (screenshot)	110
5.3	Pachetul software SAWMILL	112
5.4	Arborele de decizie pentru testul 1	114
5.5	Arborele de decizie pentru testul 2	115

LISTA DE TABELE

1.1	Aplicațiile, tehnicele și algoritmii data mining	39
2.1	Asemănări și deosebiri între E-Commerce și comerțul tradițional	49
5.1	Detaliile tehnice ale datelor de testare	112

CAPITOLUL 1

ANALIZA - DATA MINING

1.1 Elemente generale în DATA MINING

Încă de la începutul anilor 1990 se tot vorbește despre conceptul Data Mining (prescurtat DM), sau mai este întâlnit și în literatura de specialitate română *minerit* în date, în foarte multe medii, pornind de la cele academice și până la cele de afaceri sau medicale, îndeosebi.

Reprezentând o arie de cercetare cu o istorie destul de scurtă, nedepășind faza *adolescenței*, este încă disputată de câteva domenii științifice care o revendică.

Preluând afirmațiile lui Pregibon -Research Scientist Google Inc. - și anume că: *Conceptul Data Mining reprezintă un amestec de Statistică, IA (Inteligentă Artificială) și cercetare în baze de date*, sau faptul că anumiți cercetători consideră Data Mining ca fiind *un cuvânt murdar în Statistica* [122] - cel mai exact aceștia au fost statisticieni care nu considerau cu ceva timp în urmă Data Mining ca fiind ceva important pentru ei.

În ultimii ani, descoperirea de cunoștințe și instrumente de data mining au fost utilizate în principal în medii experimentale și medii de cercetare. Acum suntem într-o etapă în cazul în care instrumentele sofisticate sunt în curs de dezvoltare rapid. Grupul Meta a estimat că dimensiunea pieței pentru piata de exploatare a datelor a crescut de la cca. 50 milioane dolari în 1996 la cca. 800 milioane dolari prin 2000 [157].

În cele ce urmează vom detalia elemente fundamentale legate de conceptul Data Mining [98], cum ar fi:

- Ce este sau ce nu este conceptul Data Mining?
- De ce acest concept Data Mining?
- Cum se *sapă* sau *minerește* în date?

- Probleme rezolvabile cu metode Data Mining.
- Despre modelare și modele în conceptul Data Mining.
- Aplicații ale conceptului Data Mining.
- Exemplificarea terminologică a conceptului Data Mining.
- Confidențialitatea datelor specifice.
- Numărul tot mai mare de baze de date de pe Web ce trebuie să fie exploataate pentru a extrage informații utile.

Data mining este procesul care prezintă diferite interogări și extragerea unor informații utile, modele și tendințe (de cele mai multe ori necunoscute anterior). În esență, pentru multe organizații, obiectivele data mining includ îmbunătățirea capacitaților de marketing, detectarea modelelor anormale și estimarea viitoare bazată pe experiențele trecute și tendințele actuale. Există în mod clar o nevoie destinată acestei tehnologii ([135], [156]).

Există cantități mari de date curente și istorice. Prin urmare, bazele de date devin mai mari, tot mai dificil de a sprijini procesul decizional. În plus, datele ar putea fi din mai multe surse și mai multe domenii. Există o nevoie clară de analiză a datele pentru a sprijini planificarea și alte funcții specifice unei întreprinderi.

Data mining a fost folosit până acum în afaceri de către organizații comerciale de succes în scopul de a obține avantaje critice în competiția lor. Se bănuiește ca în viitorul apropiat acest instrument va fi folosit pentru prelucrarea bazelor de date uriașe, ca de exemplu dosarele computerizate ale pacienților, la nivel național. De fapt, chiar și în prezent, prin identificarea procedurilor medicale ce au tendința de a se grupa, prin data mining putem prezice care pacienți vor folosi noile strategii de îngrijire a sănătății, putem defini modele de comportare ale pacienților de risc, putem identifica fraudele.

