

CONTRIBUȚII LA ÎMBUNĂȚIREA TEHNOLOGIEI DE CLASIFICARE A DATELOR

Teză destinată obținerii
titlului științific de doctor inginer
la
Universitatea "Politehnica" din Timișoara
în domeniul CALCULATOARE ȘI TEHNOLOGIA
INFORMAȚIEI
de către

Ing. Claudiu-Raul ROBU

Conducător științific: prof.univ.dr.ing Vasile STOICU-TIVADAR
Referenți științifici: prof.univ.dr.ing.mat. Dumitru Dan BURDESCU
prof.univ.dr.ing. Viorel NEGRU
conf.univ.dr.ing. Ștefan HOLBAN

Ziua susținerii tezei: 27.07.2012

Seriile Teze de doctorat ale UPT sunt:

- | | |
|------------------------|---|
| 1. Automatică | 7. Inginerie Electronică și Telecomunicații |
| 2. Chimie | 8. Inginerie Industrială |
| 3. Energetică | 9. Inginerie Mecanică |
| 4. Ingineria Chimică | 10. Știința Calculatoarelor |
| 5. Inginerie Civilă | 11. Știința și Ingineria Materialelor |
| 6. Inginerie Electrică | |

Universitatea „Politehnica” din Timișoara a inițiat seriile de mai sus în scopul diseminării expertizei, cunoștințelor și rezultatelor cercetărilor întreprinse în cadrul școlii doctorale a universității. Seriile conțin, potrivit H.B.Ex.S Nr. 14 / 14.07.2006, tezele de doctorat susținute în universitate începând cu 1 octombrie 2006.

Copyright © Editura Politehnica – Timișoara, 2012

Această publicație este supusă prevederilor legii dreptului de autor. Multiplicarea acestei publicații, în mod integral sau în parte, traducerea, tipărirea, reutilizarea ilustrațiilor, expunerea, radiodifuzarea, reproducerea pe microfilme sau în orice altă formă este permisă numai cu respectarea prevederilor Legii române a dreptului de autor în vigoare și permisiunea pentru utilizare obținută în scris din partea Universității „Politehnica” din Timișoara. Toate încălcările acestor drepturi vor fi penalizate potrivit Legii române a drepturilor de autor.

România, 300159 Timișoara, Bd. Republicii 9,
tel. 0256 403823, fax. 0256 403221
e-mail: editura@edipol.upt.ro

Cuvânt înainte

Teza de doctorat a fost elaborată pe parcursul activității mele în cadrul Departamentului de Automatică și Informatică Aplicată al Universității „Politehnica” din Timișoara.

Mulțumiri deosebite se cuvin conducătorului de doctorat prof.dr.ing. Vasile STOICU-TIVADAR, care m-a sprijinit, m-a încurajat, mi-a dat forța și inspirația ca să duc lucrurile până la capăt. Îi mulțumesc domnului prof. dr.ing Ștefan HOLBAN, pentru că mi-a fost mereu aproape, cu dânsul m-am consultat începând cu septembrie 2009 și dânsul m-a direcționat către algoritmi bioinspirați. Îi mulțumesc domnului profesor și prieten Ioan FILIP, pentru sprijinul acordat în activitatea publicistică, cu dânsul am scris primele mele lucrări.

Mulțumesc familiei mele care m-a susținut în această perioadă și în mod special soției mele Mădălina care pe lângă faptul că mi-a fost alături, mi-a și citit teza și tatălui meu care mereu m-a încurajat să am încredere în mine și să am forța să finalizez teza.

Le mulțumesc tuturor celor care m-au ajutat în orice fel ca să-mi finalizez teza și consider că acesta este un nou început de drum.

Timișoara, iulie 2012

Claudiu-Raul ROBU

Robu, Claudiu-Raul

Contribuții la îmbunătățirea tehnologiei de clasificare a datelor

Teze de doctorat ale UPT, Seria 14, Nr. 6, Editura Politehnica, 2012, 182 pagini, 51 figuri, 67 tabele.

ISSN: 2069-8216

ISSN-L:2069-8216

ISBN:978-606-554-504-5

Cuvinte cheie:

clasificare, data mining, Ant Colony Optimization, algoritmi genetici, predicții, Weka

Rezumat:

Capitolul 1 prezintă definiții și terminologie, etapele procesului de data mining, domeniile de aplicare, instrumente utilizate, noțiunea de clasificare a datelor, algoritmi de clasificare, seturi de date utilizate la analiza algoritmilor, etc. Capitolul 2 tratează preprocesarea datelor și prezintă o sinteză a tehnicilor de preprocesare, aspecte legate de instrumentul conceput și implementat Arff Converter și realizează o analiză statistică a datelor unor studenți. Capitolul 3 prezintă algoritmul ACO, se analizează comportamentul lui rulând foarte multe configurații de parametrii, iar în urma analizei se emit recomandări privind alegerea parametrilor. Se prezintă diverse variante ale lui ACO și se propun 4 îmbunătățiri care se implementează în Ant-r-Miner. În capitolul 4 se prezintă aspecte legate de algoritmii genetici utilizați la clasificarea datelor și se propune un nou algoritm genetic care se implementează în Weka. Se analizează comportamentul său la diferite valori ale parametrilor de intrare și se emite o listă de recomandări privind alegerea acestora. În capitolul 5 este concepută o nouă metodă de a realiza predicții în Weka, cu ajutorul unei interfețe dinamice, funcție de setul de date analizat. În capitolul 6 sunt analizate date legate de nașterile care au avut loc la Clinica de Obstetrică și Ginecologie Bega, Timișoara în anul 2010. Se construiesc cu algoritmii propuși clasificatori în funcție de indicele apgar.

CUPRINS

1.	Introducere	7
1.1.	Noțiuni generale.....	7
1.1.1.	Definiții și terminologie	7
1.1.2.	Etapele procesului de data mining	10
1.2.	Domenii în care tehnicile de data mining se aplică cu succes.....	12
1.2.1.	Data mining în e-commerce (Web Mining)	12
1.2.2.	Data mining în medicină.....	15
1.2.3.	Data mining în educație	16
1.2.4.	Data mining în combaterea terorismului	17
1.2.5.	Sistemele de Data Mining utilizate de agențiile federale din SUA ...	19
1.3.	Clasificarea datelor	24
1.3.1.	Noțiuni generale	24
1.3.2.	Algoritmul Naive Bayes	25
1.3.3.	Algoritmul k-Nearest Neighbor	26
1.3.4.	Algoritmul C4.5.....	26
1.4.	Instrumente utilizate la realizarea de studii de data mining.....	28
1.4.1.	Weka	29
1.4.2.	R	33
1.4.3.	Comparație între Weka și R	35
1.5.	Seturi de date utilizate pentru analiza performanțelor algoritmilor	36
1.6.	Structura tezei	39
1.7.	Concluzii	41
2.	Preprocesarea datelor	42
2.1.	Noțiuni introductive	42
2.2.	Tehnici de preprocesare a datelor	43
2.2.1.	Agregarea.....	43
2.2.2.	Eșantionarea	44
2.2.3.	Reducerea dimensionalității	45
2.2.4.	Selectarea de submulțimi de caracteristici	56
2.2.5.	Crearea caracteristicilor	59
2.2.6.	Discretizarea și binarizarea	60
2.2.7.	Transformarea variabilelor.....	63
2.2.8.	Considerații privind tehnicile de preprocesare	65
2.3.	Îmbunătățirea preprocesării cu <i>Weka</i>	65
2.3.1.	Conectarea aplicației Weka la baze de date	65
2.3.2.	Preprocesarea datelor cu ajutorul lui Weka	67
2.3.3.	ARFF Converter	67
2.3.4.	Studiu de caz: Analiza datelor civice și școlare ale studenților	75
2.4.	Concluzii	80
3.	Clasificarea datelor cu ajutorul algoritmului Ant Colony Optimization	82
3.1.	Noțiuni introductive	82
3.2.	Prezentarea algoritmului.....	83
3.3.	Analiza parametrilor algoritmului	88
3.3.1.	Noțiuni generale	88
3.3.2.	Analiza empirică a parametrilor	90
3.3.3.	Analiza matematică a coeficienților algoritmului	98

3.3.4.	Considerații și recomandări cu privire la alegerea parametrilor algoritmului Ant Colony Optimization	103
3.4.	Variante ale algoritmului ACO pentru descoperirea regulilor de clasificare 105	
3.5.	Soluții de îmbunătățire a algoritmului ACO	108
3.6.	Concluzii	112
4.	Clasificarea datelor cu ajutorul algoritmului genetic AGR.....	113
4.1.	Noțiuni introductive	113
4.2.	Structura unui algoritm genetic	113
4.3.	Algoritmi genetici utilizați la clasificarea datelor	115
4.4.	AGR - Algoritmul genetic propus.....	117
4.5.	Weka extinsă cu algoritmul genetic AGR	123
4.6.	Analiza coeficienților algoritmului și rezultate experimentale	124
4.6.1.	Utilizarea tuturor regulilor descoperite de fiecare iterație.....	125
4.6.2.	Utilizarea celei mai bune reguli descoperite de fiecare iterație.....	129
4.6.3.	Analiza comportamentului algoritmului la utilizarea a diferite funcții fitness	132
4.6.4.	Recomandări cu privire la alegerea parametrilor	134
4.7.	Comparație cu alți algoritmi	134
4.8.	Concluzii	136
5.	Îmbunătățirea procesului de realizare a predicțiilor	138
5.1.	Noțiuni introductive	138
5.2.	Realizare de predicții cu Weka	138
5.3.	Îmbunătățirea procesului de realizare a predicțiilor.....	141
5.4.	Concluzii	145
6.	Analiza și clasificarea datelor legate de nașteri	146
6.1.	Noțiuni introductive	146
6.2.	Setul de date	147
6.3.	Preprocesarea datelor	148
6.4.	Analiza statistică	151
6.5.	Clasificarea.....	153
6.6.	Predicții	157
6.7.	Concluzii	160
7.	Concluzii	162
	Bibliografie	167
	Lista figurilor.....	177
	Lista tabelor	179

1. Introducere

1.1. Noțiuni generale

1.1.1. Definiții și terminologie

Oamenii au căutat tipare încă de la începutul omenirii: vânătorii au căutat tipare în comportamentul animalelor la migrare, fermierii au căutat tipare în ceea ce privește creșterea recoltelor, medicii caută tipare în simptomele pacienților. Procesul descoperirii de cunoștințe este la fel de vechi ca și omenirea, începând cu descoperirea focului și ajungând la studiile actuale privind universul. Data mining este procesul de descoperire de cunoștințe și extragere de reguli și tipare din date.

Nevoia explorării și extragerii cunoștințelor din baze de date de dimensiuni foarte mari este tot mai accentuată întrucât numărul bazelor de date de dimensiuni impresionante este în continuă creștere. Bazele de date din spitale, bănci, supermarketuri, au căpătat în timp dimensiuni impresionante. Senzorii de pe sateliți, telescoapele care scanează cerul produc seturi de date de ordinul GB / oră [TAN05]. În multe baze de date au fost stocate informații care ascund cunoștințe extrem de valoroase. Prin analizarea datelor despre clienții unui supermarket se pot determina clienții care sunt aproape pierduți și se pot iniția campanii publicitare în vederea păstrării lor. Analizarea datelor despre clienții unei bănci permite construirea profilului clientului de încredere, respectiv a clientului de neîncredere, foarte util în acordarea unui credit. Datele pacienților permit construirea profilului pacientului care suferă de o anumită afecțiune, iar acest profil poate ajuta la diagnostic sau prevenire [ATI09]. Prin intermediul *data mining* o universitate poate prezice cu o acuratețe de 85% care studenți vor absolvi și care nu. Universitatea poate folosi această informație pentru a-și concentra atenția asupra acelor studenți care sunt în pericol de a nu absolvi. Administrația Federală a Aviației Statelor Unite folosește data mining pentru a analiza datele legate de prăbușirea avioanelor, cu scopul de găsi defectele comune și a recomanda măsuri de precauție [SEI08].

Începând cu anii 90 [GOR06] se vorbește de data mining în foarte multe medii, plecând de la cele academice și până la cele de afaceri sau medicale.

Revista tehnologică online ZDNET News afirma în 8 Februarie 2001 că Data mining va deveni "una dintre cele mai revoluționare dezvoltări a deceniului următor". Recenzia MIT Technology a ales data mining ca fiind una dintre primele 10 tehnologii apărute care vor schimba lumea.

De-a lungul timpului au fost date multe definiții procesului de data mining:

Data mining este o tehnologie care permite extragerea informațiilor, utile, relevante din BD de dimensiuni mari [TWO99].

Prin data mining se înțelege procesarea datelor utilizând capabilități de căutare sofisticată și algoritmi statistici pentru a găsi șabloane și corelații în baze de date existente și de dimensiuni mari [www1].

Data mining este procesul de extragere a modelelor din date. Explorarea datelor se realizează cu scopul de a transforma aceste date în informații [PON07].

Data mining este o tehnică de a căuta modele / șabloane în baze de date de dimensiuni mari. Scopul căutării este acela de a găsi corelații între variabile (anterior necunoscute) care pot fi valoroase din punct de vedere comercial [BER04].

Data mining înseamnă extragerea de șabloane și alte informații folositoare dintr-un corp de date [www2].

Data mining este procesul care folosește tehnici statistice pentru a descoperi relațiile subtile între date, precum și construirea de modele predictive bazate pe ele [GOR06].

Data mining - un amestec de statistică, inteligență artificială și cercetare în baze de date (Prigibon – cercetător de la Google Inc.) [GOR06].

Prin data mining – se înțelege practica căutării automate de pattern-uri (modele, șabloane, tipare) în mulțimi mari de date, utilizând tehnici computaționale din statistică, învățarea automată și recunoașterea formelor [GOR06].

O definiție proprie, mai filozofică este următoarea: Data mining înseamnă învățarea din experiența trecutului reprezentată prin seturi de date, cu scopul aplicării cunoașterii dobândite în acțiunile viitoare.

Disciplinele care stau la baza data mining sunt cele din figura 1.1.

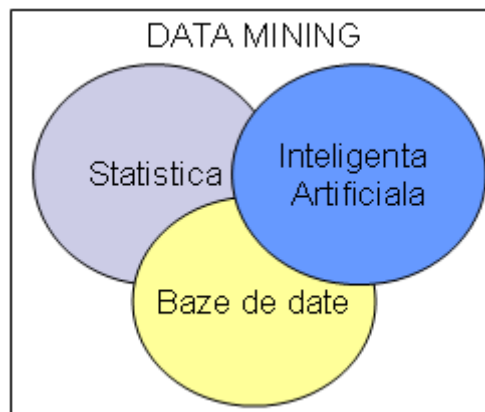


Fig. 1.1 - Disciplinele care stau la baza Data Mining

Statistica aduce în data mining tehnici ce permit identificarea de relații între variabile (regresie liniară, regresie logistică, Naive Bayes, etc).

Inteligența artificială contribuie la Data Mining prin tehnici de procesare a informației bazate pe modelul raționamentului uman.

Machine Learning este un domeniul al inteligenței artificiale care este cel mai folosit în Data Mining.

Machine Learning se ocupă cu crearea de algoritmi și tehnici care permit calculatoarelor să „învețe”.

Finalitatea Machine Learning este achiziția de descrieri structurale din exemple. Aceste descrieri structurale trebuie să poată fi folosite pentru predicții și înțelegerea fenomenului studiat.

Sistemele de baze de date reprezintă materialul care trebuie explorat.

Tehnicile de data mining cele mai utilizate sunt clasificare, regresie, clustering (grupare), reguli de asociere.

Clasificarea este o metodă predictivă care face parte din categoria metodelor de învățare supervizată. Clasele sunt valorile pe care le poate lua un atribut categorial al setului de date considerat, numit atribut clasă. Atributele diferite de atributul clasă se numesc predictorii. Scopul clasificării este de a prezice pentru o instanță dată caracterizată de un set de valori a predictorilor a clasei căreia îi aparține.

Setul de date este format din instanțe (exemple, obiecte).

Atributele sunt caracteristici sau proprietăți ale exemplurilor, exprimate prin valori numerice sau nominale (categorice). Altfel spus, atributele sunt denumirile coloanelor dintr-un tabel. Dacă într-un studiu de prognoză a vremii, dorim să prezicem cum va fi mâine (soare, înnorat, ploaie, lapoviță, grindină, ninsoare), atunci aceste valori le putem considera clase iar atributul clasă îl putem numi vreme.

Regresia este tot o metodă predictivă, similară cu clasificarea, dar valorile care se prezic nu sunt categoriale ci sunt numerice [TAN05]. În cazul regresie exemplul precedent ar presupune prezicerea temperaturii de mâine.

Clusteringul (Gruparea) este procesul de determinare a grupurilor (numite clustere) de instanțe asemănătoare dintr-un set de date. Gruparea este o tehnică de învățare nesupervizată, descriptivă, în urma utilizării ei se obțin informații cu privire la modul în care sunt organizate datele.

Detectarea asocierilor între diferite tipuri de informații care, în aparență nu au nici un fel de dependențe semantice, poate oferi concluzii interesante. Regulile de asociere au ca scop modelarea dependențelor între articole în datele tranzacționale. Un exemplu clasic de studiu utilizând regulile de asociere este analiza produselor care se cumpără împreună (cineva care își cumpără un ciocan este foarte probabil să-și cumpere și cuie).

1.1.2. Etapele procesului de data mining

Procesul de data mining presupune parcurgerea următoarelor etape (vezi figura 1.2):

- Preprocesarea datelor
- Construirea și testarea modelului (modelelor)
- Aplicarea modelului pe date noi

Preprocesarea datelor presupune citirea datelor brute, curățarea lor și prelucrarea lor în vederea optimizării procesului de *data mining*.

Datele brute pot conține valori lipsă, date introduse greșit, date expirate, date duplicat, valori aberante sau pot fi afectate de zgomot [GOR06]. În acest sens se impune o primă preprocesare a datelor brute, prin estimarea valorilor lipsă, prin eliminarea datelor duplicat sau expirate, prin eliminarea sau corectarea datelor introduse greșit, prin utilizarea unor parametrii statistici toleranți la valori extreme prin utilizarea de filtre hard și soft pentru reducerea zgomotului, etc. În urma prelucrării datelor brute se obțin date curățate.

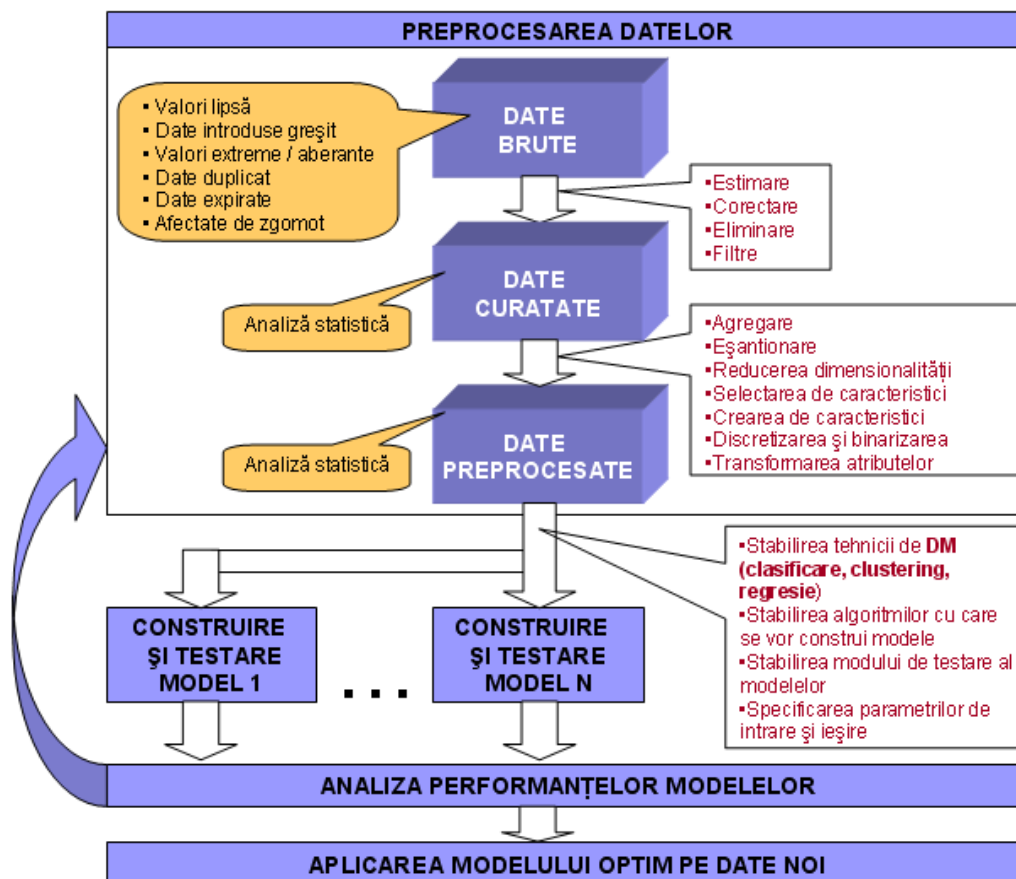


Fig. 1.2 - Etapele procesului de Data Mining

Datele curățate trebuie preprocesate în continuare, în vederea optimizării procesului de *data mining*. Câteva tehnici care pot fi utilizate în acest scop sunt prezentate în continuare:

- **Agregarea** – combinarea mai multor atribute existente într-un singur atribut, cu scopul reducerii volumului de date și al obținerii de date mai stabile (exemplu: agregarea satelor în comuna de care aparțin, a orașelor în județul de care aparțin, a vânzărilor săptămânale în vânzările lunare, etc.)
- **Eșantionarea** – se alege un eșantion reprezentativ pentru întreaga mulțime de date. O problemă importantă este alegerea optimă a volumului eșantionului astfel încât să nu se piardă date dar să nici nu fie mai multe decât este necesar.
- **Reducerea dimensionalității** – se știe că atunci când dimensiunea (numărul atributelor) crește va crește în mod implicit și împrăștierea datelor, ceea ce duce la scăderea vitezei de lucru. Acest fenomen este denumit în *data mining* ca și blestemul dimensionalității (*curse of dimensionality*). Antidotul este reducerea dimensionalității prin tehnici ca: analiza factorială și analiza componentelor principale.
- **Selectarea de submulțimi de caracteristici** – pentru eliminarea caracteristicilor redundante sau nerelevante
- **Crearea de caracteristici** - crearea de noi atribute care pot capta mai bine informațiile din date decât cele originale
- **Discretizarea și binarizarea** - trecerea de la date continue la date discrete (de exemplu trecerea de la numere reale la numere întregi)
- **Transformarea atributelor** - transformarea unor atribute vechi în atribute noi

În mod opțional, datele curățate sau preprocesate pot fi supuse unei etape de analiza statistică a valorilor fiecărui atribut. Această analiză statistică poate extrage din date informații utile, privind distribuția valorilor fiecărui atribut și ajută la înțelegerea datelor.

Construirea modelului și testarea lui. De regulă, în studiile de data mining se construiesc mai multe modele și se alege cel mai performant dintre ele. O situație rară, în care nu se justifică construirea mai multor modele, este aceea în care primul model construit este cel mai bun cu putință (de exemplu am construit un model clasificator care are o acuratețe a predicției de 100 %). De obicei însă, este necesară construirea mai multor modele cu diferiți algoritmi, compararea performanțelor obținute de toate modelele și alegerea celui mai bun model. Dacă toate modelele pe care încercăm să le construim au performanțe nesatisfăcătoare, este necesară revenirea în etapa de preprocesare, prelucrarea datelor fie prin discretizare, fie prin alte tehnici precum selectarea și crearea de caracteristici și reconstruirea modelelor, poate că în urma unor prelucrări a datelor de intrare vor obține performanțe mai bune. Construirea oricărui model se face învățând din datele disponibile. Pentru a învăța din date se pot utiliza diverse tehnici (clasificare, regresie, grupare sau reguli de asociere) și diverși algoritmi în funcție de tehnica aleasă. De exemplu pentru clasificare se pot folosi algoritmi Naive Bayes, C4.5, Suport Vector Machines, pentru regresie se pot folosi algoritmi de regresie liniară, regresie liniară multiplă, regresie logistică, pentru grupare se pot folosi algoritmi K-Means, EM, FarthestFirst, pentru a descoperi reguli de asociere algoritmul Apriori,

FPGrowth etc. Un model poate fi testat pe datele de antrenare, sau întregul set de date se împarte în date de antrenare și date de test (în acest caz va fi testat pe date noi și trebuie ales procentul de instanțe folosit la antrenare), etc. În funcție de algoritmul folosit pentru construirea modelului trebuie configurați parametrii specifici de intrare și trebuie aleși parametrii de ieșire.

Dacă cel puțin unul din modele construite obține rezultate satisfăcătoare, atunci modelul poate fi aplicat pe date noi.

1.2. Domenii în care tehnicile de data mining se aplică cu succes

Data mining este o tehnologie care are o arie foarte largă de acoperire. Data mining se poate aplica în orice domeniu în care există date salvate, ca urmare a unei experiențe și se dorește să se învețe din datele respective sau altfel spus, să se extragă cunoștințe din acele date.

Domenii în care Data Mining s-a impus mai mult decât în altele sunt e-commerce, medicină, educație, detectarea fraudelor (în special în domeniul bancar), combaterea terorismului.

1.2.1. Data mining în e-commerce (Web Mining)

Clienții sunt componenta țintă a oricărei organizații. Achiziționarea de noi clienți, satisfacerea și păstrarea celor existenți, precum și estimarea comportamentului cumpărătorilor poate îmbunătăți disponibilitatea de produse și servicii și, prin urmare, profiturile. Astfel, obiectivul oricărui exercițiu de *data mining* în e-commerce este de a îmbunătăți procesele care contribuie la furnizarea de valoare către clientul final [RAG05].

Sistemele software destinate planificării resurselor companiilor (*ERP-enterprise resource planning*) oferă deseori facilități de data mining. Unul din cele mai puternice ERP-uri este sistemul SAP [PRO09]. Modelele de data mining din SAP permit determinarea unor șabloane în comportamentul clienților și pot fi utilizate pentru a furniza răspunsuri la întrebări precum: care ofertă este cea mai potrivită, pentru care client și când ar trebui făcută acea ofertă.

Fișierele log de pe *Web* ascund informații cu privire la ce caută potențialii clienți pe un site. Își propun să cumpere sau doar să navigheze? Încearcă să-și cumpere ceva cu care sunt familiari sau ceva despre care știu puține lucruri? Încearcă să cumpere de acasă, de la servicii sau de la un hotel?

Informațiile din fișierele log sunt adesea utilizate pentru a determina ce performanță poate fi așteptată de la servere și rețea, pentru a sprijini service-ul și pentru a face interacțiunea e-business productivă.

Instrumentele de *data mining* din spatele aplicațiilor *Web*, au grijă de reprezentare cunoștințelor, de construirea profilului clientului și de modelarea predictivă a scenariilor interacțiunilor clientului. De exemplu, odată ce un client și-a cumpărat un anumit număr de servere, este foarte probabil să aibă nevoie și de anumite echipamente adiționale cum ar fi routere, switch-uri, echipamente de backup, etc. Sistemele bazate pe reguli ar putea fi folosite pentru a sugera asemenea alternative clienților [RAG05].

Data mining a fost de asemenea aplicată în detectarea modului în care clienții pot răspunde la oferte promoționale făcute de o companie e-commerce care se ocupă cu carduri de credit [ZHA03]. Tehnici precum fuzzy au fost utilizate pentru a genera reguli de tipul if-then. [RAG05].

E-commerce-ul este un domeniu foarte potrivit pentru *data mining* pentru că oferă toate ingredientele necesare pentru un *data mining* de succes: există o mulțime de înregistrări, colecțiile electronice oferă date fiabile, reflecția poate fi ușor transformată în acțiune și returnată în investiții care se pot măsura. [ANS01]

Practic fiecare companie mare care vinde cu amănuntul are un site *Web* și facilități de vânzare online poate utiliza Data Mining. Analizând fișierele log ale mai multor companii se pot identifica șabloanele de navigare ale clienților. Aceste șabloane de navigare se pot folosi pentru a sprijini un sistem de produse personalizate pentru vânzări online.

Exploatarea secvențială este procesul de aplicare a tehnicilor de *data mining* unei baze de date, în scopul de a descoperi relațiile de corespondență care există între o listă ordonată de evenimente. O aplicație importantă a tehnicilor de *data mining* secvențial este explorarea fișierelor log, pentru ca în aceste fișiere log se află înregistrată ordinea în care au fost accesate paginile unui site de către diferiți utilizatori, într-o anumită perioadă de timp [EZE05].

Autorii articolului [NIU02] prezintă o metoda de realizare a profilelor clienților în e-commerce bazată pe ierarhizarea produselor pentru o personalizare efektivă. Aceasta împarte profilul clientului în trei părți: profilul de bază aflat pe baza datelor demografice ale clientului; profilul preferențial bazat pe datele comportamentale și profilul regulă care se referă în principal la regulile de asociere. Pe baza profilelor clienților, autorii generează două tipuri de recomandări, care sunt recomandări de interes și recomandări de asociere. Aceștia propun o structură specială de date denumită arbore de profil pentru căutare și potrivire efektivă.

Web mining - este punerea în aplicare a tehnicilor de data mining pentru a descoperi modele de pe *Web*. În conformitate cu obiectivele de analiză, *Web mining* poate fi împărțit în trei tipuri diferite, care sunt explorarea utilizării *Web*-ului, explorarea conținutului *Web* și a structurii *Web*.

Explorarea structurii Web

Explorarea structurii *Web* este procesul de utilizare a teoriei grafurilor pentru a analiza nodurile și structura conexiunii într-un site *Web*. În funcție de tipul datelor *Web* structurale, explorarea structurii *Web* poate presupune extragerea de șabloane din hyperlink-uri sau explorarea tag-urilor HTML sau XML din cadrul paginilor. Se construiește o structură de arbore pentru a analiza și descrie tag-urile HTML sau XML. Un hyperlink este o componentă structurală care conectează pagina *Web* la o pagina din cadrul site-ului sau din alt site [BIN11].

Datele de structură reprezintă viziunea designerului asupra organizării conținutului în cadrul unui site. Aceasta organizare este captată prin intermediul structurii de legături între pagini, reflectată prin hyperlink-uri. Datele de structură includ de asemenea și structura intra-pagină a conținutului reprezentată de aranjarea tag-urilor HTML sau XML în cadrul paginii. Datele de structură pentru un site sunt captate în mod normal de o hartă a site-ului generată automat, hartă care reprezintă structura hyperlink a site-ului [RAG05].

Explorarea conținutului Web

Explorarea conținutului *Web* este procesul de descoperire a informațiilor utile din texte, imagini, date audio sau video de pe *Web*. Explorarea conținutului *Web* este uneori numită *Web text mining* (explorarea textului de pe *Web*) pentru că aceasta este aria de cercetare cea mai răspândită. Tehnologiile care sunt de regulă folosite în explorarea textului de pe *Web* sunt NLP (Natural Language Processing) și IR (Information Retrieval). Se analizează informațiile citite sau cu care utilizatorii unui site interacționează [www3].

Un exemplu de *Web text mining* este următorul: autorii articolului "Exploiting Wikipedia as External Knowledge for Document Clustering" [HU09] au observat că în metodele tradiționale de grupare (clustering) a textului, documentele sunt reprezentate ca fiind "bagaje de cuvinte" fără să se țină seama de informația semantică din fiecare document. De exemplu dacă două documente utilizează colecții diferite de cuvinte de bază pentru a reprezenta aceeași temă, ele ar putea fi împărțite în mod fals în două grupuri distincte. În această lucrare autorii dezvoltă două abordări "potrivirea exactă" și "potrivirea relațională" pentru a pune în corespondență documentele text cu conceptele wikipedia și ulterior cu categoriile Wikipedia. Apoi documentele text sunt grupate pe baza unei metrici de similaritate care combină informații legate de conținutul documentului, informații legate de concept și informații legate de categorii. Rezultatele experimentale arată că performanța clusteringului s-a îmbunătățit în mod semnificativ prin ralierea reprezentării documentului la conceptele și categoriile Wikipedia.

Explorarea modului de utilizare a Web-ului

Web Usage Mining descrie cum este utilizată o pagină -data și ora când a fost accesată, adresa IP a stației de pe care a fost accesată și de către cine [www3].

Mobasher [MOB04] prezintă o analiză cuprinzătoare a procesului de personalizare bazat pe explorarea modului de utilizare a *Web*-ului. Scopul acestei lucrări este de a arăta cum pot fi cântărite și utilizate efectiv, ca parte integrantă a unui sistem *Web* de personalizare, tehnicile de descoperire a șabloanelor cum sunt clustering-ul, explorarea utilizând reguli de asociere și descoperirea șabloanelor secvențiale efectuate asupra datelor de utilizare *Web*. Autorul observă că datele log colectate automat de către *Web* și de serverele de aplicații reprezintă bine comportamentul de navigare al vizitatorilor.

Mobasher [MOB04] menționează că odată ce s-a stabilit clar care sunt tipurile de date cu care urmează să se opereze (date structurale, date de conținut, date de utilizare), urmează o etapă de pregătire a datelor care se realizează prin procese cum ar fi curățirea datelor, identificarea paginilor vizualizate, identificarea utilizatorilor și identificarea tranzacțiilor. Datele curățate pot fi apoi analizate prin reguli de asociere, clustering, analiză secvențială, etc. Se poate construi o matrice de tranzacții – vizualizări care reține care pagini au fost vizitate de către utilizatori și câte secunde au stat aceștia pe fiecare pagină.

Lanțurile Markov [VAS09] pot fi utilizate pentru modelarea probabilităților de tranziție între vizualizările paginilor. Sarukkai a propus utilizarea lanțurilor Markov drept un mecanism de modelare pentru minimizarea întârzierilor sistemului, în cadrul unei analize privind utilizarea *Web*-ului [SAR00].

1.2.2. Data mining în medicină

Tehnicile de data mining se aplică cu succes în domeniul medical. La ora actuală termenul Medical Data Mining este consacrat, existând multe cărți și conferințe de Medical Data Mining.

Medicina modernă generează, aproape zilnic, cantități enorme de date eterogene. Astăzi provocarea cea mare este de a transforma aceste cantități uriașe de date în informații folositoare și cunoștințe. Datele medicale pot conține imagini (RMN), semnale (ECG), informații clinice cum ar fi temperatura, nivelul de colesterol, etc., precum și interpretarea medicului. Tot mai multe și mai multe proceduri medicale folosesc imagistica ca un instrument preferat de diagnostic, așadar există o necesitate de a dezvolta metode de exploatare eficientă în bazele de date de imagini.

Învățarea din datele unor bolnavi permite construirea profilului pacientului bolnav de o anumită boală, îmbunătățirea diagnosticării [ATI09], se pot analiza relațiile ce există între anumiți parametri medicali în vederea realizării de prognoze medicale. Tehnicile de data mining se aplică cu succes în imagistica creierului uman, genetică (pentru prezicerea structurilor proteice, pentru determinarea structurii 3D a proteinelor dată fiind secvența lor amino-acidă). Tot în domeniul geneticii se aplică tehnici de data mining pentru descoperirea unor corelații între modificările secvențelor de ADN ale diverșilor indivizi și susceptibilitatea apariției unor boli. Scopul cercetărilor din acest domeniu este îmbunătățirea diagnosticării bolilor, prevenirea mai eficientă și tratarea mai ușoară a acestora.

Tehnicile de data mining permit descoperirea de medicamente și susțin predicția pe individ (personalizarea medicamentelor).

Datorită faptului ca fișierele ce conțin date medicale pe care se aplică tehnici de data mining sunt legate de subiecți umani, referindu-se la viața lor privată, problema asigurării confidențialității datelor este foarte importantă [GAN07]. În cadrul lucrării "Anonymizing Healthcare Data: A Case Study on the Blood Transfusion Service" [MOH09] autorii expun îngrijorările legate de confidențialitatea informațiilor privind transfuziile sanguine din cadrul sistemului de sharing existent între serviciul de transfuzii de sânge al Crucii Roșii din Hong Kong și spitalele publice și identifică provocările majore care fac metodele tradiționale de anonimizare a datelor inaplicabile în acest caz. Autorii au propus un nou model de asigurare a intimității numit LKC-privacy, împreună cu un algoritm de anonimizare pentru a corespunde cerințelor de intimitate a informațiilor în cazul Serviciului de Transfuzii de Sânge. Experimentele făcute pe date reale au demonstrat că algoritmul de anonimizare propus poate reține în mod eficient informațiile esențiale în date anonime, pentru a fi analizate și este scalabil pentru anonimizarea unor seturi de date mari.

Tehnologiile transcriptomice reprezintă instrumente promițătoare pentru identificarea noilor gene implicate în îmbătrânirea cerebrală sau în bolile neurodegenerative cum ar fi boala Alzheimer. Aceste tehnologii produc date biologice masive care până acum au fost dificil de exploatat. În acest context autorii articolului "GeneMining: Identification, Visualization, and Interpretation of Brain Ageing Signatures" [SAL09] propun explorarea genelor (GeneMining), o metodologie multidisciplinară care are drept scop dezvoltarea unor strategii noi pentru analizarea datelor și proiectarea unor instrumente interactive pentru a ajuta biologia să identifice, vizualizeze și să interpreteze semnalmintele îmbătrânirii creierului. Pentru a adresa problema specifică a descoperii semnalmintelor îmbătrânirii

creierului autorii combină și aplică instrumentele existente punând accentul pe noile metode eficiente de Data Mining bazate pe șabloane secvențiale.

O aplicație de Data Mining dezvoltată pe baza contractului cu NR. 99-II CEEEX 03 – INFOSOC 4091 / 31.07.2006 de către Univ. Medicina și Farmacie "Victor Babes" Timișoara – coordonator, Universitatea Politehnica din Timișoara, Universitatea de Vest din Timișoara, Spitalul Universitar de Obstetrică-Ginecologie "Dumitru Popescu" Timișoara și Bridgeman SRL este aplicația MaternQual.

MaternQual este un sistem informatic integrat de evaluare complexă a factorilor asociați predicției riscului și calității în obstetrică. MaternQual permite culegerea și prelucrarea datelor medicale în obstetrică precum și utilizarea tehnicilor de data mining pentru procesarea acestora. Aplicația permite identificarea factorilor de risc primari și secundari în nașterea prematură, precum și predicția riscului de naștere prematură.

Un alt instrument de Data Mining în domeniul medical este aplicația Neonat, care permite înregistrarea, procesarea digitală și analiza plânsetelor nou-născuților utilizând tehnici de data mining [ROB11a]. Aplicația a fost dezvoltată în cadrul grantului CNCISIS, tip A, nr 755, cu tema "Noi direcții în diagnosticul și managementul suferinței cerebrale perinatale și consecințelor postnatale ale acesteia".

A fost propus un sistem de data mining în domeniul medical la nivel de regiune (cu adresare specială către euro-regiunea DKMT (Danube-Kris-Mures-Tisa)) [ROB11c]. Sistemul urmărește centralizarea datelor de la spitalele localizate pe teritoriul unei regiuni, într-o bază de date și apoi analiza acestor date cu ajutorul instrumentului de data mining și machine learning din *Weka*. Etapele propuse sunt următoarele: datele medicale din bazele de date ale spitalelor sunt exportate în format *XML*, potrivit standardului *HL7 CDA* [SCH06]. Acestea sunt apoi centralizate în mod automat pe un server, prin intermediul unor apeluri de servicii *web*, într-o bază de date. Datele de interes pentru studiul realizat sunt extrase din baza de date și convertite în format *ARFF* [ROB10b]. Ultima etapă constă în preprocesarea și analizarea acestora cu algoritmi din *Weka*, cu scopul de a obține concluzii medicale relevante [ROB10c].

1.2.3. Data mining în educație

În zilele noastre, la fel ca și în toate domeniile sociale, tehnologiile de informare și comunicare, au explodat în cadrul sistemului educațional. Astfel este posibilă capturarea și compilarea mai multor tipuri de informații într-o modalitate foarte simplă și necostisitoare cum ar fi datele administrative ale școlilor, registrele liceelor și universităților, înregistrările computerizate ale studenților, platformele e-learning, înregistrarea activităților, sistemele de învățare cu ajutorul calculatorului, etc.

Educational Data Mining este procesul de transformare a datelor brute, compilate de sistemele educaționale în informații folositoare care pot fi utilizate pentru luarea de decizii informate și formularea răspunsurilor la întrebările legate de cercetare.

Data mining se aplică pe date care vin din mediul educațional cu scopul de a obține o mai bună înțelegere a procesului de învățare a studenților și de a îmbunătăți calitatea și costurile din sistemul educațional.

EDM (Educational Data Mining) are 3 obiective mari:

- *Obiective pedagogice* – îmbunătățirea conceperii materialelor didactice, cu scopul de a îmbunătăți performanța academică a studenților
- *Obiective de management* – pentru a optimiza organizarea și întreținerea infrastructurii educaționale, zonele de interes, cele mai solicitate cursuri
- *Obiective comerciale* – realizarea de segmentări ale pieței și facilitarea recrutării studenților

Educational data mining este o disciplină aflată la intersecția dintre *data mining* și Pedagogie. Pedagogia contribuie cu cunoștințele intrinseci procesului de învățare. *Educational data mining* ca și *data mining* este o tehnologie care integrează tehnici multidisciplinare, de la extragerea informațiilor din baze de date prin simple comenzi SQL până la utilizarea unor algoritmi de învățare automată extrași din câmpul Inteligenței Artificiale. Printre cele mai folosite tehnici în domeniul educației se enumeră clusteringul, regulile de asociere și analiza șabloanelor.

Prin intermediul *data mining* o universitate poate prezice cu o precizie între 75 și 80 % care studenți vor absolvi și care nu. Universitatea poate folosi această informație pentru a-și concentra atenția asupra acelor studenți care sunt în pericol de a nu absolvi [DEK09].

În universitățile care oferă învățământ la distanță numărul de studenți exmatriculați este mult mai mare decât în universitățile care nu oferă această formă de învățământ. Limitarea exmatriculărilor este esențială în universități și de aceea abilitatea de a prezice care studenți vor fi exmatriculați este folositoare din mai multe motive [KOT09].

Prin intermediul regulilor de asociere se pot analiza erorile care apar împreună în timpul rezolvării exercițiilor. Scopul găsirii acestor asocieri este ca profesorul să reflecteze și să revizuiască materialele de curs sau să sublinieze aceste subtilități studenților [MAR05]. Reguli de asociere se pot aplica pentru a recomanda studenților cursurile pe care să le urmeze [BEN06].

Clustering-ul poate fi aplicat pentru a caracteriza studenții cu dificultăți. Clasificarea – arborii de decizie - se poate utiliza pentru a prezice rezultatele studenților la anumite examene [MAR05].

Un studiu de *Educational Data Mining* a fost realizat, pe datele studenților din anul universitar 2008-2009 de la Facultatea de Automatică și Calculatoare a Universității "Politehnica din Timișoara" (2100 de studenți). Studiul s-a realizat pe baza unor date civice și școlare obținute din Sistemul Informatic de Gestiune a Școlărității și a constatat în realizarea unei analize statistice la nivel de atribut, construirea unui model de regresie pentru media multianuală și gruparea datelor în două grupuri pe baza caracteristicilor comune [ROB10b].

1.2.4. Data mining în combaterea terorismului

Tragedia de la 11 septembrie a avut efecte incomensurabile și permanente atât asupra Statelor Unite cât și asupra restului lumii și a adus subiecte precum securitatea și apărarea pe prim plan. Pentru a ajuta la prevenirea unor astfel de catastrofe, un program de succes de securitate ar include dispoziții care să securizeze frontierele, sectorul de transport și infrastructura critică. Un factor critic al unui astfel de program ar fi capacitatea de sintetiză și analiză a datelor din surse multiple. În Mai 2003, a avut loc la San Francisco, un workshop la care s-a discutat

modul în care Data Mining și Machine Learning pot ajuta la analizarea unor date complexe, din surse diferite, cu scopul de a preveni activitatea teroristă viitoare. Denumirea workshopului este: "Workshop on Data Mining for Counter Terrorism and Security". Subiectele de interes de la acest workshop au inclus:

- Metode de a integra surse de date eterogene, cum ar fi text, internet, video, audio, date biometrice, precum și de vorbire
- Metode scalabile pentru a depozita surse de date răsfirate
- Identificarea tendințelor în activitățile de grup sau singulare
- Identificarea șabloanelor cu scopul de a identifica persoana
- Data Mining în domeniul securității aviației, securitatea portuară
- Web Data Mining pentru a detecta tendințele teroriste

Detectarea teroriștilor presupune împărțirea indivizilor în două clase: indivizi periculoși și indivizi nepericuloși. Acest proces este oarecum diferit de alte probleme clasice de Data Mining pentru că indivizii din prima clasă (cei periculoși) fac tot posibilul să arate ca și cei din a doua clasă [SKI03].

