

# **TEHNICI DE CLASIFICARE AUTOMATĂ A DATELOR**

**Utilizarea rețelelor neuronale în stabilirea ratingului de țară**

**Această lucrare este rezultatul proiectului „Doctorat și doctoranzi în triunghiul educație – cercetare - inovare (DOC-ECI)”, POSDRU/6/1.5/S/11, proiect cofinanțat din Fondul Social European prin Programul Operațional Sectorial Dezvoltarea Resurselor Umane 2007-2013 și coordonat de Academia de Studii Economice din București**

**ANA MARIA MIHAELA IORDACHE**

**TEHNICI DE CLASIFICARE  
AUTOMATĂ A DATELOR**

**Utilizarea rețelelor neuronale în stabilirea ratingului de țară**



**EDITURA UNIVERSITARĂ  
București**

Colecția ȘTIINȚE ECONOMICE

Referenți științifici: Prof. univ. dr. Elena-Carmen Pintilescu, Universitatea "Alexandru Ioan Cuza" din Iași, Facultatea de Economie și Administrarea Afacerilor  
Prof. univ. dr. Spircu Liliana, Academia de Studii Economice din București, Catedra de Cibernetică Economică  
Prof. Univ. Dr. Adrian Victor Bădescu, Academia de Studii Economice din București, Departamentul de Informatică și Cibernetică Economică  
Prof. univ. dr. Aivaz Kamer-Ainur, Universitatea „Ovidius” Constanța, Facultatea de Științe Economice

Redactor: Gheorghe Iovan  
Tehnoredactor: Ameluța Vișan  
Coperta: Monica Balaban

Editură recunoscută de Consiliul Național al Cercetării Științifice (C.N.C.S.) și inclusă de Consiliul Național de Atestare a Titlurilor, Diplomelor și Certificatelor Universitare (C.N.A.T.D.C.U.) în categoria editurilor de prestigiu recunoscut.

**Descrierea CIP a Bibliotecii Naționale a României**

**IORDACHE, ANA MARIA MIHAELA**

**Tehnici de clasificare automată a datelor : utilizarea rețelelor neuronale în stabilirea ratingului de țară / Ana Maria Mihaela Iordache. - București : Editura Universitară, 2022**

Conține bibliografie  
ISBN 978-606-28-1435-9

004

DOI: (Digital Object Identifier): 10.5682/9786062814359

© Toate drepturile asupra acestei lucrări sunt rezervate, nicio parte din această lucrare nu poate fi copiată fără acordul Editurii Universitare

Copyright © 2022  
Editura Universitară  
Editor: Vasile Muscalu  
B-dul. N. Bălcescu nr. 27-33, Sector 1, București  
Tel.: 021.315.32.47  
www.editurauniversitara.ro  
e-mail: redactia@editurauniversitara.ro

Distribuție: tel.: 021.315.32.47/ 0745 200 718/ 0745 200 357  
comenzi@editurauniversitara.ro  
www.editurauniversitara.ro

## CUPRINS

|  |            |
|--|------------|
| Introducere .....  | 7          |
| <b>Capitolul 1. Aspecte generale privind <i>data mining</i> și modelarea pe baza rețelelor neuronale.....</b>  | <b>11</b>  |
| 1.1 Istoric și definirea termenului de <i>data mining</i> .....  | 11         |
| 1.2 Metode și tehnici de <i>data mining</i> .....  | 14         |
| 1.3 Definiție și forma generală a unei rețele neuronale .....  | 20         |
| 1.4 Algoritmi de instruire a rețelelor neuronale .....   | 23         |
| 1.5 Rețele hebbiene. Factor de uitare.....   | 29         |
| <b>Capitolul 2. Problema clasificării obiectelor în abordarea ratingului de țară .....</b>   | <b>30</b>  |
| 2.1 Stadiul actual de cercetare al problemei clasificării obiectelor.....  | 30         |
| 2.2 Rețeaua neuronală în clasificare .....   | 31         |
| 2.3 Algoritmul de instruire al rețelei de tip perceptron multistrat.....   | 33         |
| 2.4 Algoritmul de instruire al perceptronului multistrat cu un strat ascuns de neuroni.....  | 36         |
| 2.5 Metode de calcul ale ratingului de țară utilizate de diverse agenții de rating.  | 39         |
| <b>Capitolul 3 Clasificarea țărilor din Uniunea Europeană în funcție de indicatorii economici, sociali și financiari din prisma analizei datelor .....</b> | <b>47</b>  |
| 3.1 Definiția problemei și a metodologiei de clasificare a țărilor .....   | 47         |
| 3.2 Clasificarea țărilor în funcție de grupuri diferite de indicatori .....  | 49         |
| 3.2.1 Indicatori economici .....   | 49         |
| 3.2.2 Indicatori sociali.....  | 57         |
| 3.2.3 Indicatori financiari.....   | 64         |
| <b>Capitolul 4. Clasificarea țărilor prin intermediul rețelei neuronale .....</b>  | <b>70</b>  |
| 4.1 Stabilirea (alegerea) indicatorilor pentru rețeaua neuronală .....   | 70         |
| 4.2 Asocierea unei rețele neuronale problemei de clasificare a țărilor.....  | 74         |
| 4.3 Concluzii legate de rezultatul modelării cu ajutorul rețelei neuronale.....  | 97         |
| 4.4 Recomandări pentru cercetări viitoare (algoritmi genetici).....  | 100        |
| <b>Concluzii.....</b>  | <b>105</b> |
| <b>Anexe .....</b>   | <b>113</b> |



## INTRODUCERE

O dată cu apariția calculatoarelor și cu mărirea volumului de date a apărut necesitatea de a identifica cunoștințe noi, necunoscute până în acel moment într-un timp relativ scurt.

În acest context, *data mining*, cunoscută și ca „descoperirea cunoștințelor în baze de date mari” este un instrument modern și puternic al tehnologiei informației și comunicațiilor, instrument ce poate fi folosit pentru a extrage informații utile dar încă necunoscute.

Acest lucru automatizează procesul de descoperire al unor relații și combinații în datele brute, iar rezultatele găsite ar putea fi încadrate într-un sistem automat de suport al deciziei.

Orice activitate umană, atât la nivel microeconomic, cât și la nivel macroeconomic se desfășoară în condiții de risc și incertitudine. Unele dintre acestea pot fi evitate mai ușor sau mai greu în funcție de nivelul de cunoaștere, de gradul de evaluare sau de importanța care li se acordă în fundamentarea deciziilor.