Conceptul de Data Mining [29] reprezintă *descoperirea și analiza cunoștințelor stocate în baze de date*, cu trei rădăcini generice, potrivit căror a împrumutat tehniciile de lucru și terminologia:

1. **Statistica** - reprezintă rădăcina cea mai lungă din Data Mining fără de care acest concept nu ar exista. Statistica clasică aduce în Data Mining unele tehnici definite rezumate în Analiza exploratorie a datelor (din denumirea științifică a literaturii străine EDA = Exploratory data analysis), utilizată pentru a identifica relațiile sistematice ale diferențelor variabile în momentul în care nu sunt informații suficiente

asupra naturii acestora. Dintre tehniciile analizei exploratorie a datelor clasice care sunt utilizate în Data Mining, sunt detaliate:

(a) **metodele computaționale:**

- *statistica descriptivă* (care face referire la repartiții și parametri);
- *statistici clasice* (se are în vedere media, mediana, deviația standard și.a.);
- *corelații* (reprezentând tabele multiple de frecvență);
- *metode statistice multivariate* (mai exact metoda clustering, analiza factorilor, a funcției discriminant și a corespondențelor, arbori de clasificare și de regresie, modele de regresie liniară și neliniară, serii temporale și.a.);

(b) **vizualizarea datelor**, reprezintă metoda forte de explorare a datelor prin reprezentarea informației în format vizual. Tehnicile de vizualizare sunt:

- histogramele;
- graficele rectangulare (anume box plots);
- graficele de împrăștiere (anume scatter plots);
- graficele de contur;
- graficele matriciale și.a [162].

2. **Inteligenta artificială (AI = Artificial Intelligence)** vine în ajutorul dezvoltării conceptului de Data Mining pe baza unor tehnici de procesare a informației bazate pe modelul rationamentului uman. Învățarea automată (conform literaturii de specialitate străină ML = machine learning) reprezintă o arie foarte de importantă a Inteligenței Artificiale în raport cu dezvoltarea Data Mining, prin utilizarea tehniciilor care permit computerului să poată *învăța* prin *antrenament*.

3. **Sisteme de baze de date** (conform literaturii de specilitate străină DBS = database systems), presupune cea de-a treia rădăcină a conceptului Data Mining, procurând materialul care trebuie *minerit* utilizând metodele amintite anterior.

Din perspectiva anumitor domenii necesitatea procesului de *minerit* al datelor poate fi rezumată astfel:

1. *Din punct de vedere economic* (afaceri-finanțe) - există un mare volum de date de-jă colectate în diverse domenii ca: date Web, e-commerce, super/hypermarket-uri, tranzacții finanțiar-bancare etc., care sunt pregătite pentru a fi analizate în vederea luării unor decizii optime;

2. *Din punct de vedere medical* - există la momentul de față diverse baze de date din domeniul sănătății (medico-farmaceutic), care au fost doar parțial analizate, cu mijloacele specifice medicinii și care conțin un nivel mare de informație încă neexplorată suficient;
3. *Din punct de vedere științific-cercetare* - se cunoaște o imensă bază de date din domeniul cercetării științifice din cele mai diverse domenii (mai exact astronomie, meteorologie, biologie, lingvistică și.a.) care nu pot fi explorate cu mijloacele tradiționale.

Domeniul computerelor și al Informaticii s-a dezvoltat exponential având în vedere faptul că există un imens volum de date neexplorate sistematic încă, precum și faptul că a crescut presiunea utilizării de noi metode pentru descoperirea informației *ascunse* în date, informație care este imposibil de detectat cu mijloacele tradiționale și folosind doar capacitatea umană de analiză.

Datorită *tinereții* sale, acest domeniu științific nou nu are încă o terminologie în literatură științifică străină și românească bine definită și acceptată de toți. Însă există oameni care sunt interesați de acest domeniu nou, acești oameni fără pregătirea necesară sunt dornici să cunoască mai mult pentru a-și forma o imagine cât de cât clară despre conceptul Data Mining și, mai ales, despre posibilitățile de aplicare a acestuia.

Potrivit celor relatate, unde s-a putut, s-a utilizat și terminologia internațională (în limba Engleză, în original) îndeobște folosită, pentru a ușura atât utilizarea resurselor bibliografice internaționale precum și căutarea pe Internet a termenilor de referință [58].

Având în vedere toate acestea, conceptul Data Mining nu poate fi definit ca o imagine care să ne furnizeze cât de cât fenomenul. Astfel vom încerca să redăm unele aspecte asemănătoare ale acestui concept:

- Presupune o practică a căutării automate de pattern-uri (modele, şabloane, tipare, forme etc.) în cadrul unor mulțimi mari de date, utilizând tehnici computaționale din Statistică, învățarea automată (machine learning) precum și recunoașterea formelor (pattern recognition);
- O extragere netrivială din date a informației implicate, necunoscută încă, și care poate fi folosită;
- Reprezintă știința extragerii de informație folosită din mulțimi mari de date sau baze de date;

- Presupune explorarea și analiza unor imense cantități de date, prin mijloace automate sau semi-automate, în vederea descoperirii de pattern-uri folositoare;
- Procesul descoperirii automate a informației - identificarea patternurilor ascunse și relațiilor cu datele [165], [47], [59], [60], [112].