Autorul D.B. Skilikorn a demonstrat în articolul [SKI03] că un aspect important al detectării teroriștilor este abilitatea de a detecta corelații locale la o scară mică, în contradictoriu cu corelațiile difuze la scară mare.

În contextul securității naționale, Data Mining poate fi un potențial mijloc de a identifica activități teroriste cum ar fi transferuri de bani și comunicații, precum și de a identifica și urmări teroriștii însăși, folosind date de călătorie și emigrare [SEI08].

Deși Data Mining permite extragerea de șabloane și relații, tehnologia are limitarea de a nu informa utilizatorul cu privire la valoarea sau semnificația acestor relații sau șabloane. Aceste lucruri trebuie să fie cunoscute de către utilizator [SEI08].

Câteva exemple notabile de utilizare a tehnicilor de data mining pe teritoriul Statelor Unite sunt următoarele: Departamentul de Justiție al Statelor Unite utilizează Data Mining pentru a determina șabloane de crime și a regla alocările de resurse în consecință; Administrația Federală a Aviației, folosește data mining pentru a analiza datele de prăbușire a avioanelor, cu scopul de găsi defectele comune și a recomanda măsuri de precauție [SEI08].

Data mining contribuie la combaterea terorismului prin tehnici precum reguli de asociere, analiza legăturilor, clustering, clasificare și detectarea anomaliilor. Aceste tehnici se aplică prin intermediul rețelelor neuronale, a arborilor de decizie, prin programare logică inductivă, mulțimi Rough, analiza legăturilor, care este bazată pe teoria grafurilor și k-nearest-neighbor [BHA08].

Un exemplu de utilizare a regulilor de asociere pentru combaterea criminalității este următorul: dacă John vine dintr-o țară X și este asociat cu James care are un cazier substanțial și se constată că țara X are un procent mare de atacuri teroriste atunci ar trebui ca John să fie pus sub urmărire [BHA08].

O altă tehnică de data mining folosită în acest domeniu este detectarea anomaliilor. Un exemplu în acest sens ar fi cursanții unui curs de pilotaj care se arată mult mai interesați de felul în care zboară un avion decât de felul în care decolează și aterizează [BHA08].

Analiza legăturilor este o tehnică de Data Mining folosită la determinarea șabloanelor anormale. Așa cum a fost menționat anterior, analiza legăturilor folosește multe tehnici de teorie a grafurilor. Analiza legăturilor este de asemenea folosită în *Web mining*, în special în explorarea structurii *Web*. Unele motoare de căutare cum ar fi Google folosesc o formă de analiză a legăturilor pentru a afișa rezultatele unei căutări [BHA08].

Analiza legăturilor se folosește pentru a găsi asociații interesante și pentru a determina cum se poate reduce un graf. Mulți teoreticieni ai grafurilor lucrează la probleme de reducere a acestora. O altă provocare a utilizării analizei legăturilor în combaterea terorismului este aceea de a rezolva situațiile de informații parțiale. De exemplu agenția A dispune de un graf parțial, agenția B de un alt graf parțial, și agenția C de cea de-a treia bucată. Întrebarea este cum se pot găsi asociațiile dintre grafuri când niciuna din agenții nu are o imagine completă [BHA08].

În timp ce de toate aceste aplicații de data mining pot beneficia oameni și pot fi salvate vieți, există, de asemenea, și aspecte negative ale acestei tehnologii deoarece ar putea fi o amenințare la adresa vieții private a persoanelor fizice[BHA08].

1.2.5. Sistemele de Data Mining utilizate de agențiile federale din SUA

În anul 2004 Oficiul de Contabilitate Generală al Statelor Unite (GAO - The General Accounting Office) a înaintat un raport către mai multe instituții printre care și Senatul Statelor Unite cu scopul de a prezenta și a descrie caracteristicile sistemelor de Data Mining utilizate de agențiile federale.

Misiunea GAO brațul de audit, evaluare și investigare al Congresului, este de a asista Congresul în vederea îndeplinirii responsabilităților constituționale și de a ajuta îmbunătățirea performanței și responsabilității guvernului federal pentru poporul american. GAO examinează utilizarea fondurilor publice; evaluează programele și politicile federale; furnizează analize, recomandări și alte rapoarte similare pentru a ajuta Congresul în luarea de decizii informate privind supravegherea, politica și finanțarea. Angajamentul luat de GAO pentru o guvernare bună este reflectat de valorile sale intrinseci de responsabilitate, integritate și încredere [GAO04].

Agențiile federale utilizează Data Mining cu o varietate foarte mare de scopuri, de la îmbunătățirea serviciilor și performanței până la detectarea șabloanelor și activităților teroriste. Raportul realizat de GAO analizează dacă agențiile și departamentele federale utilizează sau intenționează să utilizeze *data mining* în activitatea lor. Au fost analizate un număr de 128 de departamente federale și agenții. Dintre acestea un număr de 52 de instituții utilizează sau intenționează să utilizeze *data mining*. Aceste departamente și agenții au raportat un număr de 199 eforturi de *data mining* din care 68 erau în fază de plan iar 131 erau operaționale. Raportul este foarte amplu se întinde pe 71 de pagini, în acest subcapitol este prezentat un scurt rezumat al acestui raport, cu scopul de a ilustra puterea și importanța tehnologiei *data mining*, precum și nivelul de top la care s-a ajuns cu utilizarea ei.

O parte din cele 128 de agenții și departamente analizate pot fi vizualizate în tabelul următor:

Departamentul	Agenția
Departamentul de Agricultură	Serviciul de Inspecție al Sănătății Faunei și Florei
	Serviciile Forestiere
Departamentul de Comerț	Administrația Națională a Océanelor și Atmosferei
	Administrația Dezvoltării Economice
Departamentul de Apărare	Serviciul de Securitate al Apărării
	Agenția de Reducere a Amenințării
Departamentul Serviciilor de Sănătate și Umane	Institutes Naționale de Sănătate
	Centrele pentru Îngrijire și Servicii Medicale
Departamentul Securității Naționale	Directoratul de Știință și Tehnologie
	Biroul de Servicii pentru Cetățenie și Emigrări
Departamentul de Interne	Administrația Funciară
	Biroul de Reclamații
Departamentul de Justiție	Biroul de Alcool, Tutun, Arme de Foc și Explozibili
Departamentul Transporturilor	Administrația Federală a Aviației
	Administrația Federală a Tranzitului
Departamentul Trezoreriei	Controlul Monetar
	Biroul de Gravări și Imprimare

Tabel 1.1 - O parte din departamentele și agențiile analizate

În tabelul 1.2 se poate vizualiza care sunt scopurile eforturilor de data mining demarate și câte acțiuni au fost demarate în sprijinul acestor eforturi de către agențiile și departamentele analizate.

Scopurile eforturilor de data mining	Număr eforturi
Analizarea inteligenței și detectarea activităților teroriste	14
Detectarea activităților sau șabloanelor criminale	15
Administrarea resurselor umane	17
Analizarea informațiilor științifice și de cercetare	23
Detectarea fraudelor și abuzurilor	24
Îmbunătățirea serviciilor sau performanței	65

Tabel 1.2 - Scopul și numărul eforturilor de data mining

Departamentul Apărării a raportat cele mai multe eforturi utilizând data mining în vederea îmbunătățirii performanței, serviciilor și gestiunea resurselor umane. Apărarea și-a concentrat de asemenea cel mai frecvent eforturile în scopul analizării inteligenței și detectării activităților teroriste, pe locul doi în această ierarhie situându-se Departamentul Național de Securitate, Justiție și Educație.

Departamentul de Educație a raportat concentrarea celor mai multe eforturi în vederea detectării fraudei, deșeurilor și abuzurilor în timp ce Administrația Națională a Aeronauticii și Spațiului își canalizează majoritatea eforturilor de data mining (21 din 23) către analiza informațiilor științifice și de cercetare. Eforturile de data mining pentru detectarea activităților sau șabloanelor criminale, au fost însă împărțite în mod aproximativ egal între agențiile raportatoare.

În plus, din cele 199 de eforturi de data mining identificate, 122 utilizează informații personale. Pentru aceste eforturi, scopurile principale erau detectarea fraudei, deșeurilor și abuzurilor, detectarea activităților și șabloanelor criminale; analizarea inteligenței și detectarea activităților teroriste; și creșterea conformității fiscale.

Un exemplu al dezvoltării eforturilor la scara largă lansat în urma atacurilor de la 11 septembrie este Sistemul Antitero de schimb de informații între state, numit și MATRIX. MATRIX este utilizat la acest moment în cinci state, furnizează capacitatea de a stoca, analiza și schimba date legate de terorism precum și alte date legate de serviciile criminale între agențiile din diferite state, între state și între state și agențiile federale. Informațiile din bazele de date MATRIX includ date legate de istoricul criminalilor, înregistrări privind încarcerarea, înregistrări privind vehiculele, privind permisele de conducere, autovehicule și fotografiile digitale. Sistemul MATRIX a născut în rândul opiniei publice nemulțumiri, datorită încălcării intimității. Sistemul analizează datele personale ale foarte multor indivizi pentru a vedea care are profilul unui terorist.

Unele obiective ale data mining se concentrează pe activități umane, și prin urmare există o posibilitate destul de mare să implice informații personale. Câteva exemple de astfel de obiective sunt detectarea fraudei, a abuzurilor, detectarea activităților și șabloanelor criminale, gestiunea resurselor umane și analiza serviciilor de inteligență. În continuare sunt prezentate exemple de eforturi de data mining concentrate în acest scop:

- *Detectarea fraudei, deșeurilor și abuzurilor*. Eforturile Administrației Beneficiilor Veteranilor C & P, prin Sistemul de Analiză a Datelor de Plată, explorează datele privind pensia și compensațiile veteranilor pentru a verifica și proba eventualele fraude.
- *Detectarea activităților și șabloanelor criminale*. Departamentului de Educație prin Sistemul Inițiativa Furtului de Identitate, se concentrează către identificarea cazurilor de furt de identitate care implică împrumuturile pentru educație.
- *Gestiunea resurselor umane*. Departamentul de Resurse Umane Oracle al Forței Aeriene Americane utilizează data mining pentru a furniza informații legate de promovări, gradații, permisiile precum și alte informații relevante pentru planificarea resurselor umane.
- *Analiza serviciilor de inteligență și detectare a activităților teroriste*. Agenția de servicii inteligente în domeniul apărării Verity K2 Enterprise explorează date din comunitatea agențiilor de servicii și din surse de pe Internet pentru a identifica teroriști străini sau cetățeni americani conexați unor activități teroriste străine.

Pe de alta parte, o altă categorie de eforturi nu se concentrează în mod neapărat pe activități umane și deci nu implică informații personale. Exemple de astfel de eforturi sunt cele efectuate în scopul analizării informațiilor științifice și de cercetare. Administrația Națională de Aeronautică și Spațiu, spre exemplu, explorează seturi complexe de date științifice legate de planetă pentru a identifica șabloane și legături pentru a detecta evenimente ascunse (sistemul este denumit Mașina de Învățare și Data Mining pentru o Înțelegere mai bună a Datelor Dimensionale legate de Pământ - Machine Learning and Data Mining for Improved Data Understanding of High Dimensional Earth Sensed Data).

În tabelul de mai jos se pot vedea câteva proiecte care sunt în stadiul operațional

Organizația	Denumirea sistemului	Descrierea sistemului	Scopul sistemului
Departamentul de Agricultură	Depozitul de date Financiare	Este sistemul de raportare intern de gestiune financiară. Data mining este făcută pentru rapoarte efectuate ad-hoc și la cerere.	Gestiunea financiară
Agenția de Management al Riscului	CAE	Este parte a unui proiect mandatat de Congres pentru a ajuta Agenția de Management al Riscului în controlarea fraudei, deșeurilor și abuzurilor în programul Corporației de Asigurări Federal Crop.	Detectarea fraudei, deșeurilor și abuzurilor
Oficiul de Inventii și Mărci al SUA	Modelul Proiecției compensațiilor în Depozitul de date al companiilor	Generează și pune la dispoziție date privind proiecția compensațiilor, atât salariile cât și sporurile în cazul angajaților curenți și asupra angajărilor planificate.	Gestionarea resurselor umane
Agenția Comisariatului Apărării	Sistemul de suport al Deciziilor Corporative/ Sistemul de Management al Operațiunilor Comisariatului	Explorează date pentru a produce date analitice privind operațiunile comisariatului. Furnizează informații cum ar fi care produse se vând mai bine sau cât de cinstite sunt casierele.	Îmbunătățirea serviciilor sau performanței
Agenția Serviciilor de Apărare	Întreprinderea Verity K2	Explorează date din comunitatea serviciilor secrete și căutări pe Internet pentru identificarea teroriștilor străini sau a cetățenilor americani conexați cu activități de terorism din străinătate.	Analizarea serviciilor secrete și detectarea activităților teroriste
Agenția Serviciilor de Apărare	Chestionarea Sistemului de Date ambulatorii	Pune în corespondență diagnosticul inițial al pacienților cu rezultatul testărilor și diagnosticării ulterioare. Permite identificarea timpurie a bolilor și afecțiunilor.	Monitorizarea sănătății publice
Agenția	Modus	Este un instrument de	Detectarea

Serviciilor de Apărare	Operandi	investigație utilizat pentru identificarea și urmărirea tendințelor în comportamentul criminalistic. Leagă caracteristicile crimelor și furnizează detalii privind locul producerii crimelor precum și alți factori criminalistici.	activităților și șabloanelor criminale
Administrația Alimentelor și medicamentelor	Analiza fonetică și ortografică computațională	Este un motor de căutare care furnizează rezultate indicând cât de similare sunt denumirile a două medicamente din punct de vedere fonetic și ortografic. Scopul său este de a ajuta la siguranța evaluării denumirilor private propuse pentru a reduce confuzia legată de numele medicamentului după ce o cerere în acest sens a fost aprobată de FDA.	Îmbunătățirea siguranței
Serviciul de Management al Minereurilor	Data Mining pentru sistemul de Management al Informațiilor Tehnice Baza de date (TIMS)	Este o bază de date corporativă pentru închirierile de gaze și petrol. Baza de date este explorată pentru sprijinirea dezvoltării politicii. O zonă a data mining este identificarea închirierilor care vor fi în viitorul apropiat. Data mining a relevat faptul că închirierile cu șase sau mai multe puțuri productive într-un an aproape niciodată nu vor fi abandonate în anul următor. O altă aplicare a data mining este în siguranța operațiunilor petroliere și cu gaze. De exemplu, data mining a arătat faptul că accidentele ating o cotă maximă joia dimineața.	Îmbunătățirea serviciilor sau performanței
Administrația	Mașina de	Va găsi șabloane și relații	Analiza

Națională de Aeronautică și Spațiu	Învățare și Data Mining pentru o Înțelegere mai bună a Datelor Dimensionale legate de Pământ	În seturi de date vaste și complexe legate de știința ce studiază planeta, pentru evenimentele mici și ascunse în semnalele datelor mai mari. Va construi noi capacități de înțelegere a datelor științifice ale NASA.	informațiilor științifice și de cercetare
Administrația Națională de Aeronautică și Spațiu	Analiza Seriilor de Timp Geofizice	Va dezvolta un set de algoritmi pentru identificarea șabloanelor în cadrul activităților temporale. Datele vor constitui traiectoriile și mișcările obiectelor în cadrul imaginilor.	
Administrația Națională de Aeronautică și Spațiu	Data Mining pentru Modelul Numeric 3-D pentru Previzionarea Randamentului și Aplicarea sa în Cercetarea Atmosferică	Va automatiza analiza randamentului modelului de vreme, observarea și datele din satelit pentru a permite o înțelegere mai bună a științei privind dinamica vremii pentru a preziona evenimentele viitoare.	Analizarea informațiilor științifice și de cercetare
Comisia de Reglementare Nucleară	Datele de Raportare a Evenimentelor Titularilor (Licențiatilor)	Identifică șabloanele și tendințele siguranței nucleare în evenimentele nucleare comerciale.	Îmbunătățirea siguranței

Tabel 1.3 - Proiecte de data mining în stadiul operațional

1.3. Clasificarea datelor

1.3.1. Noțiuni generale

Clasificarea este una dintre cele mai populare tehnici de data mining. Clasificarea permite construirea unui model clasificator pe baza datelor de antrenare, testarea acestuia pe datele de test și apoi utilizarea lui pentru a realiza predicții. Construirea modelului clasificator se poate face cu ajutorul unor algoritmi specializați precum *Naive Bayes*, *k-Nearest Neighbor*, *C4.5* etc. Modelul construit poate fi utilizat pentru a prezice valoarea nominală a unui atribut clasă.

Clasificarea se poate utiliza de exemplu în domeniul medical, pentru a extrage din datele unor pacienți reguli pe care apoi să le utilizăm pentru a diagnostica noi pacienți pe baza simptomelor lor. Se pot clasifica pacienții cu probleme la inimă pe baza unor diferite tipuri de boli de inimă.

Presupunând că D este un set de date, el poate fi privit ca un set de tupluri $x_1, x_2 \dots x_n$ unde $x_1, x_2 \dots x_n$ sunt valorile unor atribute $A_1, A_2 \dots, A_n$ relevante pentru un anumit studiu. Se pot defini diferite clase sau categorii $C = \{C_1, C_2 \dots C_m\}$ de instanțe în funcție de caracteristicile lor. Problema clasificării este de a defini o funcție $f : D \rightarrow C$ în care fiecare instanță $t_i \in D$ este pusă în corespondență cu ajutorul lui $f(t_i)$ către o anumită clasă C_j .

Diversi algoritmi realizează diverse procesări cu scopul de a clasifica datele. În literatură există un număr foarte mare de algoritmi de clasificare și numărul lor continuă să crească pentru că încă nu s-a descoperit un algoritm care să funcționeze mai bine decât toți ceilalți algoritmi pe orice set de date și pentru că nici unul din algoritmii existenți nu asigură o acuratețe a predicției de 100 pe orice set de date. Capitolul prezintă sintetic trei algoritmi de clasificare foarte populari și considerați algoritmi de referință în domeniul clasificării (*C4.5*, *Naive Bayes* și *k-Nearest Neighbor*). Aceștia au fost considerați de IEEE International Conference on Data Mining (ICDM 2006) ca făcând parte din topul celor mai influenți 10 algoritmi de data mining utilizați de comunitate [WU07]. Algoritmii menționați au fost utilizați în capitolele 4, 5 și 6 pentru a construi clasificatori și pentru a compara rezultatele obținute de acești algoritmi cu cele obținute de algoritmii propuși.

1.3.2. Algoritmul Naive Bayes

Principiul algoritmului: Algoritmul urmărește construirea unui modelul predictiv bazat pe calculul unor probabilități.

Pentru fiecare instanță pe care o clasifică, algoritmul calculează care este probabilitatea să aparțină la fiecare din clasele existente. Clasa estimată va fi cea care are probabilitatea de apariție cea mai mare [WU07]. Determinarea probabilității ca o instanță să aparțină unei clase implică realizarea unui produs de probabilități în care factorii produsului sunt, în cazul atributelor nominale, probabilitățile valorilor atributelor, iar în cazul atributelor numerice, funcțiile de densitate de probabilitate. Probabilitățile valorilor tuturor atributelor nominale precum și media și deviația standard (necesare pentru calculul densității de probabilitate) în cazul atributelor numerice se vor calcula pe baza datelor din mulțimea de antrenament, deci aceste calcule vor fi făcute în etapa de antrenare. În continuare este prezentată funcționarea algoritmului pe pași.

Etapa de antrenament:

- a) Se calculează pentru fiecare clasă probabilitatea ca o instanță să aparțină clasei respective (numărul de instanțe care au clasa considerată, raportat la numărul total de instanțe)
- b) Pentru fiecare valoare a fiecărui predictor nominal, se calculează probabilitatea condiționată de fiecare clasă a valorii respective (numărul de apariții ale valorii nominale printre cele m instanțe care aparțin etichetei clasei considerate)

- c) Pentru fiecare atribut numeric se calculează media aritmetică și deviația standard cu privire la numărul de instanțe din mulțimea de antrenament care fac parte din clasa considerată

Etapă de testare:

În această etapă se parcurg instanțele din mulțimea de testare, se calculează, pentru fiecare instanță, pe baza datelor obținute în etapa de antrenament, care este probabilitatea de a aparține la fiecare din clasele existente. Pentru fiecare instanță va fi estimată clasa care are probabilitatea cea mai mare. Se compară clasele estimate cu clasele reale. Se calculează acuratețea modelului care este egală cu numărul de instanțe din mulțimea de test clasificate corect, raportat la numărul total de instanțe din mulțimea de test. Dacă rezultatul este mulțumitor, modelul poate fi utilizat pentru predicții viitoare.

1.3.3. Algoritmul *k*-Nearest Neighbor

Principiul algoritmului: Algoritmul urmărește găsirea pentru fiecare instanță din mulțimea de test a instanțelor similare din mulțimea de antrenament. Algoritmul se bazează pe principiul „Dacă merge ca o rață, măcăie ca o rață și arată ca o rață, atunci probabil că este vorba de o rață...” [TAN05].

Algoritmul ***k*-Nearest Neighbor** construiește un model clasificator, care clasifică instanțele pe baza clasei majoritare pe care o au cele mai apropiate **k** instanțe de clasa considerată. Pentru a calcula distanțele dintre instanțe se folosesc diferite măsuri cum ar distanță Euclidiană [DEZ09], distanță Mincovski, distanță Manhattan.

Etapă de antrenament:

- a) Se alege valoarea lui **K**
- b) Se alege care va fi distanța care se va utiliza ca măsură a apropierii (Euclidiană, Mincovski, Manhattan)

Etapă de testare:

- a) pentru fiecare instanță din mulțimea de testare se calculează distanța față de fiecare instanță din mulțimea de antrenament
- b) se determină care sunt cele mai apropiate **k** instanțe din mulțimea de antrenament, de instanța de testare considerată
- c) se determină care este clasa dominantă printre cele mai apropiate **k** instanțe, pe baza votului majorității
- d) instanței din mulțimea de test îi va fi atribuită clasa rezultată prin votul majorității
- e) se analizează performanțele modelului și dacă aceasta este mulțumitoare, modelul poate fi utilizat pentru predicții

1.3.4. Algoritmul C4.5

Principiul algoritmului: Algoritmul face parte din categoria arborilor de clasificare [QUI93] (sau de decizie). Algoritmul urmărește construirea pe baza mulțimii de antrenament a unui model compus din reguli care poate fi reprezentat sub forma

unui arbore de decizie. Rădăcină arborelui corespunde întregii mulțimi de antrenament. Fiecare nod al arborelui reprezintă un test după un atribut și corespunde unei submulțimi a mulțimii de antrenament. Ramurile care pornesc de la un nod reprezintă rezultatele testului. Frunzele reprezintă clasele. Structura unui arbore de decizie construit cu algoritmul C4.5 poate fi vizualizată în figura 5.1. Pentru trasarea figurii s-a considerat că atributul 1 este nominal (cu trei valori posibile), atributul 2 este numeric, atributul 3 este nominal (cu două valori posibile) iar atributul clasă poate lua două valori.

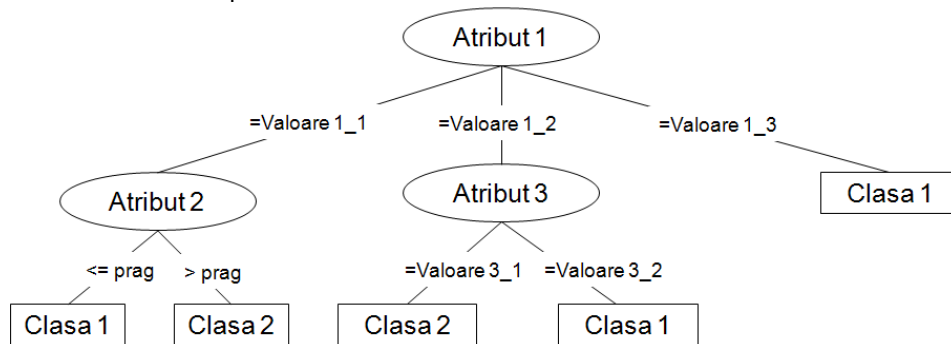


Fig. 1.3 - Structura unui arbore de decizie construit cu C4.5

Construcția arborilor este de tip inductiv [QUI86]. Inducția semnifică generalizarea unei reguli pe baza unor instanțe (sau exemple). Arborii sunt creați în mod inductiv pe baza valorilor instanțelor din mulțimea de antrenament. Pentru crearea arborelui se folosește un algoritm recursiv: divide et impera. Se va determina care este atributul care împarte cel mai bine datele, apoi instanțele se vor diviza în două subseturi. Dacă în urma divizării un subset conține mai multe clase, algoritmul se reia pentru subsetul respectiv. Dacă un subset indică o singură clasă algoritmul se oprește pentru acel subset.

Etapa de antrenament:

- în prima iterație la nivelul rădăcinii sunt prezente toate instanțele din mulțimea de antrenament. Se calculează capacitatea de separare a datelor (se calculează entropia și câștigul informațional) pentru fiecare atribut. Entropia este o măsură a dezordinii informaționale [HOL09]. *Information Gain* (câștigul informațional) oferă o măsură cantitativă a informației câștigate prin divizare (ordonare) după un anumit atribut. Atributul câștigător va fi cel care reușește să inducă cea mai multă ordine în setul de date (cel care are entropia cea mai mică și câștigul informațional cel mai mare). Acest atribut se adaugă arborelui.
- atributele pot să fie nominale sau numerice. Din fiecare atribut nominal pot porni un număr de ramuri egal cu numărul de valori nominale al atributului, iar din fiecare atribut numeric pot porni doar două ramuri: una corespunde valorilor numerice care sunt mai mici sau egale cu o valoare limită (în figura 1.3 a fost numită prag), iar cealaltă corespunde valorilor numerice care sunt peste prag (atributele numerice sunt tratate ca și atribute nominale cu două valori posibile). Pragul se determină calculând valoarea care conduce la un câștig informațional cât mai mare din mai multe puncte de împărțire posibile [QUI98]. Punctele de împărțire posibile sunt valorile aflate la jumătatea distanței dintre valorile numerice consecutive obținute prin ordonare.

- c) instanțele sunt divizate în subseturi, în funcție de valorile atributului câștigător. Pot să apară două situații:
- în unul dintre subseturi există reprezentată o singură clasă, în acest caz pentru acel subset se creează o frunză și se oprește procedura
 - în unul dintre subseturi există reprezentate cel puțin două clase, în acest caz se trece la iterația următoare cu subsetul respectiv (se reia algoritmul pentru acest subset)
- d) algoritmul se oprește când nu mai există nici un subset cu două clase reprezentate, au fost epuizate toate atributele sau au fost epuizate toate instanțele. În plus, mai pot fi specificate și alte reguli de oprire cum ar fi impunerea unei limite pentru adâncimea arborelui
- e) arborele construit în această etapă este posibil să se suprapotrivească datelor de antrenament (să nu fie suficient de general pentru a clasifica destul de corect instanțele din mulțimea de test, concept cunoscut sub numele de *overfitting*). Din această cauză se parcurge o etapă de tăiere a arborelui numită *prunning*. Tăierea se realizează parcurgând arborele și înlocuind cu noduri frunză ramurile care nu sunt relevante [WIT05].

Etapă de testare:

Pentru fiecare instanță din mulțimea de testare se parcurge arborele de la rădăcină și până la o frunză, ținând cont de valorile atributelor instanțelor testate și de testele din nodurile arborelui. Frunza la care se ajunge reprezintă clasa estimată.

Algoritmul J48 este implementarea în *Weka* a algoritmului C4.5 [MOO09].

1.4. Instrumente utilizate la realizarea de studii de data mining

În prezent există un număr foarte mare de instrumente cu care se pot realiza studii de data mining. Acestea se pot împărți în:

- instrumente comerciale (*SPSS Clementine, SAS E-Miner, MATLAB, Oracle data mining, SQL Server, SAP data mining, etc*)
- instrumente open-source (*Weka, R, Orange, Rapid Miner, Ant Miner, etc*)
- instrumente dezvoltate în urma unor prevederi contractuale pentru a rezolva anumite probleme concrete (*MaternQual* - sistem informatic integrat care utilizând tehnici de data mining permite identificarea factorilor de risc primari și secundari, pentru nașterile premature, precum și predicția riscului de naștere prematură, *NEONAT* - instrument care permite înregistrarea, procesarea digitală și analiza plânsetelor nou născuților [ROB11a]).

În cadrul tezei atenția a fost concentrată asupra instrumentelor Open Source, deoarece acestea au permis autorului să îmbunătățească tehnologia existentă prin îmbunătățirea unor algoritmi implementați în ele, prin adăugarea unor algoritmi noi, etc. Două dintre cele mai folosite instrumente la nivel mondial sunt *Weka* și *R* [HOR09]. Acestea vor fi prezentate pe scurt în continuare.

1.4.1. Weka

Weka este un program open-source de machine learning / data mining, dezvoltat în Java de către Univeristatea Waikato din Noua Zeelanda. Prima versiune internă a lui *Weka* a fost lansată în anul 1994, iar prima versiune publică a lui, versiunea 2.1 a fost lansată pe piață în anul 1996. În prezent s-a ajuns la versiunea 3.6.3 care este ultima versiune stabilă. *Weka* este un soft deosebit de util în educație, cercetare și aplicații [BOU08].

Versiunea 3.6.3 a lui *Weka* oferă următoarele facilități:

- cuprinde 70 de instrumente de preprocesare (pentru discretizare, reducerea zgomotului, selectarea atributelor, etc)
- 117 de algoritmi de clasificare și regresie (printre aceștia se găsesc J48, NaiveBayes, Random Forest)
- 11 algoritmi de clustering (cum ar fi SimpleKMeans, XMeans)
- 6 algoritmi pentru găsirea regulilor de asociere, printre care se găsește și algoritmul Apriori
- 3 interfețe grafice: Explorerul, Experimenter și KnowledgeFlow

Fișierele de date asociate cu *Weka* sunt fișierele arff (Attribute-Relation File Format), dar se pot deschide și fișiere C4.5, csv, dat, se pot importa date de pe internet sau din baze de date SQL.

Un exemplu de fișier arff este prezentat în figura 1.4.

```
@relation vreme

@attribute vremea {soare, innorat, ploios}
@attribute temperatura numeric
@attribute umiditate numeric
@attribute vant {da, nu}
@attribute jucam_golf {da, nu}

@data
soare,85,85,nu,nu
soare,80,90,da,nu
innorat,83,86,nu,da
ploios,70,96,nu,da
ploios,68,80,nu,da
ploios,65,70,da,nu
innorat,64,65,da,da
soare,72,95,nu,nu
soare,69,70,nu,da
ploios,75,80,nu,da
soare,75,70,da,da
innorat,72,90,da,da
innorat,81,75,nu,da
ploios,71,91,da,nu
```

Fig. 1.4 - Structura unui fișier *.arff

Atributele sunt precedate de eticheta @attribute și pot fi numerice, nominale (de genul soare, innorat, ploios), mai pot fi de tip *string* și *data*. În cazul tipului *data*

se va specifica formatul de dată [www4]. Datele sunt precedate de eticheta @data, instanțele se plasează pe rânduri diferite separate prin virgulă.

Interfața de start a lui *Weka* este cea din figura de mai jos:

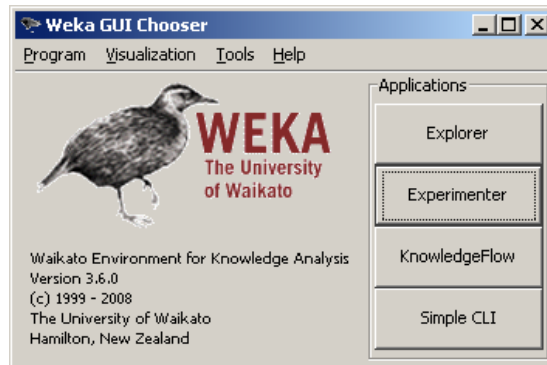


Fig.1.5 - Interfața principală a lui *Weka*

De pe interfața principală se pot rula patru aplicații: Explorer, Experimenter, KnowledgeFlow și Simple CLI. Explorer este o aplicație care permite realizarea analizei exploratorii a datelor. Interfața lui se poate vedea în figura 1.6.

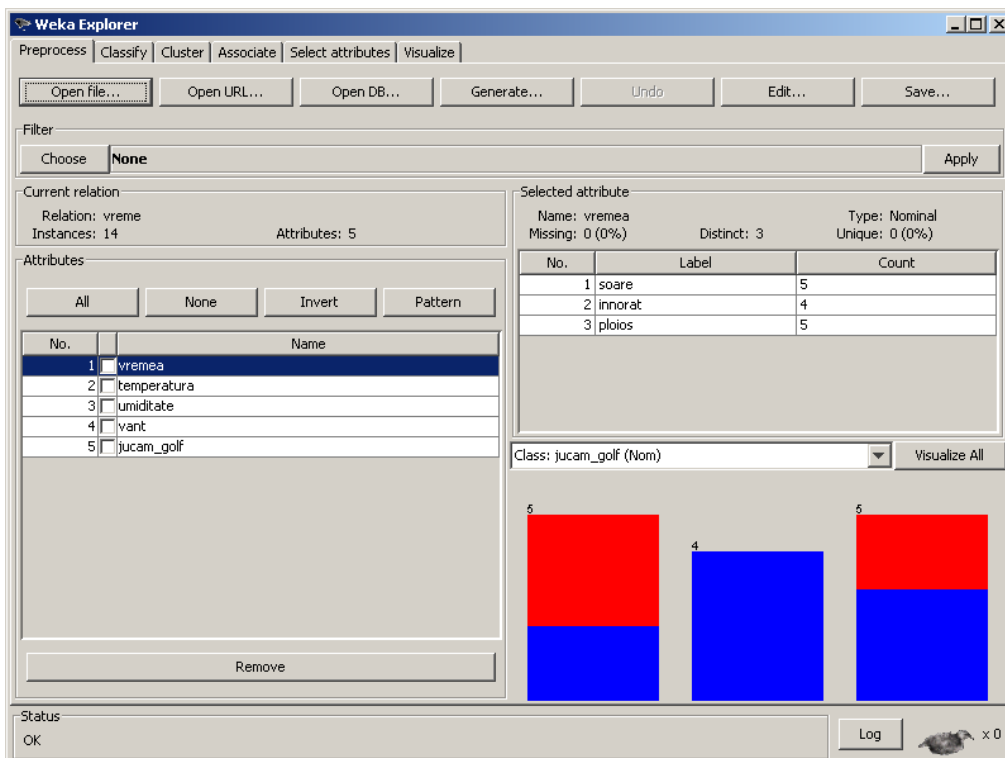


Fig.1.6 - Aplicația Explorer

Cu ajutorul comenzii *Open file...* se pot deschide fişierele de date. Odată ce acestea au fost deschise se pot vizualiza atributele din baza de date. Făcând click pe un atribut nominal se poate vizualiza care sunt valorile posibile ale acelui atribut nominal, precum și numărul de instanțe corespunzătoare fiecărei valori nominale (în exemplul considerat se poate vedea că în total avem 14 instanțe, în cinci din ele atributul vreme a avut valoarea soare, în 4 din ele a avut valoare înnorat, iar 5 a avut valoarea ploios). În cazul în care se selectează un atribut numeric se pot vizualiza automat valoarea minimă a atributului, valoarea maximă, media și deviația standard. În cadranul din partea dreaptă jos este un grafic care folosește un număr de culori egal cu numărul de clase, fiecare culoare corespunzând unei clase. Graficul ilustrează câte din instanțe aparțin unei clase și câte aparțin celeilalte clase.

Datele încărcate pot fi preprocesate, se pot selecta atributele relevante, se pot aplica filtre pentru reducerea zgomotului etc. Se poate folosi comanda Filter și apoi Apply. Pentru diverse filtre se pot stabili diverși parametri.

În continuare Tab-ul classify permite selectarea de algoritmi de clasificare și regresie, menționarea modului în care se va efectua testarea modelului construit (se poate folosi pentru testare mulțimea de antrenament, sau se poate menționa că un procent din setul inițial de date să fie utilizat pentru antrenament și cealaltă parte pentru testare, etc.). După ce se da comanda start și începe etapa de antrenare a modelului și de testare a lui, sunt prezentate rezultatele obținute.

The screenshot shows the Weka Explorer interface with the following details:

- Classifier:** J48 -C 0.25 -M 2
- Test options:**
 - Use training set:
 - Supplied test set: (Set...)
 - Cross-validation: Folds: 10
 - Percentage split: %: 56
- Classifier output:**

```

Correctly Classified Instances      9      64.2857 %
Incorrectly Classified Instances    5      35.7143 %
Kappa statistic                    0.186
Mean absolute error                 0.2857
Root mean squared error             0.4818
Relative absolute error             60 %
Root relative squared error        97.6586 %
Total Number of Instances          14

=== Detailed Accuracy By Class ===
                TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area
                0.778   0.6      0.7        0.778   0.737      0.789
                0.4     0.222   0.5        0.4     0.444      0.789
Weighted Avg.   0.643   0.465   0.629     0.643   0.632      0.789

=== Confusion Matrix ===
 a b  <-- classified as
 7 2 | a = da
 3 2 | b = nu

```

Fig. 1.7 - Model de clasificare

Se poate vedea o exprimare textuală a modelului construit și alți parametri care ne informează asupra performanțelor modelului (numărul de instanțe clasificate corect, numărul de instanțe clasificate greșit, matricea de confuzie).

Matricea de confuzie mai este numită și tabel de contingență. Numărul de elemente clasificate corect este egal cu suma elementelor de pe diagonala principală. Elementele clasificate greșit nu se găsesc pe diagonala principală. Dimensiunea matricei depinde de numărul de valori ale atributului clasă.

Alte măsuri utile pentru a compara clasificatorii sunt [BOU08]:

TP (true positive) - proporția de exemple care au fost clasificate de clasă x , printre toate exemplele care au intradevăr clasa x .

FP (false positive) – proporția de exemple care au fost clasificate de clasa x , dar de fapt aparțin unei alte clase printre toate exemplele care nu sunt de clasă x .

Precizia este numărul de exemple care au intradevăr clasa x , dintre toate care au fost clasificate ca fiind de clasă x .

În cazul în care se aplică un clasificator de tip arbore de decizie se poate vizualiza arborele de decizie creat:

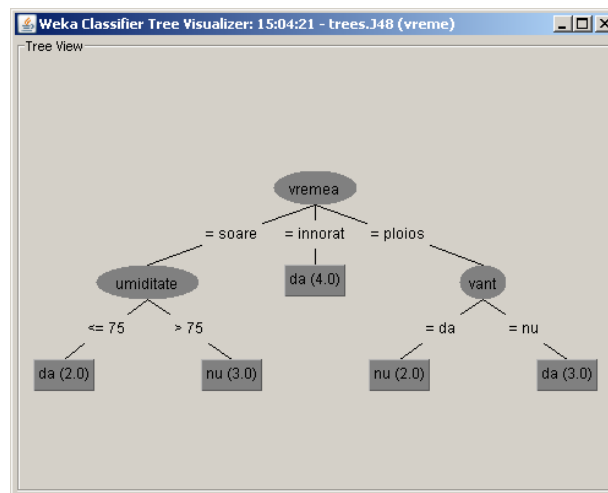


Fig.1.8 - Arbore de clasificare

Aplicația *Simple CLI* (*Command Language Interpreter*) este greu de utilizat, întrucât comenzile au foarte mulți parametri.

Aplicația *Experimenter* permite selectarea mai multor fișiere și a mai multor algoritmi, simultan. În cazul în care dorim să aplicăm mai mulți algoritmi aceluiași fișier sursă (sau mai multor fișiere sursă) putem face acest lucru manual cu explorerul sau automat utilizând *Experimenterul*. *Experimenterul* ne scutește de a repeta aceleași acțiuni.

Aplicația *KnowledgeFlow* reprezintă o alternativă la *Explorer*, care permite construirea de modele în mod grafic. Anumite funcționalități din *Explorer* nu sunt încă implementate în *KnowledgeFlow*.

Tabелul 1.4 prezintă avantajele și dezavantajele utilizării lui *Weka*. Acesta este ușor de utilizat, este flexibil, dar dacă vrem de exemplu să realizăm predicții cu un model pe care l-am construit și l-am testat trebuie să căutăm pe Internet cum se

face asta pentru că interfața nu este intuitivă. Totodată dacă vrem să înțelegem cum funcționează algoritmi, ce procesări fac, ce rezultate intermediare obțin, *Weka* nu ne este de folos, el furnizează doar rezultatul final.

Avantaje	Dezavantaje
Dispune de foarte mulți algoritmi	Unele lucruri nu sunt intuitive
Ușor de utilizat	Ascunde înțelegerea adevărată
Poate fi extins și particularizat	

Tabel 1.4 - Avantajele și dezavantajele lui *Weka*

În concluzie *Weka* este un instrument puternic, utilizat pentru a realiza diverse studii de *Data Mining*, și este folosit la scară largă. El oferă o paletă largă de algoritmi de clasificare, clustering, regresie, precum și multe instrumente foarte utile în analiza datelor. Fiind Opensource poate fi extins și personalizat.

1.4.2. R

R-ul oferă facilități pentru manipularea datelor, calcule, analiză, trasare de grafice și dezvoltare de soft. *R* este un soft și un limbaj și este considerat un dialect al limbajului *S* care a fost creat de AT & T Bell Laboratories [www5]. Cele mai multe din programele *S* (*S*-ul fiind un instrument clasic pentru comunitatea statistică) pot fi utilizate direct în *R*. Putem afirma despre *R* că este un vehicul pentru dezvoltarea de noi metode de analiză interactivă a datelor. *R* este un limbaj interpretat și orientat pe obiecte. Utilizatorul poate realiza acțiuni asupra obiectelor cu ajutorul operatorilor (aritmetici, logici și de comparare) și a funcțiilor. *R* spre deosebire de alte pachete statistice stochează toate rezultatele tot în obiecte. Sintaxa *R* este simplă, iar instrumentele de machine learning sunt distribuite în pachete [www5]. Tipurile de date din *R* pot fi numerice (numeric) sau nominale (factor). Câteva exemple de comenzi *R* sunt următoarele:

Comanda	Descriere
setwd("D:/SCOALA/doctorat/statistical aspects of data mining/Data")	Setarea directorului curent
Getwd()	Obținerea directorului curent
vector1<-c(1,10,12) vector2<-c(3,5,1)	Se creează doi vectori, unul inițializat cu valorile 1,10,12, iar celălalt cu 3,5,1
length(vector)	Returnează numărul de elemente al vectorului
My_data_set<-data.frame(attributeA=vector1,attributeB=vector2)	Crearea unui data set
mean(my_data_set[,1])	Furnizează media elementelor de pe prima linie a data set-ului creat
write.csv(my_data_set,"my_data_set_file.csv")	Salvarea datelor într-un fișier csv
Plot()	Trasarea unui grafic

Tabel 1.5 - Exemple de comenzi în *R*

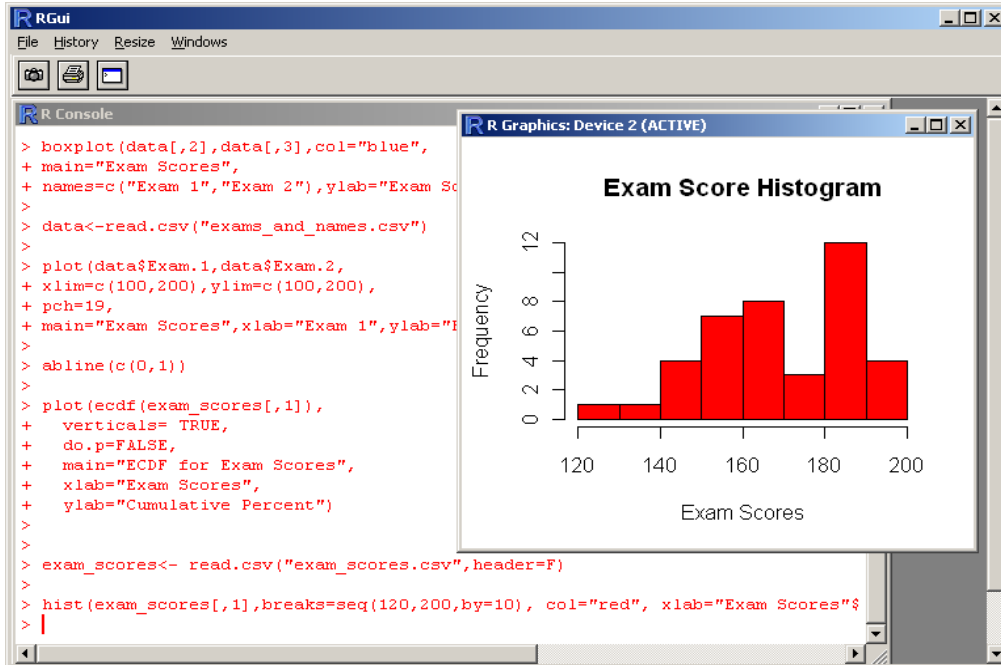


Fig. 1.9 - Interfața R – Exemplu de histogramă

Interfața R este prezentată în figura 1.9 împreună cu un exemplu de trasare a unei histograme. În figura 1.10 sunt prezentate alte două exemple de grafice trasate cu R, în cadrul unui studiu de comparare a rezultatelor la două examene cu scopul determinării care din examene este mai dificil.