Cercetarea a avut ca obiective principale: prezentarea cât mai exhaustivă a tehnicilor de clasificare automată a datelor având ca suport de decizie rețelele neuronale și elaborarea unei metodologii de evaluare a ratingului de țară în funcție de mai mulți indicatori aleși din diverse grupe: indicatori economici, indicatori financiari și indicatori sociali. Modelul de clasificare propus are aplicabilitate atât la nivel macroeconomic, cât și la nivel microeconomic.

La nivel microeconomic, în contextul globalizării, derularea de afaceri de către agenții economici pe piețele externe se va realiza numai dacă există un stimulente suficient de puternic, în măsură să motiveze companiile să-și asume riscurile implicate desfășurării activității într-o anumită țară. Ratingul de țară calculat de diverse instituții este un indicator agregat foarte important atunci când se analizează oportunitatea de a investi sau nu în acea zonă.

La nivel macroeconomic, analiza riscului de țară presupune identificarea dificultăților care pot să apară în onorarea de către un anumit stat a obligațiilor care decurg din angajamentele sale luate pe plan extern.

Lucrarea de doctorat este alcătuită din cinci capitole. În primul capitol intitulat „Aspecte generale privind data mining și modelarea pe baza rețelelor neuronale” am prezentat un istoric și o definiție a termenului data mining, o scurtă sinteză a metodelor și tehnicilor data mining, punând accent pe modelarea bazată pe rețelele neuronale.

Dezvoltarea calculatoarelor și implicit a puterii de calcul precum și a volumului mare de date au condus la necesitatea de a descoperi informații noi într-un timp relativ scurt. În acest context, începând cu anii 2000 s-a dezvoltat o nouă tehnologie denumită *data mining*. Metodele *data mining* provin din statistică, administrarea bazelor de date și din inteligența artificială (rețele neuronale, analiza datelor, procesarea imaginilor, învățarea asistată de calculator, algoritmi genetici, etc.). Printre tehnicile de *data mining* cele mai des întâlnite se numără: excluderea, clasificarea (clusterizarea), discriminarea și previziunea.

Rețelele neuronale sunt utilizate pentru soluționarea unor probleme ce nu pot fi rezolvate cu ajutorul algoritmilor convenționale, cum ar fi probleme de optimizare, probleme de clasificare, etc. O rețea neuronală este alcătuită din neuroni, grupați în straturi, fiecare

neuron dintr-un anumit strat fiind conectat la toți neuronii din stratul anterior și la toți neuronii din stratul ce urmează, cu excepția straturilor de intrare și ieșire. Cel mai important și cel mai des folosit algoritm de antrenare a unei rețele neuronale este algoritmul back-propagation.

În capitolul doi al lucrării intitulat „Problema clasificării obiectelor în abordarea ratingului de țară” am realizat o trecere în revista a stadiului actual al cercetării în domeniul clasificării obiectelor și am studiat rețeaua neuronală în problema de clasificare.

Rețelele neuronale sunt flexibile, se pot adapta cu ușurință datelor de antrenament și pot gestiona mai bine „zgomotul” și schimbările apărute în exemplele de învățare. Odată încheiat procesul de instruire (antrenare), o rețea neuronală poate genera un *output* rezonabil pentru datele de intrare care nu i-au mai fost prezentate în faza de învățare.

În domeniul clasificării obiectelor cel mai utilizat model este perceptronul multistrat. Acest model presupune existența unei rețele neuronale conectată în totalitate, având neuronii dispuși pe mai multe straturi astfel: strat de intrare, unul sau mai multe straturi ascunse și un strat de ieșire. Printre avantajele perceptronului multistrat se numără: eficiența în reamintire și eficiență mare în clasificare. Timpul îndelungat de instruire al acestui tip de rețea neuronală limitează oarecum aplicarea modelului. Mai departe am prezentat pe scurt metodologia folosită de diverse agenții pentru evaluarea ratingului de țară și a modului de acordare al calificativului pentru un anumit stat.

Pornind de la definiția conform căreia riscul reprezintă variabilitatea rezultatului posibil în funcție de un eveniment nesigur, incert, există mai multe agenții de rating, fiecare dintre acestea având propriul model de evaluare al riscului și oferind diverse produse pentru diferite segmente de utilizatori. Astfel, principalele agenții de rating existente sunt: Institutional Investor, Standard & Poor's, Political Risk Services, The Economist, Moody's și altele.

La realizarea aplicației am utilizat următoarea metodologie: alegerea grupului de țări asupra căruia se realizează aplicația (țările din Uniunea Europeană: Austria, Danemarca, Suedia, Finlanda, Germania, Cipru, Letonia, Polonia, Lituania, Olanda, Spania, Portugalia, Franța, Italia, România, Grecia, Marea Britanie, Bulgaria, Belgia, Ungaria, Malta, Slovacia, Cehia, Estonia, Slovenia, Luxemburg, Irlanda); stabilirea indicatorilor, înregistrarea valorilor acestora (pentru anul 2009) și aplicarea tehnicilor de analiză a datelor (analiza în componente principale, analiza cluster și analiza factorială) asupra fiecărei grupă de indicatori; extragerea indicatorilor care aduc aport informațional maxim din fiecare grupă; realizarea și instruirea rețelei neuronale; interpretarea rezultatelor obținute și realizarea de previziuni (simulări) asupra situației României.

Partea a treia a lucrării prezintă aplicabilitatea tehnicilor de *data mining*, mai precis extragerea informațiilor ”ascunse”, ”importante” din listele cu indicatorii țărilor înregistrați în diferite planuri (economic, social și financiar) folosind astfel tehnici de analiză a datelor (analiză în componente principale, analiză factorială și analiză cluster).

Seturile de indicatori utilizați în aplicație sunt următorii:

- Indicatori economici: consum final (% din PIB), cererea internă (% din PIB), consumul final al gospodăriilor (% din PIB), formarea brută de capital (% din PIB), formarea brută de capital fix (% din PIB), exporturi (% din PIB), importuri (% din PIB), valoarea adăugată brută (prețuri curente) (% din PIB), impozite fără subvenții pe produs (% din PIB), salarii (% din PIB), formarea brută de capital fix de către sectorul privat (% din PIB);

- Indicatori sociali: pondere angajați cu contract parțial de muncă în total angajați, ponderea populației angajate în total populație, ponderea populației în vârstă 0-14 ani în total populație, ponderea populației în vârstă 15-24 ani în total populație, ponderea populației în vârstă 25-49 ani în total populație, ponderea populației în vârstă 50-64 în total populație,



ponderea populației în vârstă 65-79 în total populație, ponderea populației în vârstă de peste 80 ani în total populație, rata șomajului, vârsta medie a populației;

- Indicatori financiari: rata anuală medie a inflației (raportată la indicii de prețuri de consum), deficit public fiscal (% din PIB), investiții publice fixe (% din PIB), datoria publică (guvernamentală) % din PIB, soldul contului curent (% PIB), împrumuturi guvernamentale externe (% din PIB).