Astfel conceptul Data Mining presupune procesul de căutare a unui *ac într-un car cu fân*, folosind un *senzor de metale*, cu scopul de a mări viteza de căutare, *automatizând* procesul respectiv.

Vom prezenta ulterior patru situații reale [112], ce ilustrează elovent ce nu este Data Mining prin comparație cu ce ar putea fi.

- *Ce anume nu poate fi reprezentat de Data Mining?*
 - Căutarea unui anumit număr de telefon într-o carte de telefoane;
 - Căutarea unei anumite informații (e.g. despre bucătărie pe Google).
- *Ce anume ar putea fi reprezentat de Data Mining?*
 - Căutarea unor nume de tipul ...escu într-un stat din SUA;
 - Gruparea laolaltă a informațiilor similare în funcție de context (e.g. despre bucătăria franceza, italiană etc., găsite pe Google).

Pentru a înțelege mai bine acest concept un exemplu mai concludent pentru sublinierea diferenței dintre ceea ce reprezintă uzual o căutare într-o bază de date și conceptul Data Mining reprezintă următorul conținut: *Cineva poate fi interesat de diferența între numărul de cumpărături de un anumit tip (e.g. electrocasnice) din cadrul unui supermarket în comparație cu un hypermarket sau, eventual, din cadrul unor supermarket-uri din două regiuni diferite.*

Astfel respectivul ia în considerație *a priori* ipoteza că există diferențe între un anumit supermarket și un hypermarket, sau între anumite vânzări dintre cele două regiuni.

În schimb în cazul Data Mining potrivit celor relatate o principală problemă poate consta de exemplu în identificarea factorilor ce influențează volumul vânzărilor, fără a avea ca bază o ipoteză considerată a priori.

În concluzie, metodele Data Mining încearcă să identifice pattern-uri și relații ascunse, care nu sunt mereu evidente. Potrivit exemplelor anterioare nu se poate pune semnul egal între o căutare/cercetare individuală a unui anumit obiect fără a se ține cont de natura sa și cercetarea de tip Data Mining care nu *caută* individualismul ci mulțimi de individualisme

care, într-un mod sau altul, pot fi grupate după diferite criterii. Diferența dintre o căutare simplă și un anumit proces Data Mining este aceea dintre căutarea unui anumit copac și identificarea unei anumite păduri.

Obiectivele Data Mining [42] pentru a desluși mai clar aria sa de aplicabilitate sunt:

1. *Obiective predictive* (presupune utilizarea unei părți din variabile pentru a prognoza una sau mai multe dintre celelalte variabile):

- Clasificarea;
- Regresia;
- Detecția deviațiilor/anomalialor.

2. *Obiective descriptive* (presupune identificarea de pattern-uri ce descriu datele și care pot fi înțelese de utilizator):

- Clustering (clusterizare);
- Descoperirea regulilor de asociere;
- Descoperirea de pattern-uri secvențiale.

1. Pentru ce utilizăm Data Mining?

Pentru a răspunde la o asemenea întrebare fără o prezentare prealabilă a tehnicielor Data Mining este necesar o prezentarea a două situații complet diferite poitrivit căreia Data Mining a fost utilizată cu succes. Pentru început se are în vedere situația privind rolul Data Mining în rezolvarea unei probleme fundamentale, din nefericire, a zilelor noastre [165].

O utilizare recentă a Datei Mining se referă la următorul fapt: *Serviciile de spionaj ale Armatei SUA ar fi descoperit cu un an înainte, pe baza utilizării tehnicielor Data Mining, indicii privind atentatele din septembrie, potrivit căreia s-ar fi identificat liderul grupului de teroriști care au pus la cale atentatul sinucigaș*. Cu părere de rău însă informațiile nu au fost luate în seamă de autorități.

Urmează o întâmplare năștimă, dar neplăcută pentru cel în cauză [59]. Pe scurt, un individ povestește că fiind plecat de casă, a fost sunat de compania telefonică la care era client pentru a-i aminti faptul că se presupunea că i-a fost furată cartela telefonică care a și fost folosită într-un mod fraudulos. Compania telefonică ajunsese la această concluzie pe baza unei analize a ultimelor con vorbiri telefonice efectuate utilizând

cartela respectivă, con vorbiri către unele locații care nu se potriveau cu pattern-ul pe care compania telefonică îl construise pentru clientul său.

Acestea sunt două exemple solide pentru a lua în considerare cu toată seriozitatea acest domeniu, pe cât de fascinant, pe atât de complex, al descoperirii de informații acolo unde cunoaștințele umane nu mai sunt de folos.