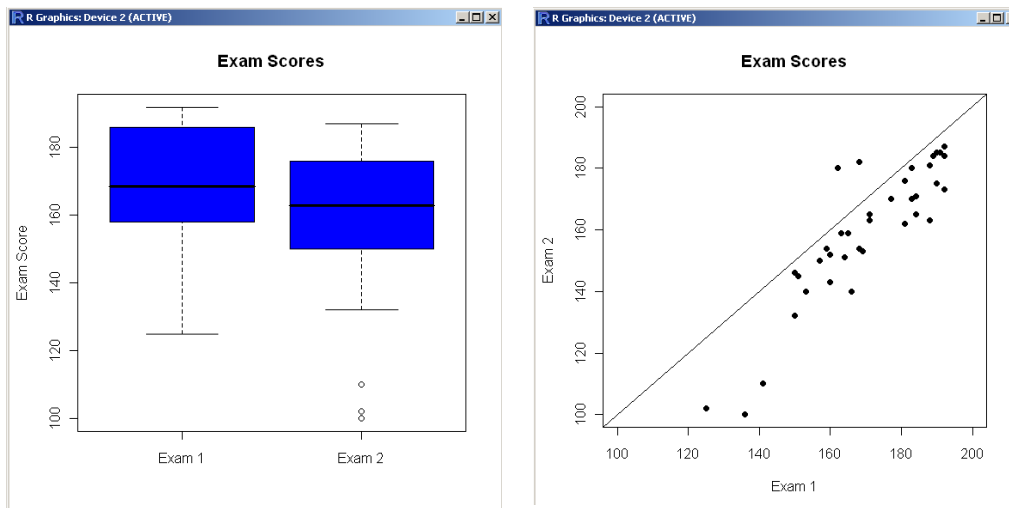


Fig.1.10 - Exemple de grafice trasate cu R

Tabelul 1.6 prezintă câteva avantaje și dezavantaje ale R-ului. R-ul are o sintaxă simplă, putem să-l programăm să facă lucrurile de care avem nevoie, cu ajutorul lui se pot crea grafice spectaculoase. Utilizatorii avansați pot să scrie cod C pentru a manipula obiectele R în mod direct. În R sunt implementate multe tehnici statistice și oferă un bun control al erorilor.

Un dezavantaj al R-ului este legat de pachete care pot să fie foarte dezorganizate, pot exista mai multe pachete cu implementări ale aceluiași algoritm, sau lucruri care ar trebui să se găsească într-un pachet pot să apară în alte pachete [www5]. Interfața grafică nu este foarte prietenoasă, iar pentru a scrie și rula un algoritm, este necesară o documentare cu privire la ce funcții trebuie utilizate și modul în care trebuie utilizate.

Pentru a face diferitele seturi de instrumente din mediile *Weka* și R accesibile într-un singur sistem unificat, autorii articolului "Open-Source Machine Learning: R Meets *Weka*" [HOR09] au sugerat un pachet R, numit RWEKA care să interfațeze funcționalitatea *Weka* către R.

Avantaje	Dezavantaje
Grafice spectaculoase	Interfață grafică nu foarte prietenoasă
Este un mediu în care sunt implementate multe tehnici statistice	Pachete destul de dezordonate
Oferă un control bun al erorilor	Necesită timp pentru utilizare
Poate executa cod C	
Se poate conecta la <i>Weka</i>	

Tabel 1.6 - Avantajele și dezavantajele lui R

1.4.3. Comparație între *Weka* și R

Atât R cât și *Weka* sunt sisteme Open Source utilizate pentru analiza datelor. Tabelul 1.7 prezintă o comparație între aceste două instrumente. Plusul arată că instrumentul corespunzător stă mai bine decât celălalt la capitolul considerat.

	<i>Weka</i>	R
GUI	+	
Algoritmii implementați	+	
Graficele trasate		+
Ușurința de utilizare	+	
Control		+
Încurajează învățarea		+
Popularitate	+	

Tabel 1.7 - Comparație între *Weka* și R

Interfața lui *Weka* este mult mai prietenoasă decât cea a R-ului, întrucât în *Weka* deschiderea unui set de date, rularea unui algoritm se face cu ajutorul unor click-uri de mouse pe opțiunile de meniu corespunzătoare, în vreme ce R-ul este un instrument dar și un limbaj, iar încărcarea datelor și execuția unui algoritm nu se poate face fără a apela din cod funcțiile corespunzătoare.

Weka are în plus față de R mult mai mulți algoritmi de clasificare, regresie, grupare și găsire a regulilor de asociere.

Ambele instrumente permit trasarea unor grafice spectaculoase, dar R-ul, în plus față de *Weka*, permite trasarea unor diagrame de tip *boxplot* (vezi figura 1.10 stânga).

Utilizatorii de *Weka* pot rula cu ușurință algoritmi, selectând-i din interfața grafică, pe când utilizatorii R-ului trebuie să apeleze algoritmi într-un mod mai greu, din cod.

În *Weka* există lucruri pe care vrem să le facem și nu putem să le facem fără să intervenim în codul său, pe când R-ul este un instrument dar și un limbaj deci poate fi programat să facă ce dorește utilizatorul.

Mulți dintre utilizatorii de *Weka* rulează algoritmi de cele mai multe ori cu parametrii implicați, fără să-și pună prea multe probleme cu privire la aceștia și în acest fel ajung să nu știe ce se întâmplă în profunzime. În R, întrucât pentru apelul oricărei funcții este necesară scrierea ei și a parametrilor, asta obligă la cunoașterea parametrilor dar complică procesul și îi cere utilizatorului să dedice mult mai mult timp.

Datorită ușurinței de utilizare și a paletii largi de algoritmi oferiți de *Weka*, acesta este mai popular decât R.

Weka este pe primul loc în preferințele autorului și de aceea a fost folosit în studiile următoare.

1.5. Seturi de date utilizate pentru analiza performanțelor algoritmilor

University of California Irvine găzduiește o arhivă numită *UCI Machine Learning Repository* [FRA10] care conține la ora actuală 211 baze de date reale, donate care sunt utilizate de comunitatea de data mining și machine learning pentru analiza empirică a comportamentului algoritmilor.

Arhiva a fost creată ca o arhivă *ftp* în anul 1987 de către David Aha și colegi absolvenți de la UC Irvine. De atunci a fost folosită de cercetători din toată lumea ca sursă primară de seturi de date. Arhiva a fost citată de mai mult de 1000 de ori, devenind astfel una din cele mai citate "lucrări" din domeniul calculatoarelor [FRA10].

În cadrul acestei teze au fost utilizate 10 seturi de date reale obținute de pe UCI Machine Learning Repository. Seturile de date utilizate sunt cele din tabelul 1.9. Capitolele în care sunt utilizate (o parte sau toate) aceste seturi de date sunt capitolele 3, 4 și 5. Coloana *Număr attribute* din tabelul 1.8 conține doar numărul de attribute predictive, nu conține și atributul clasă.

Nr	Setul de date	Număr instanțe	Număr atribute	Număr valori atribut clasa
1	Car	1728	6	4
2	Zoo	101	17	7
3	Ljubljana breast cancer	286	9	2
4	Diabetes	768	8	2
5	Cleveland heart disease	303	13	5
6	Iris	150	4	3
7	Mushrooms	8124	22	2
8	Nursurey	12960	8	5
9	Tic	958	9	2
10	Dermatology	336	34	6

Tabel 1.8 – Bazele de date utilizate pentru analiza comportamentului algoritmilor

Seturile de date conțin următoarele informații:

- **Car**, conține 1728 de instanțe și 6 atribute categoriale cu privire la evaluările unor mașini. Acestea sunt prețul de cumpărare al mașinii, prețul de întreținere a mașinii, ambele cu valorile posibile foarte mare, mare, mediu, mic, numărul de uși (2, 3, 4, 5 sau mai multe), numărul de persoane pe care poate să le transporte (2, 4, mai multe), dimensiunea portbagajului (mică, medie sau mare), siguranța (mică, medie, mare), atributul clasă cu valorile posibile (neacceptabilă, acceptabilă, bună, foarte bună).
- **Zoo** conține 101 instanțe și 17 atribute. Acestea sunt numele animalului care poate lua 100 de valori distincte, numărul de picioare poate lua valori întregi de la 0 la 9, tipul animalului este atributul clasa cu 7 valori posibile: mamifer, pasăre, reptilă, pește, amfibiu, insectă, nevertebrată, restul atributelor pot lua două valori. Acestea indică dacă animalul are păr, pene, face ouă, dă lapte, poate zbura, poate înota, este prădător, este dințat, are coloană vertebrală, respiră, este veninos, are aripioare, are coadă, este domestic.
- **Ljubljana breast cancer** – Datele au fost furnizate de medici de la Institutul de Oncologie a Centrului Medical Universitar din Ljubljana, Yugoslavia și conține 286 de instanțe și 9 atribute. O parte din acestea sunt vârsta (împărțită în intervale de 10 ani), dacă femeia este la menopauză sau nu, dimensiunea tumorii (12 intervale), gradul de malignitate (3 nivele), sânul (stâng sau drept), careul în care se găsește tumoarea, dacă sânul este iradiat sau nu. Atributul clasa indică dacă reapariția sau nu a unor evenimente specifice bolii.
- **Diabetes** – are un număr de 768 de instanțe și 8 atribute. Datele privesc un grup de amerindieni Pima, de lângă Phoenix, Arizona și sunt furnizate de Institutul Național de Diabet și Boli Digestive și de Rinichi. Atributele caracterizează femeii care au cel puțin 21 de ani și sunt următoarele: numărul de sarcini, concentrația de glucoză din plasmă la 2 ore de la

efectuarea unui test oral de toleranță la glucoză, tensiunea arterială diastolică, grosimea pielii de pe triceps, cantitatea de insulină serică la 2 ore, indicele de masă corporală, atributul clasă care indică dacă rezultatul testului este negativ sau pozitiv.

- **Cleveland heart disease** – datele au fost furnizate de Cleveland Clinic Foundation, și conțin 303 instanțe și 13 atribute. Câmpul clasă indică dacă pacientul are vreo boală de inimă. Valorile acestui câmp sunt cuprinse între 0 (nici o boala) și 4 (bolnav grav). Numele și alte informații confidențiale precum codul numeric personal au fost eliminate din baza de date. Atributele luate în considerare în studiu sunt vârsta, sexul, tipul de durere în piept (cu 4 valori posibile), tensiunea arterială în repaus, colesterolul, cantitatea de zahăr din sânge, rezultatul electrocardiogramei, ritmul cardiac maxim atins, angină pectorală indusă de exercițiu – cu valorile da sau nu, depresia ST indusă de exercițiu relativ la odihnă, panta segmentului ST extras, numărul de vase (0-3) colorate la flouoroscopie, diagnosticul bolii de inimă (stadiul din punct de vedere angiografic)
- **Iris** – Aceasta este probabil cea mai cunoscută bază de date referită în literatură care vizează descoperirea șabloanelor. Dataset-ul conține 150 de instanțe și 3 clase, de fiecare din aceste clase aparțin câte 50 de instanțe. Atributele sunt numerice și conțin informații cu privire la 3 specii de flori IRIS, Iris Setosa, Iris Versicolour și Iris Virginica. Informațiile care se memorează sunt lungimea și lățimea sepalei și lungimea și lățimea petalei. Aceste 4 atribute sunt numerice.
- **Mushrooms** conține informații despre ciuperci. Atributul clasă indică dacă acestea sunt comestibile sau otrăvitoare. Setul de date are 8124 de instanțe și 22 de atribute precum forma și culoarea pălăriei, mirosul etc.
- **Nursery** conține informații despre evaluarea aplicațiilor depuse de părinți la grădinițe. Dataset-ul are 12960 de instanțe și 8 atribute categoriale. Informațiile care se rețin privesc ocupația părinților, creșa copiilor, forma familiei, numărul de copii, condițiile de locuit, veniturile financiare ale familiei, imaginea socială și de sănătate a familiei, rezultatul evaluării (atributul clasă).
- **Tic-Tac-Toe** – este o bază de date pentru cunoscutul joc X și 0. Baza de date codifică toate situațiile posibile la sfârșitul jocului X și 0 în care X are mutarea de început. Setul de date are 958 de instanțe și 9 atribute, fiecare atribut corespunde unui pătrat de pe tabla de X și 0. Valorile posibile pentru fiecare pătrat sunt fie X, fie 0, fie B (jocul s-a încheiat și pătratul a rămas liber). Clasa are valorile posibile pozitiv și negativ (pozitiv înseamnă victorie pentru X).
- **Dermatology** – baza de date conține 34 de atribute, din care unul este întreg (vârsta) și a fost discretizat iar restul sunt nominale. Clasa reprezintă diagnosticul diferențial al bolilor *eritematoase*, *scuamoase* care este o adevărată problemă în dermatologie. Aceste boli împărtășesc toate caracteristicile clinice de eritem și scalare, cu diferențe foarte mici. Bolile din acest grup (valorile clasei) sunt psoriazis, dermatită seboreică, lichen plan, pitiriazisul rozat, dermatită cronică, și pitiriazis rubra pilaris. De obicei este necesară o biopsie pentru a diagnostica aceste boli, dar din păcate ele împart și multe caracteristici histopatologice de asemenea. Altă dificultate pentru diagnosticul diferențiat este aceea că o boală poate arata caracteristicile altei boli în stadiul de început și poate avea caracteristici

specifice în stadiile următoare. Pacienții au fost întâi evaluați din punct de vedere clinic, reținându-se 12 caracteristici. Apoi au fost luate mostre de piele pentru evaluarea a 22 de caracteristici histopatologice. Valorile acestor 22 de caracteristici s-au obținut prin analiza mostrelor de piele la microscop. În setul de date construite pentru acest domeniu, caracteristica privind istoria familiei are valoarea 1 în cazul în care oricare dintre aceste boli au fost observate în familie, și 0 altfel. Fiecare caracteristică clinică și histopatologica, a primit o valoare în intervalul de la 0 la 3. 0 indică faptul că caracteristica nu a fost prezentă, 3 indică cea mai mare prezență posibilă, și 1, 2 indică valorile relative intermediare.

1.6. Structura tezei

Capitolul 1 - prezintă definiții și terminologie, etapele procesului de data mining, domeniul în care tehnicile de data mining se aplică cu succes, algoritmi de referință pentru clasificare datelor, seturi de date utilizate pentru evaluarea performanțelor algoritmilor și două din cele mai populare instrumente utilizate la realizarea de studii de data mining (*Weka* și *R*). Cele două instrumente au fost comparate, și prezentate plusurile și minusurile fiecăruia.

Capitolul 2 - prezintă o sinteză privind tehnicile de preprocesare a datelor, aspecte legate de preprocesarea datelor cu *Weka* și aplicația concepută și implementată de autor (Arff Converter), care ușurează procesul de preluare a datelor din diferite baze de date (MySQL, Oracle, Microsoft Access, Microsoft Excel, Microsoft FoxPro) și încărcare a lor în *Weka*, etapă care ține de preprocesarea datelor. Pentru validarea soluției propuse s-a realizat o analiză a datelor unor studenți de la facultatea de Automatică și Calculatoare, preluate din aplicația GISC cu care este asigurată gestiunea școlărității la această facultate.

Capitolul 3 – prezintă aspecte legate de clasificarea datelor cu ajutorul algoritmului Ant Colony Optimization.

- A fost realizată o prezentare detaliată a algoritmului, bazată atât pe documentația existentă cât și pe studiul codului algoritmului.
- Au fost studiați parametrii algoritmului rulând 64 de configurații pe 10 seturi de date reale (prezentate în capitolul 1), obținute de pe *UCI Machine Learning Repository*. Rezultatele obținute au fost analizate atât vizual (comparativ) cât și matematic. Analiza matematică s-a realizat construind modele de regresie pe baza unor seturi de date noi construite pe baza testelor realizate.
- A fost realizată o sinteză a îmbunătățirilor aduse algoritmului din anul 2001 și până în prezent
- Au fost aduse 4 îmbunătățiri algoritmului care au fost implementate în aplicația open source *Ant Miner*. Aplicația care conține algoritmul îmbunătățit a fost numită *Ant-r-Miner*. S-a demonstrat că îmbunătățirile aduse conduc la o mai bună performanță utilizând seturi de date reale.

Capitolul 4 – prezintă aspecte legate de structura generală a unui algoritm genetic, aspecte legate de algoritmi genetici care au fost propuși în literatură pentru clasificarea datelor și propune un nou algoritm genetic în acest scop. Algoritm

genetic propus a fost numit AGR și implementat în *Weka*, contribuind în acest fel la îmbunătățirea tehnologiei de clasificare a datelor. Algoritmii genetici propuși diferă de alți algoritmi propuși în literatură pentru clasificarea datelor prin funcția fitness utilizată, facilitățile de recalculare a clasei în urma încrucișării, facilitățile de execuție repetată până se ajunge la un procent de acoperire dorit, facilitățile de utilizare a tuturor regulilor descoperite sau doar a celei mai bune reguli descoperite de fiecare iterație, facilitățile de adăugare a unei reguli implicite listei regulilor descoperite, etc. Analiza comportamentului algoritmului propus la diferite valori ale parametrilor de intrare s-a realizat rulând 14 seturi de parametri pe 5 seturi de date. Seturile de date au fost obținute de pe *UCI Machine Learning Repository* și prezentate în introducere. Analiza parametrilor s-a materializat într-o listă de recomandări privind alegerea acestor parametri. Rezultatele obținute de algoritmul pe 5 seturi de date au fost comparate cu cele obținute de algoritmi *Naive Bayes*, *k-Nearest Neighbor* și *J4.8* și au fost favorabile, în ceea ce privește acuratețea predicției, algoritmului genetic propus.

Capitolul 5 – Capitolul prezintă modul în care se pot realiza predicții cu *Weka*, soluțiile care au fost găsite pentru a îmbunătăți acest proces și îmbunătățirile care au fost aduse interfeței de clasificare a datelor. Pentru a îmbunătăți procesul de realizare al predicțiilor, au fost realizate modificări în secțiunea *Classify* a *Explorer*-ului și acestea constă în transformarea acestei interfețe într-o interfață dinamică în funcție de setul de date analizat, interfață cu care se pot realiza predicții într-un mod foarte simplu și intuitiv. Alte modificări realizate au transformat această interfață într-o interfață prietenoasă, ușor accesibilă chiar și începătorilor în *Weka*. Extensiile concepute și realizate în *Weka* contribuie la îmbunătățirea tehnologiei de clasificare a datelor.

Capitolul 6 – Capitolul prezintă o analiză a datelor legate de nașterile realizate în anul 2010 la *Clinica de Obstetrică - Ginecologie Bega, Timișoara*. Pentru realizarea studiului au fost disponibile date legate de 2326 de nașteri. Studiul a fost realizat cu ajutorul cunoștințelor și algoritmilor, respectiv instrumentelor propuse în primele 5 capitole ale tezei. Analiza datelor s-a realizat prin următoarele etape:

- Datele disponibile au fost întâi preprocesate, în Microsoft Excel prin analiza fiecărei coloane, eliminarea instanțelor cu valori lipsă sau valori introduse greșit, corectarea valorilor introduse greșit, eliminarea atributelor neimportante pentru analiză, crearea unor noi atribute, etc (vezi capitolul 2). Transformarea datelor din format *xls* în formatul *ARFF*, specific *Weka* cu ajutorul aplicației *Arff Convertor*, concepută, implementată și prezentată de autor în cadrul capitolului 2.
- Preprocesarea datelor cu ajutorul lui *Weka*, prin discreditarea atributelor numerice (vezi capitolul 2)
- Analiza statistică la nivel de atribut, care a condus la obținerea unor informații noi și interesante privind nașterile
- Dezvoltarea unor modele de clasificare cu algoritmi *Naive Bayes*, *J48*, *k-Nearest-Neighbour* din *Weka*, algoritmi care au fost prezentați în introducere
- Dezvoltarea unor modele de clasificare cu algoritmi *Ant-Miner* și *Ant-r-Miner*, prezentați în capitolul 3
- Dezvoltarea unui model de clasificare care poate estima cu o precizie de aprox 85 % unul din cele 5 intervale în care indicele apgar al nou-născutului

- Realizarea de predicții cu ajutorul noii interfețe dinamice a lui *Weka*, concepută și prezentată de autor în cadrul capitolului 5. Setul de reguli descoperit a fost transpus în format XML cu ajutorul unei aplicații concepute și implementate și prezentate de autor în capitolul 6 pentru a servi ca intrare unui sistem care le transpune în *ARDEN Syntax* [HEE05].

Capitolul 7 - prezintă concluziile și contribuțiile pe care autorul le-a adus la îmbunătățirea tehnologiei de clasificare a datelor.

1.7. Concluzii

În cadrul acestui capitol a fost realizată o sinteză amplă, folosind 47 de surse bibliografice, a celor mai importante aspecte legate de domeniul Data Mining. Au fost prezentate definiții și terminologia folosită în domeniu, etapele procesului de data mining, domenii în care tehnicile de data mining se aplică cu succes, aspecte legate de clasificarea datelor și de algoritmi folosiți pentru a construi modele clasificatoare, aspecte legate de instrumentele software utilizate în studii de data mining și de seturile de date utilizate pentru testarea și validarea algoritmilor. Structura tezei a fost de asemenea prezentată.

Contribuțiile aduse de către autor la cunoașterea în domeniu, prin intermediul acestui capitol, sunt prezentate la nivel de subcapitol în cele ce urmează:

- În subcapitolul 1.1, a fost conturat domeniul, subliniind importanța lui, au fost prezentați succint termenii existenți, un număr mare de definiții existente în domeniu, a fost introdusă o nouă definiție, au fost prezentate original etapele procesului de Data Mining, pe baza experienței practice a autorului
- În subcapitolul 1.2 – a fost realizată o sinteză amplă asupra aplicării tehnicilor de data mining în domeniile e-commerce, educație, medicină, combaterea terorismului. De asemenea a fost prezentată o sinteză a sistemelor de *data mining* utilizate de agențiile federale din SUA cu scopul de sublinia puterea și importanța tehnologiei *data mining*, precum și nivelul de vârf la care s-a ajuns cu utilizarea ei.
- În subcapitolul 1.3 a fost prezentate clar, succint aspecte generale legate de clasificarea datelor, precum și principiile de bază privind funcționarea a trei algoritmi de clasificare de referință (*Naive Bayes*, *k-Nearest Neighbor*, *C4.5*).
- În subcapitolul 1.4 au fost prezentate pe scurt aspecte legate de două dintre cele mai populare instrumente, la nivel mondial, utilizate pentru a realiza studii de data mining (*Weka* și *R*). Au fost analizate avantajele și dezavantajele pe care le are fiecare instrument și a fost realizată o comparație între cele două instrumente.
- În subcapitolul 1.5 au fost prezentate aspecte legate de cele 10 seturi de date donate, reale, obținute de pe UCI Machine Learning Repository, care au fost utilizate pe parcursul tezei, în capitolele 3,4 și 5 pentru a testa și valida algoritmi propuși și îmbunătățirile tehnologice aduse.

2. Preprocesarea datelor

2.1. Noțiuni introductive

Preprocesarea datelor este o etapă deosebit de importantă a procesului de *data mining*, practica industrială arătând că mai mult de 80 % din procesul de *data mining* se concentrează pe pregătirea datelor [ZHA05].

Datele brute pot conține valori lipsă, date introduse greșit, date expirate, date duplicat, valori aberante, etc. În acest sens se impune curățarea datelor brute, prin estimarea valorilor lipsă, prin eliminarea datelor duplicat sau expirate, prin eliminarea sau corectarea datelor introduse greșit, prin utilizarea unor parametri statistici toleranți la valori extreme, etc (vezi figura 2.1).

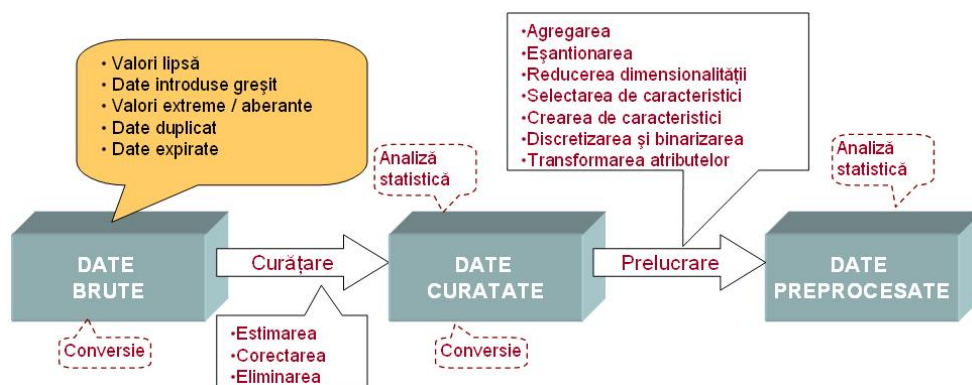


Fig. 2.1 – Preprocesarea datelor

Curățarea datelor nu este suficientă pentru ca acestea să poată fi explorate cu anumiți algoritmi de *data mining*. De foarte multe ori este necesară reducerea volumului datelor, pentru ca acestea să poată fi explorate cu algoritmi de *data mining* care sunt lenți. Reducerea volumului datelor se poate face reducând numărul de instanțe (de exemplu prin agregare și eșantionare) sau reducând numărul atributelor (prin tehnici de reducere a dimensionalității sau selectarea de submulțimi de caracteristici).

Datele inițiale pot fi nepotrivite cu anumiți algoritmi de *data mining* și în acest sens se impune o transformare a acestora prin tehnici precum discretizarea, binarizarea, crearea caracteristicilor etc.

Instrumentele de *data mining* dispun de mulți algoritmi de preprocesare a datelor care permit agregarea, eșantionarea, reducerea dimensionalității, etc. Pentru a putea aplica acești algoritmi asupra datelor ele trebuie să fie convertite în formatul acceptat de instrumentul software cu care se realizează studiul de *data*

mining. Conversia este tot o etapă a preprocesării și se poate aplica fie asupra datelor brute fie asupra datelor curățate.

Versiunea 3.6 a lui *Weka* dispune de 70 de algoritmi de preprocesare care permit discretizarea atributelor numerice, transformarea atributelor nominale în atribute binare, eșantionarea datelor, selectarea atributelor de interes pentru studiu, unirea a doua valori a unui atribut nominal într-una singură, analiza componentelor principale, etc. Fișierele de date cu care *Weka* lucrează sunt în format *ARFF*, dar *Weka* permite și încărcarea datelor din diferite baze de date utilizând fie tehnologia *JDBC* fie tehnologiile *JDBC* și *ODBC* împreună, dar acest proces este relativ complicat pentru că necesită realizarea unor configurări în *Weka* și în sistemul de operare, precum și cunoașterea limbajului *SQL*. Pentru a simplifica acest proces a fost concepută și dezvoltată de autor o aplicație numită *Arff Convertor* care permite conectarea cu ușurință la diferite tipuri de baze de date, preluarea datelor de interes din una sau mai multe tabele relaționate și exportarea lor în format *ARFF*, fără să fie necesară realizarea de configurări în Windows sau cunoașterea limbajului *SQL*. Aplicația concepută și dezvoltată simplifică procesul de preprocesare a datelor și contribuie în acest fel la îmbunătățirea tehnologiei de preprocesare a datelor cu *Weka*.

Capitolul prezintă o sinteză amplă și riguroasă a principalelor tehnici de preprocesare a datelor. De asemenea este prezentată aplicația *Arff Convertor* și pentru a demonstra utilitatea ei un studiu de caz realizat cu *Arff Convertor* și cu *Weka* pe datele studenților din anul universitar 2008-2009 de la Facultatea de Automatică și Calculatoare a Universității "Politehnica" din Timișoara.

2.2. Tehnici de preprocesare a datelor

2.2.1. Agregarea

Agregarea este procesul de combinare a două sau mai multe obiecte (instanțe) într-un singur obiect (instanță), pe baza caracteristicilor comune ale acestora. Presupunând că dispunem de o bază de date cu vânzările realizate în fiecare zi, de mai multe magazine din locații diferite pe parcursul unui an de zile, un exemplu de agregare ar fi înlocuirea tuturor tranzacțiilor corespunzătoare unui magazin cu o singură tranzacție, care să reprezinte totalul vânzărilor realizate de magazinul respectiv pe parcursul unui an de zile. În acest exemplu agregarea se realizează prin însumare, în alte cazuri se poate realiza prin medie. Un alt exemplu de agregare ar fi reducerea datelor de la 365 de zile la 12 luni.

Agregarea oferă câteva avantaje:

- seturile de date mai mici care rezultă în urma agregării necesită pentru analiză mai puțină memorie și mai puțin timp de procesare, așadar se pot folosi algoritmi mai consumatori de timp
- Agregarea permite o privire de sus asupra datelor în locul unei priviri de jos
- Comportamentul grupurilor de obiecte este mai stabil decât cel al indivizilor.

Un dezavantaj al agregării este posibila pierdere a amănuntelor interesante [TAN05]. În cazul în care se face o agregare la nivel de luni ale anului nu se mai pot determina zilele cu vânzările cele mai mari.

2.2.2. Eșantionarea

Eșantionarea este necesară pentru că de multe ori procesarea întregului set de date este mare consumatoare de timp.

În unele cazuri utilizarea unui algoritm de eșantionare poate reduce numărul instanțelor până în punctul în care se poate utiliza pentru analiza datelor un algoritm mai bun dar mai mare consumator de timp.

Eșantionul ales trebuie să fie reprezentativ. Un eșantion este reprezentativ dacă are aceleași proprietăți ca și setul de date original. Dacă proprietatea de interes este media obiectelor se va alege un eșantion la care obiectele au aceeași medie ca și cea a setului de date original [TAN05].

Exemplul cel mai familiar de eșantionare ne este oferit de sondajele de opinie (cu scop de informare politică sau socială) în care doar o proporție foarte mică a populației (câteva mii din câteva milioane) este interviuată.

Există mai multe tehnici de eșantionare, în continuare vor fi prezentate pe scurt cele mai cunoscute și anume:

- eșantionarea aleatoare
- eșantionarea prin stratificare
- eșantionarea adaptivă

Cea mai simplă tehnică de eșantionare este **eșantionarea aleatoare (simple random sampling)**, în cadrul acestei tehnici toate instanțele din populație au probabilități de selecție egale. Această tehnică poate fi implementată în două moduri:

- eșantionarea fără înlocuire (o instanță aleasă în mod întâmplător este lăsată în continuare în populație putând ca apoi să fie aleasă din nou)
- eșantionarea cu înlocuire (o instanță aleasă în mod întâmplător este eliminată din populație ca să nu mai poată fi aleasă din nou)

Când populația este formată din diferite tipuri de obiecte și fiecare tip de obiect apare de un anumit număr de ori, atunci prin tehnica descrisă mai sus este posibil ca obiectele care au o frecvență scăzută să nu fie reprezentate adecvat (unele tipuri de obiecte pot să nu fie selectate deloc). Acest aspect poate conduce la probleme atunci când este necesară o reprezentare corespunzătoare a tuturor tipurilor de obiecte. De exemplu când se construiesc modele de clasificare, o cerință critică este ca clasele rare să fie reprezentate în mod adecvat în eșantion. De aceea este necesară o schemă de eșantionare care este capabilă să se adapteze la diferite frecvențe ale item-urilor de interes. O astfel de tehnică este **eșantionarea prin stratificare**. Varianta cea mai simplă a acestei tehnici este extragerea unui număr egal de obiecte din fiecare grup, chiar dacă grupurile au numere de obiecte diferite. O altă variantă extrage din fiecare grup un număr de obiecte proporțional cu dimensiunea grupului [TAN05].

Odată selectată tehnica de eșantionare trebuie aleasă dimensiunea eșantionului astfel încât să nu se piardă informație relevantă dar nici să nu avem un eșantion prea mare. Determinarea dimensiunii potrivite a eșantionului este de multe ori dificilă așa că se folosesc tehnici de **eșantionare adaptivă**.

Eșantionarea adaptivă (progresivă) începe cu un eșantion mic pe care îl tot mărește până se obține un eșantion suficient de mare. Eșantionarea adaptivă

elimină nevoia de a determina inițial dimensiunea eșantionului, însă necesită o metodă de a determina dacă eșantionul este suficient de mare. Dacă se utilizează un eșantion adaptiv pentru a învăța un model predictiv, deși acuratețea modelului predictiv crește odată cu creșterea volumului eșantionului, la un moment dat, câștigul de acuratețe obținut prin creșterea dimensiunii eșantionului este neglijabil. În acel moment se poate opri creșterea volumului eșantionului.

2.2.3. Reducerea dimensionalității

Termenul "Reducerea dimensionalității" este adeseori rezervat pentru acele tehnici care reduc dimensionalitatea prin faptul că creează noi atribute care sunt combinații ale vechilor atribute. Dimensionalitatea poate fi redusă prin tehnici de algebră liniară cum sunt Analiza Componentelor Principale (**Principal Component Analysis**), Descompunerea Valorilor Singulare (**Singular Values Decomposition**), etc.

Avantajele reducerii dimensionalității sunt:

- mulți algoritmi de *data mining* lucrează mai bine dacă dimensionalitatea este redusă
- modelul rezultat poate să fie mai ușor înțeles dacă sunt mai puține atribute implicate
- reducerea dimensionalității poate să elimine atributele nerelevante și să reducă zgomotul
- datele pot să fie mai ușor vizualizate

Creșterea numărului de atribute al unui set de date, determină necesitatea creșterii numărului de instanțe în mod exponențial, dacă dorim să se păstreze un anumit nivel de acuratețe a datelor. Acest fenomen este cunoscut sub numele de "Blestemul Dimensionalității".

Blestemul dimensionalității este problema cauzată de creșterea exponențială în volum, asociată cu adăugarea unor extra-dimensiuni unui spațiu matematic [VER05].

Pentru a ilustra fenomenul "curse dimensionality", se consideră că spațiul complet al probabilității pentru o variabilă este reprezentat de intervalul unitate (0,1) și se consideră generarea a 10 instanțe în cadrul acestui interval. Fiecare instanță va reprezenta 10 % din spațiul probabilității (în medie). Apoi considerăm o a doua variabilă definită pe un alt interval ortogonal (0,1), de asemenea fiind reprezentată de 10 instanțe. Se vor produce 10 puncte în planul definit de liniile ortogonale x_1 și x_2 , reprezentând un nou spațiu de probabilitate. Dar noul spațiu are $10 * 10 = 100$ de unități, așa că fiecare dintre aceste 10 puncte reprezintă doar 1 % din spațiul de probabilitate. Ar fi necesare 100 de puncte pentru ca fiecare punct să reprezinte același 10 % din spațiul de probabilitate care era reprezentat de 10 puncte în spațiul cu o singură dimensiune .

Acest lucru este ilustrat în figura 2.2 în care cele 10 puncte nu sunt generate aleator pentru a ilustra acoperirea diminuată. Numărul de puncte ar trebui să crească exponențial pentru a menține o anumită acuratețe.

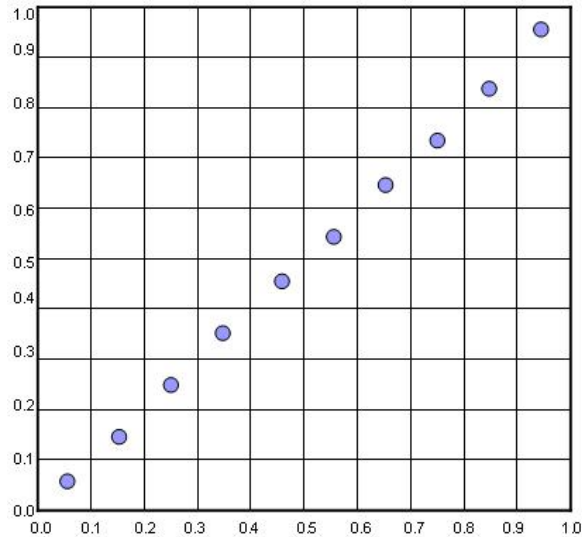


Fig. 2.2 – puncte în spațiul bidimensional

Tehnicile de bază folosite în reducerea dimensionalității sunt următoarele:

- Analiza componentelor principale (*PCA*)
- Descompunerea valorilor singulare (*SVD*)
- Analiza factorială
- Încapsularea locală lineară (Locally linear embedding - *LLE*)
- Scalarea multidimensională (Multidimensional Scaling - *MDS*)

Analiza componentelor principale

Analiza componentelor principale (PCA) este o tehnică de algebră liniară, pentru atribute continue, care găsește noi atribute (componente principale) care sunt :

- combinații liniare ale vechilor atribute
- ortogonale (perpendiculare) unele pe altele
- captează maximul de variație din date

Analiza componentelor principale caută să explice structura corelațiilor unui set de variabile (atribute), folosind un set mai mic de combinații liniare a acestor variabile. Aceste combinații liniare sunt numite componente. Variabilitatea totală a unui set de date, produsă de un set complet de n variabile, poate fi adesea surprinsă printr-un set mai mic de k combinații liniare a acestor variabile, ceea ce înseamnă că în cele k componente există aproape la fel de multă informație ca și în cele n variabile inițiale. Dacă se dorește, analiza poate apoi să înlocuiască cele n variabile originale cu cele $k < n$ componente, în așa fel încât setul de date să fie format din m instanțe cu k componente în loc de m instanțe cu n variabile [LAR06].

Scopul Analizei Componentelor Principale este de a găsi noi atribute care captează mai bine variabilitatea datelor decât cele vechi [JOL02]. Mai explicit, prima dimensiune (atribut) este aleasă pentru a capta o variabilitate cât mai mare. A doua dimensiune este ortogonală față de prima și captează cât de mult poate din variabilitatea rămasă, ș.a.m.d.

PCA are câteva caracteristici:

- tinde să identifice șabloanele cele mai puternice din date. De aceea ea poate fi utilizată ca o tehnică de găsire a șabloanelor.
- adeseori cel mai mult din variabilitatea datelor poate fi captat printr-o mică parte din setul total de dimensiuni. Ca un rezultat, reducerea dimensionalității utilizând PCA, poate determina date de dimensiuni relativ mici și este posibilă aplicarea tehnicilor care nu funcționează bine cu date de dimensiuni mari.
- zgomotul din date este mai slab decât tiparele, prin urmare reducerea dimensionalității poate să reducă mult din zgomot [TAN05].

Fiind dată o matrice D cu m linii și n coloane (m linii reprezintă obiectele sau instanțele și n coloane reprezintă atributele) matricea de covarianță a matricei D se notează cu S și este formată din elementele [LAR06]:

$$s_{ij} = \text{covarianța}(d_{*i}, d_{*j}) = \frac{\sum_{k=1}^m (d_{ki} - \mu_i)(d_{kj} - \mu_j)}{m} \quad (2.1)$$

În cuvinte s_{ij} este covarianță atributelor i și j . μ_i este valoarea medie a atributului i . Covarianța a două atribute este măsură a cât de mult două atribute variază împreună. Valoarea pozitivă a unei covarianțe indică că atunci când o variabilă crește și cealaltă tinde să crească. Valoarea negativă a unei covarianțe indică că atunci când o variabilă crește cealaltă tinde să scadă. Dacă o covarianță are valoarea 0 atunci cele două variabile sunt independente [LAR06]. Dacă $i=j$ atributele sunt aceleași și atunci covarianța este de fapt varianța atributului. Dacă matricea D este preprocesată așa încât media fiecărui atribut este 0, atunci $S = D^T D$. Scopul PCA este de a găsi o transformare a datelor care satisface următoarele proprietăți [JOL02]:

1. Fiecare pereche de noi atribute are covarianță 0 (pentru atribute distincte)
2. Atributele sunt ordonate în acord cu cât de mult surprinde fiecare atribut din varianța datelor
3. Primul atribut captează cât de mult se poate din varianța datelor
4. Subiect al cerințelor de ortogonalitate, fiecare atribut succesiv captează cât de mult se poate din varianța rămasă [TAN05]

O transformare a datelor care are aceste proprietăți poate să fie obținută analizând valorile proprii (eigenvalues) ale matricei de covarianță.

Valorile proprii și vectorii proprii sunt definiți doar pentru matrice pătratice. Matricea de covarianță exprimă varianța atributelor și este întotdeauna o matrice pătratică.

Valorile proprii măsoară cantitatea de varianță „explicată” de fiecare componentă principală și se calculează după cum este prezentat în continuare.

Fie S o matrice pătratică cu dimensiunea de $n \times n$. Numărul λ este o valoare proprie a matricei S dacă există un vector v diferit de 0 astfel încât:

$$Sv = \lambda v \quad (2.2)$$

În acest caz vectorul v este numit vector propriu al matricei S corespunzător valorii proprii λ . Relația (2.2) poate fi scrisă sub forma :

$$(S - \lambda I_n)v = 0 \quad (2.3)$$

Unde I_n este matricea unitate de dimensiune $n \times n$. Dacă matricea $S - \lambda I_n$ este inversabilă, se poate înmulți relația (2.3) cu inversa și se obține:

$$(S - \lambda I_n)^{-1}(S - \lambda I_n)v = (S - \lambda I_n)^{-1}0 \Rightarrow v = 0 \quad (2.4)$$

Dar noi căutăm un vector v diferit de 0, deci este necesar ca matricea $S - \lambda I_n$ să fie neinversabilă, adică să aibă determinantul egal cu 0. Calculând determinantul matricei $S - \lambda I_n$ se obține un număr maxim de n valori proprii

$\lambda_1, \lambda_2, \dots, \lambda_n$. Pentru a determina un vector propriu $v = \begin{bmatrix} v_1 \\ v_2 \\ \dots \\ v_n \end{bmatrix}$ diferit de 0 care

corespunde unei valori proprii λ , se rezolvă sistemul de ecuații liniare din relația (2.3).

Valorile proprii ale matricelor de covarianță S sunt nenegative și pot să fie ordonate sub forma $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{n-1} \geq \lambda_n$. Suma valorilor proprii este egală cu numărul variabilelor inițiale. Fie $U = [u_1, u_2, \dots, u_n]$ matricea vectorilor proprii ai matricei de covarianță S . Acești vectori proprii sunt ordonați în așa fel încât vectorul de pe poziția i corespunde valorii proprii de pe poziția i (valorile proprii la rândul lor sunt ordonate descrescător). Matricea D a fost preprocesată în așa fel încât media fiecărui atribut este 0.

În continuare se calculează matricea $D' = DU$. Fiecare nou atribut (din D') este o combinație liniară a atributelor originale. În mod specific, ponderile combinației liniare pentru atributul de pe poziția i sunt componentele vectorului propriu de pe poziția i . Varianța atributului de pe poziția i este λ_i . Suma varianței atributelor originale este egală cu suma varianțelor noilor attribute. Noile attribute sunt numite componente principale [TAN05]. Primul nou atribut este prima componentă principală, al doilea nou atribut este a doua componentă principală, etc. Vectorul propriu care este asociat cu cea mai mare valoare proprie indică care este direcția în care datele au cea mai multă varianță.

Vectorii proprii ai lui S definesc un nou set de axe. Analiza componentelor principale poate fi văzută ca o rotație a axelor originale către un nou set de axe care

sunt aliniate cu variabilitatea din date. Variabilitatea totală a datelor se păstrează dar noile atribute sunt acum necorelate.