Pentru fiecare set de indicatori am aplicat tehnici de analiză a componentelor principale, tehnici de analiză cluster și analiză factorială pentru o primă clasificare a țărilor în funcție de indicatorii din fiecare grupă.

Ultimul capitol al lucrării intitulat "Clasificarea țărilor prin intermediul rețelei neuronale" prezintă asocierea unei rețele neuronale pentru problema de clasificare a țărilor în funcție de indicatorii care aduc maximum de informație din cele trei grupe de indicatori analizate în capitolul patru.

Pentru a asocia o rețea neuronală de tip perceptron multistrat problemei analizate este necesar să se cunoască mai întâi care sunt indicatorii care aduc maximum de informație din fiecare grupă analizată și apoi în câte clase se împart țările în funcție de acești indicatori.

În urma aplicării tehnicilor de analiză a datelor indicatorii care aduc aport informațional maxim din fiecare grupă sunt: cererea internă, formarea brută de capital, exporturi (din domeniul economiei reale); ponderea populației angajate în total populație, ponderea populației în vârstă de 15-24 ani în total populație, vârsta medie a populației (din plan social); deficit public fiscal, investiții publice fixe, împrumuturi guvernamentale externe (din domeniul financiar).

După instruirea rețelei neuronale s-a constatat că sunt țări care și-au schimbat clasa în care erau inițial. Acest lucru a fost posibil datorită faptului că statele din Uniunea Europeană au fost împărțite în trei părți: date de antrenare (80% din totalul țărilor), date de validare și date de testare.

Tot în acest capitol se prezintă anumite previziuni (scenarii) de îmbunătățire a performanțelor economice, sociale și financiare ale României cu scopul de a accede în clasa țărilor dezvoltate. Previziunile au fost obținute pe baza scorurilor și a funcțiilor de transfer prin rețeaua neuronală deja antrenată.

La nivelul Uniunii Europene este importantă stabilirea de politici economice, sociale și/sau financiare cu scopul evitării dezechilibrelor macroeconomice și urmărirea realizării unei creșteri sustenabile sprijinită, acolo unde este necesar, de adoptarea unor reforme structurale. Pentru realizarea acestui lucru, în ultimii ani se vorbește tot mai mult de un nou concept numit macroprudențialitate. Dar oricât de eficiente ar fi politicile macroprudențiale, acestea nu pot înlocui politicile macroeconomice, sociale sau financiare adoptate la nivelul fiecărui stat.



## CAPITOLUL 1

### ASPECTE GENERALE PRIVIND *DATA MINING* ȘI MODELAREA PE BAZA REȚELELOR NEURONALE

Termenul *data mining* presupune analiza datelor din diverse puncte de vedere (aspecte) cu scopul de a extrage cunoștințe pentru a le folosi mai departe în fundamentarea deciziilor la nivel micro sau macroeconomic. Din cele mai vechi timpuri oamenii au fost atrași de funcționarea creierului uman și de modul de adoptare a deciziilor de către acesta. În acest sens au început să studieze rețelele neuronale, să elaboreze diverși algoritmi de instruire a acestora. Datorită volumului mare de informație rezultat pe baza modelării neuronale au fost dezvoltate și rețelele hebbiene, principala lor caracteristică fiind introducerea în calcul a unui factor de uitare cu scopul de a păstra mai departe informația cu adevărat importantă.

#### 1.1 Istoric și definirea termenului de *data mining*

*Data mining* s-a dezvoltat ca o consecință a dezvoltării bazelor de date. Colectarea datelor în diverse formate electronice a început în anii '60 permițând o analiză retrospectivă a datelor prin intermediul calculatorului. Bazele de date relaționale au apărut în anii '80 împreună cu Structured Query Language (SQL). Anii '90 sunt caracterizați de o explozie a informațiilor. Pentru stocarea lor au început să se folosească depozitele de date (data warehouses). *Data mining* a apărut ca răspuns la provocările cu care s-au confruntat comunitatea specialiștilor în baze de date, care se ocupau cu cantități masive de date, aplicarea analizei statistice și aplicarea tehnicilor de căutare, specifice inteligenței artificiale, asupra datelor.<sup>1</sup>

Principalele domenii care pot fi considerate surse de date sunt:

- mediul economic;
- telecomunicații (în telefonie se generează anual aproape 20 exabytes de date);
- sateliți; internet – ordin de mărime: terabytes;
- biblioteci: (Library of Congress 20tb-3pb date);
- agenții legale (baza de date cu amprente FBI: 1 petabyte);
- experimente științifice.

*Data mining* este aplicat într-o varietate de domenii, începând cu managementul de investiții până la astronomie. Importanța și potențialul de aplicare al *data mining* a fost recunoscut în marketing, domeniul bancar, asigurarea sănătății, telecomunicații, pentru aplicații cum ar fi analiza coșului de piață, pentru promovarea eficienței, analiza vulnerabilității clienților, managementul relațiilor cu clienții, crearea de portofoliu, detectarea fraudei în telefonia celulară etc.

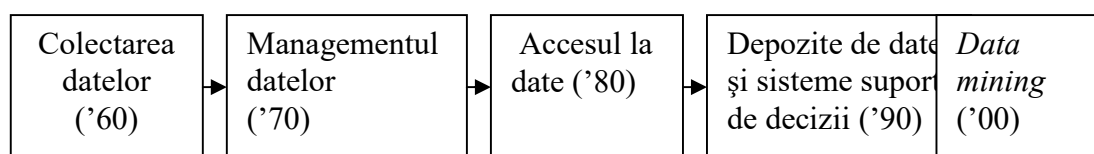
În 1995, Gartner Group Advanced Technology Research Note a poziționat *data mining* și inteligența artificială pe primul loc între cele cinci zone tehnologice cheie care "vor

---

<sup>1</sup> Ionita Angela, Asupra termenului de minerit de date (data mining), [http://ria.ici.ro/ria2005\\_2/art03.html](http://ria.ici.ro/ria2005_2/art03.html)

avea în mod clar impact peste un spectru larg de domenii industriale în următorii trei până la cinci ani".<sup>2</sup>

**Data mining** este un pas în procesul KDD (knowledge discovery in data – descoperirea cunoștințelor din date), care constă în aplicarea de algoritmi de descoperire și analiză de date, care sub limitele eficienței de calcul vor produce o mulțime de modele pe datele studiate (figura 1).