Există un număr mare de companii care au ca obiect principal de activitate Data mining [59]. Acest lucru se datorează mai ales cererii tot mai mari de servicii de Data Mining de către piața finanțier-economică (a se vedea, spre exemplu, domeniul ca: Business intelligence (BI), Business Performance Management (BPM), Customer Relationship Management (CRM) etc.), și piața de medicină (Health Informatics, e-Health etc.), fără a neglija și alte domenii de interes precum telecomunicațiile, meteorologia, biologia și.a.

Pornind de la pro gnoza în domeniul marketingului potrivit marilor companii transnaționale și trecând prin analiza tendințelor de tranzacționare a acțiunilor de la Bursă, realizarea profilului clientului fidel, modelarea cererii de produse farmaceutice, automatizarea diagnosticării cancerului, precum și detectarea fraudelor bancare și a uraganelor, clasificarea stelelor și galaxiilor și.a, se observă o paletă diversă de domenii potrivit căreia tehniciile Data Mining sunt folosite eficient, fapt ce dă un răspuns clar întrebării: de ce se utilizează Data Mining?

Însă nu trebuie considerat faptul că Data Mining poate rezolva orice problemă de descoperire de informație utilă în date. Identic ca și mineritul originar, și în cazul Data Mining este posibil să se *sape în muntele* de date chiar dacă nu se dă peste filonul de aur al cunoașterii. Descoperirea de cunoștințe sau de informații utile depinde foarte mult de anumiți factori, începând cu *muntele* de date în care se *minerește* și finalizând cu *uneltele* de Data Mining utilizate și, desigur, să nu uităm priceperea *minerului*. Dacă nu există *aur* în munte, degeaba se sapă. Filonul de aur, dacă există, este necesar să fie identificat și evaluat corect și concret, iar apoi, dacă este eficientă exploatarea sa, este necesar să fie executată cu unelte de minerit corespunzătoare.

2. Cum se *minerește* în date?

Procesul de *minerit* în date, are la bază trei etape definitorii ai procesului de Data Mining:

- Explorarea datelor, presupune *curățarea* și transformarea datelor, selectarea de sub mulțimi de date și a caracteristicilor, acolo unde avem un mare număr de

variabile §.a.;

- Construirea modelului și validarea acestuia face referire la considerarea diferențelor modele și alegerea aceluia cu cea mai bună performanță în prognoză - evaluarea competitivă a modelelor (conform literaturii străine denumit *competitive evaluation of models*);
- Aplicarea modelului la date noi, în vederea producerii de programe sau estimării corecte pentru problemele cercetate.

Conform specialistului R. Groth [59], se pot identifica cinci etape ale *mineriei datelor* în date:

- (a) **Pregătirea datelor (data pre-processing).** Înainte de utilizarea unei tehnici de Data Mining, este absolut necesară pregătirea datelor *brute* (raw data) în vederea analizării eficiente. Se cunosc mai multe trăsături ale pregătirii inițiale a datelor înaintea procesării propriu-zise cu ajutorul tehniciilor Data Mining [112]. Se are în vedere pentru început problema calității datelor (și anume existența zgomotului (noise), a valorilor extreme sau aberante (outlier/anomaly), a valorilor lipsă (missing values), a datelor dupicate, introduse incorect sau expirate și. a.).

În concordanță cu problemele detectate privind calitatea datelor, se are în vedere rezolvarea acestora cu metode specifice. Luând ca exemplu existența *zgomotului*, mai exact acele distorsiuni ale valorilor (măsurătorilor) originale produse de perturbări aleatoare, se utilizează diferite tehnici de *filtrare*, care au ca scop îndepărțarea sau reducerea efectului distorsiunilor.

Un alt exemplu este cel al procesării semnalelor (signal processing), în afara filtrelor electronice (hard), filtrele *matematice* (soft), adică algoritmi matematici care sunt folosiți în vederea modificării componente armonice a semnalului. În concluzie, dacă există unele valori extreme, adică cele care se abat foarte mult de la media celorlalte valori ale datelor, se are în vedere îndepărțarea acestora, sau utilizarea unor parametri (statistici) care să nu fie aşa de sensibili la aceste valori extreme (utilizarea medianei în locul mediei, care este *sensibilă* la valorile extreme).

Cazul valorilor lipsă este foarte des întâlnit în practica Data Mining și are o multitudine de cauze. Astfel, se utilizează diferite metode, ca: eliminarea obiectelor sau instanțelor care au atribută lipsă, estimarea valorilor lipsă și înlocuirea