Determinarea numărului optim de componente principale

O problemă care apare este *determinarea numărului optim de componente extrase*. Criteriile pentru această problemă sunt următoarele [LAR06]:

- *Criteriul valorilor proprii* - Valorile proprii exprimă importanța componentelor principale. Criteriul valorilor proprii prevede că doar componentele cu valori proprii mai mari de 1 trebuie reținute
- *Criteriul proporției varianței* – analistul trebuie să specifice cât de mult din variabilitatea totală dorește ca principalele componente să ia în calcul. Apoi analistul selectează componentele una după alta până la obținerea proporției dorite a variabilității.
- *Criteriul comunalității minime* - Comunalitatea reprezintă proporția de varianță a unei variabile particulare care este folosită în comun cu alte variabile. Comunalitățile reprezintă importanța globală a fiecărei variabile în PCA, ca și un întreg. De exemplu o variabilă cu o comunalitate mult mai mică decât celelalte variabile indică că această variabilă împarte mult mai puțină variabilitate comună cu celelalte variabile și contribuie mai puțin la soluția PCA. Comunalitățile care sunt foarte scăzute pentru o variabilă anume ar trebui să fie o indicație pentru analist că acea variabilă s-ar putea să nu participe la soluția PCA (spre exemplu s-ar putea să nu fie un membru al niciuneia dintre componentele principale). Per global, valori mari ale comunalității indică că principalele componente au extras cu succes o proporție mare a variabilității din variabilele originale. Valori scăzute ale comunalității indică faptul că mai este multă variație în setul de date care nu a fost luată în calcul de principalele componente. Valorile comunalității sunt calculate ca fiind suma ponderilor componentelor pentru o variabilă dată. Comunalitățile mai mici de 0.5 pot fi considerate ca fiind prea scăzute având în vedere că aceasta ar însemna că variabila are în comun mai puțin de jumătate din variabilitatea sa în comun cu celelalte variabile.
- *Criteriul Scree Plot* – Un scree plot este o reprezentare grafică a valorilor proprii în funcție de numărul pe care îl au componentele. Scree plot-urile sunt utile pentru găsirea pragului maxim în ceea ce privește numărul componentelor care trebuie reținute. Cele mai multe scree plot-uri sunt destul de asemănătoare ca și formă, începând de sus în partea stângă, căzând destul de rapid, aplatizându-se apoi la un anumit punct. Acesta datorită faptului că prima componentă explică de obicei cea mai mare parte din variabilitate, următoarele câteva componente explică într-o măsură moderată variabilitatea, în timp ce ultimele componente într-o măsură mică. Criteriul scree plot este următorul: numărul maxim de componente care ar trebui extrase este imediat anterior punctului în care graficul, începe pentru prima oară să se îndrepte formând o linie orizontală.

Descompunerea valorilor singulare (SVD)

SVD este o metodă pentru transformarea unor variabile corelate în variabile necorelate, care redau mai bine relațiile dintre datele originale. În același timp SVD este o metodă pentru identificarea și ordonarea dimensiunilor de-a lungul cărora datele au cea mai mare variație. Odată identificat unde este cea mai multă variație, este posibilă găsirea celei mai bune aproximări, a datelor originale, care folosește mai puține dimensiuni. Prin urmare SVD poate fi văzută ca o metodă pentru reducerea dimensionalității. SVD se bazează pe o teoremă din algebra liniară, care spune că o matrice dreptunghiulară A cu m linii și n coloane, poate fi împărțită într-un produs de 3 matrice în felul următor:

$$A = U\Sigma V^T \quad (2.5)$$

unde matricea U are dimensiunea $m \times m$, Σ are dimensiunea $m \times n$, iar V este de dimensiune $n \times n$. U și V sunt matrice ortogonale, așadar $UU^T = I_m$ și $VV^T = I_n$ [TAN05]. Matricea Σ este o matrice diagonală, entitățile de pe diagonală ei sunt ne-negative și sunt ordonate în ordine descrescătoare: $\sigma_1 \geq \sigma_2 \geq \dots \sigma_n \geq 0$. Valorile $\sigma_1, \sigma_2, \dots, \sigma_n$ se numesc **valori singulare** și se calculează extrăgând radicalul din valorile proprii. Vectorii coloană V, v_1, v_2, \dots, v_n se numesc **vectori singulari dreapta**, iar coloanele lui U se numesc **vectori singulari la stânga**.

Vectorii proprii ai matricei $A^T A$ sunt vectorii singulari dreapta (coloanele matricei V), iar vectorii proprii ai matricei AA^T sunt vectorii singulari stânga (coloanele matricei U). Valorile proprii nenule ale matricelor $A^T A$ sau AA^T sunt pătratele valorilor singulare.

Algoritmul SVD urmărește descompunerea matricei A în produs de trei matrice, etapele algoritmului sunt următoarele:

- Se calculează produsul matricelor $A^T A$
- Se determină valorile proprii ale matricei $A^T A$, extrăgând radical din acestea se obțin valorile singulare, se construiește matricea Σ amplasând valorile proprii pe diagonală principală în ordine descrescătoare
- Se determină vectorii proprii ai matricei $A^T A$ care corespund valorilor proprii calculate. Acești vectori proprii se amplasează pe coloanele matricei V , și apoi se calculează V^T
- Se calculează produsul matricelor AA^T , apoi se determină vectorii proprii ai acestui produs, care reprezintă coloanele matricei U

După parcurgerea etapelor de mai sus matricea A a fost descompusă ca produs a matricelor U, Σ și V^T . Reducerea efectivă a dimensionalității se realizează punând pe 0 valorile singulare cele mai mici [LES06] ale matricei Σ (luând în considerare doar valorile proprii cele mai mari - cele care surprind cea mai multă varianță). Punerea pe 0 a anumitor valori din matricea Σ , care este o matrice diagonală, determină obținerea unor coloane nule în această matrice. Coloanele nule

din matricea Σ vor determina anularea unor coloane din matricea A , atunci când se calculează produsul $U\Sigma V^T$. Calcularea produsului $U\Sigma V^T$, după ce au fost realizate prelucrările menționate, determină obținerea unei aproximări a matricei inițiale, care va avea câteva coloane nule. Colonele nule vor fi eliminate în acest fel reducându-se dimensionalitatea.

Analiza factorială

Analiza factorială este o metodă care permite comasarea informației dintr-un număr de variabile originale într-un set mai mic de dimensiuni (factori) cu o pierdere minimă de informație.

Tehnicile *PCA* și *SVD* produc noi atribute care sunt combinații ale atributelor originale. Analiza factorială exprimă atribute existente ca și combinații liniare a unui mic număr de atribute ascunse (factori ascunși sau factori latenți). Motivația este următoarea: deseori există caracteristici ale obiectelor, care sunt greu de măsurat în mod direct dar care par a fi legate de caracteristici măsurabile.

Un exemplu de analiză factorială este următorul [GOR06]: să presupunem că staff-ul unui lanț de supermarket-uri dorește să măsoare gradul de satisfacție a cumpărătorilor cu privire la serviciile oferite. Se consideră că prezintă interes următorii doi factori: a) gradul de satisfacție privind modul de servire și b) gradul de satisfacție privind calitatea produselor comercializate. Pentru măsurarea satisfacției clienților cu privire la serviciile menționate se realizează un sondaj de opinie printre 1000 clienți cărora li se adresează 10 întrebări. Fiecare răspuns a unui client reprezintă un scor. Scorul răspunsului dat de un client la o întrebare poate fi pus sub forma $7 * \text{satisfacție serviciu} + 5 * \text{satisfacție produs}$ (7 și 5 reprezintă încărcările factorilor iar factorii ascunși sunt satisfacție serviciu și satisfacție client).

Fie f_1, f_2, \dots, f_p atributele ascunse. Aceste atribute sunt atribute noi și au valori pentru fiecare obiect. Considerăm că matricea de date originale are dimensiunea de $m \times n$ și se numește D , iar noua matrice se numește $F = f_1, f_2, \dots, f_p$ și are dimensiunea de $m \times p$. Modelul de analiză factorială standard presupune că există următoarea relație între noile și vechile obiecte [TAN05]:

$$d_{i*}^T = \Lambda f_{i*}^T + \varepsilon \quad (2.6)$$

$$d_{ij} = \lambda_{j1} f_{i1} + \lambda_{j2} f_{i2} + \dots + \lambda_{jp} f_{ip} + \varepsilon_i \quad (2.7)$$

Λ care are intrările λ_{ki} este o matrice cu dimensiunea de $n \times p$ a încărcărilor factorilor (factor loadings) care indică, pentru fiecare din atributele originale, modul în care valorile originale depind de factori ascunși, cum ar fi noile atribute.

Factorii vor fi ortogonali unul pe celălalt ceea ce înseamnă că vor fi necorelați.

Încapsularea locală lineară

Încapsularea locală lineară (LLE - *Locally linear embedding*) este un algoritm de învățare nesupervizată care realizează maparea datelor de dimensiuni mari în mod neliniar spre un spațiu de dimensiune mică [SAU00]. Algoritmul este cel de mai jos:

1. Se găsesc vecinii cei mai apropiați pentru fiecare punct de date x_i din spațiul n dimensional. Pentru a determina distanța dintre puncte se poate calcula distanța Euclidiană și vecinii cei mai apropiați se pot determina cu tehnica k-Nearest- Neighbours
2. Se calculează ponderile w_{ij} care reconstruiesc cel mai bine fiecare punct x_i pe baza vecinilor săi, x_j . Se exprimă fiecare punct x_i ca și o combinație liniară a altor puncte $x_i = \sum_j w_{ij} x_j$, unde $\sum_j w_{ij} = 1$ și $w_{ij} = 0$ dacă nu este un vecin apropiat a lui x_i
3. Se găsesc coordonatele fiecărui punct din spațiul dimensional cu mai puține dimensiuni, al dimensiunii specificate p utilizând ponderile găsite în pasul 2 [TAN05].

În etapa a doua, matricea de ponderi W ale cărei intrări sunt w_{ij} , se calculează minimizând eroarea de aproximare la pătrat așa cum este măsurată în ecuația (2.8). Matricea W poate fi calculată rezolvând problema celor mai mici pătrate [SAU00].

$$error(W) = \sum_i \left(x_i - \sum_j w_{ij} x_j \right)^2 \quad (2.8)$$

Erorile de reconstruire sunt măsurate cu funcția de cost din relația 2.8 care adună pătratele distanțelor dintre toate punctele de date și reconstrucțiile lor.

Etapa 3 efectuează reducerea efectivă a dimensionalității. Fiind date o matrice ponderi și un număr de dimensiuni p specificate de utilizator, algoritmul construiește o vecinătate păstrând integrarea datelor în spațiul dimensional cu mai puține dimensiuni. Dacă y_i este vectorul din spațiul dimensional cu mai puține dimensiuni care corespunde lui x_i și Y este noua matrice de date al cărei rând i este y_i atunci aceasta poate fi realizată găsind un Y care minimizează ecuația următoare:

$$error(Y) = \sum_i \left(y_i - \sum_j w_{ij} y_j \right)^2 \quad (2.9)$$

Scalarea multidimensională

Tehnica standard *MDS*

Fiind date M obiecte în spațiul n dimensional și informația de distanță dintre ele, cum ar fi matricea de distanțe de dimensiune $M \times M$, ceea ce se urmărește este găsirea a M puncte în spațiul p – dimensional ($p < n$) astfel încât distanțele dintre obiecte să fie menținute atât de bine pe cât este posibil [FAL95].

Un exemplu clasic în domeniul scalării multidimensionale este următorul: având o matrice care conține distanțele dintre cele mai mari M orașe ale unei țări, scalarea multidimensională încearcă să găsească coordonatele corecte pentru aceste orașe.

În continuare se va descrie tehnica *MDS* standard pentru proiectarea datelor într-un spațiu p – dimensional. Întâi se calculează matricea de distanțe D , în care fiecare poziție d_{ij} reprezintă distanța de la obiectul i la obiectul j [TAN05]. Cele mai uzuale măsuri pentru calculul distanței dintre obiecte sunt *distanța Minkovski*, *distanța Euclidiană* și *distanța Manhattan* [DEZ09]. *Distanța Minkovski* [STE02] între două obiecte i și j care au n atribute, se calculează cu formula:

$$d_{ij} = \sqrt[r]{\sum_{k=1}^n (x_{ik} - x_{jk})^r} \quad (2.10)$$

Dacă r are valoarea 2, distanța este *Euclidiană*, iar dacă are valoarea 1, este *Minkovski*.

Fie d'_{ij} distanța dintre obiecte după ce acestea au fost transformate în spațiul p dimensional. Tehnica *MDS* clasică încearcă să asocieze fiecare obiect unui punct p - dimensional în așa fel încât o cantitate numită *stress* să fie minimizată [TAN05]. Cantitatea numită *stress* este definită după cum urmează:

$$stress = \sqrt{\frac{\sum_{ij} (d'_{ij} - d_{ij})^2}{\sum_{ij} d_{ij}^2}} \quad (2.11)$$

Versiunea clasică a lui *MDS* este un exemplu de tehnică *MDS metrică*, care presupune că disimilaritățile sunt variabile continue. Tehnicile *MDS non-metrică* presupun că disimilaritățile sunt categoriale. Tehnicile încearcă să asocieze obiectele punctelor p -dimensionale și apoi să modifice punctele pentru a reduce mărimea *stress*.

MDS începe cu o ghicire și o îmbunătățește în mod iterativ până când nu mai este posibilă nici o îmbunătățire. În versiunile sale cele mai simple algoritmul asociază fiecare obiect unui punct p dimensional (folosind o euristică sau chiar aleator). Apoi algoritmul examinează fiecare punct, calculează distanțele dintre el și celelalte puncte și mută punctul pentru a fi cât mai aproape de situația inițială.

FastMap

Algoritmul original *MDS* nu este potrivit pentru aplicații pe scară largă pentru că poate să necesite stocarea în memorie a unei matrice de distanțe de dimensiune $N \times N$ și poate avea o complexitate de $O(N^3)$ [ATA10]. De-a lungul timpului au fost propuse diverse îmbunătățiri pentru acest algoritm, printre care și *FastMap*. *FastMap* determină o coordonată la un moment dat, examinând un număr constant de rânduri din matricea de distanțe [PLA10].

Scopul algoritmului este de a găsi M puncte în spațiul p -dimensional, a căror distanțe Euclidiene să se potrivească cu distanțele Euclidiene din matricea de distanțe de dimensiune $M \times M$. Algoritmul realizează proiectarea punctelor pe k direcții mutual ortogonale. Se aleg două obiecte depărtate, numite pivoți și se alege linia care trece prin acești pivoți, în spațiul n dimensional. Ideea de bază a algoritmului este de a proiecta obiectele pe această linie. Determinarea proiecțiilor obiectelor pe această linie se realizează cu ajutorul *legii cosinusului*.

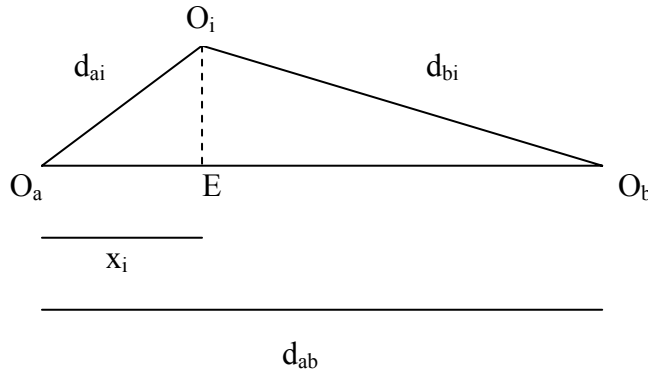


Fig. 2.3 – Proiectarea unui punct pe o dreaptă

Exprimând înălțimea triunghiului $O_a O_i O_b$ prin intermediul teoremei lui Pitagora în triunghiurile $O_a O_i E$ și $E O_i O_b$, se obțin relațiile de mai jos.

$$\begin{aligned}
 O_i E^2 &= d_{ai}^2 - x_i^2 \\
 O_i E^2 &= d_{bi}^2 - (d_{ab} - x_i)^2 \\
 d_{ai}^2 - x_i^2 &= d_{bi}^2 - (d_{ab} - x_i)^2 \\
 d_{ai}^2 - x_i^2 &= d_{bi}^2 - d_{ab}^2 + 2d_{ab}x_i - x_i^2 \\
 d_{bi}^2 &= d_{ai}^2 + d_{ab}^2 - 2d_{ab}x_i
 \end{aligned} \tag{2.12}$$

Legea cosinusului (numită și teorema lui Pitagora generalizată) spune că în orice triunghi este valabilă relația (2.12). Rezolvând ecuația pentru x_i se obține prima coordonată a obiectului O_i .

$$x_i = \frac{d_{ai}^2 + d_{ab}^2 - d_{bi}^2}{2d_{ab}} \quad (2.13)$$

Se observă că datorită ecuației (2.13) se pot mapa obiectele ca și puncte pe o linie păstrând o parte din informațiile legate de distanță: de exemplu dacă O_i este aproape de pivotul O_a atunci x_i va fi mic [FAL95]. În felul acesta problema este rezolvată pentru $p=1$ (pentru spațiul unidimensional). Maparea obiectelor ca și puncte în spațiul 2-dimensional sau p dimensional se bazează pe următoarea observație: considerăm un hiperplan de dimensiune $n-1$ care este perpendicular pe linia (O_a, O_b) . Fiecare punct din cele M puncte din spațiul n dimensional va fi proiectat pe hiperplanul de dimensiune $n-1$. Se vor calcula distanțele dintre proiecțiile punctelor și algoritmul se va relua recursiv. Fie O'_i și O'_j proiecțiile obiectelor i și j pe hiperplanul de dimensiune $n-1$. Distanța $D'(\)$ dintre O'_i și O'_j poate fi calculată pe baza distanțelor dintre punctele O_i și O_j în felul următor:

$$\left(D'(O'_i, O'_j)\right)^2 = \left(D(O_i, O_j)\right)^2 - (x_i - x_j)^2 \quad (2.14)$$

Abilitatea de a calcula distanța $D'(\)$ ne permite să calculăm proiecțiile pe o a doua dreaptă ortogonală pe prima. Astfel se poate rezolva problema pentru spațiul bidimensional și repetând același algoritm recursiv pentru spațiul p - dimensional [FAL95].

ISOMAP

MDS nu este indicat pentru reducerea dimensionalității atunci când punctele au o relație non lineară complicată unele față de altele. ISOMAP care este o extensie a tradiționalului MDS a fost dezvoltat pentru a lucra cu astfel de seturi de date. Algoritmul ISOMAP:

1. Se găsesc cei mai apropiați vecini ai fiecărui punct de date și se creează un grafic de ponderi conectând un punct la cei mai apropiați vecini. Nodurile sunt punctele de date iar ponderile legăturilor sunt distanțele dintre puncte.
2. Se redefinesc distanțele dintre puncte pentru a fi lungimea celei mai scurte căi dintre 2 puncte în graficul de vecinătăți
3. Se aplică clasicul *MDS* asupra noii matrice de distanțe.

Cei mai apropiați vecini pot fi definiți fie prin luarea punctelor de proximitate k , fie prin luarea tuturor punctelor pe o rază specificată de la un punct. Scopul etapei a doua este de a calcula **distanța geodezică** mai degrabă decât distanța Euclidiană. Distanța Euclidiană între 2 orașe aflate în colțuri opuse ale pământului este lungimea unui segment de linie care trece prin pământ, în timp ce distanța geodezică între 2 orașe este lungimea celui mai scurt arc pe suprafața pământului [TAN05] .

Probleme uzuale

O problemă comună a tehnicilor de reducere a dimensionalității este *calitatea rezultatului*. Se pune problema dacă o tehnică poate să producă o reprezentare a datelor de o încredere rezonabilă într-un spațiu cu mai puține dimensiuni. O reprezentare a datelor în spațiul cu mai puține dimensiuni trebuie să capteze caracteristicile importante ale datelor și să elimine aspectele care sunt irelevante sau chiar potrivnice. Reușita acestei operații depinde de tipul de date și de distribuția datelor care pot fi analizate prin abordarea de reducere a dimensionalității.

Tehnicile precum *PCA*, *SVD* și *analiza factorială* presupun o relație lineară între vechile și noile seturi de atribute. Deși această prezumție poate fi adevărată în multe cazuri, în multe alte cazuri este necesară abordarea neliniară. Mai exact, algoritmi cum sunt *ISOMAP* și *LLE* au fost dezvoltati pentru a trata relațiile neliniare [TAN05].

Complexitatea spațială și temporală a algoritmilor de reducere a dimensionalității este o problemă cheie. Cei mai mulți algoritmi de care am discutat au complexitate temporală sau spațială de $O(M^2)$ sau mai mare, unde M este numărul de obiecte. Aceasta limitează aplicabilitatea lor asupra unor seturi mari de date deși eșantionarea poate fi utilizată destul de eficient uneori. *FastMap* este singurul algoritm prezentat aici care are complexitate liniară, temporală și spațială.

Un alt aspect important al algoritmilor de reducere a dimensionalității se referă la *producerea aceluiși rezultat de fiecare dată când rulează*. *PCA*, *SVD* și *LLE* produc același rezultat. *Analiza factorială* și *tehnicele MDS* pot produce rezultate diferite la rulări diferite.

O altă problemă este determinarea numărului de dimensiuni pentru reducerea dimensionalității. Tehnicile pe care le-am luat în considerare pot realiza o reducere a dimensionalității în aproape orice număr de dimensiuni. În multe situații trebuie făcută o alegere între un număr mic de dimensiuni și o eroare de aproximare mai mare și o eroare de aproximare mai mică și mai multe dimensiuni.

2.2.4. Selectarea de submulțimi de caracteristici

Selectarea submulțimilor de caracteristici este procesul de identificare și de eliminare a cât mai multor caracteristici redundante sau irelevante [KOT06]. Caracteristicile pot fi clasificate în felul următor:

- *Relevante* – acestea au influență asupra ieșirii și rolul lor nu poate fi asumat de celelalte caracteristici
- *Irelevante* – nu au nicio influență asupra ieșirii
- *Redundante* – există o redundanță dacă o caracteristică poate prelua rolul alteia

Caracteristicile redundante dublează multe sau toate informațiile conținute în una sau mai multe alte atribute. Un exemplu de redundanță este memorarea datei nașterii și a CNP-ului în același tabel (CNP-ul include și data nașterii). Caracteristicile irelevante conțin informații care sunt nefolositoare pentru procesul de data mining derulat. De exemplu numărul de ordine al pacienților este irrelevant pentru analiza dacă suferă de o anumită boală sau nu. Caracteristicile redundante și relevante pot reduce acuratețea clasificării.

Deși anumite atribute irelevante și redundante pot fi eliminate imediat utilizând bunul simț și cunoștințele din domeniu, pentru selectarea celui mai bun subset de caracteristici este necesară o abordare sistemică [TAN05]. Abordarea ideală pentru selectarea caracteristicilor, implică încercarea tuturor subseturilor posibile de caracteristici ca date de intrare pentru algoritmul de *data mining* care ne interesează și selectarea subsetului care produce cele mai bune rezultate. Această metodă prezintă avantajul că reflectă obiectivul în funcție de algoritmul de *data mining* care va fi folosit. Din păcate fiindcă numărul de subseturi care implică n atribute este 2^n , o asemenea abordare este impracticabilă în cele mai multe situații fiind necesare strategii alternative. Există 3 abordări standard pentru selectarea subseturilor de caracteristici și anume: integrat, filtrat și împachetat.

Abordarea integrată

Selectarea caracteristicilor are loc în mod natural ca parte a algoritmului de *data mining*. În mod specific în timpul operării algoritmului de *data mining*, chiar algoritmul decide atributele pe care le va folosi și cele pe care le va ignora. Algoritmii utilizați pentru construirea clasificatorilor arborilor decizionali, operează deseori în această manieră.

Abordarea filtru

Caracteristicile sunt selectate înainte de rularea algoritmului de *data mining* utilizând o abordare care este independentă de sarcina de *data mining*. De exemplu putem selecta seturi de atribute a căror corelație aferentă este minimă.

Funcțiile de evaluare filtru pot fi împărțite în următoarele categorii [KOT06]:

- *Distanță* – Pentru o problemă cu două clase, o caracteristică X este preferată față de altă caracteristică Y dacă X induce o mai mare diferență între probabilitățile condiționate a celor două clase decât Y
- *Informație* – caracteristica X este preferată față de caracteristica Y dacă câștigul de informație de la caracteristica X este mai mare decât cel de la caracteristica Y
- *Dependență* – coeficientul este o măsură clasică de dependență și poate fi utilizat pentru a găsi corelația dintre o caracteristică și o clasă. Dacă corelația caracteristicii X cu clasa C este mai mare decât cea a caracteristicii Z cu clasa C atunci caracteristica X este cea preferată
- *Consistență* – două exemple sunt în conflict dacă au aceleași valori pentru un subset de caracteristici dar aparțin unor clase diferite

Abordarea împachetată

Aceste metode utilizează algoritmul de *data mining* țintă, ca și o cutie neagră, pentru găsirea celui mai bun subset de atribute, într-o modalitate similară algoritmului ideal descris mai sus, dar fără enumerarea tuturor subseturilor posibile.

Procesul de selectare al caracteristicilor constă în 4 părți (vezi figura 2.4): un algoritm de selectare care generează subseturile propuse de caracteristici și încearcă să găsească un subset optim, o măsură pentru evaluarea subsetului care determină cât de bun este subsetul propus, un criteriu de oprire și o procedură de validare.

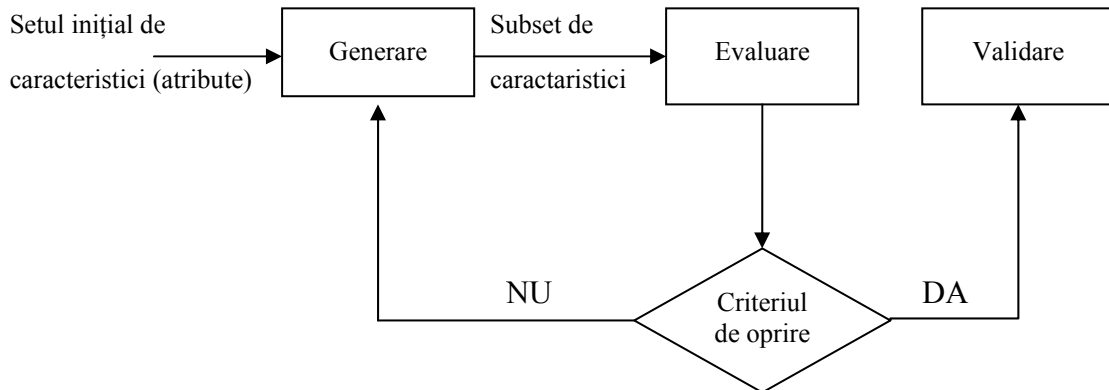


Fig. 2.4 – Fluxul procesului de selectare a unui subset de caracteristici

Selectarea caracteristicilor reprezintă o căutare a tuturor subseturilor de caracteristici. Diferite tipuri de strategii de căutare pot fi utilizate, însă strategia de căutare ar trebui să fie rapidă și ar trebui să găsească seturi de caracteristici optime sau aproape optime. De obicei nu este posibilă satisfacerea ambelor cerințe, astfel că sunt necesare compromisuri [TAN05].

Pentru că numărul de subseturi poate fi enorm și examinarea tuturor acestora nu este practică, este necesar un gen de criteriu de oprire. Criteriul de oprire poate fi unul din următoarele:

- adăugarea sau ștergerea oricărei caracteristici nu produce un subset mai bun
- s-a obținut un subset optim în acord cu o anumită funcție de evaluare
- s-a executat un număr maxim de iterații, se alege cel mai bun subset găsit
- s-a obținut un subset de o dimensiune dorită

În final, odată ce un subset de caracteristici a fost selectat, rezultatele algoritmului de *data mining* ținută pe subsetul selectat, trebuie validate. O abordare de evaluare directă este rularea algoritmului cu întreg setul de caracteristici și compararea rezultatelor complete cu rezultatele obținute utilizând subsetul de caracteristici. În cel mai bun caz subsetul de caracteristici va produce rezultate care sunt mai bune sau aproape la fel de bune ca și cele produse utilizând toate caracteristicile. O altă abordare de validare este utilizarea unui număr de algoritmi de selectare a caracteristicilor diferiți, pentru a obține subseturi de caracteristici și a compara rezultatele rulării algoritmului de *data mining* pe fiecare subset.

Ponderea caracteristicilor

Ponderea caracteristicilor este o alternativă la păstrarea sau eliminarea caracteristicilor. Caracteristicilor mai importante li se desemnează o pondere mai mare, în timp ce caracteristicilor mai puțin importante li se dă o pondere mai mică. Aceste ponderi sunt conferite uneori pe baza cunoștințelor în domeniu privind importanța relativă a caracteristicilor. În mod alternativ, ele pot fi determinate automat. De exemplu anumite scheme de clasificare cum sunt *SVM*, produc modele de clasificare în care fiecare caracteristică are o anumită pondere. Caracteristicile cu pondere mai mare joacă un rol mai important în cadrul modelului. Normalizarea

obiectelor, atunci când se calculează similaritatea cosinusoidală poate fi privită ca un tip de pondere a caracteristicilor [TAN05].

2.2.5. Crearea caracteristicilor

Crearea caracteristicilor este un proces care descoperă informații lipsă despre relațiile dintre caracteristici și mărește spațiul caracteristicilor deducând sau creând noi caracteristici. Presupunând că inițial dispunem de n caracteristici A_1, A_2, \dots, A_n după construirea de caracteristici putem avea încă m caracteristici $A_{n+1}, A_{n+2}, \dots, A_{n+m}$. De exemplu o nouă caracteristică A_k ($n < k \leq n + m$), poate fi obținută realizând o operație logică între atributele A_i și A_j . De exemplu o problemă bidimensională cu $A_1 = \text{lățimea}$ și $A_2 = \text{lungimea}$, poate fi transformată într-o problemă unidimensională, după ce se construiește aria $A_3 = A_1 * A_2$ [LIU03b]. Prin crearea de noi caracteristici se poate introduce redundanță. Aceasta poate fi apoi eliminată prin selectarea unor subseturi de caracteristici.

Deseori este posibilă crearea din atributele originale a unui nou set de atribute care captează în mod mai eficient informațiile dintr-un set de date. Mai mult, numărul de noi atribute poate fi mai mic decât numărul de atribute originale, permițându-ne să recoltăm toate beneficiile reducerii dimensionalității menționate mai sus. Crearea unui nou set de caracteristici din datele brute, originale este cunoscută drept *extragerea caracteristicilor*.

Maparea caracteristicilor

Procesul de transformare al caracteristicilor, poate extrage un set de noi caracteristici din setul original de caracteristici printr-o mapare funcțională. După extragerea caracteristicilor vom avea b_1, b_2, \dots, b_m unde $m < n$ și $b_i = f_i(a_1, a_2, \dots, a_n)$ și f_i este o funcție de mapare. De exemplu, pentru fiecare obiect o se pot defini: $b_1(x) = c_1 * a_1(x) + c_2 * a_2(x)$, unde c_1 și c_2 sunt constante [KOT06].

Construcția caracteristicilor

Uneori caracteristicile din seturile originale de date au informațiile necesare dar nu au forma potrivită pentru algoritmul de data mining. În această situație una sau mai multe caracteristici construite din caracteristicile originale pot fi mai folositoare decât acestea. Un exemplu în acest sens este următorul: presupunem că dispunem de un set de informații despre artefactele istorice, printre care și volumul și masa fiecărui artefact. Dacă considerăm că aceste artefacte sunt realizate dintr-un număr mic de materiale cum ar fi lemn, bronz, lut, aur și vrem să clasificăm aceste artefacte din perspectiva materialului din care sunt construite, cel mai ușor ar fi să construim o nouă caracteristică densitatea = masa / volum și să ne folosim de ea pentru a genera o clasificare corectă. Deși au existat tentative de a efectua automat construcția caracteristicilor (de exemplu algoritmul *Gala*), prin explorarea unor combinații matematice simple ale atributelor existente, cea mai utilizată abordare este construcția caracteristicilor utilizând expertiza în domeniu [TAN05].

Un algoritm care realizează construcția de caracteristici este algoritmul *Gala*. Acesta execută construcția caracteristicilor în timpul construirii unui arbore de

clasificare a deciziilor. Noi caracteristici sunt construite la fiecare nod al arborelui efectuând o căutare *branch and bound* în spațiul caracteristicilor. Căutarea este executată prin combinarea iterativă a caracteristicii cu cea mai mare valoare a câștigului de informație cu o caracteristică originală de bază care întrunește anumite criterii de filtrare. Gala construiește noi caracteristici binare utilizând operatori logici cum sunt conjuncția, negația [KOT06].

2.2.6. Discretizarea și binarizarea

Unii algoritmi de data mining în special anumiți algoritmi de clasificare, necesită ca datele să fie sub forma atributelor categoricale [ZHA05]. Algoritmii care găsesc șabloane de asociere necesită ca datele să fie sub forma atributelor binare, astfel deseori este necesară transformarea unui atribut continuu într-un atribut categorial (discretizare), iar atât atributele continue cât și cele discrete pot necesita a fi transformate în unul sau mai multe atribute binare (binarizare). Mai mult dacă un atribut categorial are un număr mare de valori (categorii) sau anumite valori nu au loc frecvent atunci poate fi benefic pentru anumite sarcini de data mining să reducă numărul de categorii prin combinarea unora dintre valori [TAN05].

La fel ca și în cazul selectării caracteristicilor, cea mai bună abordare de discretizare și binarizare este cea care produce cel mai bun rezultat pentru algoritmul de data mining care va fi utilizat pentru analiza datelor. De obicei aplicarea unui astfel de criteriu în mod direct nu este practică. Drept urmare discretizarea și binarizarea sunt efectuate într-o modalitate în care satisfac un criteriu care are o bună performanță pentru sarcina de data mining avută în vedere.

Discretizarea are multe avantaje: atributele discrete sunt mai ușor de utilizat, înțeles și explicat decât cele continue. Discretizarea face învățarea mai clară și mai rapidă. În general rezultatele obținute utilizând caracteristici discrete sunt mai compacte și mai clare, putând fi mai ușor examinate, utilizate și reutilizate.

Discretizarea atributelor continue

Discretizarea este folosită în general în problemele de clasificare și cele de analiză a asocierilor. Procesul de inducție al arborilor de decizie este mult mai complicat când se folosesc atribute continue decât atunci când se folosesc atribute discrete. Este necesară discretizarea atributelor continue fie înainte de inducția arborelui de decizie, fie în timpul procesului de construire al arborelui. Algoritmii C4.5 și CART, care sunt folosiți pe scară largă utilizează diverse tehnici pentru a evita lucrul direct cu valori continue [LIU02b]. În general cea mai bună discretizare depinde de algoritmul folosit și de celelalte atribute considerate.

Transformarea unui atribut continuu într-un atribut categorial presupune două acțiuni: determinarea numărului de categorii în care se va face transformarea și deciderea modului în care se va face maparea valorilor continue în aceste categorii. În primul pas, după ce valorile atributului continuu sunt ordonate, ele sunt divizate în n intervale, specificând $n-1$ puncte de împărțire (split point). În continuare, toate valorile dintr-un interval sunt mapate la aceeași valoare categorială. De aceea problema discretizării este de a găsi câte puncte de împărțire vor fi și unde vor fi ele plasate. Rezultatul poate fi ilustrat fie sub forma unor seturi

de intervale $\{(x_0, x_1], (x_1, x_2], \dots, (x_{n-1}, x_n]\}$ sau sub forma unor inegalități $x_0 < x \leq x_1, \dots, x_{n-1} < x \leq x_n$ [TAN05].

Metodele de discretizare pot să fie împărțite după mai multe criterii în metode supervizate versus nesupervizate, metode locale versus globale, metode statice respectiv dinamice și, metode sus – jos respectiv metode jos- sus [DOU95].

Cel mai comun criteriu de clasificare a algoritmilor de discretizare este clasificarea lor în algoritmi nesupervizați, care discretizează atributele fără să ia în considerare etichetele de clasă și algoritmi supervizați care discretizează atributele ținând cont de atributul clasă.

Metodele locale produc partiții care sunt aplicate pentru a localiza regiuni ale spațiului de instanțe. Metodele globale produc o mască peste întregul spațiu n dimensional al instanțelor continue, în care fiecare caracteristică este împărțită în

regiuni independente de celelalte atribute. Maska conține $\prod_{i=1}^n k_i$ regiuni, unde k_i este numărul de partiții ale atributului al i -lea [DOU95].

Multe tehnici de discretizare necesită un parametru k care indică numărul maxim de intervale care va fi produs în procesul de discretizare. Metodele statice realizează o trecere prin date pentru fiecare caracteristică și determină valoarea lui k pentru fiecare caracteristică, într-un mod independent de celelalte caracteristici. Tehnicile dinamice realizează o căutare prin spațiul valorilor k posibile, pentru toate caracteristicile simultan, astfel captând interdependențele în discretizarea caracteristicilor [DOU95]. O metodă dinamică discretizează valorile continue când clasificatorul se construiește (precum în algoritmul C4.5) în timp ce discretizarea statică se realizează înainte de procesul de clasificare [LIU02b].

Metodele de discretizare pot fi împărțite în metode sus-jos și în metode jos-sus. Metodele sus-jos încep cu o listă goală de puncte de tăiere și adaugă noi puncte la această listă împărțind intervalele pe măsură ce procesul de discretizare avansează. Metodele jos-sus încep cu o listă care conține toate valorile continue, ca și puncte de tăiere și se elimină treptat din aceste puncte, unind intervalele pe măsură ce procesul de discretizare avansează [KOT06].

Etapele unui proces de discretizare tipic sunt următoarele: (1) se sortează valorile continue ale caracteristicii care urmează să se discretizeze, (2) se evaluează puncte de împărțire sau intervale adiacente pentru unire, (3) în acord cu un anumit criteriu se împart sau se unesc intervalele de valori continue, (4) se oprește discretizarea.

Discretizarea nesupervizată

Cea mai simplă metodă de discretizare este metoda de discretizare nesupervizată numită discretizarea de dimensiune egală. Aceasta calculează maximul și minimul caracteristicii care urmează să fie discretizate și împarte domeniul observat în k intervale de dimensiuni egale. O altă metodă nesupervizată este discretizarea bazată pe frecvență egală. Aceasta numără câte valori sunt în atributul care urmează să fie discretizat și îl împarte în intervale conținând același număr de instanțe [KOT06]. Prin acest tip de discretizare se pierde informația legată

de clasă pentru că se combină în același interval valori care sunt legate de etichete de clase diferite.

Discretizarea supervizată

Discretizarea supervizată urmărește împărțirea caracteristicilor în intervale cât mai pure din punct de vedere al împrăștierii etichetelor de clasă. În acest sens este necesară localizarea punctelor de împărțire într-un mod care maximizează puritatea intervalelor. În continuare este prezentată o abordare simplă, bazată pe entropie.

Fie k numărul de etichete de clasă diferite, m_i numărul valorilor în intervalul al i -lea al unei partiții și m_{ij} numărul valorilor clasei j în intervalul i . Atunci entropia e_i a intervalului al i -lea este dată de ecuația:

$$e_i = \sum_{j=1}^n p_{ij} \log_2 p_{ij} \quad (2.15)$$

Unde $p_{ij} = \frac{m_{ij}}{m_i}$ este probabilitatea apariției clasei cu valoarea j în al i -lea

interval.

Entropia totală a partiției este media ponderată a entropiilor intervalelor individuale:

$$e = \sum_{i=1}^n w_i e_i \quad (2.16)$$

Unde m este numărul valorilor, $w_i = \frac{m_i}{m}$ este fracția valorilor în al i -lea

interval și n este numărul de intervale. În mod intuitiv entropia unui interval este o măsură a dezordinii acestuia. Dacă un interval conține doar valorile unei clase (este perfect pur) atunci entropia este 0 și nu contribuie la entropia totală. Dacă clasele valorilor într-un interval apar regulat (intervalul este cât de impur posibil) atunci entropia este maximă.

O abordare simplă pentru împărțirea unui atribut continuu începe prin împărțirea valorilor inițiale astfel încât cele două intervale care rezultă să aibă entropii minime. Această tehnică necesită considerarea fiecărei valori ca un posibil punct de împărțire pentru că se prezumă că fiecare interval conține seturi ordonate de valori. Procesul de împărțire este repetat apoi într-un alt interval alegând de obicei intervalul cu cea mai mare entropie până ce este atins un număr de intervale specificat de utilizator sau este satisfăcut un criteriu de oprire [TAN05].

Binarizarea

O tehnică simplă pentru binarizarea atributelor categoricale este următoarea: Dacă există m atribute categoricale, atunci se atribuie fiecare valoare categoricală unui întreg din intervalul $[0, m-1]$. Dacă atributul este ordinal atunci ordinea

trebuie să se mențină prin atribuire. Chiar dacă atributul este reprezentat original utilizând întregi, această reprezentare este necesară dacă întregii nu fac parte din intervalul $[0, m - 1]$. În continuare fiecare din acești întregi se convertesc la numărul binar corespunzător. Pentru a reprezenta în binar întregii sunt necesari n biți unde $n = \lceil \log_2(m) \rceil$. De exemplu pentru o variabilă categorială $\{insuficient, suficient, bine, foarte_bine\}$ sunt necesare două variabile binare, precum în tabelul de mai jos:

Valoare categorială	Valoare întregă	x1	x2
Insuficient	0	0	0
Suficient	1	0	1
Bine	2	1	0
Foarte bine	3	1	1

Tabel 2.1 – binarizare calificative

O astfel de transformare poate cauza complicații cum ar fi crearea de relații neintenționate între atributele create. De exemplu atributele x_1 și x_2 sunt corelate pentru că codificarea valorii categoriale *foarte_bine* necesită ambele atribute. Mai mult, analiza asociațiilor necesită atribute binare asimetrice, unde doar prezența valorii 1 este importantă [TAN05]. Pentru problemele de asociere este necesară introducerea unui atribut binar pentru fiecare atribut categorial, după cum se vede în tabelul 2.2.

Valoare categorială	Valoare întregă	x1	x2	x3	x4
Insuficient	0	1	0	0	0
Suficient	1	0	1	0	0
Bine	2	0	0	1	0
foarte bine	3	0	0	0	1

Tabel 2.2 – binarizare în probleme de asociere

Dacă numărul atributelor categoriale care rezultă în acest fel este prea mare atunci tehnicile prezentate în continuare pot să fie utilizate pentru a reduce numărul valorilor categoriale înainte de binarizare.

2.2.7. Transformarea variabilelor

Transformarea variabilelor se referă la o transformare aplicată tuturor valorilor unei variabile. Transformarea se poate realiza prin funcții simple sau normalizare.

Funcții simple

Pentru acest tip de transformare a variabilei se aplică o funcție matematică simplă în mod individual, fiecărei valori. Dacă x este o variabilă atunci exemple a unor astfel de transformări includ x^k , $\log(x)$, e^x , \sqrt{x} , $\frac{1}{x}$, $\sin(x)$. În statistică

transformările variabilelor în special radicalul, logaritmul și $1/x$ sunt deseori utilizate pentru a transforma datele care nu au o distribuție gaussiană, în date care au o astfel de distribuție.

Transformările variabilelor trebuie aplicate cu prudență pentru că modifică natura datelor. Deși se dorește acest lucru pot apărea probleme dacă natura transformării nu este evaluată complet. De exemplu, transformarea $\frac{1}{x}$, inversează ordinea elementelor variabilei. Pentru a clarifica efectul unei transformări este important să se urmărească dacă transformarea menține ordinea, dacă se poate aplica tuturor valorilor (în special celor negative și valorii 0) și care este efectul acesteia asupra valorilor cuprinse între 0 și 1 [TAN05].

Normalizarea și Standardizarea

Un alt tip uzual de transformare a variabilelor este **standardizarea și normalizarea** unei variabile.

Scopul standardizării sau normalizării este de a face ca un set întreg de valori să aibă o proprietate anume. Un exemplu tradițional este cel al "Standardizării unei variabile în statistică". Dacă \bar{x} este media valorilor atributelor și s_x este

deviația lor standard, atunci transformarea $x' = \frac{(x - \bar{x})}{s_x}$ creează o nouă variabilă

care are o medie de 0 și o deviație standard de 1. Normalizarea menționată este cunoscută sub denumirea *normalizarea z-score*. Dacă am combina într-un anumit fel variabile diferite, atunci o astfel de transformare este deseori necesară pentru a evita ca o variabilă cu valori mari să domine rezultatele calculului. De exemplu compararea oamenilor pe baza a două variabile: vârsta și venitul. Pentru oricare doi oameni diferența de venit va fi mult mai mare în termeni absoluți (milioane de lei) decât diferența de vârstă (mai puțin de 150 de ani). Dacă diferențele în spectrul valorilor vârstei și venitul nu sunt luate în considerare, atunci comparația dintre oameni va fi dominată de diferențele de venit [TAN05]. Dacă similaritatea sau nesimilaritatea a doi oameni este calculată utilizând măsuri specifice, atunci în multe cazuri cum ar fi distanța euclidiană, veniturile vor domina calculele.

Media și deviația standard sunt puternic afectate de valori extreme, așa că transformarea menționată mai sus este adeseori modificată. Mai întâi media este înlocuită de mediană și anume valoarea de mijloc. În al doilea rând, deviația standard este înlocuită de deviația standard absolută. Mai specific, dacă x este o variabilă, atunci deviația standard absolută a x este dată de $\sigma_A = \sum_{i=1}^m |x_i - \mu|$, unde x_i este a i -a valoare a variabilei, m este numărul de obiecte iar μ este fie media fie mediana.

O altă tehnică de normalizare este *normalizarea min-max* [KOT06]:

$$x' = \frac{x - \min}{\max - \min} (\text{noul_max} - \text{noul_min}) + \text{noul_min} \quad (2.17)$$

în care x' reprezintă noua valoare iar x reprezintă vechea valoare.