**Figura 1.** Evoluția tehnologiei spre *data mining*

**Data mining** reprezintă o arie de intersecție între învățarea asistată de calculator, statistică și baze de date.<sup>3</sup>

**Data mining** reprezintă procesul de selecție, explorare și modelare a unor seturi mari de date pentru a descoperi modele necunoscute, utile în dezvoltarea afacerilor.<sup>4</sup>

Fayyad U., Piatetsky-Shapiro G., Smyth P afirmă că „**data mining** este procesul de identificare de cunoștințe valide, noi, potențial utile și, în final, inteligibile din baze de date care sunt folosite în luarea deciziilor hotărâtoare în domeniul afacerilor, mai precis, „**data mining** reprezintă aplicarea unor algoritmi specifici pentru extragerea modelelor din date”.<sup>5</sup>

Din punct de vedere al dicționarului explicativ al limbii române, **data mining** este o expresie improprie care asociază procesul de căutare și găsimă a cunoștințelor dintr-o cantitate mare de date cu procesul de extragere a mineralelor din roci. Ambele procese se referă la extragerea esențialului din domenii diferite.

În conformitate cu Carta Alba a soluțiilor de management de date de la IBM, **data mining** este procesul de extragere de informații valide, necunoscute anterior și, în final, inteligibile din baze mari de date, informații folosite în luarea deciziilor hotărâtoare în domeniul afacerilor. Extragerea informațiilor poate fi utilizată la formarea de modele de predicție sau de clasificare, pentru a identifica relații între înregistrările din bazele de date sau pentru a furniza un sumar al bazelor de date care sunt minerite. **Data mining** constă dintr-un număr de operații, fiecare fiind suportată de o varietate de tehnici cum sunt reguli de inducție, rețele neuronale, clusterizare conceptuală, descoperire asociativă etc. În multe domenii din lumea reală cum ar fi analiza de marketing, analiza financiară, detectarea fraudei, informațiile extrase necesită utilizarea cooperativă a mai multor operații și tehnici de minerit de date.<sup>6</sup>

Toate cele trei aspecte provin din tehnologiile de **data mining**, specifice diferitelor domenii, și se ocupă de problema ușurinței în utilizare prin ascunderea complexității algoritmilor fundamentali de **data mining**.

„...Pentru a concluziona, instrumentele de **data mining** sunt orientate către a obține abilități de analiză multimedia și capacități de analizare simultană a numeroase tipuri de baze de date. Utilizatorii finali sau specialiștii în **data mining** vor fi persoane care operează cu

<sup>2</sup> <http://www.thearling.com/text/dmwhite/dmwhite.htm>

<sup>3</sup> Marcel Holsheimer: *Data Mining by Business Users: Integrating Data Mining in Business Processes*. KDD Tutorial Notes, 1999

<sup>4</sup> SAS Institute

<sup>5</sup> Fayyad U.M., Piatetsky-Shapiro G., Smyth P., *Knowledge Discovery in Database*, AAAI Press, 1996

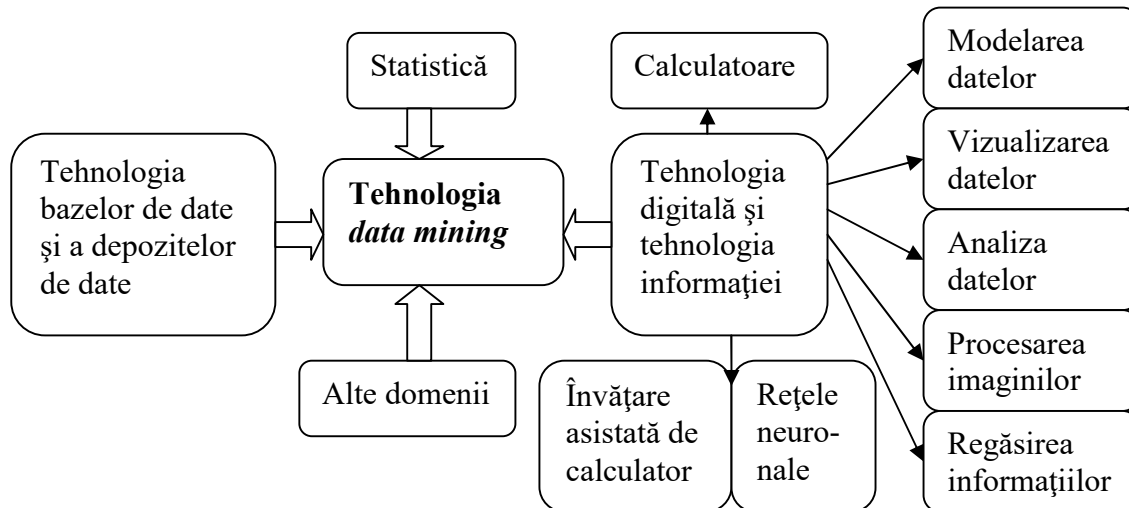
<sup>6</sup> IBM's Data Mining Technology - White Paper data Management Solutions, April, 1996

instrumentele care sunt înglobate în pachetele software standard, cel mai probabil utilizând tehnologii de rețele neuronale."<sup>7</sup>

Alte definiții ar putea fi:

- „*data mining* este un proces de descoperire de relații, șabloane și cunoștințe din date;
- „*data mining* este procesul prin care informațiile și cunoștințele sunt extrase din volumele mari de date folosind tehnici care sunt mai mult decât o simplă căutare în date”.

Principalele discipline care se intersectează cu domeniul *data mining* sunt: statistica, programarea calculatoarelor, învățarea asistată de calculator, tehnologia bazelor de date, tehnologia digitală, tehnologia informației, etc (figura 2).



**Figura 2.** Discipline cu care se intersectează *data mining*

În concluzie, *data mining* a evoluat ca o tehnologie datorită a două fenomene complementare:

- volumul de date în continuă expansiune în urma dezvoltării tehnologiilor de baze de date și a instrumentelor de colecționare a datelor;
- nevoia de cunoaștere concretizată prin necesitatea de a filtra și interpreta toate aceste volume de date stocate în baze de date, depozite de date sau bănci de date.

Procesul de descoperire a cunoștințelor din cantități mari de date este un proces complex format, în principal, din următoarele **etape**, care se succed din punct de vedere cronologic:

- *curățirea datelor* (data cleaning) se realizează prin înlăturarea datelor inutile și a celor inconsistente;
- *integrarea datelor* (data integration) constă în combinarea datelor provenite din mai multe surse de date diferite;
- *selectarea datelor* (data selection) reprezintă extragerea datelor relevante pentru analiză din depozitele de date disponibile (baze de date, depozite de date, www etc.);
- *transformarea datelor* (data transformation) constă în punerea datelor în formate unitare, corespunzătoare pentru *data mining* (analiza în vederea descoperirii de cunoștințe) prin realizarea unor operații rezumative și/sau de agregare (de unificare);
- *data mining* este definită prin extragerea unor modele de date aplicând asupra datelor rezultate după parcurgerea etapelor anterioare și metode inteligente, denumite generic metode

<sup>7</sup> ATEL, L., S. PERRY, D. TAYLOR, A. BROWN, S. TIKE: The Future of Data Mining

*data mining*. Acest proces este esențial pentru descoperirea cunoștințelor utile „ascunse” în depozitele de date;

- *evaluarea modelelor* (pattern evaluation) se realizează prin evaluarea modelelor de date extrase pentru identificarea celor care reprezintă cunoștințele care interesează în mod real;

- *prezentarea cunoștințelor* (knowledge presentation) obținute utilizatorilor lor, prin folosirea unor tehnici de vizualizare și de reprezentare adecvate.