Normalizarea se mai poate realiza utilizând tehnica *normalizarea prin scalare zecimală*:

$$x' = \frac{x}{10^i} \quad (2.18)$$

unde i este cel mai mic întreg astfel încât $\max(x') < 1$.

2.2.8. Considerații privind tehnicile de preprocesare

Datele brute pot fi afectate de zgomot, pot exista valori aberante, valori lipsă, date duplicat, date introduse greșit sau date expirate [GOR06]. În acest sens se impune o *prelucrare a datelor brute*, prin utilizarea de filtre hard și soft pentru reducerea zgomotului, prin utilizarea unor parametri statistici toleranți la valori extreme, prin estimarea valorilor lipsă, prin eliminarea datelor duplicat, a datelor introduse greșit sau a datelor expirate. În urma prelucrării datelor brute se obțin date curățate.

Algoritmii de data mining care sunt lenți necesită foarte mult timp de rulare dacă setul de date este foarte mare. De aceea este necesară reducerea volumului datelor. Aceasta se poate face reducând numărul de instanțe (de exemplu prin agregare și eșantionare) sau reducând numărul atributelor (prin tehnici de reducere a dimensionalității sau selectarea de submulțimi de caracteristici).

Datele inițiale pot fi nepotrivite cu anumiți algoritmi de data mining și în acest sens se impune o transformare a acestora prin tehnici precum discretizarea, binarizarea, etc. De exemplu algoritmul de clasificare ID3 nu funcționează decât cu atribute nominale. Un set de date care conține atribute numerice poate fi clasificat cu ID3 doar dacă atributele numerice se discretizează.

2.3. Îmbunătățirea preprocesării cu Weka

2.3.1. Conectarea aplicației Weka la baze de date

Încărcarea datelor din diferite baze de date (citirea datelor) în instrumentele software cu care se va realiza studiul de data mining este tot o etapă a preprocesării. De multe ori și această etapă întâmpină greutăți datorită multitudinii de formate existente și datorită tehnologiilor de lucru care de multe ori nu sunt ușor accesibile utilizatorilor care nu sunt specialiști în calculatoare.

Conectarea lui *Weka* la baze de date, deși este posibilă este greoaie. Ea se poate realiza în două moduri: utilizând doar tehnologia Java Database Connectivity (JDBC) sau utilizând atât tehnologia JDBC cât și tehnologia Open Database Connectivity (ODBC) [www6].

Prima variantă de a lega *Weka* la diferite baze de date este utilizând *JDBC* driver-ul pentru baza de date la care se dorește conectarea. Fișierul **.jar* ce conține *JDBC* driverul trebuie adăugat în variabila de mediu care îi spune lui *JAVA* unde să caute după clase și care se numește *CLASSPATH*. Dacă nu există această variabilă de mediu atunci ea se va crea. Următoarea etapă este cea de navigare în directorul *weka/experiment/* a arhivei *weka.jar* și de setare a fișierului cu proprietăți *DatabaseUtils.props* corespunzător bazei de date (*DatabaseUtils.props.mysql*, *DatabaseUtils.props.oracle*, etc). În mod normal în acest fișier trebuie setat doar *jdbcDriver* și *jdbcURL*. Un exemplu de setări pentru o conectare la o bază de date *MySQL* locală este următorul:

```
jdbcDriver=com.mysql.jdbc.Driver
jdbcURL = jdbc:mysql://localhost:3306/db_name
```

Fișierul configurat trebuie redenumit în *DatabaseUtils.props* și amplasat în directorul de start a lui *Weka* pentru a putea fi recunoscut. Ultima etapă presupune deschiderea bazei de date cu ajutorul comenzii *OpenDB...a Explorerului* din *Weka* și construirea și rularea interogării *SQL* cu scopul aducerii datelor din baza de date în *Weka*.

Cea de-a doua variantă de conectare utilizează *JDBC-ODBC bridge* și presupune parcurgerea următoarelor trei etape: în prima etapă se utilizează *ODBC* pentru a crea un nou *Data Source Name (DSN)*, (se alege driverul corespunzător bazei de date la care se dorește conectarea, se menționează numele *DSN* -ului și datele necesare realizării conexiunii la baza de date. Dacă acel driver nu se găsește în *Windows*, el trebuie descărcat și instalat). A doua etapă este cea de completare a fișierului cu proprietăți *DatabaseUtils.props* cu datele necesare realizării conexiunii *JDBC*:

```
jdbcDriver=sun.jdbc.odbc.JdbcOdbcDriver
jdbcURL=jdbc:odbc:dbname
```

unde, *dbname* reprezintă numele *DSN*-ului creat. Acest fișier se găsește în directorul *weka/experiment/* al arhivei *weka*. Fișierul *DatabaseUtils.props* trebuie să fie recunoscut de *Weka*, pentru aceasta cea mai simplă variantă este mutarea lui în directorul lui *Weka*. Următoarea etapă este cea de deschidere a bazei de date, din meniul *Open DB...a Explorerului* din *Weka*. Dacă deschiderea bazei de date a reușit, urmează scrierea interogării *SQL* pentru aducerea datelor din baza de date în *Weka*. Datele odată aduse în *Weka* pot fi salvate în format *ARFF*.

Preluarea datelor dintr-o bază de date *Access* sau *FoxPro* cu ajutorul tehnologiilor prezentate mai sus nu este accesibilă utilizatorilor de *Weka* care nu sunt specialiști în calculatoare (necesită realizarea unor configurări în *Windows* și în *Weka* precum și cunoașterea limbajului *SQL*). Acesta este motivul pentru care a fost concepută și dezvoltată aplicația *ARFF Converter* care are rolul de a simplifica acest proces făcându-l accesibil oricui.

2.3.2. Preprocesarea datelor cu ajutorul lui Weka

Weka oferă foarte multe facilități pentru preprocesarea datelor. Instrumentele de preprocesare în *Weka* se numesc filtre. *Weka* conține filtre pentru analiza componentelor principale, discretizare, binarizare, selectarea de atribute, crearea de atribute etc. În figura 2.5 pot fi vizualizate o parte din filtrele pe care *Weka* le oferă. Atributele de tip șir de caractere pot fi convertite în atribute nominale, iar atributele nominale pot fi convertite în atribute binare sau șir de caractere. Atributele numerice pot fi transformate în atribute binare sau nominale. Selectarea unor submulțimi de atribute (caracteristici) se poate realiza din tab-ul *Select attributes*. Utilizatorul trebuie să menționeze metoda de căutare a unor submulțimi de atribute și metoda de evaluare a lor. Printre metodele de căutare a unor submulțimi de caracteristici se află metodele: *BestFirst* și *ExhaustiveSearch*. Două metode de evaluare a subseturilor sunt: *CfsSubsetEval* și *InfoGainAttributeEval*. Evaluatorul *CfsSubsetEval* evaluează valoarea unui subset de atribute considerând abilitatea de prezicere individuală a fiecărei caracteristici împreună cu gradul de redundanță dintre ele. Sunt preferate subseturile de caracteristici care sunt foarte corelate cu clasa și mai puțin corelate între ele. Evaluatorul *InfoGainAttributeEval* – evaluează valoarea unui atribut măsurând câștigul de informație în raport cu clasa.

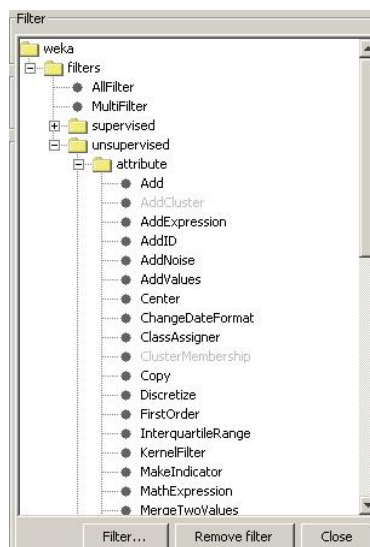


Fig. 2.5 - Preprocesarea datelor în *Weka*

2.3.3. ARFF Convertor

Noțiuni introductive

Weka este un program de machine learning și *data mining* foarte puternic și folosit pe scară largă. Fișierele de intrare pentru *Weka* sunt fișierele *ARFF* (*Attribute-Relation File Format*). Conectarea lui *Weka* la baze de date (vezi subcapitolul 2.10.1) este dificilă și necesită cunoștințe de *SQL*. O alternativă la acest proces complicat, o

reprezintă exportarea datelor din aceste baze de date în format *ARFF*, și apoi încărcarea acestor date în *Weka*[ROM08]. O soluție posibilă în acest sens este utilizarea soft-ului *Cahit Arf*, care se poate conecta la diferite baze de date și poate realiza conversii în *ARFF* din diferite formate, dar utilizarea acestui instrument necesită cunoștințe de SQL și restrânge în acest fel numărul utilizatorilor [ALK10].

Ca o soluție la problema menționată, autorul a conceput și dezvoltat un instrument (*Arff Convertor*), a cărui utilizare nu necesită cunoașterea limbajului *SQL*, care permite conectarea cu ușurință la baze de date (în prezent *Oracle*, *MySQL*, *Microsoft Access*, *Microsoft Visual FoxPro*, *Microsoft Excel*) și generarea unui fișier de intrare pentru *Weka* (în format *ARFF*) cu datele dorite din una sau mai multe tabele relaționate ale bazei de date (vezi figura 2.6).

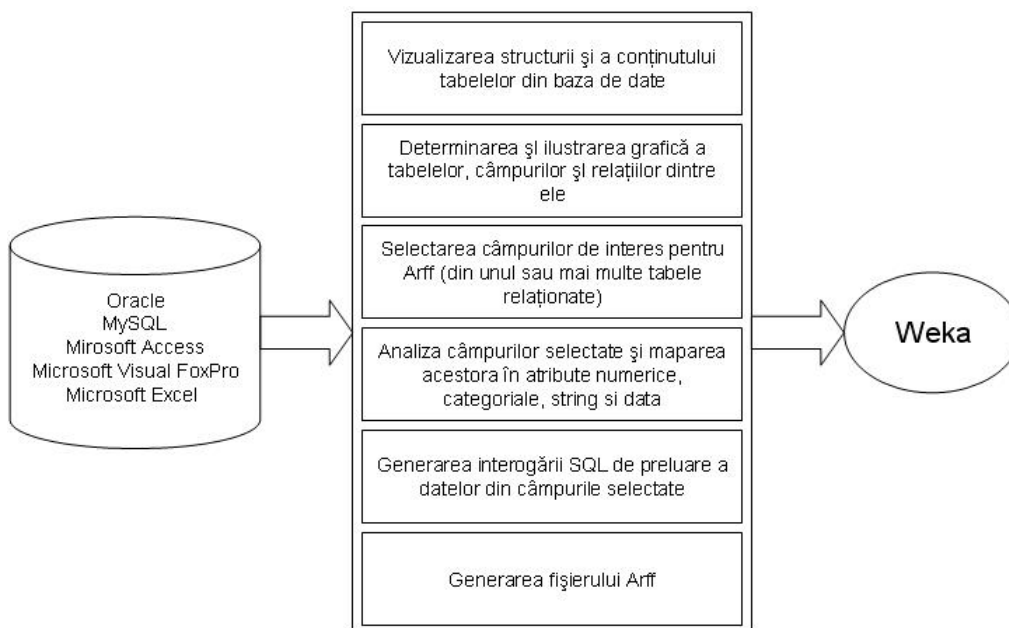


Fig. 2.6 - Funcționalitățile lui *ARFF* Convertor

Arff Convertor permite vizualizarea structurii și a conținutului tabelor din baza de date, ilustrarea grafică a tabelor (un tabel va fi reprezentat printr-un control de tip *listbox* iar câmpurile lui ca și item-uri în *listbox*) și a relațiilor dintre ele, selectarea cu ajutorul unor dublu-click-uri a câmpurilor de interes pentru generarea *ARFF*-ului, determinarea și alegerea tipurilor de atribute care vor deveni câmpurile selectate, etc. Aplicația analizează tipul și conținutul câmpurilor selectate din diferitele baze de date la care se conectează și sugerează tipul pe care îl vor avea câmpurile respective în format *ARFF*. De exemplu câmpurile de tip *tinyint*, *smallint*, *int*, *integer*, *bigint*, *real*, *double*, *float*, *decimal*, *numeric* dintr-o baza de date *MySQL* vor putea fi convertite în formatul *ARFF* fie în tip numeric fie într-un tip categorial. Variabilele de tip *varchar*, *char* din *MySQL* vor putea fi convertite fie în format *String*, fie în format categorial. Utilizatorul trebuie să aleagă tipul atributelor, iar dacă se determină că un atribut are mai puțin de 10 valori care se repetă atunci tipul implicit (sugerat) este categorial. Altfel poate fi numeric sau *String* în funcție de tipul inițial. Variabilele de tip boolean din toate bazele de date menționate vor

putea fi convertite automat în variabile categoriale în *ARFF*, categoriile posibile fiind: $\{Yes, No\}$, $\{Y, N\}$, $\{True, False\}$, $\{T, F\}$, $\{Da, Nu\}$, $\{D, N\}$.

Câmpurile selectate sunt analizate și dacă se constată că provin din tabele diferite se încearcă găsirea unei relații fizice între tabele. Dacă se determină o astfel de legătură, aceasta este utilizată pentru a genera o comandă SQL care va prelua datele din ambele tabele, (comandă ce va conține clauza *left join*). Dacă nu se determină o legătură între tabele (este posibil ca legătura între ele să fie una logică) atunci se va afișa o interfață care va conține tabelele selectate și câmpurile lor, utilizatorul trebuind să menționeze câmpurile de legătură. Dacă au fost selectate doar câmpuri dintr-un tabel, se generează comanda SQL care preia valorile câmpurilor selectate din tabelul ales. Aplicația dispune de interfețe dinamice care se adaptează la conținutul bazei de date deschise.

Pentru a demonstra buna funcționare și utilitatea soft-ului dezvoltat, autorul a preluat datele studenților de la Facultatea de Automatică și Calculatoare, Universitatea Politehnică din Timișoara, anul universitar 2008-2009 din Sistemul informatic de gestiune a școlarității, în format Microsoft Visual FoxPro. Datele preluate, sub forma a două tabele relaționate în format *dbf*, au fost convertite cu ajutorul aplicației propuse în format *ARFF*. Pentru a demonstra utilitatea acestei acțiuni s-a realizat o analiză statistică cu ajutorul lui *Weka* pe fiecare din atributele alese. Analiza statistică a permis extragerea unor concluzii interesante cu privire la nivelul de pregătire al studenților de la AC, cu privire la proveniența lor, numărul de credite obținute după primul semestru, etc. Detalii despre acest studiu se găsesc în subcapitolul 2.10.4.

Detalii tehnice și de implementare

Aplicația realizează următoarele prelucrări:

- explorarea structurii și conținutului bazei de date deschise
- selectarea tabelor de interes
- afișarea dinamică a tabelor, câmpurilor și a legăturilor fizice dintre tabele
- selectarea câmpurilor de interes pentru *ARFF*

```
while (mai sunt câmpuri selectate){ //parcurea listei câmpurilor selectate
    Adaugă un combobox pentru acest câmp
    if (câmpul curent este numeric){
        - adaugă în combobox opțiunea de atribut numeric
        - determină categoriile posibile
        - adaugă în combobox opțiunea de atribut categorial
    }
    else
    if (câmpul curent este sir de caractere){
        - adaugă în combobox opțiunea de atribut de tip string
        - determină categoriile posibile
    }
}
```

```

        - adaugă în combobox opțiunea de atribut categorial
    }
    else
        if (câmpul curent este boolean){
            - adaugă în combobox opțiunea de atribut categorial cu categoriile
posibile {Yes,No}, {Y,N}, {True,False}, {T,F},
{Da,Nu} și {D,N}
        }
        else
            if (campul curent este de tip data)
                - adaugă în combobox opțiunea Data
    }

- determinarea numărului de tabele din care provin câmpurile selectate

if (câmpurile selectate provin dintr-un singur tabel){
    - generează comandă select pentru preluarea datelor selectate
    - scrie datele în ARFF
}
else{
    determină relațiile dintre tabelele de care aparțin câmpurile selectate
    if (dacă au fost determinate relații)
        - generează o comandă SQL cu clauza tip left join care preia date
din tabele
    else{ //utilizatorul trebuie să menționează legăturile
        - generează o interfață cu tabele selectate și toate câmpurile lor
        - utilizatorul menționează legăturile dintre câmpuri
        - generează o comandă SQL cu clauza tip left join care preia date
din tabele
    }
}

```

Proiectul a fost realizat în mediul *Microsoft Visual Studio 2008*, *Visual C#* și permite conectarea la baze de date *Oracle*, *MySql*, *Microsoft Access*, *Microsoft Excel*, *Microsoft Visual FoxPro* [FIL08]. Tehnologia utilizată pentru conectarea la aceste baze de date este *OLEDB* (Object Linking and Embedding, Database). Pentru conectarea la bazele de date menționate au fost utilizați următorii furnizori: *MSDAORA.1* pentru *Oracle*, *MySQLProv* pentru *MySQL* (a fost instalat conectorul *MyOLEDB.3*), *Microsoft.Jet.OLEDB.4.0* pentru *Microsoft Access* și *Microsoft Excel*, *VFPoledb.1* pentru *Microsoft Visual Fox Pro* (a fost instalat conectorul *Microsoft Visual FoxPro Oledb Provider*). Figura 2.7 prezintă interfața principală a aplicației.

Tabelele din baza de date deschisă sunt afișate sub formă arborescentă într-un control *TreeView*. Declanșarea evenimentului *AfterSelect* al controlului *TreeView* (determinată de un click pe numele unui tabel) determină afișarea în controlul *ListView* a structurii tabelului selectat iar în controlul *DataGridView* a conținutului tabelului.

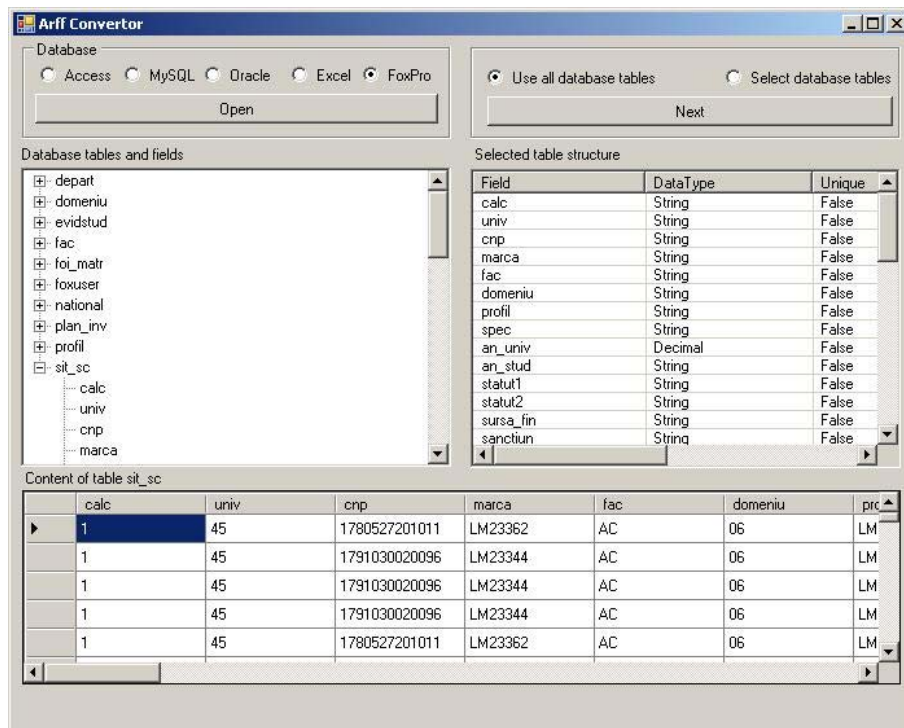


Fig. 2.7 - Fereastra principală a aplicației

Obținerea tuturor tabelelor din bazele de date (*Access*, *MySQL*, *Excel* și *FoxPro*) a fost făcută utilizând clasa *CatalogClass* din spațiul de nume *ADOX* și interfața *Tables*.

```
CatalogClass cat = new CatalogClass();
cat.set_ActiveConnection(string_conectare);
Tables tb = cat.Tables;
```

În cazul bazei de date *Oracle* pentru obținerea denumirilor tabelelor și a câmpurilor din fiecare tabel au fost rulate următoarele două interogări: *select table_name from use_tables* și *select column_name from user_tab_columns where table_name = '' + tabel + '' order by column_id* [FIL09].

Obținerea structurii tabelelor s-a realizat parcurgând colecția de coloane a *DataTable*-ului și accesând proprietățile: *ColumnName*, *DataType*, *Unique*, *AllowDBNull*, *AutoIncrement*, *DefaultValue*, *MaxLength*.

În situația în care numărul de tabele din baza de date este foarte mare se pot alege doar tabelele de interes în generarea *ARFF*-ului (vezi figura 2.8).

Toate tabelele sau doar tabelele de interes (cele care au fost selectate prin interfața din figura 2.8) sunt afișate sub forma de *label*-uri iar câmpurile lor sunt

afișate în *listbox*-uri (vezi figura 2.9). Relațiile dintre tabele sunt ilustrate prin linii de culori generate aleator. La producerea evenimentului *DoubleClick* al *ListBox*-ului asociat unui tabel, în *ListBox*-ul *Fields in ARFF* va fi adăugată denumirea tabelului apoi punct și apoi denumirea câmpului pe care s-a făcut dublu click-ul.

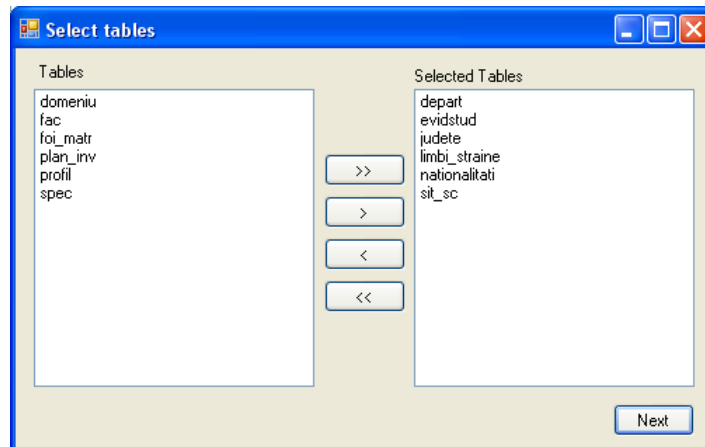


Fig. 2.8 - Selectarea tabelelor de interes

Atât label-urile cât și *listbox*-urile au fost amplasate într-un control panel și au fost create dinamic, din cod, în acord cu tabelele selectate din baza de date deschisă.

Ca și structuri ajutătoare la trasarea legăturilor dintre tabele, s-a utilizat un *ArrayList*, format dintr-o colecție de obiecte, fiecare obiect corespunzând unui tabel. Un obiect din *ArrayList* este de clasă *clsTabel* și are proprietăți care păstrează numele tabelului, coordonatele colțului stânga sus, numărul de câmpuri din tabel, denumirile câmpurilor.

Pentru determinarea legăturilor dintre tabele, s-a utilizat în cazul bazei de date Access, interfața *Tables*, clasa *CatalogClass* din spațiul de nume *ADOX*. Au fost parcurse toate tabelele selectate care făceau parte din colecția *Tables* a obiectului *CatalogClass*. În cadrul fiecărui *Tabel* au fost parcurse toate *Key*-urile, determinând tabelele de legătură cu ajutorul proprietății *Key.RelatedTable*. Pentru determinarea coloanelor de legătură s-au folosit *Key.Columns[i].Name* și *k.Columns[i].RelatedColumn*. Relațiile determinate au fost scrise într-un control *listbox*, dar a fost implementată și facilitatea de trasare grafică a acestora (vezi figura 2.9).

Trasarea liniilor care ilustrează legăturile dintre tabele s-a implementat cu ajutorul unei metode care determină indexul coloanei primită ca și parametru. Această informație servește la determinarea din dreptul cărui *Item* din cadrul *Listbox*-ului trebuie să pornească linia.

O altă metodă implementată returnează al câtelea tabel de pe interfață este tabelul primit ca și parametru. Această metodă a fost utilizată pentru a determina în ce poziții se află tabele pe care trebuie să le legăm, unul față de celălalt. Dacă indexul tabelului 1 modulo 5 este egal cu indexul tabelului 2 modulo 5 înseamnă că

tabelele sunt pe aceeași coloană (5 este numărul de tabele de pe un rând), dacă indexul tabelului 1 % 5 > indexul tabelului 2 % 5 înseamnă că tabelul 1 este în stânga tabelului 2, altfel înseamnă că tabelul 2 este în stânga tabelului 1.

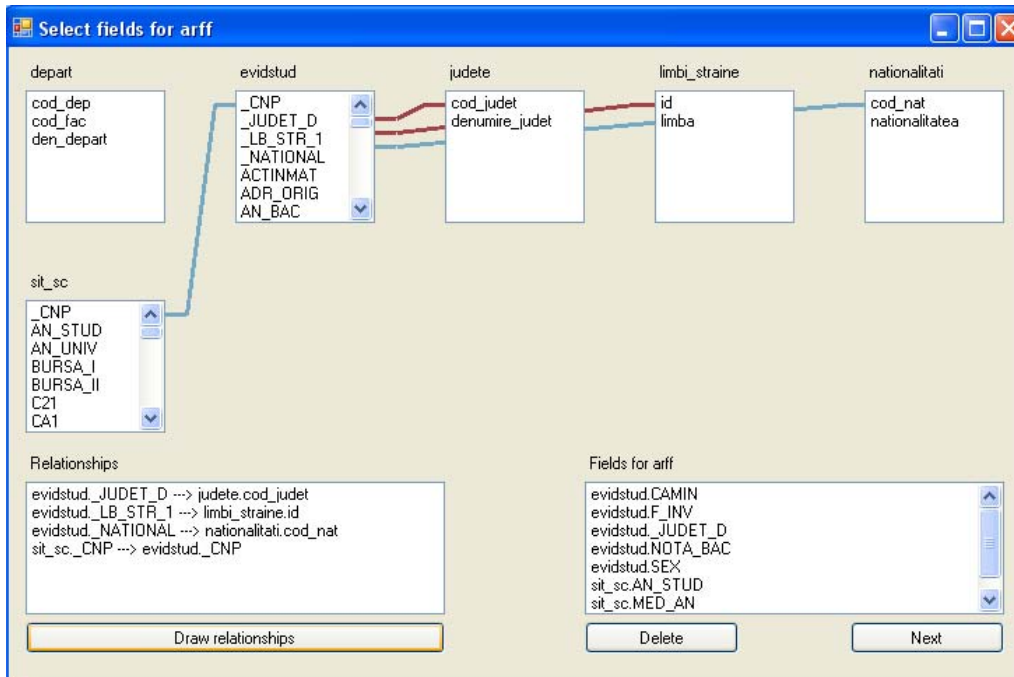


Fig.2.9 - Selectarea câmpurilor din tabele diferite pentru ARFF

Pozițiile relative ale tabelelor contează pentru a ști dacă linia trasată pleacă de pe partea stângă a listboxului corespunzător tabelului 1 sau de pe partea dreaptă a lui, respectiv dacă ajunge pe partea stângă sau pe partea dreaptă a listboxului corespunzător tabelului 2.

Stabilirea tipului pe care îl vor avea câmpurile care au fost selectate, în formatul ARFF (numerice, categoriale, string sau date) se realizează cu ajutorul unor controale de tip ComboBox (se adaugă dinamic câte un control pentru fiecare câmp selectat) (vezi figura 2.10). Conținutul combobox-urilor este următorul: în cazul unui câmp de tip șir de caractere, în combobox-ul corespunzător se va adăuga un item cu denumirea atributului și tipul string și un alt item cu denumirea atributului și categoriile determinate prin analiza conținutului câmpului. În cazul unui câmp de tip numeric (*Int32*, *Double*, *Decimal*, etc) combobox-ul corespunzător va conține un atribut numeric și un atribut categorial cu toate valorile atributului (nu vor apărea duplicate).

În cazul câmpurilor de tip boolean, ele vor deveni attribute categoriale cu valorile posibile $\{Yes, No\}$, $\{Y, N\}$, $\{True, False\}$, $\{T, F\}$, $\{Da, Nu\}$, $\{D, N\}$. Dacă pentru un atribut de tip șir de caractere sau numeric se determină că are mai puțin de 10 categorii, atunci valoarea implicită a atributului va fi cea categorială. Dacă se determină că are mai mult de 10 categorii valoarea implicită va fi string sau numeric. Câmpurile de tip DateTime pot deveni doar attribute Date.

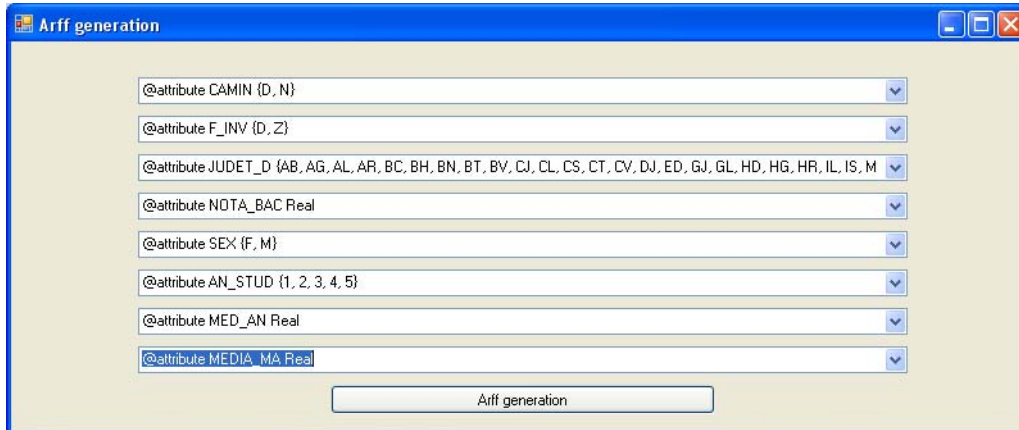


Fig. 2.10 - Stabilirea tipului atributelor

Înainte de generarea *ARFF*-ului se va analiza dacă câmpurile care au fost alese să devină atribute în *ARFF* aparțin unui singur tabel sau aparțin mai multor tabele. Dacă aceste câmpuri fac parte din tabele diferite, înseamnă că acele tabele sunt relaționate și interogarea care se generează va fi de tipul left - join. Se analizează dacă există o legătură fizică între tabelele selectate, iar dacă aceasta nu există, i se solicită utilizatorului să menționeze tabelul părinte, tabelul copil și câmpurile de legătură (vezi figura 2.11). Relația (fie introdusă de utilizator, fie determinată automat) este utilizată pentru a genera interogarea ce va prelua datele din câmpurile menționate și care este de tip left join. Un exemplu de interogare generată este următorul:

```
select evidstud.f_inv, evidstud.st_civ,evidstud.camin, evidstud.std_abs,
evidstud.sex
,evidstud.judet_d,sit_sc.an_stud,sit_sc.co1,sit_sc.r11,sit_sc.co2,sit_sc.r22,sit_sc.c2
1,sit_sc.r21 from (evidstud left join sit_sc on evidstud.cnp=sit_sc.cnp)
```

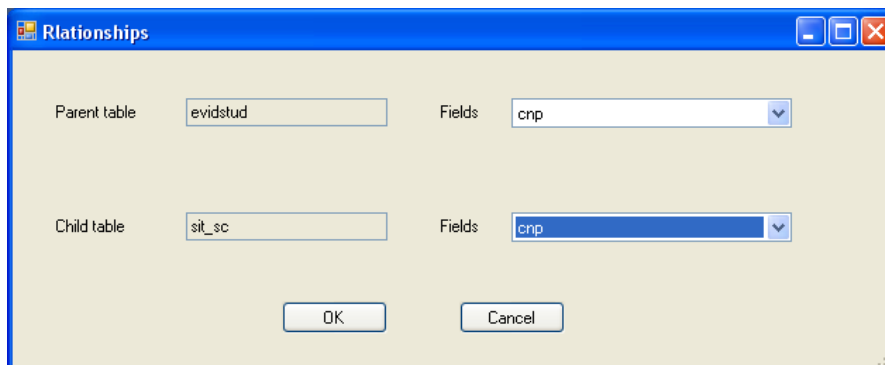


Fig. 2.11 – Specificarea legăturilor dintre câmpuri

În figura 2.12 poate fi vizualizat un exemplu de fișier *ARFF* generat.

```

gisc.arff - Notepad
File Edit Format View Help
@relation studenti_AC

@attribute CAMIN {D, N}
@attribute F_INV {D, Z}
@attribute JUDET_D {AB, AG, AR, BC, BH, Bh, BH, BN, BT, BV, C}
@attribute NOTA_BAC Real
@attribute SEX {F, M}
@attribute AN_STUD {1, 2, 3, 4, 5}
@attribute MED_AN Real
@attribute MEDIA_MA Real

@data
N, D, TM, 9.2, M, 2, 9.154, 8.577
N, D, TM, 9.2, M, 1, 8.8
N, D, TM, 9.2, M, 3, 8.727, 8.627
D, Z, ?, 0, M, 1, 4, 0
N, D, TM, 7.66, M, 3, 5.083, 7.371
N, D, TM, 7.66, M, 2, 7.385, 7.371
N, D, TM, 7.66, M, 1, 7.357, 0
N, Z, TM, 9.45, M, 1, 8.875, 8.875
N, Z, TM, 9.45, M, 2, 9.5, 9.187
N, D, TM, 0, M, 3, 7.917, 7.965
N, D, TM, 0, M, 3, 7.917, 7.981
N, D, TM, 0, M, 1, 7.643, 0
N, D, TM, 0, M, 2, 8.385, 0
N, D, TM, 6.35, M, 2, 4, 7.071
N, D, TM, 6.35, M, 1, 7.071, 7.071
N, D, TM, 6.25, M, 1, 4, 0

```

Fig. 2.12 - Fișierul ARFF generat

2.3.4. Studiu de caz: Analiza datelor civice și școlare ale studenților

Pentru a demonstra utilitatea aplicației *Arff Convertor*, autorul a realizat un studiu de caz cu ajutorul aplicației dezvoltate și cu *Weka* pe datele studenților înscriși în anul universitar 2008-2009 la Facultatea de Automatică și Calculatoare din cadrul Universității Politehnica din Timișoara. Bazele de date pe care a fost făcut studiul au fost preluate din softul de gestiune al școlarității utilizat de către secretariatul facultății, în format *dbf*. Au fost preluate două tabele: *evidstud.dbf* – care conține datele civice ale studenților și are 70 de câmpuri și *sit_sc.dbf* – care conține datele școlare ale studenților și are 43 de câmpuri. Cele două tabele sunt relaționate prin *cnp*-ul studenților, relația fiind de unul la mai mulți (unui rând din *evidstud.dbf*, îi corespund mai multe rânduri din *sit_sc.dbf*). Un student are asociate în fișierul *sit_sc.dbf* un număr de rânduri egal cu numărul de ani universitari în care a fost student. Într-o primă fază au fost selectați în fișierul *sit_sc.dbf* doar studenții din anul universitar 2008-2009 (din toți anii de studiu și de la toate specializările și ciclurile atât batchelor cât și master). În urma selecției au rezultat un număr de 2112 studenți. În a doua fază a fost folosită aplicația *Arff Convertor*, pentru a genera un fișier ARFF care conține doar atributele de interes: din cele 113 de câmpuri care sunt în tabela *evidstud* și *sit_sc* s-au ales pentru analiză un număr de 17 câmpuri, 11 din fișierul *evidstud.dbf* și 6 din fișierul *sit_sc.dbf*.

Din fișierul *evidstud.dbf* au fost alese pentru analiză următoarele: forma de învățământ starea civilă, cămin, sex, județul de domiciliu, nota bacalaureat, media multianuală în liceu, media de admitere, cetățenia, limba străină 1, limba străină 2.

Din fișierul *sit_sc.dbf* au fost alese pentru analiză următoarele câmpuri: numărul de credite la sfârșitul semestrului 1, numărul de restanțe la sfârșitul semestrului 1, numărul de credite la sfârșitul semestrului 2, numărul de restanțe la sfârșitul semestrului 2, media anuală, media multianuală. Selectarea câmpurilor de interes din cele două tabele *dbf* a fost făcută cu ajutorul aplicației *Arff Convertor*.

Studiul a continuat prin realizarea cu ajutorul lui *Weka* a unei analize statistice pe fiecare din atributele alese[WIT05].

Rezultatele obținute se referă la studenții de la AC din 2008-2009 și sunt următoarele:

- 96 % din studenții (2027 studenți) de la AC urmează învățământul la zi și doar 4 % (85 studenți) învățământul la distanță
- Majoritatea studenților 89 % (1883 studenți) sunt necăsătoriți, 1% sunt căsătoriți (22 de studenți), 1.5 % sunt divorțați (33 studenți) și 0.15 % (3 studenți) sunt văduvi
- 64 % din studenți stau în cămin (vezi figura 2.13)

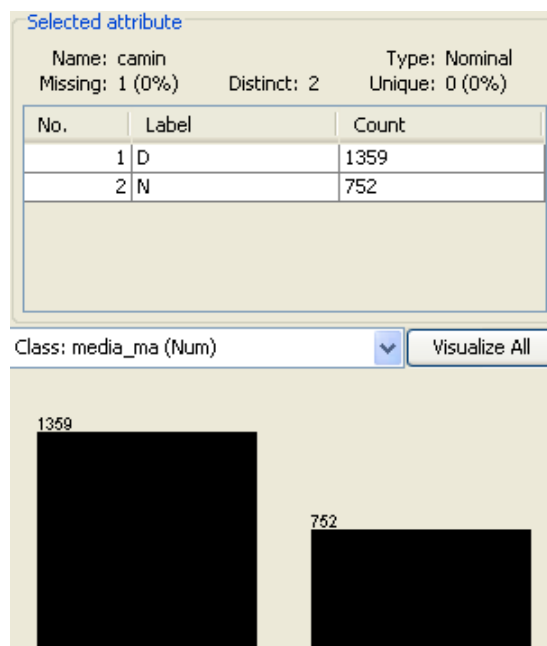


Fig. 2.13 - Numărul de studenți care stau în cămin

- 70 % din studenți sunt băieți
- Cei mai mulți studenți de la Facultatea de Automatică și Calculatoare sunt din Timiș 39 %, pe locul 2 se află Hunedoara 11.3 %, apoi Caraș- Severin 10.1 %, Arad 9.3 %, Bihor 5.9 %, Gorj 5.5 %, Mehedinți 4.7%, etc. (vezi figura 2.14)

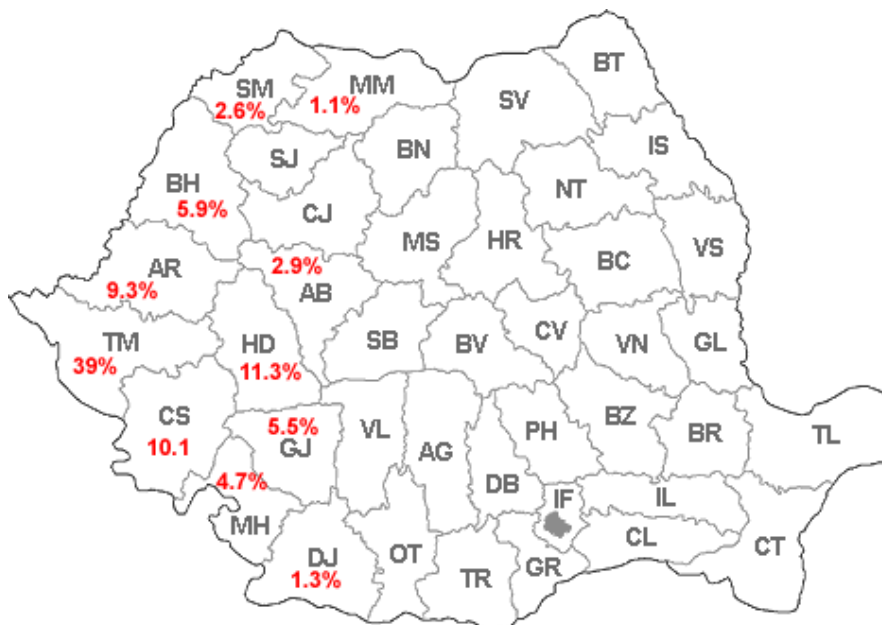


Fig. 2.14 - Domiciliul studenților de la Facultatea de Automatică și Calculatoare

- Nota la examenul de bacalaureat este de tip numeric și a fost discretizată cu ajutorul lui *Weka* în 10 intervale (vezi figura 2.15) . După cum se observă 74.5 % din studenți au nota la bacalaureat cuprinsă între 9 și 10, 15.5 % din studenți au nota între 8 și 9, 2.5 % din studenți au nota între 7 și 8 și doar 0.7 % din studenți au nota între 6 și 7. În schimb pentru 6.6 % din studenți nu a fost introdusă nota la bacalaureat în baza de date.

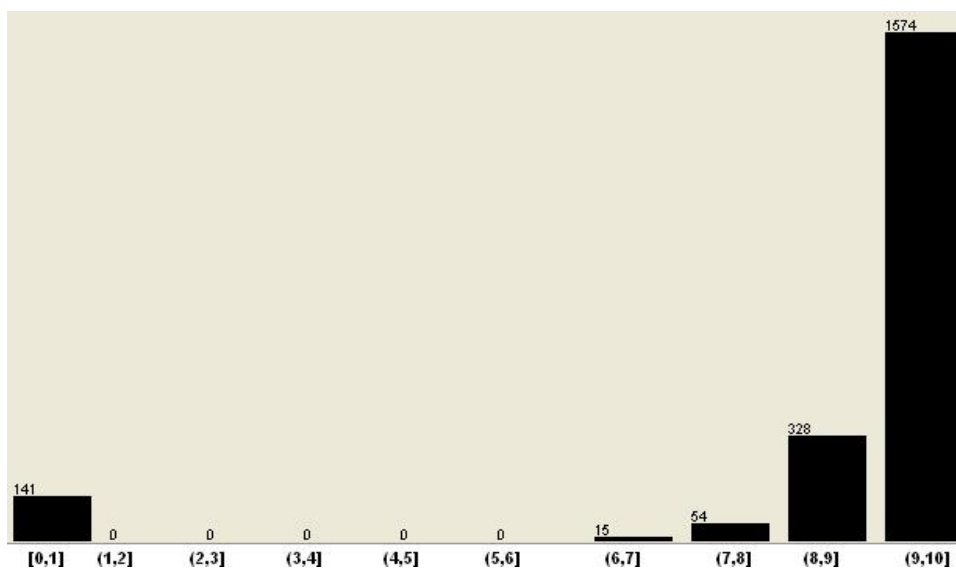


Fig. 2.15 - Nota la bacalaureat

- Media de admitere este destul de ridicată, dar scăzută față de cea de la bacalaureat: 40 % dintre studenți au intrat cu medii de admitere între 9 și 10, 30% cu medii între 8 și 9, 13.8 % cu medii între 7 și 8, 8.2 % cu medii între 6 și 7 și 2.6 % cu medii între 5 și 6. Pentru 5.6 % dintre studenți nu a fost completat câmpul media de admitere în baza de date

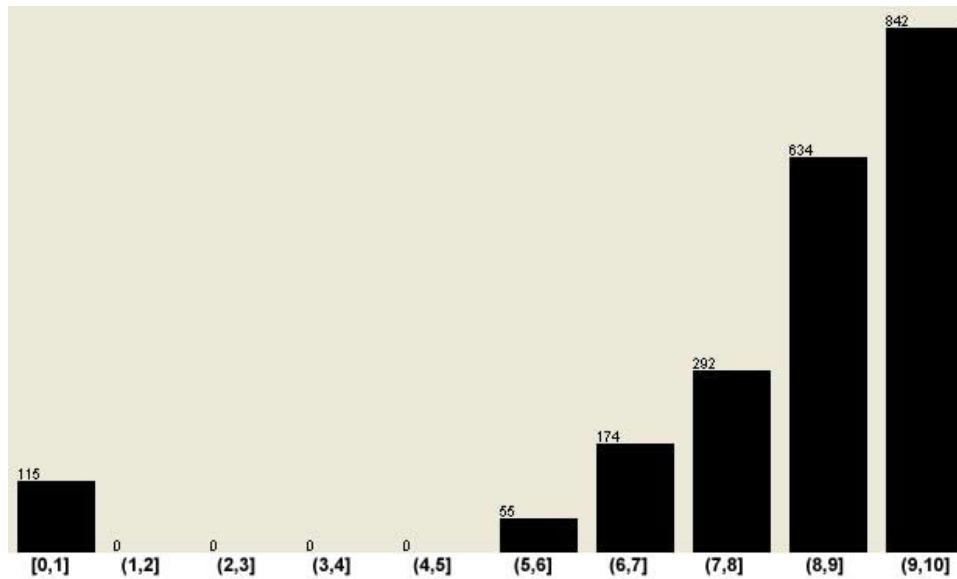


Fig. 2.16 - Media de admitere

- Media multianuală în liceu: 42.2 % de studenți au avut media multianuală cuprinsă între 9 și 10, 21.3 % au avut media între 8 și 9, 8.71 % au avut media multianuală între 7 și 8, 2.51 % au avut media între 6 și 7, 1 % au avut media între 5 și 6 și pentru 1 % a fost introdus greșit că au avut media multianuală între 4 și 5.
- Din punct de vedere al cetățeniei 99.2 % din studenți au cetățenie română, dar există și 6 studenți moldoveni, 7 slovaci, un albanez, un sârb
- Prima limbă studiată în liceu de studenți este engleza pentru 56.9 % studenți, franceza pentru 13.4 % studenți, germana pentru 4.1 % și rusa pentru 0.2 % studenți
- Limba străină 2, este franceză pentru 41 % studenți, engleză pentru 15.3 %, germană pentru 14.4 %
- Numărul de credite la sfârșitul semestrului 1 este cuprins între 0 și 32 (vezi figura 2.17). În medie studenții au strâns 25.4 credite la sfârșitul semestrului 1. Valoarea a fost discretizată, împărțită în 2 intervale și analizată. (89.9 % din studenți au obținut peste 16 credite, restul au obținut sub 16 credite în primul semestru și sunt în pericol de exmatriculare)

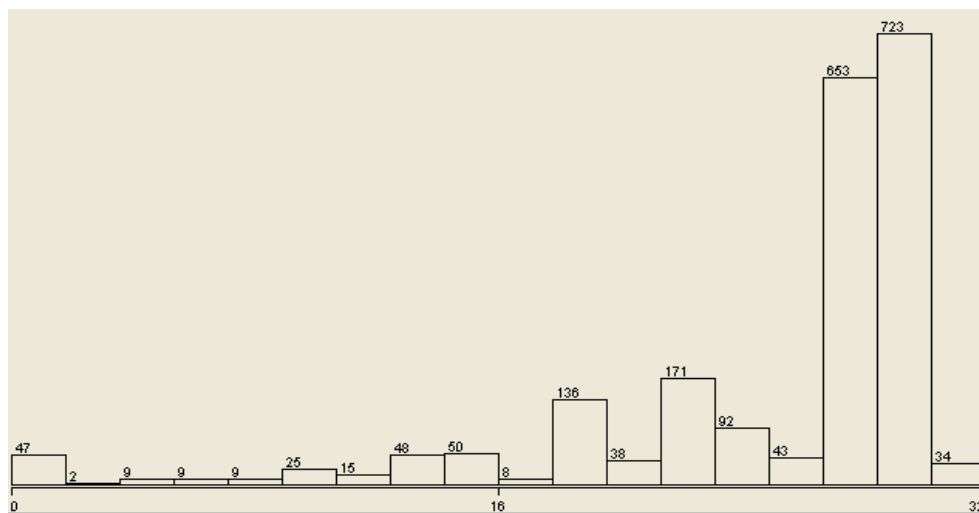


Fig. 2.17 - Numărul de credite la sfârșitul semestrului 1

- Numărul de credite la sfârșitul semestrului 2 a fost introdus doar pentru 946 de studenți. Dintre aceștia 82.1 % au obținut peste 15 credite
- Numărul de restanțe la sfârșitul semestrului 1 este 0 pentru 67.1 % studenți. Numărul de restanțe la sfârșitul semestrului 2 este 0 pentru 71.26 % studenți dar nu este un număr credibil pentru că aici nu au fost completate date despre creditele restante și numărul de restanțe pentru mulți studenți

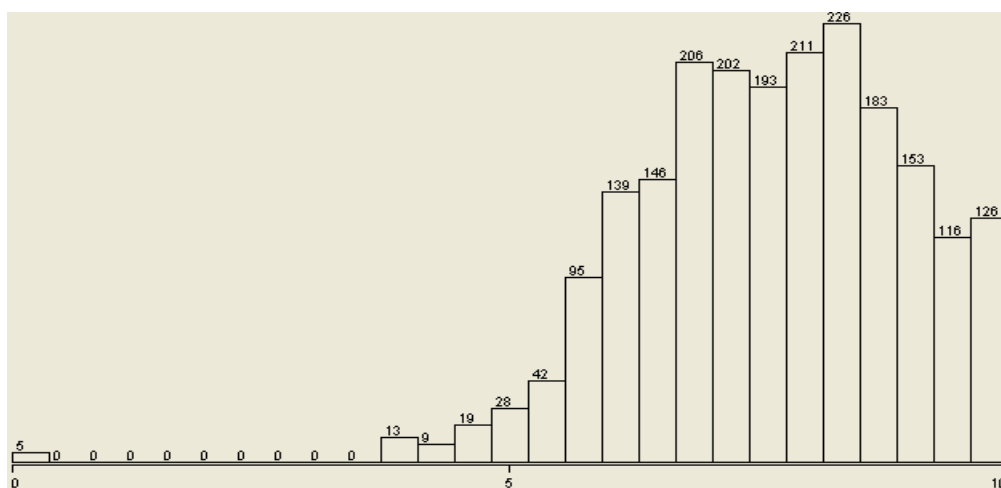


Fig. 2.18 - Media anuală

- Media anuală are o distribuție gaussiană, aspect care poate fi vizualizat în figura 2.18. Valoarea medie a mediei anuale este de 7.6 cu o deviație standard de 1.3
- Media multianuală are și ea o distribuție gaussiană, valoarea medie este de 7.5 cu o deviație standard de 1.5.

2.4. Concluzii

Preprocesarea datelor este etapă care consumă cel mai mult timp în studiile de data mining [ZHA05].

Algoritmii de data mining extrag cunoștințe din date. Succesul extragerii acestor cunoștințe depinde de calitatea datelor pe care operează acești algoritmi. Dacă datele sunt nepotrivite sau conțin informație nerelevantă, atunci algoritmii pot produce rezultate de o calitate mai slabă sau pot să nu descopere nimic util în date. De aceea etapa de preprocesare a datelor este o etapă deosebit de importantă.

Când setul de date este prea mare, este posibil să nu poată fi folosiți anumiți algoritmi de data mining. Pentru a obține mai puține date se poate reduce numărul instanțelor (prin agregare sau eșantionare) respectiv numărul atributelor (prin tehnici de reducere a dimensionalității, selectare de submulțimi de caracteristici, creare de caracteristici, etc).

Datele inițiale pot să fie nepotrivite cu anumiți algoritmi, în acest sens se impune o transformare a acestora. Transformarea datelor se poate realiza prin discretizare, operație prin care attributele continue sunt transformate în attribute categoriale, prin binarizare, prin aplicarea unor funcții matematice simple respectiv prin normalizare.

Weka este un program de data mining foarte puternic, el oferind foarte multe facilități de preprocesare (pentru discretizare, reducerea zgomotului, selectarea de submulțimi de caracteristici, etc) dar și numeroși algoritmi de clasificare, regresie, clustering, reguli de asociere, etc. Înainte de preprocesarea datelor cu *Weka* prin tehnici precum cele menționate este necesară parcurgerea altor etape ale preprocesării care constau în curățarea datelor prin eliminarea duplicatelor și a datelor expirate, corectarea sau eliminarea datelor introduse greșit, dar și conversia datelor în formatul de intrare cerut de *Weka*. Deși *Weka* permite preluarea datelor din baze de date cu ajutorul unor tehnologii, utilizarea acestora cere utilizatorului să facă configurări în Windows, în *Weka* și să cunoască limbajul SQL pentru a prelua datele. Aceste cerințe restrâng aria utilizatorilor de *Weka* la specialiști în domeniul calculatoarelor, deși *Weka* ar putea fi utilizat de către oricine dorește să analizeze date cărora le cunoaște înțelesul. Ca o soluție la această problemă autorul a conceput și dezvoltat o aplicație care se poate lega foarte ușor la baze de date și poate genera un fișier de intrare pentru *Weka*, în format *ARFF* cu datele din una sau mai multe tabele relaționate, fără a fi necesare cunoștințe de SQL și fără să fie necesară realizarea de configurări în sistemul de operare sau în *Weka*. Aplicația dezvoltată se numește *Arff Convertor* și permite conectarea la baze de date *Oracle*, *MySQL*, *Microsoft Access*, *Microsoft Excel*, *Microsoft Visual FoxPro*. Printre facilitățile aplicației se numără următoarele: selectarea câmpurilor de interes din unul sau mai multe tabele relaționate (fără a fi nevoie de cunoașterea limbajului SQL, prin dubluclick pe numele câmpului), vizualizarea grafică prin intermediul unor interfețe dinamice a structurii, relațiilor dintre tabele și a conținutului acestor baze de date, precum și generarea de fișiere de intrare pentru *Weka* (fișiere *ARFF*) care conțin datele din unul sau mai multe tabele ale bazelor de date menționate. Aplicația determină attributele numerice, string și date și propune attribute categoriale. Modul de vizualizare al structurii și conținutului tabelelor, precum și al relațiilor dintre ele,

facilitează înțelegerea datelor, lucru care este foarte important pentru procesul de data mining [BER09].

Pentru validarea aplicației *Arff Convertor*, autorul a realizat cu aplicația dezvoltată și cu *Weka* un studiu de caz pe o bază de date Microsoft Visual Fox Pro, cu datele a 2100 de studenți, care a fost preluată din aplicația de gestiune a școlarității care funcționează în cadrul Universității "Politehnica" din Timișoara, facultatea Automatică și Calculatoare. Studiul a constatat în realizarea unei analize statistice a datelor studenților. O parte din concluziile desprinse în urma analizei sunt:

- din Timiș provin doar 39% dintre studenții de la AC restul sunt din județele învecinate
- 64 % din studenții de la AC stau în cămin
- 70 % sunt băieți
- 74.5 % au media la bacalaureat între 9 și 10
- 40 % dintre studenți au media de admitere între 9 și 10
- În medie studenții obțin 25.4 credite la sfârșitul primului semestru

Studiul a dovedit că unealta propusă funcționează bine și este utilă, simplificând mult legătura dintre bazele de date sursă și programul *Weka*, foarte utilizat în data mining.

Contribuțiile care au fost aduse prin intermediul acestui capitol sunt:

- Realizarea unei sinteze ample și riguroase cu privire la principalele tehnici de preprocesare a datelor, bazată pe studiul a 30 de referințe
- Îmbunătățirea tehnologiei de preprocesare a datelor prin concepția și dezvoltarea unui instrument software permite conversia datelor din 5 formate de baze de date în formatul de date folosit de *Weka* (*ARFF*). Instrumentul propus dispune de interfețe dinamice, în acord cu datele care se prelucrează, iar utilizarea lui este simplă și nu necesită realizarea de configurări în sistemul de operare, în *Weka* sau cunoștințe de *SQL*
- Realizarea unei analize a datelor studenților de la facultatea de Automatică și Calculatoare din anul universitar 2008-2009 cu ajutorul lui *Arff Convertor* și *Weka*, analiză care oferă o imagine de ansamblu asupra provenienței studenților, pregătirii lor, etc.

3. Clasificarea datelor cu ajutorul algoritmului Ant Colony Optimization

3.1. Noțiuni introductive

Algoritmul *Ant Colony Optimization* face parte din rândul algoritmilor bioinspirați, el este utilizat la optimizarea funcțiilor din anul 1996, a fost propus de Dorigo et al [DOR96], iar utilizarea lui pentru descoperirea regulilor de clasificare a fost propusă de Parpinelli et al în anul 2001 [PAR01].

Autorii au demonstrat că acest algoritm poate descoperi în date reguli a căror acuratețe este comparabilă cu cea a unor algoritmi considerați clasici, chiar și fără niciun fel de încercare de optimizare a parametrilor algoritmului. Din Anul 2001 și până în prezent au fost făcute multe studii cu acest algoritm și au fost aduse multe îmbunătățiri în ceea ce privește calculul funcției euristice, modul de actualizare al feromonului, modul de selectare a termenilor pentru reguli, etc. Cu toate acestea, la fel ca în cazul tuturor algoritmilor de clasificare, nu a fost găsită o variantă care să fie mai bună decât altele pe orice set de date. Probleme precum alegerea parametrilor pentru algoritm în așa fel încât să se obțină cele mai bune rezultate posibile precum și îmbunătățirea procesărilor algoritmului sunt deschise în continuare.

Algoritmul propus în [PAR01] și [PAR02] a fost inițial implementat în limbajul C și testat pe diverse seturi de date. Ulterior algoritmul a fost implementat în Java și distribuit sub forma unui sistem Open Source numit *Ant Miner*. Ultima versiune este versiunea 1.2.1, dezvoltată în anul 2007 de către Fernando Meyer în colaborare cu Parpinelli [PAR02].

În cadrul acestui capitol au fost tratate următoarele probleme:

- A fost realizată o prezentare detaliată a algoritmului, bazată atât pe documentația existentă cât și pe studiul codului algoritmului.
- Au fost studiați parametrii algoritmului rulând câte 64 de configurații pe parametri pe fiecare din cele 10 seturi de date reale, obținute de pe *UCI Machine Learning Repository*. Rezultatele obținute au fost analizate atât empiric, comparativ cât și matematic construind modele de regresie și matrice de corelare. Pe baza analizei făcute a fost emis un set de recomandări cu privire la alegerea acestor parametri
- A fost realizată o sinteză a îmbunătățirilor aduse algoritmului din anul 2001 și până în prezent
- Au fost aduse 4 îmbunătățiri algoritmului care au fost implementate în aplicația open source *Ant Miner*. Aplicația care conține algoritmul îmbunătățit a fost numită *Ant-r-Miner*. S-a demonstrat că îmbunătățirile aduse conduc la o mai bună performanță utilizând date reale.

3.2. Prezentarea algoritmului

Algoritmul *Ant Colony Optimization* este inspirat din comportamentul furnicilor. Acestea, deși sunt aproape oarbe, reușesc să găsească întotdeauna drumul cel mai scurt dintre mușuroi și sursa de hrană. Studiile au arătat că acest lucru este posibil datorită unei substanțe, numită feromon cu care furnicile marchează traseul pe care îl urmează. O furnică lasă o urmă de feromon pe unde circulă și este influențată de mirosul de feromon lăsat de alte furnici [DOR04]. Feromonul este supus, odată cu trecerea timpului, fenomenului de evaporare. Studiile au arătat că dacă pe traseul dintre mușuroi și sursa de hrană se plasează un obstacol, după un timp, toate furnicile îl vor ocoli pe drumul cel mai scurt. Explicația este următoarea: când primele furnici ajung la obstacol ele aleg să îl ocolească prin partea stângă sau prin partea dreaptă cu o probabilitate egală. La întoarcere, mirosul de feromon de pe traseul mai lung este mai slab decât cel de pe traseul mai scurt, datorită fenomenului de evaporare. Prin urmare, mai multe furnici vor urma calea mai scurtă. După un anumit număr de ture efectuate între cele două destinații, toate furnicile se vor plimba pe traseul cel mai scurt [MAR11].

Regulile descoperite sunt de forma:

IF <conditions> THEN <class>

Partea condițiilor se numește antecedentul regulii și conține o combinație logică a atributelor predictoare în forma $T_{1,3}$ și $T_{2,1}$ și ... (sau *atributul1=valoarea3* și *atributul2=valoarea1*). Clasa reprezintă consecventul regulii și este de forma *atribut_clasa=valoarea2*.

Pentru antrenarea și testarea algoritmului se folosește tehnica numită *cross validation*. Setul inițial de date se împarte într-un număr de subseturi (mulțimi) pe care utilizatorul îl dorește și îl specifică prin parametrul *Folds number*. Mulțimile obținute au un număr aproximativ egal de elemente. Algoritmul va fi rulat de un număr de ori egal cu numărul de mulțimi. La fiecare rulare una din mulțimi va fi folosită pentru test și celelalte pentru antrenare.

Parametrii de intrare ai algoritmului sunt numărul de subseturi în care se va împărți setul de date (*NS*), numărul de furnici (*NF*), numărul minim de instanțe acoperite de orice regulă (*NMIR*), numărul maxim de instanțe care pot rămâne neacoperite de reguli (*NMINR*), numărul de reguli pentru convergență (*NC*), numărul de iterații (*NI*). Înainte de rulare algoritmului acești parametri de intrare trebuie să primească valori.

Algoritmul în pseudocod:

```
- alege aleator pentru fiecare instanță din setul de date în care din cele NS subseturi face parte
ns=1; //index pentru numărul de subseturi
sumaAP =0;
while ns <= NS // NS= numărul de subseturi
    mulțimea de test= subsetul ns
    mulțimea de antrenare =subseturile diferite de subsetul ns

    lista regulilor = 0 // inițial lista regulilor este goală
```

while (numărul de instanțe din mulțimea de antrenare > NMINR)

f=1; //indexul furnicii
 c=1; //index folosit pentru testul de convergența
 - inițializează toți T_{ij} cu aceeași cantitate de feromon

repetă

- furnica Ant_f începe cu o regulă goală și construiește regula adăugând rând pe rând câte un T_{ij} la regulă
 - calculează clasa regulii construite
 - aplică elagajul asupra regulii construite (rule pruning)
 - actualizează cantitatea de feromon în toți T_{ij}
 if (regula construită este identică cu regula precedentă)
 atunci c=c+1
 altfel c=1
 end if
 f=f+1; // se trece la următoarea furnică

până când (f >= NF sau (c >= NC)

- selectează cea mai bună regulă construită de toate furnicile
 - adaugă regula găsită la lista regulilor descoperite
 mulțimea_de_antrenare = mulțimea_de_antrenare - instanțele acoperite de regula descoperită

end while

- adaugă la sfârșitul listei de reguli o regulă care nu are nici o condiție în antecedent, și care are în consecvență clasa majoritară dintre clasele mulțimii de antrenament

t=0; //index pentru instanțele din mulțimea de testare

while (t < nr_inst_mulțimea_testare)

nr=0; //index pentru lista de reguli

while (nr < nr_reguli)

if (antecedent regulii acoperă instanța(t))
 dacă consecvență regula coincide cu clasa instanței
 nr_clasificate_corect = nr_clasificate_corect + 1
 altfel
 nr_clasificate_incorect = nr_clasificate_incorect + 1
 end if

end while

t=t+1;

end while

acuratețea_predicției = nr_clasificate_corect / nr_inst_mulțimea_testare
 sumaAP = sumaAP + acuratețea_predicției;
 ns=ns+1

end while

medieAP = sumaAP / nr_mulțimi

Împărțirea setului de date în NS subseturi

Împărțirea instanțelor în subseturi (mulțimi) se face parcurgând toate instanțele și generând pentru fiecare instanță un număr aleator de la 0 și până la numărul de subseturi -1. Numărul generat reprezintă subsetul de care aparține instanța. Întrucât toate subseturile trebuie să aibă un număr aproximativ egal de elemente, inițial se calculează cât este numărul de elemente maxim pe care trebuie să-l aibă fiecare subset = numărul de elemente din setul de date împărțit la numărul de subseturi +1. Dacă restul acestei împărțiri este 0 atunci se mai scade un element din dimensiunea grupului. O anumită instanță este atribuită unei mulțimi obținută prin generarea aleatoare a numărului dacă acea mulțime mai are loc să primească instanțe (nu a depășit deja numărul de instanțe permis). Dacă l-a depășit, reia generarea aleatoare.

Lista de reguli

Inițial lista regulilor este goală. Din regulile pe care le-au descoperit maxim NF furnici se alege cea mai bună și se adaugă listei de reguli. Din mulțimea de antrenare se elimină instanțele acoperite de regula adăugată listei de reguli. Procesul continuă până când în mulțimea de antrenare rămân un număr de instanțe neacoperite de nici o regulă mai mic decât valoarea $NMINR$ stabilită de utilizator.

Construcția regulii

O furnică pornește de la o regulă fără antecedent și consecvent. Furnica construiește întâi antecedentul regulii și apoi alege valoarea consecventului. Construirea antecedentului se face adăugând câte un termen ($T_{i,j}$) la un moment dat. Un termen $T_{i,j}$ este o pereche atribut-valoare. Termenul $T_{i,j}$ are semnificația: atributul i are valoarea j . Alegerea termenului care va fi adăugat la regulă depinde de cantitatea de feromon asociată lui $T_{i,j}$ și de o funcție euristică. Fiecare termen $T_{i,j}$ are asociată o cantitate de feromon. Termenii din regula descoperită de furnică corespund traseului urmat de aceasta. Inițial toți $T_{i,j}$ din setul de date (toate perechile atribut valoare) au asociată o cantitate egală de feromon (cea din relația 3.1). După construcția fiecărei reguli termenii folosiți în regulă primesc o cantitate de feromon mai mare decât restul.

$$\tau_{ij}(t=0) = \frac{1}{\sum_{i=1}^a b_i} \quad (3.1)$$

O furnică continuă să adauge termeni la regula curentă până când, orice termen ar adăuga la regulă, ar face ca regula să acopere un număr de instanțe mai mic decât parametrul $NMIR$ specificat de utilizator sau până când au fost epuizate toate atributele disponibile.

Elagajul regulilor (rule pruning)

După ce regula a fost construită se aplică asupra ei o funcție care încearcă să o îmbunătățească. Această îmbunătățire încearcă pe de o parte să facă regula mai performantă din punct de vedere calitativ dar și mai simplă. Performanța unei reguli este calculată cu formula 3.2.

$$Q = \frac{TP}{TP + FN} \times \frac{TN}{FP + TN} \quad (3.2)$$

Unde,

- TP - este numărul de instanțe acoperite care satisfac și antecedentul și consecventul regulii
- FP - numărul de instanțe care satisfac antecedentul regulii dar nu satisfac consecventul
- FN - numărul de instanțe care nu satisfac antecedentul regulii dar satisfac consecventul
- TN - numărul de instanțe care nu satisfac nici antecedentul regulii și nici consecventul ei

Îmbunătățirea regulii se face eliminând acei termeni din regulă care o fac să fie mai puțin performantă (principiu cunoscut sub denumirea elagajul regulilor). Eliminarea termenilor se face eliminând rând pe rând toți termenii din regulă și calculând calitatea regulii după ce se elimină fiecare termen, cu ajutorul formulei 3.2. Dacă eliminarea oricărui termen nu duce la obținerea unei reguli mai bune, elagajul e încheiat. Altfel se elimină din regulă termenul care prin eliminare îmbunătățește cel mai mult regula. În urma eliminării unui termen, este posibil ca consecventul regulii să trebuiască schimbat. După eliminarea oricărui termen se recalculează consecventul regulii.

Actualizarea feromonului

După ce regula a fost construită trebuie actualizată cantitatea de feromon asociată cu toți termenii i, j existenți, indiferent dacă au făcut parte sau nu din regulă. Cantitatea de feromon a termenilor care fac parte din regulă trebuie mărită, pentru că ea corespunde traseului urmat de furnică, iar cantitatea de feromon din ceilalți termeni trebuie scăzută pentru a simula fenomenul de evaporare a feromonului pe traseele neumblate [MAT08]. Creșterea cantității de feromon corespunzătoare termenilor care fac parte din regulă se face în mod proporțional cu performanța regulii, după formula 3.3.

$$\tau_{ij}(t+1) = \tau_{ij}(t) + \tau_{ij}(t) \cdot Q \quad (3.3)$$

Scăderea cantității de feromon se face indirect, prin normalizare. Normalizarea se realizează împărțind cantitatea de feromon asociată fiecărui termen la suma cantității de feromon a tuturor termenilor. Normalizarea se face după ce feromonul asociat termenilor folosiți de regulă a fost crescut.

Analiza regulii descoperite

Se verifică dacă regula curentă este identică cu regula descoperită de furnica precedentă. Dacă regula curentă este identică cu regula descoperită de NC furnici sau numărul de furnici curent a ajuns egal cu parametrul NF setat de utilizator atunci se selectează cea mai bună regulă descoperită de furnici, se adaugă aceasta la lista de reguli și se elimină instanțele care sunt acoperite de regula din setul de antrenare. Condiția de oprire legată de parametrul NC , corespunde în realitate la situația când toate furnicile urmează aceeași cale. Procesul se repetă până când numărul de instanțe din mulțimea de antrenare devine mai mic decât parametrul $NMINR$ stabilit de utilizator. Dacă condiția de oprire legată de numărul de furnici sau cea legată de testul de convergență nu sunt îndeplinite, atunci se

trece la următoarea furnică care va descoperi și ea o regulă, ghidată fiind de traseul parcurs de furnicile precedente (cantitățile de feromon actualizate).

Regula implicită

În urma etapei de antrenare de mai sus, a fost descoperită o listă de reguli. Această listă este posibil să nu acopere toate instanțele din mulțimea de test. De aceea la sfârșitul ei se adaugă o instanță care nu are nici o condiție în antecedent, iar în consecvență are clasa majoritară din setul de antrenare. Dacă se întâmplă ca o instanță să nu poată fi clasificată cu regulile descoperite de furnici, atunci regula de final va atribui instanței clasa majoritară.

Testarea regulilor descoperite

În continuare se parcurg rând pe rând instanțele din mulțimea de test. Pentru fiecare instanță se parcurge lista de reguli, în ordinea în care au fost descoperite, căutând o regulă care să poată clasifica instanța (să o acopere prin antecedent). Odată găsită aceasta regulă se compară clasa prezisă de regulă cu clasa instanței, și dacă ele coincid, se mărește numărul instanțelor clasificate corect, iar dacă nu, se mărește numărul instanțelor clasificate greșit.

După epuizarea numărului de instanțe din mulțimea de test se calculează acuratețea predicției care este raportul dintre numărul de instanțe clasificate corect și numărul total de instanțe.

Tot procesul descris mai sus se reia de un număr de ori egal cu parametrul NS menționat de utilizator. La fiecare iterație mulțimea de test este un alt subset din cele NS subseturi iar mulțimea de antrenament este formată din cele $NS-1$ subseturi.

Se va calcula media acurateței predicției obținută la cele NS rulări.

Probabilitatea de alegere a termenilor

Probabilitatea ca un termen $T_{i,j}$ să fie adăugat la regula curentă depinde de cantitatea de feromon a termenului și de valoarea unei funcții euristice după cum se poate vedea în formula 3.4.

$$P_{ij}(t) = \frac{\tau_{ij}(t)\eta_{ij}}{\sum_i \sum_j \tau_{ij}(t)\eta_{ij}}, \forall i \in I \quad (3.4)$$

Unde,

- η_{ij} este valoarea funcției euristice și reprezintă o măsură a puterii predictive a lui $T_{i,j}$
- τ_{ij} este cantitatea de feromon disponibilă curent, la momentul t în poziția i , j care corespunde traseului urmat de furnică
- a este numărul total de atribute,
- b_i este numărul total de valori ale atributului i
- I este mulțimea atributelor care nu au fost încă folosite de furnica curentă

Cu cât este mai mare valoarea lui η_{ij} cu atât mai relevant pentru clasificare este $T_{i,j}$ și deci cu atât mai mare este probabilitatea de a-l alege. Cantitatea de

feromon este independentă de termenii care apar deja în regula curentă parțială, dar este dependentă total de calea urmată de furnicile anterioare.

Alegerea termenului i, j care să fie adăugat regulii curente se face în mod aleatoriu. Cu cât valoarea lui P_{ij} este mai mare cu atât $T_{i,j}$ are mai multe șanse să fie ales.

Funcția euristică ilustrează abilitatea pe care o are $T_{i,j}$ de a îmbunătăți acuratețea predicției [PAR01]. Valoarea acestei funcții euristice pentru un termen implică o măsură a entropiei asociată cu acel termen (vezi formula 3.5). Pentru fiecare $T_{i,j}$ entropia lui este dată de formula 3.6.

$$\eta_{ij} = \frac{\log_2(k) - Entropie(T_{ij})}{\sum_i^a \sum_j^{b_j} \log_2(k) - Entropie(T_{ij})} \quad (3.5)$$

$$Entropie(T_{ij}) = - \sum_{w=1}^k \left(\frac{freq T_{ij}^w}{|T_{ij}|} \right) * \log_2 \left(\frac{freq T_{ij}^w}{|T_{ij}|} \right) \quad (3.6)$$

Unde,

- k este numărul de clase
- $|T_{ij}|$ este numărul total de instanțe în partiția $|T_{ij}|$ (numărul de instanțele în care atributul i are valoarea j)
- $freq T_{ij}^w$ este numărul de instanțe din partiția $|T_{ij}|$ care au clasa w

Entropia redă dezordinea informațională. Cu cât este mai mare valoarea entropiei cu atât mai uniform distribuite sunt clasele și cu atât mai puțin puternică va fi puterea de predicție a lui $T_{i,j}$ [PAR02]. În mod evident se urmărește să se adauge termeni care dau o mare putere predictivă regulii curente. Cu cât este mai mare valoarea lui $Entropie(T_{ij})$ cu atât este mai mică probabilitatea de a adăuga acest termen regulii curente.

Entropia termenilor nu se schimbă pe parcursul rulării și de aceea se calculează ca un pas de preprocesare.

3.3. Analiza parametrilor algoritmului

3.3.1. Noțiuni generale

Înainte de utilizarea algoritmului *Ant Colony Optimization* este necesară configurarea parametrilor algoritmului care sunt: numărul de furnici (NF), numărul minim de instanțe acoperite de orice regulă ($NMIR$), numărul maxim de instanțe care pot rămâne neacoperite de reguli ($NMINR$), numărul de reguli pentru convergență (NC), numărul de iterații (NI).

Un algoritm de clasificare trebuie să ruleze rapid, să descopere reguli cât mai simple și cel mai important, trebuie să aibă o acuratețe a predicției cât mai apropiată de 100 %.

Pentru a determina în ce mod trebuie aleși parametrii algoritmului astfel încât obiectivele de mai sus să fie atinse, algoritmul a fost testat pe 10 dataset-uri din *UCI Machine Learning Repository*, pentru diverse valori ale parametrilor.

Informații detaliate cu privire la dataset-urile utilizate sunt prezentate în subcapitolul 1.5. Atributele continue din dataset-urile *Ljubljana breast cancer*, *Diabetes*, *Cleveland heart disease*, *Iris*, *Dermatology* au fost discretizate cu ajutorul lui *Weka*.

Din 2001 până în prezent s-au realizat în literatura de specialitate diverse studii asupra utilizării algoritmului *Ant Colony Optimization* la clasificarea datelor. În aceste studii, autorii au utilizat diferite seturi de date obținute de pe *UCI Machine Learning Repository* și parametrii algoritmului au fost aleși de multe ori fără vreo încercare de optimizare [PAR01][PAR02][SMA06]. Problema optimizării parametrilor, respectiv a analizei sensibilității algoritmului la parametrii săi, a rămas o direcție de cercetare în multe studii [PAR01][PAR02][LIU02a][LIU03a]. Tabelul 3.1 prezintă valorile alese pentru parametrii algoritmului în diverse studii.

Autorii articolelor din tabelul 3.1 au demonstrat că rezultatele care se obțin cu algoritmul *Ant Colony Optimization*, respectiv cu diverse variante ale acestuia sunt competitive cu cele obținute cu alți algoritmi chiar și fără a încerca să optimizeze acești parametri. Analiza parametrilor algoritmului și alegerea lor astfel încât să se obțină rezultate optime este o problemă de cercetare deschisă în continuare.

Articolul	NF	NMIR	NMINR	NC
[PAR01]	3000	10	10	10
[PAR02]	1500			
[SMA06]	3000	5	10	5
[SAL11]	1500	5	10	10
[WAN04]	200	8	15	8
[NEJ08]	7000	5	10	10

Tabel 3.1 – parametrii aleși pentru algoritm în diverse studii

După cum se vede în tabelul de mai sus, valorile alese pentru *NF* variază de la 200 la 7000, valorile pentru parametrii *NMIR* și *NC* de la 5 la 10. Se pun probleme precum: care sunt avantajele alegerii unui număr mare de furnici? Se obțin reguli cu o acuratețe mai bună? Se obțin reguli mai puține, mai ușor de înțeles? Timpul de execuție este mai mic? Aceleași probleme se pun și pentru restul parametrilor și lucrarea de față încearcă să găsească răspuns la aceste întrebări.

Pentru ca seturile de date alese să poată fi analizate cu *Ant Miner* versiunea 1.2.1 (software Open-Source care implementează algoritmul *Ant Colony Optimization* în problema descoperirii regulilor de clasificare – vezi 3.1) a fost necesară o etapă de preprocesare care a constat în discreditarea atributelor continue, deoarece *Ant Miner* nu funcționează decât pentru attribute categoriale. Discretizarea datelor s-a făcut cu ajutorul lui *Weka*.

Testele au fost făcute rulând aplicația în paralel pe un număr de 30 de calculatoare dintr-un laborator al Universității "Politehnica" din Timișoara. Toate

calculatoarele au configurație identică Intel Core 2 Quad CPU, 2,4 GHz, 2,4 GHz , 4 GB RAM.

Acuratețea predicției s-a obținut în felul următor:

- fiecare set de date a fost împărțit în mod aleatoriu în 10 subseturi (mulțimi) care conțin un număr de instanțe aproximativ egal.
- algoritmul a fost rulat de 10 ori pe cele 10 mulțimi, la fiecare rulare o mulțime din cele 10 a fost folosită pentru testare și celelalte 9 mulțimi pentru învățare.
- valorile medii ale acurateței predicției sunt cele prezentate în continuare

Au fost rulate în total 64 de configurații de parametri pe fiecare din cele 10 seturi de date alese.

Analiza parametrilor algoritmului s-a realizat în două moduri: empiric și matematic. Analiza empirică a constat în compararea rezultatelor obținute cu 31 de configurații de parametri, s-a analizat influența pe care o are modificarea fiecărui parametru de intrare asupra acurateței predicției și asupra timpului de execuție a algoritmului.

Analiza matematică a încercat descoperirea unor relații matematice între parametrii de intrare și acuratețea predicției (s-au construit modele matematice de regresie) și s-au construit matrice de corelare, s-au analizat relațiile dintre parametrii de intrare și cel de ieșire cu ajutorul coeficienților din modelele de regresie și din matricele de corelare. Pentru a realiza analiza matematică s-au construit noi seturi de date cu valorile alese pentru parametrii de intrare și cu acuratețea predicției obținută rulând algoritmul cu parametrii aleși. Pentru a avea un volum de date mai mare pentru analiză s-au utilizat toate rezultatele obținute cu cele 64 de configurații de parametri. Utilitarul cu care au fost construite modelele de regresie și matricele de corelație este *DataFit* [www7].

3.3.2. Analiza empirică a parametrilor

3.3.2.1. Analiza parametrului Număr de furnici (NF)

Pentru a analiza efectul pe care îl are alegerea numărului de furnici asupra acurateței predicției s-au realizat teste cu 13 numere de furnici care iau valori cuprinse între 5 și 1500, în vreme ce ceilalți parametri au fost păstrați la valori fixe (vezi tabelul 3.2).

Parametru	Varianta 1
NF	5, 10, 15, 20, 25, 30, 40, 50, 100, 200, 500, 1000, 1500
NMIR	5
NMINR	10
NC	10
NI	100

Tabel 3.2 – Valorile parametrilor algoritmului

S-a determinat pentru fiecare set de date din tabelul 3.1 care este cea mai bună valoare pentru acuratețea predicției și cu ce număr de furnici s-a obținut cea mai bună valoare.

Rezultatele obținute în ceea ce privește acuratețea predicției pentru seturile de date alese și parametrii din tabelul 3.2 sunt prezentate în tabelul 3.3. Valorile maxime au fost subliniate cu **bold**, *italic* și underline, iar cele minime doar cu **bold**.

Nr. furnici	5	10	15	20	25	30	40	50	100	200	500	1000	1500
Datasetul	5	10	15	20	25	30	40	50	100	200	500	1000	1500
Car	84.7	84.3	85	85.4	86.11	86.1	84.2	85.3	85.9	86.9	85.36	85.07	86.4
Zoo	92.47	87.62	91.23	90.95	93.66	92.53	92.78	90.69	91.73	92.05	91.44	83.21	91.63
Ljubljana breast cancer	74.79	74.16	72.74	73.24	75.13	73.53	74.04	76.58	71.97	72.37	76.75	72.43	75.93
Diabetes	74.85	73.69	73.56	74.32	73.7	71.74	73.3	72.39	71.61	73.42	72.4	75.37	72
Cleveland heart disease	78.99	78.24	81.24	76.85	75.15	73.19	75.83	74.41	77.54	78.08	77.96	77.41	76.29
Iris	94	95.3	94.6	95.3	96	95.3	94	96	95.33	94	94.67	96	94.67
Mushrooms	98	97.97	96.15	96.24	97.27	96.17	97.91	95.65	94.41	93.88	94.56		
Nursurey	86.38	86.45	86.52	86.49	86.52	86.47	86.39	86.51	86.48	86.45	86.5		
Tic	72.13	70.35	72.12	71.39	70.67	70.66	70.77	71.09	69.32	71.93	71.2	71.17	70.16
Dermatology	94.29	94.56	95.59	94.56	94.53	94.8	95.1	94.55	94.8	95.1	95.58	94.81	95

Tabel 3.3 – acuratețea predicției pentru varianta 1 de parametri

După cum se observă valorile predicției pentru seturile de date *Mushrooms* și *Nursurey* în cazul a 1000 respectiv 1500 de furnici au rămas necompletate, datorită faptului că la trei rulări consecutive în care aplicația a fost lăsată să meargă vreme de 4 ore la fiecare rulare, aplicația nu a putut verifica acuratețea predicției pentru nici măcar o mulțime din cele 10.

Tabelul 3.4 arată valorile maxime și minime pentru acuratețea predicției, diferența dintre valoarea minimă și cea maximă, precum și numărul de furnici pentru care s-au obținut aceste valori.

Dataset	Valoare maximă acuratețe	Nr furnici acuratețe maximă	Valoare minimă acuratețe	Nr furnici acuratețe minimă	Diferența între acuratețea maximă și cea minimă
Car	86.9	200	84.2	40	2.7
Zoo	93.66	25	83.21	1000	10.45
Ljubljana breast cancer	76.75	500	71.97	100	4.78
Diabetes	75.37	1000	71.61	100	3.76
Cleveland heart disease	81.24	15	73.19	30	8.05
Iris	96	50,1000	94	5,40,200	2
Mushrooms	98	5	93.88	200	4.12
Nursurey	86.52	15,25	86.38	5	0.14
Tic	72.13	5	69.32	100	2.81
Dermatology	95.59	15	94.29	5	1.3

Tabel 3.4 – valorile maxime și minime ale acurateții predicției și numărul de furnici cu care s-au obținut acestea, pentru varianta 1 de parametri

După cum se observă din tabelul de mai sus, alegerea numărului de furnici poate influența acuratețea predicției între 0.14 (cazul setului de date *Nursurey*) și

10.45 procente (cazul setului de date Zoo). Tabelul 3.5 prezintă numărul de furnici și pentru fiecare furnică, numărul de seturi de date în care ea a produs cea mai bună acuratețe a predicției și cea mai slabă acuratețe a predicției.

Nr furnici	5	10	15	20	25	30	40	50	100	200	500	1000	1500
Nr dataseturi acuratețe maximă	2	0	3	0	2	0	0	1	0	1	1	2	0
Nr dataseturi acuratețe minimă	3	0	0	0	0	1	2	0	3	2	0	1	0

Tabel 3.5 – numerele de furnici și la câte valori maxime și minime au condus în ceea ce privește acuratețea predicției în varianta 1 de parametri

Analizând tabelul 3.5 putem afirma că numărul de furnici care are cea mai mare șansă să conducă la cea mai mare acuratețe a predicției este 15, întrucât cu acest număr de furnici s-a obținut de 3 ori valoarea maximă a predicției și niciodată valoarea minimă. Următorul număr de furnici care a produs rezultate bune este 25 de 2 ori maxim și niciodată minim. Cel mai slab rezultat se poate spune că l-a produs numărul de 100 de furnici, care a dus la 3 valori minime și niciun maxim, urmat de numărul 40 de furnici, care a produs 2 valori minime și niciun maxim, la egalitate cu numărul 1500 care deși nu a produs valori minime și maxime nu a putut fi rulat pe seturi de date voluminoase precum *mushrooms* și *nursurey*.

Evoluția acurateții predicției la creșterea numărului de furnici se observă și în figura de mai jos. Pe axa x este reprezentat vectorul de furnici cu $v[0]=5$, $v[1]=10, \dots, v[12]=1500$ furnici.

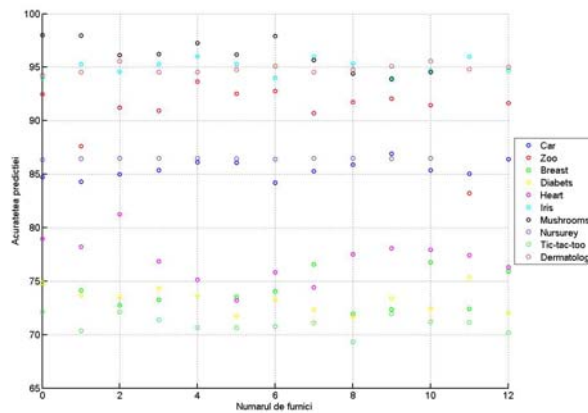


Fig. 3. 1 – Acuratețea predicției pentru varianta 1 de parametri

În ceea ce privește timpul de execuție în care s-au obținut rezultatele din tabelul 3.4, acesta este exprimat în secunde în tabelul 3.6.

Nr. furnici	5	10	15	20	25	30	40	50	100	200	500	1000	1500
Datasetul	7	12	19	28	32	41	54	72	144	274	717	1414	2011
Car	9	19	28	37	47	58	79	93	191	381	977	1991	3074
Zoo													
Ljubljana breast cancer	2	4	6	9	10	14	16	22	40	81	227	420	633

Diabetes	3	7	8	11	14	16	22	20	37	86	277	414	779
Cleveland heart disease	1	3	4	5	6	7	11	13	24	51	130	263	396
Iris	0	0	0	0	1	1	1	2	4	8	24	47	64
Mushrooms	457	885	1297	1637	2320	3033	4358	5062	9485	18572	42351	∞	∞
Nursurey	72	127	186	247	309	376	513	644	1267	2569	6662	∞	∞
Tic	8	11	23	29	42	51	61	78	158	229	734	1405	2112
Dermatology	74	119	169	203	230	266	317	428	768	1320	3320	6312	9679

Tabel 3.6 – timpul de execuție în varianta 1 de parametri

Analizând timpii de execuție putem afirma că întotdeauna un număr mic de furnici este de preferat față de un număr mare de furnici, întrucât timpul de execuție crește direct proporțional cu numărul de furnici. În cadrul setului de date *Mushrooms* timpul de execuție ajunge la 11 ore pentru 500 de furnici. Figura de mai jos arată evoluția timpului de execuție la creșterea numărului de furnici:

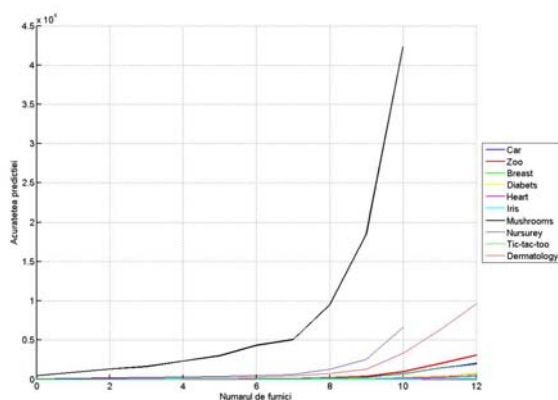


Fig 3.2 – Timpul de execuție pentru varianta 1 de parametri

3.3.2.2. Analiza parametrului Numărul minim de instanțe pe care orice regulă trebuie să le acopere (NMIR)

Pentru a analiza influența parametrului *NMIR* asupra funcționării algoritmului s-a considerat numărul de 15 furnici care a produs cele mai bune rezultate în testele anterioare. Pentru *NMIR* s-au luat valorile 5, 10 și 15. Ceilalți parametri au fost păstrați la valorile din tabelul 3.7.