Metodele *data mining* provin din calculul statistic clasic, din administrarea bazelor de date și din inteligența artificială. Ele nu înlocuiesc metodele tradiționale ale statisticii, ci sunt considerate a fi extinderi ale tehnicilor grafice și statistice. Deoarece softului îi lipsește intuiția umană, rezultatele metodelor *data mining* vor trebui supuse în mod sistematic unei supravegheri umane.

Structura tipică de date potrivită pentru *data mining* conține observațiile plasate pe linii iar variabilele plasate pe coloane. Domeniile sau intervalele de valori pentru fiecare variabilă vor trebui să fie definite precis, evitându-se cât mai mult posibil exprimările vagi.

## 1.2 Metode și tehnici de *data mining*

Rezultatele obținute cu ajutorul tehnicilor de *data mining* pot varia și sunt specifice fiecărui tip de utilizator. În general, tehnicile de *data mining* se aplică într-o bază de date din două motive:

- validarea unei ipoteze privind datele, corectitudinea acestora, dar și ipoteze statistice;  
- descoperirea de noi caracteristici din date. Descoperirea, la rândul ei, poate fi împărțită în descriere și predicție. Descrierea datelor se realizează, fie prin calcularea de diverși indicatori statistici (elementari) din datele brute, cum ar fi media, dispersia, abaterea medie pătratică, etc., fie prin aplicarea asupra datelor a tehnicilor avansate de analiză, cum ar fi analiza cluster, analiza componentelor principale, analiza discriminantă, etc. Predicția are ca scop exprimarea cât mai exactă a unor valori viitoare pentru datele aflate în analiză.

Cele mai folosite tehnici de *data mining* sunt:

1. **Excluderea** ce presupune tratarea din punct de vedere informatic a datelor. Aceasta presupune identificarea și eliminarea înregistrărilor care conțin anomalii în comparație cu celelalte date (date aberante, înregistrări cu date lipsa, etc).

2. **Clasificarea (clusterizarea)** este operația prin care obiectele dintr-o mulțime dată sunt repartizate în submulțimi numite „clase” în funcție de asemănările și deosebirile dintre ele.

3. **Discriminarea** este diferită de clasificare din cauza necesității, în aplicații, a unor cunoștințe anterioare asociate claselor. De obicei este necesar un eșantion din fiecare clasă pentru a realiza o analiză discriminantă. Analiza discriminantă are următoarele obiective: investigarea diferențelor între grupuri, separarea eficientă a grupurilor, identificarea variabilelor importante discriminatorii, testarea ipotezelor legate de diferențele dintre grupuri, clasificarea de noi observații în grupuri pre-existente.

4. **Previzionarea** ce se realizează pe baza trendului.

Dacă ne referim la metodele „de învățare” (*data mining*) acestea pot fi grupate în două categorii: nesupervizate, respectiv supervizate.

Metodele de învățare nesupervizate includ următoarele:

1) Analiza componentelor principale (ACP). Scopul său este de a reduce dimensionalitatea datelor multivariate prin „integrarea” variabilelor corelate, transformând liniar variabilele inițiale în variabile necorelate între ele.

2) Analiza factorială (Factor Analysis) permite extragerea unui număr mic de factori ascunși care explică cea mai mare parte a variabilității comune și determină corelațiile observate între datele inițiale.

3) Analiza clasificării (Cluster Analysis) are ca scop de a grupa cazurile (observațiile) în clustere (grupuri, categorii).

Printre metodele de învățare supervizate se numără:

1) Regresia liniară multiplă având scopul de a descrie variațiile unei variabile  $y$  în funcție de valorile altor variabile prin intermediul unor funcții liniare sau liniarizabile. Variabila  $y$  se numește dependentă, iar celelalte poartă denumirea de variabile independente sau explicative.

2) Regresia logistică este un caz particular al regresiei liniare multiple. În acest tip de regresie „răspunsul” este o variabilă binară de tipul „da sau nu”.

3) Analiza discriminării este o tehnică de statistică multivariată frecvent utilizată pentru a construi un model descriptiv/predictiv al separării („discriminării”) pe grupuri bazat pe variabilele predictor observate și de a clasifica fiecare observație într-unul din grupuri. În analiza discriminantă sunt utilizate mai multe atribute cantitative pentru a separa o singură variabilă de clasificare.

4) Rețelele neuronale asociază unui input vectorial un anumit output vectorial. Rețelele neuronale oferă modele matematice de grupare neliniare în special pentru clasificarea unor obiecte având la bază conceptul de rețea bioelectrică din creierul uman formată de neuroni și sinapsele acestora. Principala trăsătură a acestor rețele este capacitatea de a învăța pe bază de exemple, folosindu-se de experiența anterioară pentru a-și îmbunătăți performanțele.

În scopul obținerii beneficii maxime din aplicarea algoritmilor de *data mining* asupra unui set de date este necesară asigurarea integrității și valabilității datelor. Principalele acțiuni care se au în vedere sunt<sup>8</sup>: curățirea, transformarea, sintetizarea, reducerea și discretizarea datelor. (Kennedy 1998, Pzle 1999)

**Curățirea datelor.** Datele din lumea reală sunt de obicei incomplete și afectate de zgomot. Curățirea lor presupune adăugarea valorilor lipsă, reducerea zgomotului și înlăturarea valorilor aberante. Valorile lipsă provin cel mai adesea din: disfuncționalități ale echipamentelor de colectare, neconcordanțe între datele într-un set de date sau omisiunea acestora, atunci când datele nu sunt înțelese sau sunt considerate banale. Practicile cele mai folosite pentru tratarea acestor cazuri sunt: ignorarea întregului tuplu din baza de date, completarea manuală a datelor lipsă, folosirea unui caracter special care arată o valoare lipsă, folosirea unor aproximări ale valorilor lipsă.

**Sintetizarea datelor** presupune efectuarea unei sinteze exhaustive a bazelor de date multidimensionale, băncilor de date sau a fișierelor într-o schemă unitară. Sintetizarea datelor se poate realiza pe trei niveluri (Han și Kamber, 2001): sintetizarea schemei bazei de date, detectarea și rezolvarea neconcordanțelor dintre valorile datelor (de exemplu să se păstreze aceeași scală a unității de măsură), managementul datelor redundante.

**Transformarea datelor.** Cele mai utilizate tehnici de transformare a datelor sunt: eliminarea zgomotului din date, agregarea, generalizarea, construirea de caracteristici (atribute) noi și normalizarea. Cele mai utilizate tehnici de normalizare, potrivit cercetătorilor Weiss și Indurkha sunt: (Weiss and Indurkha, 1998)

---

<sup>8</sup> Negnevitsky Michael, *Artificial Intelligence – A guide to intelligent systems, second edition*, Editura Pearson Education, England, 2005, ISBN 0321-20466-2

- normalizarea min-max presupune aplicarea unei transformări liniare asupra datelor:

$$x' = \frac{x - \min}{\max - \min}, x' \in [0, 1].$$

- normalizare după scorul z. Un atribut A este normalizat în funcție de medie și dispersie  $x' = \frac{x - \text{media}_x}{\text{dispersia}_x}$ .

- normalizare prin scalare zecimală se realizează prin aplicarea formulei:  $x' = \frac{x}{10^j}$ , unde j este cel mai mic întreg care satisface relația  $\max(|x'|) < 1$ .

Și să revenim asupra metodei de **clasificare**. Aceasta reprezintă procesul de găsim a proprietăților comune dintr-un set de date și clasificarea acestora în diferite clase conform unui model de clasificare. (Chen, 1996). Clasificarea presupune parcurgerea a două etape:

- existența unui set de învățare sau de antrenare folosit pentru a construi un model care clasifică datele în anumite grupuri (clase);
- folosirea modelului de clasificare construit pentru un nou set de date sau pentru a extrage reguli de clasificare noi.

Cele mai utilizate tehnici de clasificare sunt: clasificarea bayesiană, arbori de decizie și rețele neuronale.

**Clasificarea Bayesiană** presupune asignarea unui element x la una din clasele  $c_1, c_2, \dots, c_N$ , date, folosind un model de probabilitate definit conform teoremei lui Bayes.

O rețea bayesiană este un graf orientat ale cărui noduri reprezintă variabilele și ale cărui muchii sunt legăturile cauzale sau de influență dintre acestea. Fiecărui nod i se asociază un tabel de probabilități ce descriu relațiile dintre acesta și părintele său.

Strategiile de decizie Bayes se aplică în probleme de clasificare cu număr mare de clase cu scopul a minimiza erorile, având la bază probabilitatea condiționată. Probabilitatea condiționată este folosită pentru a măsura încrederea ca un eveniment aleator să aibă loc știind că alt eveniment aleator a avut loc. Fie două evenimente A și B, probabilitatea condiționată ca

evenimentul A să aibă loc știind că evenimentul B a avut loc este  $P(A|B) = \frac{P(A \cap B)}{P(B)}$ . Fie

$(\Omega, \Sigma, P)$  un spațiu de probabilitate, B un eveniment arbitrar din  $\Sigma$  și  $\{A_1, \dots, A_n\}$  o partiție a spațiului  $\Omega$ . Atunci,  $P(A_i|B) = \frac{P(B|A_i) \cdot P(A_i)}{P(B)} = \frac{P(B|A_i) \cdot P(A_i)}{\sum_{i=1}^n P(B|A_i) \cdot P(A_i)}$ , unde  $P(B) > 0$  și

$P(A_i) > 0, i=1, \dots, n$ .  $P(A_i|B)$  se numește probabilitate posterioară,  $P(A_i)$  este probabilitatea apriorică,  $P(B|A_i)$  poartă denumirea de verosimilitate. Scopul este de a minimiza probabilitatea de a greși.

Fie  $D_j$  regula de decizie referitoare la clasa  $\Omega_i$ . Fiind dat un vector X, eroarea relativă la clasa  $\Omega_i$  este definită de  $P\{\text{eroare}/x\} = q - P(\Omega_i|x)$ . Se minimizează probabilitatea de a greși. Conform regulii bayesiene de decizie se alege  $D_j$  dacă  $P(\Omega_j|x) > P(\Omega_i|x), i \in \{1, \dots, j-1, j+1, \dots, r\}$ .

Pentru un nou set de date care urmează a fi clasificate utilizând un clasificator bayesian, presupunem că fiecare atribut este o variabilă aleatoare. Fiind dat un obiect cu atributele  $\{A_1, A_2, \dots, A_p\}$  se dorește asignarea sa în clasa  $\Omega_i$ . Clasificarea este corectă atunci când probabilitatea condiționată  $P(\Omega_i|A_1, A_2, \dots, A_p)$  este maximă. Se calculează probabilitățile posterioare  $P(\Omega_i|A_1, A_2, \dots, A_p)$  pentru toate clasele  $\Omega_i$ . Se alege apoi clasa care are probabilitatea maximă.



**Arborele de decizie** este un arbore care conține în noduri câte un test pentru o anumită proprietate, fiecare arc având o valoare a proprietății testate în nodul din care pleacă arcul respectiv, iar în fiecare nod terminal o clasă. Arborii de decizie sunt folosiți pentru a selecta cea mai bună direcție de acțiune în situațiile în care apare incertitudinea. Nodurile frunză specifică decizia la care se ajunge pentru calea care duce la acel nod.

ID3 (**Iterative Dichotomiser 3**)<sup>9</sup> este un algoritm de construire inductivă a unor arbori de decizie. Euristica algoritmului ID3 măsoară câștigul informațional pe care îl aduce fiecare atribut și alege ca test acel atribut care maximizează acest câștig. Pentru o mulțime de mesaje  $M = \{m_1, m_2, \dots, m_n\}$  și o probabilitate  $p(m_i)$  de apariție a fiecărui mesaj, conținutul informațional al unui mesaj din  $M$  se definește astfel:  $I(M) = \sum_{i=1}^n -p(m_i) \log_2(p(m_i))$ .

Algoritmul ID3 folosește teoria informației pentru a selecta atributul care oferă cel mai mare câștig informațional în clasificarea exemplilor de învățare. ID3 alege atributul considerat cel mai important, și împarte exemplele în subseturi corespunzătoare valorilor posibile ale atributului. În continuare se apelează ID3 recursiv pentru subseturile obținute și lista de atribute rămase. Procedul recursiv se termină când toate exemplele dintr-un subset au aceeași clasificare. Dacă se termină toate atributele și subsetul conține încă exemple de clasificare diferite, atunci înseamnă că sunt exemple cu aceeași descriere dar cu clasificare diferită. Acesta se poate datora mai multor motive:

- o parte din date sunt incorecte;
- datele sunt corecte dar atributele sunt insuficiente;
- datele sunt corecte dar clasificarea implică un anumit grad de nedeterminism.