Parametru	Varianta 2	Varianta 3	Varianta 4
NF	15	15	15
NMIR	5	10	15
NMINR	10	10	10
NC	10	10	10
NI	100	100	100

Tabel 3.7 – Valorile parametrilor algoritmului

Rezultatele obținute în ceea ce privește timpul de execuție și acuratețea predicției pentru cele 3 variante sunt prezentate în tabelul 3.8. Cu *AV_n* s-a notat

acuratețea predicției obținută cu varianta n de parametri, iar cu $TEVn$ s-a notat timpul de execuție obținut cu varianta n de parametri. În cazul variantei 4 de parametri, au fost 4 seturi de date pentru care nu s-a obținut un rezultat (datorită intrării aplicației în buclă infinită). În cazul în care nu s-a putut ajunge la un rezultat, au mai fost încercate încă 2 rulări cu același set de parametri pe același set de date. Dacă în urma acestor 3 încercări s-a ajuns la un rezultat, atunci acuratețea predicției a fost completată în tabelul 3.8 iar dacă nu s-a ajuns la rezultat atunci numărul maxim de mulțimi, care a putut fi testat la una din cele 3 rulări, a fost completat în tabel în formatul cifră de la 0 la 9 și apoi litera f (ex. 5f= s-au încercat 3 rulări, niciuna nu a putut fi finalizată din cauza faptului că aplicația a intrat în buclă infinită, performanța maximă obținută este că la una din aceste 3 rulări s-au putut verifica 5 folds).

Setul de date	AV2	AV3	AV4	TEV2	TEV3	TEV4
Car	85	83.68	3f	19	17	∞
Zoo	91.23	84.91	80.51	28	28	26
Ljubljana breast cancer	72.74	74.89	3f	6	6	∞
Diabetes	73.56	73.04	72.5	8	9	7
Cleveland heart disease	81.24	79.22	8f	4	4	∞
Iris	94.6	95.33	94	0	0	0
Mushrooms	96.15	96.95	97.06	1297	1212	1427
Nursurey	86.52	86.51	86.5	186	185	188
Tic	72.12	72.44	9f	23	28	∞
Dermatology	95.59	93.72	94.02	169	160	151

Tabel 3.8 - acuratețea predicției și timpul de execuție pentru variantele 2, 3 și 4 de parametri

Analizând rezultatele din tabelul de mai sus observăm că varianta 2 de parametri conduce la acuratețea maximă în 6 cazuri, varianta 3 de parametri în 3 cazuri și varianta 4 într-un singur caz. Deci creșterea $NMIR$ conduce la o scădere a acurateții predicției și conduce la multe situații de blocare. Cele mai bune rezultate s-au obținut pentru varianta 2, în care $NMIR$ ia valoarea 5.

3.3.2.3. Analiza parametrului Numărul maxim de instanțe neacoperite de nici o regulă ($NMINR$)

Pentru a analiza comportamentul algoritmului, a fost considerată varianta 5 de parametri, în care au fost păstrați parametrii din varianta 4, mai puțin numărul de instanțe neacoperite de nicio regulă care în loc să primească valoarea 10 primește valoarea 20.

Parametru	Varianta 4	Varianta 5
NF	15	15
NMIR	15	15
NMINR	10	20
NC	10	10
NI	100	100

Tabel 3.9 - Valorile parametrilor algoritmului în variantele 4 și 5 de parametri

Rezultatele obținute sunt prezentate comparativ în tabelul 3.10. Un prim lucru care se observă este că creșterea numărului de instanțe neacoperite de nici o regulă face ca programul să găsească soluții în cazul seturilor de date pentru care acest lucru nu a fost posibil cu varianta 4 de parametri.

Varianta 5 de parametri este mai bună pentru că a găsit valori maxime în cazul a 6 seturi de date.

Setul de date	AV4	AV5	TV4	TV5
Car	3f	84.66	∞	16
Zoo	80.51	79.42	26	25
Ljubljana breast cancer	3f	75.53	∞	6
Diabetes	72.5	73.44	7	7
Cleveland heart disease	8f	73.92	∞	3
Iris	94	93.33	0	0
Mushrooms	97.06	97.7	1427	1338
Nursurey	86.5	86.47	188	180
Tic	9f	73.29	∞	22
Dermatology	94.02	93.74	151	153

Tabel 3.10 – comparație între valorile maxime obținute cu variantele 4 și 5 de parametri

Deși ne așteptam ca creșterea numărului de instanțe neacoperite de nici o regulă să determine o scădere a performanței, în acest caz lucrurile nu stau așa. În 6 cazuri din 10 s-a obținut o acuratețe mai bună a predicției dacă se lasă maxim 20 de instanțe neacoperite, în 4 cazuri s-a obținut o acuratețe a predicției mai bună dacă se lasă maxim 10 instanțe neacoperite. O posibilă explicație pentru acest rezultat este numărul mare de rulări nereușite, în cazul cu fiecare regulă trebuie să acopere minim 15 instanțe și pot rămâne 10 instanțe neacoperite de nicio regulă. Pentru a lămurii acest aspect s-a realizat o verificare suplimentară, cea din tabelul 3.11.

Parametru	Varianta 6	Varianta 7
NF	15	15
NMIR	5	5
NMINR	10	30
NC	10	10
NI	100	100

Tabel 3.11 – Variantele 6 și 7 de parametri

Setul de date	AV 6	AV 7	TE 6	TE 7
Car	85	84.83	19	21
Zoo	91.23	81.19	28	24
Ljubljana breast cancer	72.74	73.86	6	5
Diabetes	73.56	72.63	8	9
Cleveland heart disease	81.24	78.43	4	3
Iris	94.6	95.33	0	0
Mushrooms	96.15	97.12	1297	1583
Nursurey	86.52	86.54	186	176
Tic	72.12	70.23	23	19
Dermatology	95.59	95.07	169	148

Tabel 3.12 – acuratețea predicției și timpul de execuție pentru variantele 5 și 6 de parametri

În varianta 6 se obțin în cazul a 6 seturi de date, valori ale acurateții predicției mai mari decât în cazul variantei 7. În cazul setului de date zoo diferența este semnificativă (10 %). Deci acuratețea predicției este în general mai bună dacă configurăm numărul maxim de instanțe neacoperite de nici o regulă la o valoare cât mai mică, excepție de la acest fapt fac situațiile în care numărul de instanțe care trebuie acoperit de o regula este mare și conduce la multe situații de blocare. În acest caz este mai bine să mărim *NMINR*, vom obține rezultate mai bune.

3.3.2.4. Analiza parametrului Numărul de reguli pentru convergență (NC)

Pentru a analiza acest parametru considerăm valorile 2, 5, 10, 20. Cele 4 variante de parametri analizați sunt prezentate în tabelul 3.13.

Parametru	Varianta 8	Varianta 9	Varianta 10	Varianta 11
NF	15	15	15	15
NMIR	5	5	5	5
NMINR	10	10	10	10
NC	2	5	10	20
NI	100	100	100	100

Tabel 3.13 - variantele 7, 8, 9 și 10 de parametri

Tabelul 3.14 prezintă rezultatele obținute cu variantele de parametri 8, 9, 10 și 11.

Setul de date	AV8	AV9	AV10	AV11	TEV8	TEV9	TEV10	TEV11
Car	85.71	85.59	85	85.07	4	10	19	36
Zoo	84.57	93.11	91.23	90.55	6	14	28	58
Ljubljana breast cancer	74.89	73.81	72.74	74.79	1	3	6	13
Diabetes	73.7	72.28	73.56	74.48	1	2	8	19
Cleveland heart disease	78.92	77.03	81.24	77.04	0	1	4	8
Iris	94.67	93.33	94.6	95.33	0	0	0	1
Mushrooms	97.4	96.47	96.15	96.83	284	586	1297	2411
Nursurey	86.36	86.47	86.52	86.45	40	94	186	353
Tic	72.34	71.7	72.12	72.85	3	13	28	41
Dermatology	91.56	93.48	95.59	94.56		94	157	256

Tabel 3.14 - acuratețea predicției și timpul de execuție pentru variantele 8, 9, 10 și 11 de parametri

Analizând rezultatele din tabelul 3.14 putem afirma că pentru variantele de parametri 8, 10 și 11 s-au obținut maxime ale acurateții predicției de 3 ori, iar pentru varianta 9 s-a obținut maximum acurateții predicției o singură dată. Creșterea numărului de reguli pentru convergență determină creșterea timpului de execuție. Deci cel mai bun rezultat, obținut judecând prin prisma acurateții predicției respectiv timpului de execuție, este pentru varianta 8 de parametri.

3.3.2.5. Analiza parametrului Numărul de iterații (NI)

O regulă este descoperită fie dacă furnica curentă a descoperit aceeași regulă ca și NC-1 furnici, fie dacă furnica curentă nu a descoperit aceeași regulă ca și NC-1 furnici, dar s-a executat numărul de iterații stabilit.

Influența pe care o are acest parametru NI este analizată în continuare. S-au considerat 4 variante ale parametrilor:

Parametru	Varianta 12	Varianta 13	Varianta 14	Varianta 15
NF	15	15	15	15
NMIR	5	5	5	5
NMINR	10	10	10	10
NC	2	2	2	2
NI	10	50	100	200

Tabel 3.15 - variantele 12, 13, 14 și 15 de parametri

Setul de date	AV12	AV13	AV14	AV15	TEV12	TEV13	TEV14	TEV15
Car	85.47	84.72	85.71	85.13	3	4	4	4
Zoo	93.77	90.18	84.57	91.83	6	6	6	5
Ljubljana breast cancer	73.73	74.9	74.89	74.16	1	1	1	1
Diabetes	72.65	74.48	73.7	72.64	1	1	1	1
Cleveland heart disease	77.91	77.46	78.92	78.33	0	0	0	1
Iris	94.67	96	94.67	94	0	0	0	0
Mushrooms	97.06	97.12	97.4	96.77	238	241	284	265
Nursurey	86.42	86.06	86.36	86.15	39	37	40	36
Tic	74	73.36	72.34	71.49	4	2	3	5
Dermatology	94.02	92.39	91.56	92.82	50	47	46	51

Tabel 3.16 - acuratețea predicției și timpul de execuție pentru variantele 12, 13, 14 și 15 de parametri

În ceea ce privește acuratețea predicției, cel mai bun rezultat se obține pentru varianta 12 de parametri (10 iterații) și anume 4 valori maxime.

Analizând timpul de execuție în cazul celor 4 variante, observăm că se obțin valori asemănătoare, cu toate că numărul de iterații ia valori într-o plajă mare de la 10 la 200. O posibilă explicație este că numărul mic de reguli pentru convergență face ca valoarea acestui parametru să nu mai conteze, bucla repetitivă oprindu-se de fiecare dată datorită neîndeplinirii condiției legate de numărul de reguli pentru convergență și nu datorită depășirii numărului de iterații. Pentru a verifica această teorie s-a ales un număr mare de reguli pentru convergență (30) și valorile din tabelul 3.17:

Parametru	Varianta 16	Varianta 17	Varianta 18	Varianta 19
NF	15	15	15	15
NMIR	5	5	5	5
NMINR	10	10	10	10
NC	30	30	30	30
NI	10	50	100	200

Tabel 3.17 - variantele 16, 17, 18 și 19 de parametri

Tabelul 3.18 prezintă acuratețea predicției pentru variantele 16, 17, 18 și 19 de parametri, precum și timpii de execuție.

Setul de date	AV16	AV17	AV18	AV19	TV16	TV17	TV18	TV19
Car	85.77	85.25	86.46	84.96	14	55	56	64
Zoo	91.12	91.7	92.73	91.62	28	92	89	84
Ljubljana breast cancer	75.56	73.74	72.77	73.46	5	19	20	20
Diabetes	73.57	74.22	73.84	73.83	4	18	22	28
Cleveland heart disease	76.96	80.1	78.99	77.32	3	12	12	13
Iris	96	94	96	94	0	1	1	1
Mushrooms	96.86	96.48	95.38	96.76	1094	3825	4152	5709
Nursurey	86.42	86.47	86.46	86.52	169	528	536	541
Tic	70.12	72.75	69.29	73.39	15	59	61	54
Dermatology	94.53	95.38	95.59	94.55	97	346	347	360

Tabel 3.18 - acuratețea predicției și timpul de execuție pentru variantele 16, 17, 18 și 19 de parametri

Comparând evoluția timpilor de execuție din tabelul 3.18 cu cea din tabelul 3.16, observăm că timpul de execuție din tabelul 3.18 are o alură crescătoare iar cel din tabelul 3.16 este mai degrabă constant, ceea ce înseamnă că în cazul 3.16, valorile alese pentru numărul de instanțe nu au prea mare importanță datorită valorii mici a lui NC . Analizând timpii de execuție din tabelul 3.18 putem spune că ei sunt comparabili pentru variantele 17, 18 și 19 de parametri și sunt mult mai mici în cazul variantei 16 de parametri (excepție face setul de date *mushrooms* în care timpul de execuție este direct proporțional cu numărul de parametri). Deci un număr mic de iterații determină execuția mai rapidă a algoritmului, iar creșterea numărului de iterații peste o anumită valoare (în cazul experimentelor noastre, valoarea 50) poate să nu modifice decât puțin timpul de execuție. În ceea ce privește acuratețea predicției, cel mai bun rezultat se obține pentru varianta 18 de parametri (100 de iterații) și anume 4 valori maxime.

3.3.3. Analiza matematică a coeficienților algoritmului

În a doua parte a studiului s-a încercat descoperirea unor relații matematice între parametrii de intrare și parametrul de ieșire al algoritmului (acuratețea predicției). S-au construit modele matematice de regresie și s-au construit matrice de corelare, s-au analizat relațiile dintre parametrii de intrare și parametrul de ieșire cu ajutorul coeficienților din modelele de regresie și din matricele de corelare.

Modelele matematice de regresie construite permit estimarea ieșirii (acurateței predicției) în funcție de parametri de intrare ai algoritmului (numărul de furnici - NF , numărul minim de instanțe acoperite de orice regulă - $NMIR$, numărul maxim de instanțe care pot rămâne neacoperite de reguli - $NMINR$, numărul de reguli pentru convergență - NC , numărul de iterații - NI).

Pentru determinarea acestor modele și pentru construirea matricelor de corelare s-au construit 10 noi seturi de date. Fiecare set de date conține valorile alese pentru parametrii de intrare și acuratețea predicției obținute rulând algoritmul cu parametrii aleși. Toate seturile de date construite conțin 64 de instanțe, pentru că s-au realizat 64 de rezultate experimentale pentru fiecare set de date (s-au

analizat 64 de variante de parametri pentru fiecare set de date). Variantele de parametri analizate sunt prezentate în tabelul 3.20. Alegerea acestor variante de parametri a fost influențată de valorile implicite pe care acești parametri le au în *Ant Miner* [PAR02], de valorile care au fost alese de alți cercetători în studiile pe care le-au realizat (vezi tabelul 3.2), de rezultatele care s-au obținut pentru acuratețea predicției pe parcursul execuției (vezi studiul comparativ din subcapitolul 3.3.2) și de semnificația mărimilor. În cazul configurațiilor de parametri de intrare care au determinat algoritmul să ruleze la infinit (acestea sunt de obicei configurațiile de parametri în care *NMIR* are valoarea 15) în coloana aferentă variabilei dependente (acuratețea predicție) a fost completată valoarea 100 împărțită la numărul de valori ale atributului clasă (s-a completat probabilitatea de a alege clasa corectă fără niciun fel de informație suplimentară).

NF	NMIR	NMINR	NC	NI
5, 10, 15, 20, 25, 30, 40, 50, 100, 200, 500, 1000, 1500	5	10	10	100
5, 10, 15, 20, 25, 30, 40, 50, 100, 200, 500, 1000, 1500	10	10	10	100
5, 10, 15, 20, 25, 30, 40, 50, 100, 200, 500, 1000, 1500	15	10	10	100
5, 10, 15, 20, 25, 30, 40, 50, 100, 200, 500, 1000, 1500	15	20	10	100
15	5	30	10	100
15	5	10	2, 5, 20	100
15	5	10	2	10, 50, 100, 200
15	5	10	30	10, 50, 100, 200

Tabel 3.19 – Valorile parametrilor de intrare ai algoritmului

Instrumentul software utilizat pentru a construi modelele de regresie și pentru a determina matricele de corelare este *DataFit* [www7]. Inițial s-a încercat construirea unor modele de regresie folosind un total de 10 variabile independente pe lângă cei 5 parametri de intrare menționați, încă 5 parametri care caracterizează setul de date. Aceștia sunt: numărul de instanțe, numărul de atribute, numărul de valori ale atributului clasă, numărul de valori lipsă, numărul total de valori ale atributelor. Pe baza acestor 10 variabile independente utilitarul *DataFit* nu a putut construi modele de regresie.

Analiza modelelor de regresie

Modelele de regresie obținute utilizând ca și variabile independente cei 5 parametri de intrare ai algoritmului sunt prezentate în tabelul 3.20.

Setul de date	Modelul	R ²
Car	$AP = 0.013 * NF - 3.876 * NMIR + 2.638 * NMINR - 0.002 * NC - 0.001 * NI + 78.518$	0.55
Zoo	$AP = \exp(-0.011 * NMIR - 0.004 * NMINR + 0.001 * NC + 4.591)$	0.83
	$AP = -0.918 * NMIR + -0.293 * NMINR + 0.075 * NC - 0.003 * NI + 97.48$	0.83
Ljubljana breast cancer	$AP = 0.001 * NF - 1.744 * NMIR + 1.269 * NMINR - 0.034 * NC - 0.003 * NI + 71.674$	0.58
Diabetes	$AP = 0.001 * NF - 1.220 * NMIR + 0.758 * NMINR + 0.013 * NC - 0.001 * NI + 72.378$	0.40
Cleveland heart disease	$AP = 0.004 * NF - 2.923 * NMIR + 2.101 * NMINR - 0.012 * NC + 0.001 * NI + 72.349$	0.39
Iris	$AP = 0.024 * NMIR - 0.006 * NMINR + 0.009 * NC - 0.006 * NI + 95.68$	0.08
Mushrooms	$AP = -0.035 * NF + 0.055 * NMIR + 0.002 * NMINR - 0.039 * NC + 98.136$	0.90

	AP = $\exp(0.001 \cdot \text{NMIR} + 4.589)$	0.89
Nursurey	AP = $-0.049 \cdot \text{NF} + 0.125 \cdot \text{NMIR} - 0.038 \cdot \text{NMINR} - 0.026 \cdot \text{NC} + 0.002 \cdot \text{NI} + 88.341$	0.87
	AP = $\exp(-0.001 \cdot \text{NF} + 0.002 \cdot \text{NMIR} - 0.001 \cdot \text{NMINR} + 4.487)$	0.83
Tic	AP = $-1.579 \cdot \text{NMIR} + 0.998 \cdot \text{NMINR} - 0.042 \cdot \text{NC} + 70.804$	0.60
Detmatology	AP = $0.001 \cdot \text{NF} + 0.035 \cdot \text{NMIR} + 0.002 \cdot \text{NMINR} + 0.07 \cdot \text{NC} - 0.001 \cdot \text{NI} + 93.395$	0.29

Tabel 3.20 – Modelele de regresie determinate și performanțele lor

R^2 este coeficientul determinării multiple și se calculează cu formula 3.7.

$$R^2 = 1 - \frac{RSS}{TSS} \quad (3.7)$$

unde RSS este suma reziduală a pătratelor și se calculează cu formula 3.8, iar TSS este suma totală a pătratelor și se calculează cu formula 3.9.

$$RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (3.8)$$

$$TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad (3.9)$$

Unde, Y_i este valoarea reală a variabilei dependente (a acurateții predicției) pentru instanța i și \hat{Y}_i este valoarea estimată de modelul de regresie pentru variabila dependentă și \bar{Y} este valoarea medie a variabilei dependente.

Mărimea RSS reprezintă suma reziduală a pătratelor (*Residual Sum of Squares*) și se calculează ca și sumă a pătratelor diferențelor dintre valorile reale ale variabilei dependente și valorile estimate de modelul de regresie, cu formula 3.8. O potrivire perfectă între acestea ar conduce la o sumă reziduală a pătratelor de valoare 0.

Mărimea R^2 ia valori cuprinse între 0 și 1 și măsoară cât de bine surprinde modelul de regresie variația variabilei dependente. De exemplu dacă R^2 are valoarea 0.8, atunci 80 % din variația variabilei dependente este explicată de modelul de regresie.

Utilitarul *DataFit* a determinat pentru fiecare set de date câte 3 modele. În tabelul 3.20 au fost prezentate doar modelele de regresie care au obținut cele mai bune performanțe pentru fiecare set de date. După cum se observă unele din aceste modele surprind bine variația variabilei dependente (cele care au mărimea R^2 apropiată de valoarea 1) iar altele au o performanță slabă.

Analizând cele 10 seturi de date alese pentru analiză se observă că 4 dintre ele conțin date medicale (breast, diabetes, heart, dermatology), 3 conțin date care caracterizează elemente din natură (animale sau plante: zoo, iris, mushrooms), unul conține evaluările unor mașini (car), unul conține evaluarea aplicațiilor depuse la grădinițe și unul codifică mutările posibile la jocul X și 0.

Modelele de regresie care s-au construit pe seturile de date medicale descriu între 29 % și 58 % variația acurateții predicției, deci sunt foarte slabe. Două din modelele de regresie construite pe seturi de date care caracterizează elemente din natură, zoo și mushrooms, au obținut o performanță mai bună de 83 % respectiv 90% iar setul de date iris a obținut cea mai slabă performanță dintre toate seturile de date, de 8 %. Modelul de regresie construit pe setul de date car, are o performanță slabă de 55 %. Pe baza setului de date ce conține evaluarea aplicațiilor la grădinițe s-au construit modele de regresie bune de 83% și 89 %. Modelul de regresie construit pe baza setului de date tic surprinde 60 % din variația mării dependente.

În continuare au fost analizați coeficienții parametrilor din modelele de regresie determinate. Valorile acestora arată cât de mult influențează acești parametri variabila dependentă. Valorile acestor parametri pot fi vizualizate în tabelul 3.21. S-au luat în considerare doar modelele care au o performanță medie sau bună (pentru care parametrul R^2 este mai mare de 0.5). S-a determinat modulul pentru fiecare parametru și maximum dintre valorile absolute pentru fiecare model. Valorile maxime au fost subliniate. Toate valorile maxime în valoare absolută au fost asociate coeficienților parametrului *NMIR*, ceea ce înseamnă că acuratețea predicției depinde într-o bună măsură de modul în care se alege acest parametru. Valorilor cele mai mari în modul ale acestui parametru le corespund valori negative, ceea ce înseamnă că, de obicei, creșterea acestui parametru determină o scădere a acurateții predicției.

NF	NMIR	NMINR	NC	NI	NF	NMIR	NMINR	NC	NI
0.013	-3.876	2.638	-0.002	-0.001	0.013	3.876	2.638	0.002	0.001
0	-0.011	-0.004	0.001	0	0	0.011	0.004	0.001	0
0	-0.918	-0.293	0.075	-0.003	0	0.918	0.293	0.075	0.003
0.001	-1.744	1.269	-0.034	-0.003	0.001	1.744	1.269	0.034	0.003
-0.035	0.055	0.002	-0.039	0	0.035	0.055	0.002	0.039	0
0	0.001	0	0	0	0	0.001	0	0	0
-0.049	0.125	-0.038	-0.026	0.002	0.049	0.125	0.038	0.026	0.002
-0.001	0.002	-0.001	0	0	0.001	0.002	0.001	0	0
0	-1.579	0.998	-0.042	0	0	1.579	0.998	0.042	0

Tabel 3.21 – Coeficienții parametrilor din modelele de regresie determinate

S-a determinat maximum valorilor absolute, exceptând parametrul *NMIR* (următorul parametru ca și importanță). În 6 din cele 9 modele, valoarea maximă a corespuns parametrilor lui *NMINR*, în unul dintre ele lui *NC* și în unul dintre ele lui *NF*. Deci următorul parametru ca și importanță este *NMINR*. Creșterea lui *NMINR* poate determina creșterea sau scăderea performanței în funcție de setul de date (pentru 4 seturi de date determină creșterea performanței și pentru 4 scăderea ei).

Analiza coeficienților de corelare

Următoarea etapă a analizei a constat în analiza coeficienților din matricele de corelare aferente fiecărui set de date. Tabelul 3.22 prezintă matricea de corelare aferentă setului de date car. Au fost analizați parametrii care influențează cel mai mult parametrul de ieșire potrivit matricei de corelare. În cazul setului de date car,

se poate observa că, creșterea *NMIR* determină o scădere substanțială a performanței, iar creșterea lui *NMINR* determină o creștere a performanței.

	NF	NMIR	NMINR	NC	NI	AP
NF	1.00	0.13	0.02	-0.06	0.02	0.14
NMIR	0.13	1.00	0.41	-0.14	0.06	-0.51
NMINR	0.02	0.41	1.00	-0.06	0.03	0.22
NC	-0.06	-0.14	-0.06	1.00	-0.05	0.06
NI	0.02	0.06	0.03	-0.05	1.00	-0.03
AP	0.14	-0.51	0.22	0.06	-0.03	1.00

Tabel 3.22 – Matricea de corelare aferentă setului de date car

Tabelul 3.23 prezintă coeficienții de corelare dintre variabila dependentă și variabilele independente potrivit matricelor de corelare aferente fiecărui set de date.

Setul de date	NF	NMIR	NMINR	NC	NI
Car	0.14	-0.51	0.22	0.06	-0.03
Zoo	-0.01	-0.14	0.04	0.30	0.36
Ljubljana breast cancer	-0.06	-0.54	0.27	0.05	-0.04
Diabetes	-0.02	-0.49	0.15	0.07	-0.03
Cleveland heart disease	0.01	-0.44	0.21	0.05	-0.02
Iris	-0.07	-0.18	-0.1	0.1	-0.22
Mushrooms	-0.95	-0.11	-0.01	0.04	-0.02
Nursurey	-0.93	-0.10	-0.01	0.05	-0.02
Tic	-0.11	-0.60	0.20	0.06	-0.03
Detmatology	0.30	0.15	0.06	0.39	-0.04

Tabel 3.23 – Coeficienții de corelare dintre variabila dependentă și variabilele independente

Potrivit coeficienților de corelare, parametrul *NMIR* are o influență importantă asupra ieșirii și creșterea lui determină de obicei o scădere a performanței (excepție face setul de date dermatology). Parametrul *NF* ia în general valori în jurul lui 0, deci nu are o influență foarte mare asupra ieșirii, excepție de la această regulă o constituie seturile de date mari mushrooms și nursurey, în care se observă că creșterea lui *NF* determină o scădere puternică a ieșirii.

Creșterea parametrului *NMINR* determină o creștere ușoară a ieșirii, dar sunt și câteva situații în care ieșirea este influențată slab, într-un mod negativ de creșterea acestui parametru. În cazul tuturor seturilor de date creșterea lui *NC* determină creșterea valorii variabilei dependente. Parametrul *NI* are de regulă o influență negativă asupra ieșirii, excepție face doar setul de date zoo.

Valorile coeficienților din modelele de regresie și din matricea de corelare demonstrează o parte din observațiile făcute în analiza empirică și anume:

- creșterea numărului de furnici nu conduce la o acuratețe mai bună și din contră, dacă setul de date este mare, duce la o scădere a acurateței prin faptul că algoritmul nu ajunge la rezultat
- creșterea valorii parametrului număr minim de instanțe acoperite de fiecare regulă descoperită (*NMIR*) determină o scădere a acurateței

3.3.4. Considerații și recomandări cu privire la alegerea parametrilor algoritmului Ant Colony Optimization

Rulările au fost realizate rulând aplicația în paralel pe un număr de 30 de calculatoare dintr-un laborator al Universității "Politehnica" din Timișoara. Toate calculatoarele au o configurație identică Intel Core 2 Quad CPU, 2,4 GHz, 2,4 GHz, 4 GB RAM. Au fost rulate și analizate 64 de configurații de parametri pe fiecare din cele 10 seturi de date alese.

Analiza făcută a demonstrat că nu există o configurație de parametri universală care să conducă la rezultate mai bune decât altele pe orice set de date, dar a facilitat analiza comportamentului algoritmului pe diferite seturi de date și extragerea unor concluzii și recomandări utile cu privire la alegerea acestor parametri.

Concluziile cu privire la influența pe care o au parametrii algoritmului asupra acurateței predicției și asupra timpului de execuție și recomandările cu privire la alegerea lor sunt prezentate în continuare.

- Numărul de furnici este numărul maxim de reguli candidate, construite și tăiate (sau fasonate) în timpul unei iterații a buclei repetitive a algoritmului. O furnică este asociată unei reguli. Cu cât este mai mare numărul de furnici cu atât mai multe reguli candidate sunt evaluate pe iterație și cu atât mai lent este sistemul. Experimentele au arătat că, creșterea numărului de furnici determină creșterea direct proporțională a timpului de execuție, fără să garanteze o creștere a performanței. Dacă numărul de furnici este ales foarte mare (de genul 1000, 1500) și setul de date este mare (comparabil ca dimensiune cu *mushrooms* sau *nursurey*) există o șansă mare ca aplicația să facă utilizatorul să aștepte mult fără să ajungă la un rezultat. Rezultate bune s-au obținut pentru numărul de 15 furnici. Alegerea unui număr de furnici mare duce cu siguranță la un timp de execuție mare dar nu se poate spune ca duce și la o creștere a acurateței predicției, așa ca nu recomandăm alegerea unui număr mai mare de 500 de furnici
- Creșterea *NMIR* conduce la o scădere a acurateței predicției și conduce la multe situații de blocare. Cele mai bune rezultate s-au obținut pentru varianta 2, în care *NMIR* ia valoarea 5. Creșterea *NMIR* poate determina algoritmul să intre în buclă infinită în multe situații. Ieșirea din aceste situații de blocare se poate realiza prin creșterea *NMINR*. Datorita faptului că testele în care acest parametru a luat valoarea 15 au făcut ca algoritmul să nu poată fi rulat pe 4 seturi de date din 10, nu recomandăm alegerea unei valori mai mari de 10 pentru acest parametru. Dacă este necesară alegerea unei valori mai mari din diverse motive, atunci creșterea *NMINR* poate mări șansele de a obține rezultate bune.
- Creșterea *NMINR* în general duce la o scădere a performanței, (lucru confirmat și de rezultatele din tabelul 3.12) excepție fac cazurile în care sunt puse condiții restrictive pentru *NMIR*, care conduc la bucle infinite, care pot fi rezolvate prin creșterea *NMINR* (lucru confirmat de rezultatele din tabelul 3.10). Nu recomandăm alegerea unui număr mai mare de 10 pentru acest parametru decât dacă *NMIR* este mai mare de 10.
- Cu cât numărul de reguli pentru convergență este mai mare cu atât timpul de execuție este mai mare. Testele făcute au dus la un număr egal de 3

maxime pentru acuratețea predicției în cazul în care NC ia pe rând valorile 2, 10 și 20. În această situație varianta preferată este cea care are valoarea cea mai mică pentru NC deoarece ea se execută mai repede. Întrucât în bucla repetitivă sunt puse condiții de oprire legate de numărul de reguli pentru convergență dar și de numărul de iterații, alegerea unei valori foarte mici pentru numărul de reguli pentru convergență face ca parametrul NI să nu mai conteze. De aceea influența parametrului NI a fost analizată și în cazul în care NC are valoarea 2 dar și în cazul în care are valoarea 30. Pe de altă parte numărul de reguli pentru convergență reprezintă o garanție că regula găsită este bună. Cu cât mai multe furnici converg către o regulă cu atât mai multe șanse sunt ca regula să fie bună. Ținând cont de rezultatele obținute dar și de semnificațiile mărimilor credem că o valoare între 5 și 10 este potrivită pentru acest parametru.

- Un număr mic de iterații determină execuția mai rapidă a algoritmului, iar creșterea numărului de iterații peste o anumită valoare (în cazul experimentelor noastre, valoarea 50) poate să nu modifice decât puțin timpul de execuție (datorită faptului că bucla *while* se va opri nu datorită condiției legate de numărul de iterații ci condiției legate de NC). În ceea ce privește acuratețea predicției cel mai bun rezultat s-a obținut pentru $NI = 10$ dacă $NC=2$, și $NI=100$, dacă $NC=30$. Întrucât NI și NC sunt condiții de oprire a aceleași bucle repetitive, când se alege o valoare pentru NI trebuie să se țină seama de valoarea lui NC .

Studiul empiric realizat a arătat că acuratețea predicției poate fi influențată cu 10 % de valorile alese pentru parametrii algoritmului, iar timpul de execuție poate fi de la câteva secunde la câteva ore în funcție de valorile acestor parametri. Analiza unui set de date cu ajutorul aplicației *Ant Miner* necesită testarea mai multor configurații de parametri și alegerea celei care duce la cea mai bună performanță. Recomandările de mai sus cresc șansele de a obține un raport bun între acuratețea predicției și timpul de execuție.

Analiza matematică a coeficienților a permis construirea unor modele de regresie pe baza rezultatelor experimentale obținute efectuând câte 64 de măsurători pe fiecare set de date și construirea unor matrice de corelare care redau legătura dintre parametrii de intrare și parametrul de ieșire ales (acuratețea predicției). Analiza matematică a demonstrat prin valorile obținute pentru coeficienții de regresie respectiv prin valorile termenilor din matricea de corelare că:

- Un număr mare de furnici nu duce la o acuratețe mai bună, ba chiar poate duce la o acuratețe slabă dacă setul de date analizat este mare și NF este foarte mare
- Creșterea $NMIR$ determină o scădere a performanței.

O direcție de dezvoltare ulterioară este analiza influenței pe care o au parametrii de mai sus asupra numărului de reguli și numărului de condiții descoperite și găsirea de secvențe de parametri care se potrivesc pentru anumite tipare de seturi de date.

3.4. Variante ale algoritmului ACO pentru descoperirea regulilor de clasificare

În [PAR02] autorii au comparat rezultatele obținute cu *Ant Miner* cu cele obținute cu algoritmul clasic *CN2*, pe 6 seturi de date. Pe 4 din acestea a fost mai bun *Ant Miner*, pe unul a fost mai bun *CN2* și în unul au fost la egalitate. De aceea au fost considerați echivalenți în ceea ce privește acuratețea predicției, dar *Ant Miner* s-a dovedit net superior în ceea ce privește numărul de reguli descoperite și numărul de condiții din aceste reguli. Cu cât regulile sunt mai puține și mai simple cu atât sunt mai ușor de folosit. Direcțiile de dezvoltare sunt modificarea *Ant Miner* în așa fel încât să poată lucra cu atribute continue, ca să nu mai fie necesară discretizarea lor într-un pas de preprocesare precum și analiza modului în care se comportă algoritmul cu alte funcții euristice sau de actualizare a feromonului.

În [PAR01] autorii au comparat algoritmul cu algoritmul clasic *C4.5*, pe 6 seturi de date. În ceea ce privește acuratețea predicției rezultatele au fost comparabile, dar *Ant Miner* a fost net superior în ceea ce privește numărul de reguli și numărul de condiții din fiecare regulă. O direcție de cercetare din [PAR01] este analiza sensibilității algoritmului la parametrii săi. Alte direcții de cercetare se referă la funcțiile de actualizare a feromonului, funcția euristică, și lucrul cu parametrii continui.

În [LIU02a] se subliniază faptul că funcția euristică din *Ant Miner* se bazează pe calculul entropiei, care este destul de costisitor din punct de vedere computațional. Autorii au propus o altă funcție euristică, mult mai simplă care se bazează pe densitate, așa cum se vede în formula 3.10. $Majority_classT_{ij}$ reprezintă numărul de instanțe care aparțin clasei majoritare din partiția în care atributul i ia valoarea j . $|T_{ij}|$ reprezintă numărul de instanțe pentru care atributul i are valoarea j . S-a demonstrat pe un set de date că rezultatele obținute sunt aceleași în ceea ce privește acuratețea predicției și numărul de reguli și de condiții ca și în *Ant miner*.

$$\eta_{ij} = \frac{majority_classT_{ij}}{|T_{ij}|} \quad (3.10)$$

În [LIU03a] autorii afirmă că, datorită modului de actualizare a cantității de feromon din *Ant Miner*, furnicile converg prea repede către o singură regulă construită și au propus o metodă de actualizare a cantității de feromon din formula 3.11. Cu ajutorul parametrului ρ se poate controla cât de repede se evaporă feromonul. Cu cât ρ este mai mare cu atât mai repede se evaporă feromonul.

$$\tau_{ij}(t) = (1 - \rho) \cdot \tau_{ij}(t-1) + \left(1 - \frac{1}{1+Q}\right) \cdot \tau_{ij}(t-1) \quad (3.11)$$

Pentru a mări probabilitatea de a alege termeni care nu au fost aleși în regulile precedente, au fost introduse două numere aleatoare subunitare q_1 și q_2 . Dacă numărul q_1 generat aleatoriu este mai mic decât un prag φ care poate fi

configurat, atunci se alege în mod aleatoriu, proporțional cu performanța, un termen ij , iar dacă este mai mare atunci se alege termenul care are cea mai mare valoare pentru P_{ij} .

```

If (  $q1 \leq \varphi$  )
Loop
    If (  $q2 \leq \sum_{j \in J_i} P_{ij}$  ) atunci alege  $T_{ij}$ 
Endloop
Else
    Alege termenul  $T_{ij}$  cu max  $P_{ij}$ 

```

A fost demonstrat pe 2 seturi de date că *Ant Miner 3* dă rezultate mai bune pentru acuratețea predicției și numărul de reguli este mai mare în *Ant Miner 3*. Acuratețea a fost calculată cu formula :

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN} \quad (3.12)$$

În [SAL11] autorii au propus 5 extensii ale algoritmului și au testat 12 combinații ale acestor extensii pe 23 de seturi de date. De asemenea au fost evaluate performanțele a 3 algoritmi performanți *Ripper*, *PART*, *C4.5* pe aceste 23 de seturi de date. S-a obținut că *C4.5* are cea mai bună performanță. Apoi s-au comparat rezultatele celei mai bune variante de *Ant Miner* din cele 12 încercate cu performanțele obținute de *C4.5* și s-au obținut rezultate comparabile în ceea ce privește acuratețea predicției, iar în ceea ce privește numărul de reguli s-a obținut că cea mai bună versiune de *Ant Miner* produce reguli mai puține decât *C4.5*.

Cele 5 extensii propuse sunt următoarele:

- Utilizarea unor multiple tipuri de feromon, câte un tip pentru fiecare clasă; O furnică întâi selectează clasa regulii și apoi depozitează tipul corespunzător de feromon. În *Ant Miner*, versiunea originală, întâi se construiește antecedentul regulii și apoi regula este atribuită la clasa pe care antecedentul îl acoperă cel mai bine. În versiunea modificată întâi se alege clasa regulii și apoi termenii care se vor adăuga în antecedentul regulii, alegerea termenilor se face ținând seamă de clasa aleasă. O furnică este influențată nu de toate furnicile ci doar de acelea care au construit reguli cu același consecvent (ca și când ar fi diferite tipuri de feromon și o furnică nu este influențată de orice tip). Efectul acestei schimbări a fost nevoia de a modifica structurile bidimensionale (atribut, valoare) ce păstrează cantitatea de feromon și valorile euristice în structuri tridimensionale (atribut, valoare, clasa)
- S-a utilizat un intensificator al calității pentru a premia regulile care au o calitate bună și pentru a penaliza regulile care au o calitate slabă în termeni de feromon.

$$\Delta \tau(t)_k = \begin{cases} 2 \cdot Q(R_t), & \text{if } (confidence(R_t) \geq \phi_1) \\ Q(R_t), & \text{if } (\phi_1 \geq confidence(R_t) \geq \phi_2) \\ Q(R_t) - 2, & \text{if } (\phi_2 \geq confidence(R_t)) \end{cases} \quad (3.13)$$

Unde $\Delta\tau(t)_k$ este cantitatea de feromon de tip k care va fi depozitată de furnica t pe fiecare pereche atribut-valoare care aparține antecedentului regulii R_t precum și clasei k . $Q(R_t)$ este calitatea regulii, iar ϕ_1 și ϕ_2 sunt praguri de confidență la care calitatea este contrastată. Calitatea regulii este evaluată cu altă funcție decât în *Ant Miner*, cea din formula:

$$Q(R_r) = \frac{TP}{\underbrace{TrainingSet}_{Support(R_r)}} + \frac{TP}{\underbrace{Matches}_{Confidence(R_r)}} \quad (3.14)$$

- S-a permis utilizarea operatorului de negare logică în antecedentul regulii, regulile descoperite fiind de forma de mai jos:
If($A^i=V_{ij}$) *and* ($A_k NOT=V_{kl}$) *and*...
Noduri de negare sunt adăugate pentru un atribut dacă are cel puțin 3 valori. Regulile construite în acest fel au o mai puternică acoperire a setului de date decât celelalte, deci se va produce o lista de reguli mai scurtă
- Au fost incluse furnicile încăpățânate care țin cont de propria istorie atunci când adaugă termeni unei reguli. Fiecare furnică memorează cea mai bună regulă pe care a construit-o mergând pe urmele proprii. Fiecare furnică lasă un număr de urme care este configurabil.
- S-au utilizat furnici cu personalitate care au propriile valori pentru α și β (coeficienți pentru feromon și pentru funcția euristică). Personalitatea se manifestă prin faptul că unele furnici vor da mai multă importanță informației legate de feromon iar altele informației legate de funcția euristică. Valorile pentru α și β sunt întregi și cuprinse între 1 și 3. Ele sunt selectate probabilistic de către fiecare furnică înainte de construcția regulii, folosind informația legată de feromon

În [LAL11] autorii propun modificări importante lui *Ant Miner* care îi permit să descopere reguli care în antecedent pot avea operatorii AND, OR, XOR, NOT. Antecedentul fiecărei reguli descoperite poate fi reprezentat sub forma unui arbore binar, în care nodurile interioare corespund operatorilor acceptați AND, OR, XOR, iar frunzele corespund perechilor *atribut=valoare* sau *atribut not =valoare*. Metoda propusă se numește *TREE MINER*. Pentru a construi o regulă o furnică trebuie să traverseze un graf de reguli. A altă diferență față de *Ant Miner* este că nu mai este necesar elagajul regulilor. Pentru actualizarea feromonului a fost introdus un nou parametru numit *conv_rate* care poate fi configurat astfel încât mărirea cantității de feromon pentru termenii adăugați în regulă să se poată face mai mult sau mai puțin. Dacă cantitatea de feromon este mai mare atunci regulile converg mai repede, altfel mai târziu, asta favorizând explorarea. Aplicația a fost testată pe 4 seturi de date și comparată cu *Ant Miner*, în ceea ce privește acuratețea predicției rezultatele sunt comparabile iar numărul de reguli descoperite de *Tree Miner* este mai mic.

În [SAI11] a fost modificată funcția care asigură elagajul regulilor și a fost introdusă o funcție care asigură evaporarea dinamică a feromonului. Algoritmul a fost rulat pe 5 seturi de date și comparat cu *Ant Miner*, iar rezultatele arată că modificările propuse au dus la îmbunătățirea acurateții predicției, regulile descoperite sunt mai ușor de înțeles și timpul de execuție este mai scurt.

În [SMA06] a fost propus algoritmul *MulO-AntMiner* care permite predicția a mai multe atribute clasa. Ca un efect consecventul unei reguli poate sa conțină diferite atribute.

În [WAN04] a fost propusă o nouă versiune a lui *Ant Miner*, versiune care poate determina un set de reguli neordonate (În ant miner regulile descoperite trebuie aplicate instanțelor în ordinea descoperirii). Regulile descoperite cu versiunea de *Ant Miner* din [WAN04] pot fi aplicate instanțelor în orice ordine.

3.5. Soluții de îmbunătățire a algoritmului ACO

Utilizarea algoritmului *Ant Colony Optimisation* la descoperirea regulilor de clasificare a fost propusa de *Parpinelli et al.* în anul 2001. Algoritmul propus de ei a fost întâi implementat în limbajul C și testat pe diverse seturi de date. Ulterior algoritmul a fost implementat în Java și distribuit sub forma unui sistem Open Source numit *Ant Miner*. Ultima versiune este versiunea 1.2.1, dezvoltata în anul 2007 de catre Fernando Meyer în colaborare cu lui Parpinelli. În fig. 3.3 este prezentată interfața principală a aplicației. În cadrul acestui subcapitol au fost propuse soluții de îmbunătățire a algoritmului *Ant Colony Optimisation* și a modului de utilizare a acestuia la descoperirea regulilor de clasificare. Îmbunătățirile propuse vizează modul de calcul a cantității inițiale de feromon asociat fiecărui termen, modul de actualizare a cantității de feromon pentru termenii ce fac parte din regula descoperită, modul de evaluare a calității regulilor, modul de alegere a termenilor care vor face parte din noua regulă. Modificările au fost implementate în *Ant Miner*, și aplicația care le conține a fost numită *Ant-r-Miner*. În figura 3.4 este prezentată interfața grafică a aplicației extinsă.