Pentru a rezolva problema valorilor și a mări viteza de execuție se poate aplica algoritmul C4.5. Acesta se bazează pe algoritmul ID3, ajutând la îmbunătățirea adoptării deciziei prin care derivă regulile de clasificare. Algoritmul C4.5 generează un arbore de decizie prin partiționarea recursivă a mulțimii de date, printr-o strategie de parcurgere în adâncime (“depth-first”). Algoritmul ia în considerare toate testele posibile pentru împărțirea mulțimii de date și selectează testul care aduce cel mai mare aport (câștig) informațional. Acuratețea clasificării și timpul de rezolvare pot fi acceptabile pentru baze de date mici, dar în cazul bazelor de date de dimensiuni mari ori timpul de execuție crește, ori acuratețea clasificării descrește considerabil.

Din punct de vedere informatic clasificarea presupune împărțirea unei baze de date în baze de date mai mici prin algoritmi de tipul „divide-et-impera” (împarte și domină), pe baza conceptului de similitudine. În procesul de clasificare se poate preciza sau nu de la început numărul de clase.

Pentru ca un proces de clasificare să fie corect trebuie îndeplinite următoarele condiții (Han și Kamber, 2001): similaritate între grupuri mare și similaritate mică în interiorul grupurilor.

Tehnicile de clasificare pot fi împărțite în funcție de următoarele trei criterii (Jain, 1999):

- în funcție de metoda identificării clusterilor pot fi clasificate în tehnici de partiționare ierahice, bazate pe densitate și pe rețele;
- în funcție de tipul de date cu care se lucrează, pot fi algoritmi statistici și conceptuali;
- în funcție de teoria folosită pentru a extrage clusterii, putem avea clusterizarea fuzzy, clusterizare „clară” și algoritmi care se bazează pe teoria rețelelor neuronale (Kohonen).<sup>10</sup>

<sup>9</sup> Negnevitsky Michael, Artificial Intelligence – A guide to intelligent systems, first edition, Editura Pearson Education, England, 2002, ISBN 0201-71159-1

<sup>10</sup> Negnevitsky Michael, Artificial Intelligence – A guide to intelligent systems, second edition, Editura Pearson Education, England, 2005, ISBN 0321-20466-2

În continuare voi prezenta pe scurt cațiva algoritmi de clasificare.

**Algoritmul *k-means*** este un instrument bine cunoscut și frecvent folosit în clusterizare. Presupunând  $n$  vectori de atribute  $x_1, x_2, \dots, x_n$  și  $k$  clustere,  $k < n$ , fixat. Fie  $m_i$  vectorul valorilor medii a vectorilor din clusterul  $i$ . Pentru separarea clusterelor se poate utiliza un clasificator bazat pe distanța minimă. Se poate spune că un *item* (obiect)  $x$  aparține clusterului  $i$  dacă  $\|x - m_i\|$  este minimă dintre toate cele  $k$  distanțe posibile.

S-a constatat că rezultatele obținute prin *k-means* depind de alegerea valorilor medii inițiale ale clusterelor. De aceea, o variantă pentru îmbunătățirea performanțelor algoritmului este rulara sa multiplă și alegerea partiționării celei mai bune. O măsură a calității clusterizării este dată de calcularea entropiei claselor rezultate. Din mai multe variante de clusterizare, varianta optimă are cea mai mare sumă a entropiilor, ceea ce arată faptul că

datele din clustere sunt cât mai omogene  $H_c = -\sum_{i=1}^k p_i \log p_i$ ,  $K^* = \arg \max_c H_c$ , unde  $H_c$  este

entropia unei clusterizări,  $p_i$  sunt probabilitățile de apartenență ale valorilor la un anumit cluster ( $c$ ), iar  $K^*$  este strategia optimă de clusterizare. Printre dezavantajele pe care le are algoritmul *k-means* se numără următoarele: se aplică numai seturilor de date numerice, necesită specificarea numărului de clustere în avans, are capacitate redusă de a lucra cu date aberante sau date care sunt afectate de zgomot, precum și incapacitatea de a descoperi clustere cu forme neconvexe. Pentru a rezolva aceste inconveniente ale lui *k-means* s-au dezvoltat alți algoritmi. Algoritmul CLARANS este un algoritm tipic care folosește medoidii și care se bazează pe algoritmi PAM și CLARA. Un medoid este elementul central dintr-o clasă. Acest algoritm se folosește pentru a identifica  $k$  clase optime. O dată ce s-au identificat cei  $k$  medoidi, datele rămase sunt asignate grupului corespunzând celui mai aproape medoid. Apoi se verifică dacă nu se poate înlocui un medoid cu oricare dintre vecini. Dacă orice alt punct îmbunătățește calitatea grupului, atunci acesta devine noul medoid și procesul se repetă. Dacă nu, se alege un optim local și se aleg aleator alte date pentru a găsi un optim global. Complexitatea algoritmului CLARANS crește liniar cu numărul de obiecte.

Pentru a manipula baze de date spațiale au fost creați alți doi algoritmi, care au la bază algoritmul CLARANS: SD (CLARANS) și NSD (CLARANS), utilizatorul având posibilitatea de a alege tipul de reguli extrase în procesul de învățare.

Algoritmii ierarhici construiesc arbori în care fiecare nod reprezintă un cluster. Mai departe ei se pot împărți în algoritmi cumulativi și algoritmi separatori, urmărind o abordare de jos în sus sau de sus în jos. Cei mai utilizați algoritmi sunt:

- CURE: Un algoritm robust care lucrează cu valori aberante și poate identifica clusteri din forme nesferice (Ghula, 1998);

- ROCK: un algoritm de clusterizare robust folosit pentru date booleene și categoriale. El introduce două noi concepte: vecinii unui punct și legăturile lor (Guha, 2000);

- **Birch** este un algoritm care se utilizează pentru a clasifica seturi mari de date și are la bază două elemente: caracteristica după care se face gruparea și arborele de grupare. Acest algoritm construiește o dendrogramă numită arbore de clasificare în timp ce inspectează setul de date. Există două etape cheie în aplicarea algoritmului: scanarea (parcurgerea) bazei de date și construirea unui arbore și apoi aplicarea algoritmului de grupare pentru a clasifica noduri «frunză».<sup>11</sup>

---

<sup>11</sup> Holsheimer Marcel, Data Mining by Business Users: Integrating Data Mining in Business Processes, San Diego, California, USA, 1999, ISBN 1-58113-171-2

Fiind dat un număr de  $N$  obiecte  $\{Q_i\}$  se calculează următorii indicatori pentru a determina similaritatea:<sup>12</sup>

- centroid-ul (centrul de greutate)  $X_o = \frac{\sum_{i=1}^N X_i}{N}$  ;

- raza (distanța medie de la puncte la centrul de greutate):  $R = \left( \frac{\sum_{i=1}^N (X_i - X_o)^2}{N} \right)^{\frac{1}{2}}$  ;

- diametrul (distanța medie dintre oricare două elemente  $X_i$  și  $X_j$  ale unui cluster):

$$D = \left( \frac{\sum_{i=1}^N \sum_{j=1}^N (X_i - X_j)^2}{N(N-1)} \right)^{\frac{1}{2}} ;$$

- măsurarea gradului de asemănare între două obiecte  $X_i$  și  $X_j$ . În general se utilizează metrica lui Minkowski:  $d_p(X_i - X_j) = \left( \sum_{k=1}^N |X_{ik} - X_{jk}|^p \right)^{\frac{1}{p}}$ , unde  $N$  este dimensiunea obiectului (lungimea vectorului, etc). Distanța euclidiană se obține pentru  $p=2$ , iar distanța Manhattan pentru  $p=1$ ;

- media între două cluster:  $D2 = \left( \frac{\sum_{i=1}^{N_1} \sum_{j=N_1+1}^{N_1+N_2} (X_i - X_j)^2}{N_1 N_2} \right)^{\frac{1}{2}}$ , unde  $N_1$  și  $N_2$  reprezintă

dimensiunea primului cluster, respectiv celui de-al doilea cluster;

- media în interiorul clusterului:  $D3 = \left( \frac{\sum_{i=1}^N \sum_{j=1}^N (X_i - X_j)^2}{N(N-1)} \right)^{\frac{1}{2}}$ , unde  $N$  reprezintă dimensiunea

clusterului;

- creșterea varianței:

$$\left( \sum_{k=1}^{N_1+N_2} \left( X_k - \frac{\sum_{l=1}^{N_1+N_2} X_l}{N_1 + N_2} \right)^2 - \sum_{i=1}^{N_1} \left( X_i - \frac{\sum_{l=1}^{N_1} X_l}{N_1} \right)^2 - \sum_{j=N_1+1}^{N_1+N_2} \left( X_j - \frac{\sum_{l=N_1+1}^{N_1+N_2} X_l}{N_2} \right)^2 \right)^{\frac{1}{2}} .$$

<sup>12</sup> Anghelache Constantin, Tratat de statistică teoretică și economică, Editura Economică, București, 2008, ISBN 9789737093806

Fiecare nod al arborelui reprezintă un grup de obiecte ce se caracterizează prin trei indicatori:  $(N, LS, SS)$ , unde  $N$  este numărul de obiecte din cluster,  $LS$  și  $SS$  au următoarele formule de calcul:

$$\text{- pentru un nod frunză, } N: LS = \sum_{P_i \in N} P_i \text{ și } SS = \sum_{P_i \in N} |P_i|^2 ;$$

$$\text{- pentru un nod părinte (non-frunză), care are nodurile copil } N_1, N_2, \dots, N_k: LS = \sum_{i=1}^k LS \text{ al}$$

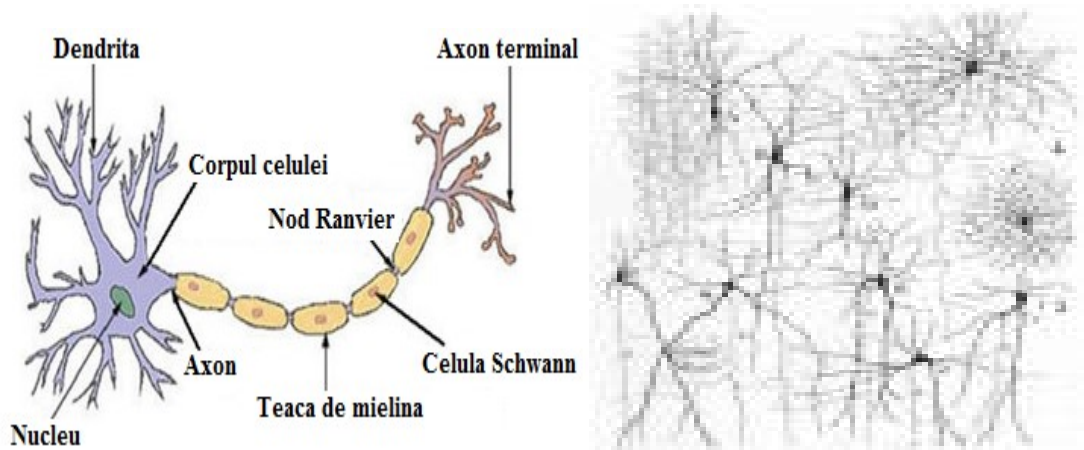
nodului  $N_i$  și  $SS = \sum_{i=1}^k SS$  al nodului  $N_i$ .

Fiecare nod non-frunză conține un anumit număr de copii. Numărul de copii pe care un nod non-frunză poate să îi aibă este limitat de un prag numit factor de adăugare. Diametrul unui subcluster dintr-un nod nu poate depăși un anumit prag.

Inserarea unui obiect în arbore se realizează prin traversarea arborelui de sus în jos începând cu rădăcina în funcție de distanțele euclidiană și Manhattan. Obiectul este introdus în cel mai apropiat cluster sub un nod frunză. În cazul în care introducerea unui obiect într-un subcluster determină depășirea diametrului și implicit a pragului stabilit, atunci se creează un nou subcluster. Împărțirea unui nod frunză este efectuată de prima identificare a perechii de subgrupuri de sub nodul frunză care sunt separate de cea mai mare distanța inter-cluster. Apoi toate celelalte subcluster sunt împărțite între cele două grupuri, având ca suport apropierea față de cele două subgrupuri.

### 1.3 Definiție și forma generală a unei rețele neuronale

La fel ca un creier uman o rețea neuronală este alcătuită din neuroni și din conexiunile dintre ei (figura 3). În rețelele neuronale aceste conexiuni poartă denumirea de ponderi (weights). Informația electrică este simulată prin valori specifice stocate în aceste ponderi. Prin simpla schimbare a acestora se poate simula schimbarea structurii de conectare a rețelei.<sup>13</sup>



**Figura 3.** Reprezentarea neuronilor în creierul uman

<sup>13</sup> Căleanu Catalin- Daniel, Tîponut Virgil, Rețele Neuronale. Arhitecturi și algoritmi, Editura Politehnica Timișoara, 2002, ISBN 973-9389-66-X