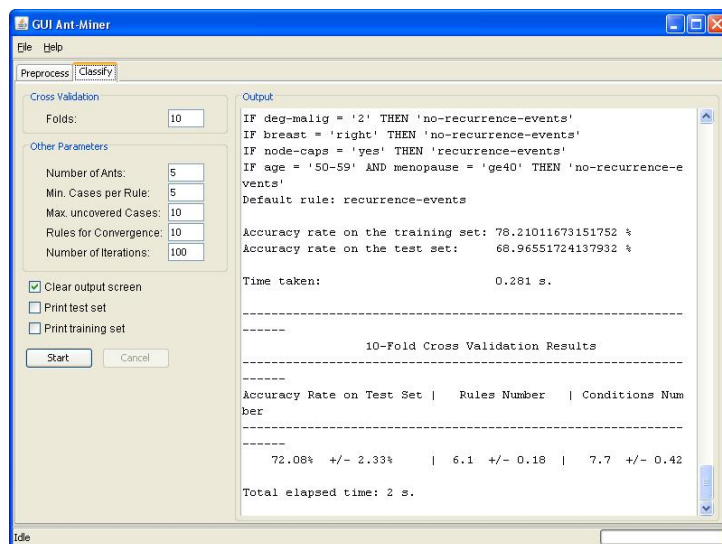


Fig. 3.3 - Ant Miner 1.2.1

Funcția propusă pentru a evalua calitatea unei reguli în *Ant-r-Miner* realizează o sumă între acuratețea predicției și comprehensibilitatea regulii. Acuratețea regulii este exprimată ca un raport între numărul de instanțe care au

Întrădevăr clasa prezisă de regulă raportat la toate instanțele acoperite de regulă (din care unele au clasa acoperită de regulă, iar altele nu au clasa acoperită de regulă). Comprehensibilitatea unei reguli a fost exprimată ținând cont de numărul de condiții din antecedentul regulii. Se consideră că o regulă este cu atât mai valoroasă cu cât este mai ușor de înțeles (are comprehensibilitatea mai mare). Comprehensibilitatea este subunitară și este cu atât mai mare cu cât regula are mai puține condiții în antecedent. Funcția propusă este cea din formula 3.15. Termenul *nar* reprezintă numărul de condiții din antecedentul regulii (numărul de atribute din regulă), unite prin AND (un atribut nu poate apărea decât o dată într-o regulă) iar termenul *nta* reprezintă numărul total de atribute.

$$Q = tp / (tp + fp) + (nta - nar) / nta \quad (3.15)$$

În *Ant-Miner* fiecărei perechi atribut-valoare *i* se atribuie o cantitate inițială de feromon calculată cu formula 3.1. În tabelul de mai jos cu FI se notează feromonul inițial iar cu AP acuratețea predicției. Cantitatea de feromon din coloana FI s-a obținut cu formula 3.1. Acuratețea predicției pentru valorile inițiale este în coloana AP.

Setul de date	FI	AP (FI)	AP (FI= 0.2)	AP (FI= 0.5)	AP (FI= 1)	AP (FI= 5)	AP (FI= 10)
Car	0.08	84.89	85.24	84.5	85.59	85.48	84.66
Zoo	0.019	88.15	90.12	91.41	89.12	90.1	90.03
Ljubljana breast cancer	0.018	73.86	74.68	73.14	72.69	73.46	72.32
Diabetes	0.052	73.96	72.26	73.31	73.55	71.34	72.13
Cleveland heart disease	0.068	77.31	77.54	78.18	75.49	78.54	74.85
Iris	0.105	93.33	94.67	94	95.33	94.67	94.67
Mushrooms	0.007	96.38	97.93	96.15	96.05	96.16	97.01
Nursurey	0.072	84.02	85.52	85.41	85.63	84.88	83.57
Tic	0.034	70.14	72.33	72.02	74.2	74.21	74.1
Dermatology	0.018	92.64	89.64	88.83	90.75	90.04	91.59

Tabel 3.24 - acuratețea predicției pentru diferite valori ale cantitatii inițiale de feromon

Analizând rezultatele din tabelul 3.24 se poate constata că se obțin valori comparabile pentru acuratețea predicției chiar dacă feromonul inițial se calculează cu formula 3.1, ca în *Ant Miner* sau se aleg pentru feromonul inițial valori constante. Atribuirea unor valori constante feromonului inițial este de preferat pentru că este mai rapidă computațional decât calculul sumei numărului de valori pe care îl are fiecare atribut categorial și împărțirea lui $\log_2(k)$ la acesta suma (unde *k* reprezintă numărul de clase). După cum se observă din tabelul 3.25 cele mai bune valori s-au obținut dacă feromonul inițial ia valoarea 1. În *Ant-r-Miner* inițial fiecare termen are asociată cantitatea 1 de feromon.

În urma descoperirii unei reguli este necesară actualizarea cantității de feromon a tuturor termenilor. Termenii care au fost aleși în regulă corespund traseelor urmate de furnică, deci cantitatea de feromon asociată lor se va mări, iar cantitatea de feromon a termenilor care nu fac parte din regulă trebuie diminuată pentru a simula fenomenul de evaporare. Creșterea cantității de feromon a termenilor din regulă se face în *Ant-r-Miner* pe baza principiului de recompensare a

regulilor bune respectiv de penalizare a celor slabe din [SAL11]. Diferența față de [SAL11] este ca se utilizează un singur factor care asigură intensificarea, respectiv diminuarea calității. Acest factor a fost numit *Quality factor* (factor de calitate) și utilizatorul poate stabili care este valoarea lui, precum și a limitelor în care el acționează (vezi formulele 3.16 și 3.17). Factorul de calitate a fost notat cu n și poate lua doar valori mai ≥ 1 . Dacă calitatea regulii este sub o limita inferioară l_1 stabilită de utilizator atunci cantitatea de feromon va crește cu $\tau_{ij} \cdot Q/n$, pentru a penaliza regula slabă. Dacă calitatea regulii este între cele 2 limite stabilite, atunci cantitatea de feromon crește exact ca în Ant Miner cu $\tau_{ij} \cdot Q$. Dacă calitatea regulii este peste o limită superioară stabilită (l_2), atunci cantitatea de feromon se mărește cu $\tau_{ij} \cdot Q \cdot n$, pentru a premia regula bună. În cazul în care se dorește utilizarea acestei facilității se poate fie să se aleagă pentru l_1 valoarea 0 și pentru l_2 valoarea 1 fie să se aleagă valoarea 1 pentru n . l_1 și l_2 trebuie să ia valori reale subunitare.

$$\tau_{ij}(t+1) = \tau_{ij}(t) + \tau_{ij}(t) \cdot \lambda, \quad (3.16)$$

$$\lambda = \begin{cases} Q/n, \forall Q < l_1 \\ Q, \forall Q, l_1 \leq Q \leq l_2 \\ Q \cdot n, \forall Q > l_2 \end{cases} \quad (3.17)$$

În [LIU02a] autorii au demonstrat că se pot obține rezultate la fel de bune și cu o funcție euristică mai simplă, în care în loc de calculul entropiei, care este destul de costisitor, din punct de vedere computațional se calculează o densitate. De aceea funcția euristică utilizată în *Ant-r-Miner* este cea din formula 3.10, lucrarea [LIU02a].

În Ant Miner furnicile converg repede către regula construită, de aceea în *Ant-r-Miner* a fost introdus factorul de convergență ($0 \leq cf \leq 1$), configurabil de utilizator. Dacă utilizatorul alege valoarea 1 pentru acest factor, atunci se va alege termenul cu probabilitatea maximă. O valoare apropiată de 0 a factorului de convergență face ca alegerea termenilor să fie aleatoare.

```
double r=(random.nextInt() << 1 >>> 1) % 101 ;
r=r / 100; // r numar aleator intre 0 si 1

if ( r <= cf )
    alege  $T_{ij}$  cu  $P_{ij}$  maxim
else
    alege  $T_{ij}$  prin selecția de tip ruletă
```

Figura 3.4 prezintă interfața aplicației *Ant-r-Miner/*. După cum se vede, în interfață a fost adăugat un nou *JPanel*, denumit *New Parameters* în care utilizatorul poate introduce valori pentru parametrii noi introduși (*Quality Factor*, *Lower limit*, *Upper Limit*, *Convergence factor*).

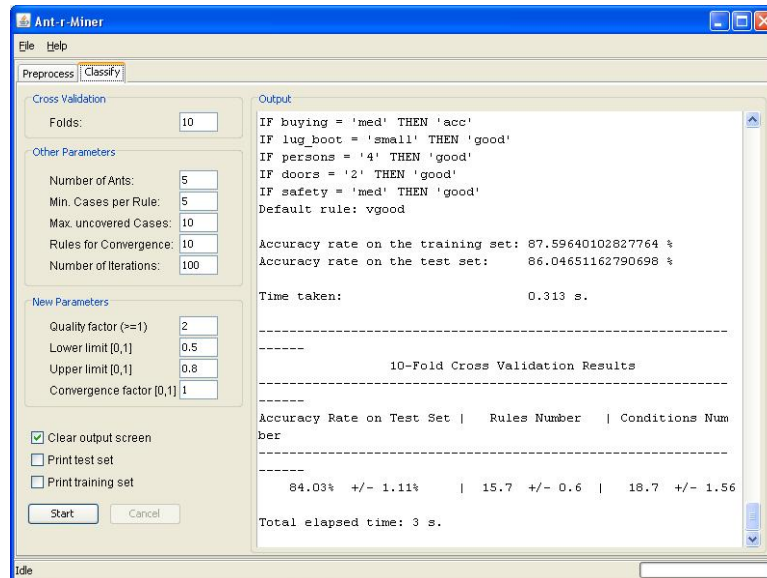


Fig 3.4 - Aplicația propusă, Ant-r-Miner

Algoritmul propus a fost testat pe 10 seturi de date și rezultatele lui au fost comparate cu cele obținute cu Ant Miner. Valorile alese pentru parametrii algoritmului au fost NF=15, NMIR=5, NMINR=10, NC=10, NI=100.

Pentru parametrii noi introduși în Ant-r-Miner s-au ales valorile: *Quality Factor*=2, *Lower limit*=0.5, *Upper Limit*=0.8, *convergence factor*=0.7,

Tabelul 3.25 prezintă comparativ acuratețea predicției și timpul de execuție obținute testând Ant Miner și Ant-r-Miner pe seturile de date alese.

Setul de date	Ant Miner (AP)	Ant-r-Miner (AP)	Ant Miner (TE)	Ant-r-Miner (TE)
Ljubljana breast cancer	72.74	72.95	7	6
Car	85	89.7	19	10
Dermatology	95.59	89.9	158	86
Diabetes	73.56	77.35	10	4
Cleveland heart disease	81.24	78.51	4	9
Iris	94.6	95.33	0	0
Mushrooms	96.15	99	1566	1002
Nursurey	86.52	88.68	183	107
Tic	72.12	81.83	23	8
Zoo	91.23	86.26	31	19

Tabel 3.25 – acuratețea predicției cu Ant-Miner și Ant-r-Miner

Rezultatele din tabelul 3.25 arată că se obțin rezultate puțin mai bune cu *Ant-r-Miner* decât cu *Ant Miner* în ceea ce privește acuratețea predicției, iar în ceea ce privește timpul de execuție rezultatele obținute cu *Ant-r-Miner* sunt mult mai bune. În ceea ce privește numărul de reguli și de condiții *Ant Miner* este superior.

3.6. Concluzii

Algoritmul *Ant Colony Optimization* este folosit cu succes la clasificarea datelor. Din anul 2001 când a fost propusă utilizarea lui pentru descoperirea regulilor de clasificare, au fost realizate multe studii cu ajutorul acestui algoritm și problema alegerii unei secvențe de parametri optime a fost lăsată ca o direcție de cercetare de foarte mulți cercetători. În cadrul acestui capitol a fost studiat codul algoritmului și prezentată în detaliu funcționarea lui, au fost analizate foarte multe configurații de parametri și s-a emis o listă de recomandări privind alegerea acestora. De asemenea problema îmbunătățirii algoritmului a fost atinsă de mulți cercetători. Au fost prezentate sintetic îmbunătățirile pe care o parte din aceștia le-au adus și au fost propuse și operate noi îmbunătățiri. Acest lucru a fost posibil datorită faptului că aplicația care implementează algoritmul se numește *Ant Miner* și este Open Source.

Contribuțiile aduse în cadrul acestui capitol sunt următoarele:

- A fost realizată o prezentare detaliată a algoritmului, bazată atât pe documentația existentă cât și pe studiul codului algoritmului.
- Au fost studiate parametrii algoritmului rulând câte 64 de configurații pe fiecare din cele 10 seturi de date alese și analizând rezultatele. Analiza rezultatelor s-a făcut atât vizual, comparativ cât și matematic construind și analizând modele de regresie și matrice de corelare
- Analiza făcută a condus către o listă de recomandări privind modul de alegere al parametrilor, utile pentru toți utilizatorii acestui algoritm
- A fost realizată o sinteză a îmbunătățirilor aduse algoritmului din anul 2001 și până în prezent
- Au fost aduse îmbunătățiri algoritmului care vizează modul de calcul al cantități inițiale de feromon a termenilor, modul de actualizare a feromonului, modul de evaluare a calității regulilor descoperite și modul de selecție a termenilor. Modificările propuse au fost implementate în aplicația open source *Ant Miner*. Aplicația care conține algoritmul îmbunătățit a fost numită *Ant-r-Miner*. S-a demonstrat că îmbunătățirile aduse conduc la o mai bună performanță utilizând date reale.

Pe viitor se dorește extinderea aplicației oferindu-i utilizatorului posibilitatea să aleagă dintr-o listă de funcții, funcția pe care dorește să o folosească pentru a evalua calitatea regulilor. O altă direcție de dezvoltare este de oferire a posibilității utilizatorului să-și construiască în mod dinamic funcția de calitate, cu ajutorul termenilor *TP*, *FP*, *TN*, *FN* și a unor operații elementare gen adunare, scădere, înmulțire, împărțire. Analiza influenței pe care o au valorile parametrilor asupra numărului de reguli și asupra numărului de condiții este de asemenea o direcție de dezvoltare ulterioară.

4. Clasificarea datelor cu ajutorul algoritmului genetic *AGR*

4.1. Noțiuni introductive

Algoritmii genetici sunt tehnici adaptive care pot fi folosite cu succes în rezolvarea unor probleme de căutare și optimizare complexe [MAU00]. Ei sunt bazați pe principiile geneticii și ale selecției naturale enunțate de Darwin ("supraviețuiește cel mai bine adaptat"). Algoritmii genetici au fost utilizați cu succes în *data mining*, pentru a determina reguli de clasificare [DEH06], pentru a căuta centrele de clustere adecvate [MAU00], pentru a selecta atributele de interes în prezicerea valorii unui atribut țintă [FRE02], etc. Clasificarea instanțelor a fost realizată și prin utilizarea unor algoritmi hibrizi, bazați pe algoritmi genetici și *Particle Swarm Optimization* [DIN09], respectiv *Naïve Bayes* și *k-Nearest Neighbor* [AIC10]. Câteva aplicații în care algoritmii genetici au fost aplicați cu succes pentru a rezolva probleme de clasificare sunt: clasificarea amprentelor, clasificarea bolilor de inimă, clasificarea emoțiilor de pe fața umană, clasificarea imaginilor de la sateliți [YAN06], etc.

Capitolul prezintă sintetic aspecte privind algoritmii genetici utilizați la clasificarea datelor, funcțiile *fitness* ale acestora și propune un nou algoritm genetic pentru clasificarea datelor. Algoritmii genetici propus a fost implementat în Open Source-ul *Weka*, utilizând limbajul *Java*. Mediul de dezvoltare utilizat a fost *Eclipse*. Sunt prezentate aspecte legate de procesările algoritmului, de implementare și de rezultatele experimentale obținute prin aplicarea algoritmului pe 5 seturi de date obținute din *UCI Machine Learning Repository*. S-a analizat comportamentul algoritmului utilizând 14 variante de parametri. A fost analizat comportamentul algoritmului atât în cazul în care se adaugă la lista regulilor descoperite toate regulile descoperite după epuizarea numărului de generații cât și în cazul în care se adaugă doar cea mai bună regulă descoperită după rularea numărului de generații stabilit. Sunt comparate rezultatele obținute rulând algoritmul cu 3 funcții *fitness* diferite. Funcția *fitness* care a condus la cele mai bune rezultate a fost inclusă în algoritm.

Rezultatele obținute de *AGR* pe seturile de date alese au fost comparate cu cele obținute de algoritmii *Naive Bayes*, *k-Nearest-Neighbour* și *J48* pe aceleași seturi de date și în aceleași condiții. Concluzia a fost că algoritmul propus este competitiv, poate fi utilizat la clasificarea datelor putând să conducă la rezultate mai bune pe anumite seturi de date decât alți algoritmi clasici.

4.2. Structura unui algoritm genetic

Forma generală a unui algoritm genetic este următoarea [MUK09]:

```

Determină populația inițială;
WHILE a condiție de oprire nu este îndeplinită DO
BEGIN
    Selectează indivizii pentru reproducere;
    Creează moștenitori prin încrucișarea indivizilor selectați;
    Eventual aplică mutația unor indivizi;
    Calculează următoarea generație
END

```

Populația inițială este aleasă aleatoriu. Selecția este procedura prin care sunt aleși indivizii ce vor participa la formarea generației următoare. Indivizilor mai puternici, mai bine adaptați li se vor da mai multe șanse de a face parte din populația următoare. Performanța indivizilor este măsurată cu ajutorul unei funcții *fitness* care depinde de problema de rezolvat. Selecția indivizilor pentru reproducere se poate face în mai multe moduri. Cea mai populară variantă este selecția de tip ruletă. Fiecare cromozom are asociată o felie a ruletei. "Feliile ruletei sunt cu atât mai late cu cât *fitness*-ul cromozomului asociat feliei este mai mare".

O altă tehnică de selecție este selecția bazată pe rang. Aceasta previne convergența prematură ce apare la utilizarea schemelor de selecție proporțională cu *fitness*-ul. Indivizii sunt ordonați în funcție de valoarea *fitness*-ului iar probabilitatea de selecție este proporțională cu rangul ocupat.

Selecția prin trunchiere, elimină o parte din cromozomii cu cea mai slabă performanță, iar în locul lor se generează alții, după diferite scheme posibile. De exemplu se alege un nou cromozom x , folosind principiul ruletei. Dacă x diferă de toți ceilalți cromozomi ai populației actuale, atunci el este inclus în populație; în caz contrar, este supus operatorului de "mutație" până ce devine diferit de ceilalți cromozomi și este inclus în populație.

Selecția turneu alege în mod aleatoriu k indivizi iar dintre aceștia doar cei mai buni j sunt selectați pentru supraviețuire. Procedura se repetă până se obține numărul dorit de indivizi. Este cea mai eficientă din punct de vedere al complexității timp.

Selecția de tip elitist păstrează în generația următoare, cei mai buni k indivizi ai generației curente. Dacă nu se rețin acești indivizi ei pot să dispară din cauză că nu au fost selectați pentru supraviețuire sau fiindcă au fost distruși prin aplicarea operatorilor. Acest lucru garantează faptul că cel mai puternic membru al populației din generația următoare nu va fi mai slab decât cel din generația curentă.

Încrucișarea a doi indivizi se poate face într-un singur punct sau în mai multe puncte. Noii indivizi ai populației sunt obținuți prin combinarea unor părți alternative de material genetic provenind de la părinți.

Mutația permite algoritmului genetic să găsească noi soluții în cadrul populației și îl protejează împotriva pierderii de informație în cazul unor încrucișări nepotrivite. Dacă operatorul de selecție reduce diversitatea în populație, cel de mutație determină o nouă creștere a diversității. Cu cât probabilitatea mutației este mai mare, cu atât mai redus este riscul convergenței premature, dar apare un nou risc datorită faptului că o rată mare de mutație va transforma algoritmul genetic într-un algoritm de căutare aleatoare.

4.3. Algoritmi genetici utilizați la clasificarea datelor

Algoritmii genetici propuși în literatură pentru clasificarea datelor, au structura din subcapitolul 4.2, dar utilizează diverse moduri de reprezentare a cromozomilor, de calcul al funcțiilor *fitness*, de încrucișare a indivizilor, etc.

În [DEH06] un cromozom e reprezentat ca o înșiruire de m biți (vezi formula 4.1).

$$m = \sum_{i=1}^n nr(a_i) + n - 1 \quad (4.1)$$

Unde,

- n reprezintă numărul de atribute,
- a_i este atributul i
- $nr(a_i)$ reprezintă numărul de valori ale atributului i

Pentru reprezentarea fiecărei perechi atribut-valoare, $a_i = v_{ij}$, unde a_i este atributul i și v_{ij} este valoarea j a atributului i , sunt necesari $nr(a_i) + 1$ biți, câte un bit pentru fiecare valoare posibilă a atributului și un bit a cărui valoare indică dacă atributul face parte din regulă sau nu.

Funcția *fitness* propusă în [DEH06] este cea din relația 4.2. Funcția calculează performanța regulii luând în considerare acuratețea predicției, claritatea regulii dar și gradul de interes al regulii.

$$Fitness = \frac{w_1 C(R) + w_2 predAcc + w_3 RInt}{w_1 + w_2 + w_3} \quad (4.2)$$

Unde $C(R)$ este comprehensibilitatea, $predAcc$ este acuratețea predicției și $RInt$ este gradul de interes al regulii, iar w_1, w_2, w_3 sunt ponderi definite de utilizator. Comprehensibilitatea se calculează în (4.3) unde M este numărul maxim de condiții al unei reguli.

$$C(R) = M - \text{number of condition } (R) \quad (4.3)$$

Acuratețea predicției este calculată cu relația (4.4).

$$predAcc = \frac{|A \& C|}{|A|} \quad (4.4)$$

Gradul de interes al regulii se calculează în (4.5) unde $InfoGain(A_i)$ este câștigul informațional adus de atributul i , și $|\text{dom}(G)|$ este cardinalitatea domeniului pentru atributul clasă G .

$$RInt = 1 - \frac{\sum_{i=1}^{n-1} InfoGain(A_i)}{\log_2(|dom(G)|)} \quad (4.5)$$

În [VIV10a] reprezentarea cromozomilor se face la fel ca în [DEH06]. După încrucișare sunt realizate unele prelucrări ale cromozomilor rezultați.

În primul rând, se determină pentru fiecare atribut care este valoarea care apare cel mai des în cromozomii rezultați. În următorul pas se înlocuiește valoarea determinată cu o valoare obținută aleatoriu. Dacă valoarea *fitness-ului* crește în urma acestei înlocuiri atunci se păstrează valoarea obținută aleatoriu.

A doua modificare constă în alegerea aleatoare a unui atribut din fiecare cromozom. Dacă el face parte din regulă, se elimină din ea și se evaluează *fitness-ul* regulii. Dacă acesta a crescut atunci atributul rămâne eliminat. Selecția utilizată este selecția de tip ruletă, iar încrucișarea utilizată este încrucișarea în două puncte. Mutația se realizează asupra unui singur individ din populație ales aleatoriu și basculând un bit al său ales aleatoriu. Cei mai buni n cromozomi sunt selectați din părinți și copii ca să facă parte din noua generație. Restul cromozomilor sunt aleși din cromozomii copii. Algoritmul se poate opri dacă nu se constată nici o îmbunătățire în valoarea *fitness-ului* vreme de m generații sau dacă încă nu a fost generat numărul de generații dorit. *Fitness-ul* se calculează cu formula 4.6.

$$Fitness = \frac{tp}{tp + A \cdot Fn} \cdot \frac{tn}{tn + B \cdot fp} \quad (4.6)$$

Unde $0.2 < A < 2$ și $1 < B < 20$. A și B sunt parametri definiți de utilizator care determină rata de convergență a soluției și sunt determinați experimental în acord cu cerințele.

În [VIV10b] este propus algoritmul PGA (*Partitioning Genetic Algorithm*), care împarte setul de date în partiții, algoritmul genetic este aplicat fiecărei partiții și un set de reguli este extras de pe fiecare partiție. Apoi seturile de reguli se combină pentru a forma o listă de reguli candidat. În al doilea pas suportul și confidența se calculează pentru fiecare regulă din regulile candidate folosind întregul set de date. Cromozomii se reprezintă utilizând câte un bit pentru fiecare valoare posibilă a unui atribut și încă un bit care indică dacă perechea atribut=valoare considerată, face parte din regulă sau nu. Mutația face ca un atribut să primească o altă valoare în mod aleatoriu. Au fost dezvoltati 2 operatori care controlează dimensiunea regulilor adăugând sau ștergând termeni din reguli. Operatorul de inserare activează o gena punând pe 1 bitul care indică dacă perechea atribut valoare face parte din regulă, iar operatorul de ștergere pune pe 0 acest bit. Cei mai buni cromozomi din generația curentă sunt copiați în noua generație nealterați. Funcția *fitness* propusă este cea din relația 4.7:

$$Fitness = w_1 \cdot \frac{|A \& C - 0.5|}{|A|} + w_2 \cdot \frac{L - n}{L - 1} \quad (4.7)$$

Primul termen al funcției măsoară acuratețea predicției iar cel de-al doilea măsoară comprehensibilitatea regulii.

- $|A \& C|$ este numărul de instanțe care satisfac și antecedentul și consecventul regulii
- $|A|$ este numărul de instanțe care satisfac antecedentul regulii.
- L este numărul maxim de condiții ale regulii
- n este lungimea regulii.

Articolul [FRE02] este unul de sinteză și prezintă două metode de codificare a cromozomilor, prima dintre ele este cea utilizată în [DEH06] [VIV10a] [VIV10b] iar cea de-a doua este prezentată în continuare. Metoda propusă utilizează pentru fiecare atribut un număr de biți egal cu numărul de valori posibile ale atributului. Un atribut lipsește dintr-o regulă dacă toți biții alocați atributului au fie valoarea 1 fie valoarea 0. Dacă reprezentarea valorii unui atribut conține mai mulți de 1, înseamnă că în regulă apare atributul respectiv cu mai multe valori unite prin SAU logic. În ceea ce privește reprezentarea consecventului, articolul prezintă mai multe variante:

- consecventul se include în genomul individului, făcându-l astfel subiect al evoluției
- valorile consecventului se asociază pe rând cu toți indivizii din populație. Pentru a determina reguli care prezic diferite clase se rulează algoritmul de un număr de ori egal cu numărul de clase.
- se alege clasa cea mai potrivită pentru regulă, imediat ce antecedentul acesteia a fost determinat. Alegerea clasei se poate face după diverse criterii: fie se alege clasa care are cei mai mulți reprezentanți care satisfac antecedentul și consecventul, fie se alege clasa care face ca *fitness*-ul regulii să fie maxim.

Funcția *fitness* propusă în [FRE02] este cea din formula 4.8.

$$Fitness = w_1 \times (CF \times Comp) + w_2 \times Simp \quad (4.8)$$

Primul termen din relația 4.8 surprinde acuratețea predicției, iar al doilea surprinde cât de ușor de înțeles este regula. $CF = tp / (tp + fp)$ este factorul de confidență iar $Comp = tp / (tp + fn)$ este factorul de complitudine și redă care este proporția de instanțe acoperite de antecedentul regulii din cele care au clasa prezisă de regulă [FRE02]. *Simp* reprezintă simplitatea regulii (este invers proporțională cu numărul de condiții din antecedentul regulii). w_1 și w_2 sunt ponderi definite de utilizator.

4.4. AGR - Algoritmul genetic propus

Potrivit rezultatelor experimentale prezentate în articolele referite în subcapitolul 4.3, algoritmi genetici au un mare potențial de aplicare la clasificarea datelor. Cu toate acestea, versiunea 3.6 a lui *Weka*, care este ultima versiune stabilă la ora actuală, nu dispune de un algoritm genetic. Algoritmii genetici din literatură prezentați în cadrul subcapitolului 4.3 diferă în ceea ce privește codificarea cromozomilor, funcțiile *fitness* care le utilizează, procesările care le efectuează, etc. Problema modului în care aceștia asigură că regulile descoperite pot clasifica orice instanță nu este lămurită clar. Așa cum se va vedea în capitolul 4.5 regulile

descoperite după rularea algoritmului pentru numărul de generații stabilit nu acoperă toate instanțele. Pe de altă parte câtă vreme în literatură nu există un algoritm genetic clasificator care să asigure o acuratețe a predicției de 100 % pe orice set de date, problema analizei și propunerii unor noi algoritmi genetici în acest scop, rămâne una deschisă.

Pornind de la observațiile de mai sus în cadrul acestui capitol a fost propus un nou algoritm genetic numit *AGR*, algoritm genetic care diferă după cunoștințele autorului de cei propuși în literatură prin următoarele:

- *permite recalcularea clasei în urma încrucișării.* Un cromozom codifică o regulă. În urma încrucișării cromozomilor, clasa cromozomilor rezultați este posibil să nu fie cea mai potrivită. De aceea a fost dezvoltată o facilitate care permite recalcularea clasei în urma încrucișării.
- *permite stabilirea unui procent minim de acoperire pentru regulile descoperite.* Algoritmul genetic rulează pentru numărul de generații stabilit și descoperă un set de reguli. Se determină acuratețea regulilor dar și acoperirea lor. Dacă aceste reguli descoperite au o acoperire mai mică decât cea impusa atunci ele se adaugă unei liste de reguli, instanțele din mulțimea de antrenament acoperite de ele se elimină din aceasta și tot algoritmul se reia pana se ajunge la acoperirea dorită.
- *permite determinarea și adăugarea unei reguli implicite listei regulilor descoperite.* Regula care se determină nu are condiții în antecedent, poate acoperi orice instanță și o clasifică ca aparținând de clasa majoritară. Aceasta este o metodă rapidă de a crește acoperirea listei regulilor descoperite.
- *permite adăugarea la lista de reguli descoperite doar a celei mai bune reguli, obținută după rularea algoritmului pentru numărul de generații stabilit sau a tuturor regulilor descoperite.* După rularea algoritmului pentru numărul de generații stabilit se descoperă mai multe reguli. Utilizatorul poate opta, dacă le consideră valide, pe toate sau doar pe cea mai bună dintre ele
- *utilizează o funcție fitness diferită de cele utilizate de alți algoritmi.* Au fost testate 3 funcții *fitness* și a fost aleasă cea cu care s-au obținut cele mai bune rezultate.

Pseudocod-ul algoritmului genetic propus este prezentat în continuare:

```

if (a fost selectat clasificatorul algoritmului genetic) then
  - împarte aleatoriu setul de date ales în mulțime de antrenament și mulțime de test
  - alege aleatoriu instanțele care formează populația
  - codifică cromozomii, instanțele din mulțimea de antrenament și cele din mulțimea de test

  lista_tuturor_regulilor_descoperite={ }
  contor_iterații=1;

  do{ //iterația n
    contor_generații=0;

    do{

```

```

- determină valorile fitness ale cromozomilor și suma acestora
- selectează o populație intermediară de cromozomi după principiul ruletei
- încrucișează unii cromozomi din populația intermediară
- aplică operatorul de mutație populației
- recalculează clasa noilor cromozomi (opțiune a utilizatorului)

    contor_generații+=1;
}while (contor_generații < nr_generații);

- afișează populația finală
- elimină regulile duplicate (cromozomii)
- adaugă regulile descoperite sau doar cea mai bună dintre regulile descoperite (opțiune a utilizatorului) la lista_tuturor_regulilor_descoperite
- decodifică cromozomii și afișează regulile codificate de aceștia
- determină și afișează gradul de acoperire al regulilor descoperite pe mulțimea de antrenament
- calculează și afișează acuratețea predicției pe mulțimea de antrenament
- elimină din mulțimea de antrenament instanțele acoperite de regulile descoperite
- generează o nouă populație de cromozomi din mulțimea de antrenament
    contor_iterații+=1;
}while (gradul de acoperire al regulilor descoperite < gradul de acoperire impus de utilizator)

- determină și adaugă o regulă implicită listei de reguli (opțiune a utilizatorului)
- afișează toată lista de reguli descoperite
- testează regulile pe mulțimea de test și se afișează acoperirea și acuratețea predicției
- afișează timpul de execuție al algoritmului

```

Împărțirea setului de date în mulțime de antrenament și mulțime de test se face în mod aleatoriu. Procentul de instanțe din mulțimea de test este preluat din interfață din caseta de text `m_PercentText`.

Instanțele din populația inițială sunt alese în mod aleatoriu din setul de date.

Codificarea cromozomilor

Cromozomii se obțin prin codificarea binară a instanțelor din populație. O genă este reprezentată de valoarea unui atribut. Fiecare cromozom codifică o regulă. Regulile din populația inițială sunt exprimate sub forma:

```
IF atribut1=valoare1 și atribut2=valore2 și... THEN clasa=c1
```

Fiecare regulă (cromozom) se codifică sub forma unui șir de biți [VIV10a]. Un atribut care poate lua m valori va fi codificat pe m biți (câte un bit pentru fiecare valoare posibilă). În cazul în care avem 3 atribute A1 cu 4 valori posibile, A2 cu 2 valori posibile și A3 cu 3 valori posibile și un atribut clasa cu 3 valori posibile, codificarea unei reguli se poate face ca în tabelul 1 (biții sunt numerotați de la

dreapta la stânga). În cazul în care setul de date conține atribute numerice acestea se discretizează.

Atribut	A1	A2	A3	C
Cod	1000	01	001	100
Regulă	If A1=valoare1 and A2=valoare2 and A3=valoare3 then C=valoare1			

Tabel 4.1 - Codificarea cromozomilor

Instanțele din mulțimea de antrenament și cele din mulțimea de test se codifică și ele binar.

Determinare performanță cromozomi

În continuare începe o secvență repetitivă de cod, care se execută de un număr de ori egal cu numărul de generații stabilit de utilizator. Se determină valorile *fitness* pentru toți cromozomii din populație. Funcția *fitness* folosită este cea din formula 4.9. Primul factor reprezintă acuratețea predicției (vezi formula 4.10) iar cel de-al doilea sensibilitatea regulii (vezi formula 4.11). S-a analizat comportamentul algoritmului la diferite funcții *fitness* (vezi subcapitolul 4.5.3) și funcția *fitness* cu care s-au obținut cele mai bune rezultate este cea din formula 4.9, care a fost inclusă în algoritm.

$$Fitness = PA \cdot Senz \quad (4.9)$$

$$PA = \frac{tp}{tp + fp} \quad (4.10)$$

$$Senz = \frac{tp}{tp + fn} \quad (4.11)$$

Selecție cromozomi

Algoritmul creează o populație intermediară în care sunt puși cromozomii selectați în mod proporțional cu performanța lor, după principiul ruletei. Pseudocodul algoritmului este prezentat mai jos.

- determină suma valorilor *fitness* corespunzătoare tuturor cromozomilor
- r = un număr aleatoriu în intervalul $[0, \text{suma}]$

```

suma_partială=0
For i=0 to nr_indivizi_populație
    suma_partială+=fitness(i);
    if(suma_partială>r)
        Break;
Return i

```


Încrucișare cromozomi

Se parcurge populația intermediară și pentru fiecare cromozom se generează un număr aleatoriu, cu distribuție uniformă între 0 și 1. Dacă acest număr este mai mare decât probabilitatea de încrucișare, stabilită de utilizator, atunci cromozomul va fi copiat în noua populație, dacă nu, el se va încrucișa cu altul pentru a forma doi descendenți în noua populație. Încrucișarea se realizează într-un punct, iar punctul de încrucișare este ales aleatoriu (între 0 și numărul de biți din antecedent). Tabelul 4.2 exemplifică încrucișarea cromozomilor în punctul de taiere 3, obținut aleatoriu. Spațiile dintre grupurile de biți nu fac parte din reprezentarea cromozomilor dar au fost puse pentru a ajuta utilizatorul să delimiteze atributele.

Cromozom1	1000 10 001 100
Cromozom2	0001 01 010 001
Cromozom1'	1001 01 010 001
Cromozom2'	0000 10 001 100

Tabel 4.2 - Încrucișarea cromozomilor

Mutație cromozomi

În următoarea etapă se parcurg toți cromozomii din noua populație și pentru fiecare individ se generează aleatoriu un număr între 0 și 1. Dacă acest număr este mai mic decât probabilitatea de mutație (parametru stabilit de utilizator), atunci asupra cromozomului respectiv se va aplica operatorul de mutație. Operatorul de mutație generează un număr aleatoriu cuprins între 0 și numărul de biți din antecedent. Bitul astfel obținut va fi basculat (dacă are valoarea 1 devine 0 și invers).

Recalculare clasă

În reprezentarea cromozomilor a fost inclusă și clasa. Este posibil însă ca în urma operațiilor de încrucișare și mutație cromozomii obținuți să prezică mai bine o altă clasă decât cea a părinților. În acest sens, dacă utilizatorul dorește, bifează o opțiune care determină pentru fiecare cromozom din noua populație care este clasa cea mai potrivită. S-a considerat că, clasa cea mai potrivită este cea care conduce la cele mai multe instanțe clasificate corect.

Algoritmul prezentat mai sus se repetă pentru un număr de generații stabilit de utilizator. Indivizii din populația finală sunt afișați.

Convergența cromozomilor

În urma execuției numărului de generații stabilit se produce fenomenul de convergență către cel mai bun individ. Astfel populația finală obținută după numărul de generații stabilit, va avea foarte multe duplicate de cromozomi (spre exemplu, în cazul setului de date car.arff populația finală are între 1 și 3 cromozomi unici). Metoda `elimina_cromozomi_duplicat` va lăsa în populația finală doar cromozomii unici. Cromozomii din populația finală codifică regulile descoperite.

Reține doar cea mai bună regulă din fiecare iterație

Dacă este bifată opțiunea *Reține doar cea mai bună regulă din fiecare iterație*, se determină regula (cromozomul) care are cel mai mare *fitness* din cele descoperite. Doar această regulă va fi considerată validă. Ea se va adăuga listei de

reguli descoperite și va fi decodificată și afișată. Dacă nu este bifată opțiunea menționată atunci toate regulile descoperite după execuția numărului de generații stabilit sunt considerate valide și vor fi adăugate listei regulilor descoperite, respectiv vor fi decodificate și afișate.

Decodificarea regulilor

Regulile codificate de cromozomi sunt decodificate și afișate în caseta de text `m_OutTextGA`. Operația de decodificare poate întâlni câteva probleme specifice:

- *Atribute care au toți biții cu valoarea 0.* În urma operațiilor de încrucișare și mutație repetate este posibil ca unele atribute să aibă toți biții puși pe 0. Acele atribute nu vor apărea în regulă.
- *Atribute care au toți biții cu valoarea 1.* Astfel de atribute vor fi eliminate din regulă pentru că nu influențează cu nimic rezultatul clasificării.
- *Atribute care au mai multe valori de 1 (dar nu toate).* Dacă un atribut are mai multe valori de 1 atunci el va apărea în regula de mai multe ori, operatorul SAU va lega valorile posibile ale atributului (de ex. codul 1001 al unui atribut n se decodifică Atribut n =valoare 1 SAU atribut n =valoare 4)

Acoperirea regulilor

Se calculează și se afișează acoperirea pe care o au regulile obținute asupra mulțimii de antrenament, precum și acuratețea predicției. Acoperirea regulilor s-a calculat numărând câte instanțe din mulțimea de antrenament sunt acoperite de antecedentul regulilor descoperite. Acuratețea predicției s-a calculat împărțind numărul de instanțe care respectă și antecedentul și consecventul regulilor la numărul de instanțe care respectă antecedentul regulilor.

Algoritmul a fost testat pe toată durata implementării pe setul de date *car.arff*. S-a constatat că regulile descoperite după rularea a 50 de generații pe o populație de 100 de cromozomi au o acuratețe foarte bună, dar acoperirea acestor reguli este slabă. În acest sens au fost concepute și implementate două soluții de creștere a acoperirii regulilor.

O primă soluție de creștere a acoperirii regulilor este prezentată în continuare. Utilizatorul are posibilitatea să impună ca algoritmul să continue căutarea de reguli până când se găsește un număr de reguli care asigură procentul de acoperire al instanțelor din mulțimea de antrenament dorit de utilizator. Dacă regulile descoperite acoperă un procent de instanțe mai mic decât cel dorit de utilizator, atunci ele se adaugă listei de reguli, instanțele din mulțimea de antrenament acoperite de aceste reguli sunt eliminate, se reinițializează populația cu instanțe alese aleatoriu din mulțimea de antrenament și algoritmul începe o nouă iterație. Dacă acoperirea regulilor descoperite este peste valoarea minimă impusă de utilizator, regulile descoperite sunt testate pe mulțimea de test și rezultatele sunt afișate. Se determină timpul de execuție al algoritmului și se afișează.

O altă soluție de a crește acoperirea regulilor este de a determina și adăuga la lista regulilor descoperite o regulă implicită. Această regulă implicită nu are nici o condiție în antecedent, în felul acesta ea va cuprinde toate instanțele care nu au putut fi clasificate cu regulile descoperite. Clasa regulii implicite este întotdeauna clasa cea mai populară printre instanțele rămase în setul de antrenare. Utilizatorul poate opta prin intermediul unui control din interfață dacă dorește sau nu determinarea unei regulii implicite și adăugarea ei la lista de reguli descoperite prin rularea repetată a algoritmului genetic.

4.5. Weka extinsă cu algoritmul genetic AGR

Algoritmul genetic propus a fost denumit *AGR* și a fost implementat în versiunea 3.6.4 a lui *Weka* (limbajul Java). Figura 4.1 prezintă interfața aplicației extinse.

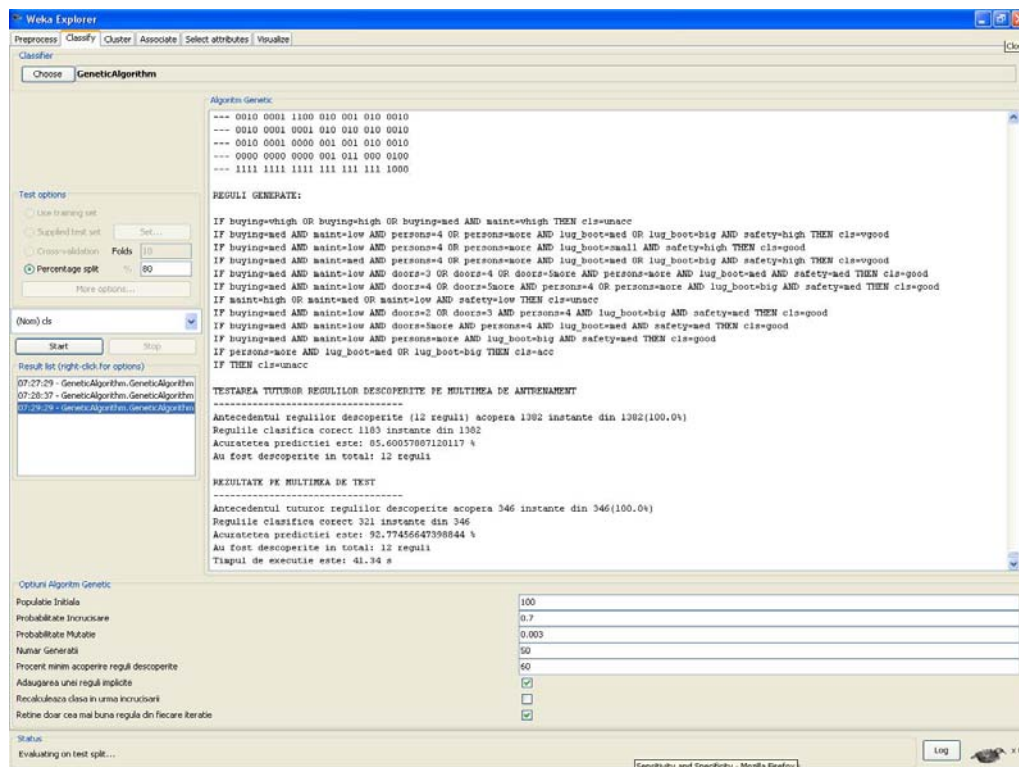


Fig. 4.1 - Weka cu algoritmul genetic

Codul sursă al algoritmului a fost implementat în fișierul *ClassifierPanel.java*. Pentru a adăuga clasificatorul *GeneticAlgorithm* în lista de clasificatori (vezi figura 4.2) a fost creat un nou pachet *weka.classifiers.GeneticAlgorithm* și a fost completat numele pachetului în lista clasificatorilor din fișierul *GenericPropertiesCreator.props* din pachetul *weka.gui*.

În pachetul *weka.classifiers.GeneticAlgorithm* adăugat a fost creată clasa *GeneticAlgorithm*, derivată din *Classifier*, care implementează interfețele *UpdateableClassifier*, *TechnicalInformationHandler*.

Codul sursă al algoritmului, precum și codul aferent modificării interfeței grafice a fost scris în fișierul *ClassifierPanel.java*. În interfață au fost adăugate o serie de controale de tip *JPanel*, *JLabel*, *JTextField*, *JCheckBox*, pentru a permite utilizatorului să precizeze valorile unor parametri precum: dimensiunea populației, probabilitatea de încrucișare, probabilitatea de mutație, procentul minim permis de

acoperire a regulilor descoperite, numărul de generații, dacă dorește sau nu adăugarea unei reguli implicite la lista regulilor descoperite, dacă dorește sau nu recalcularea clasei în urma încrucișării și mutației.

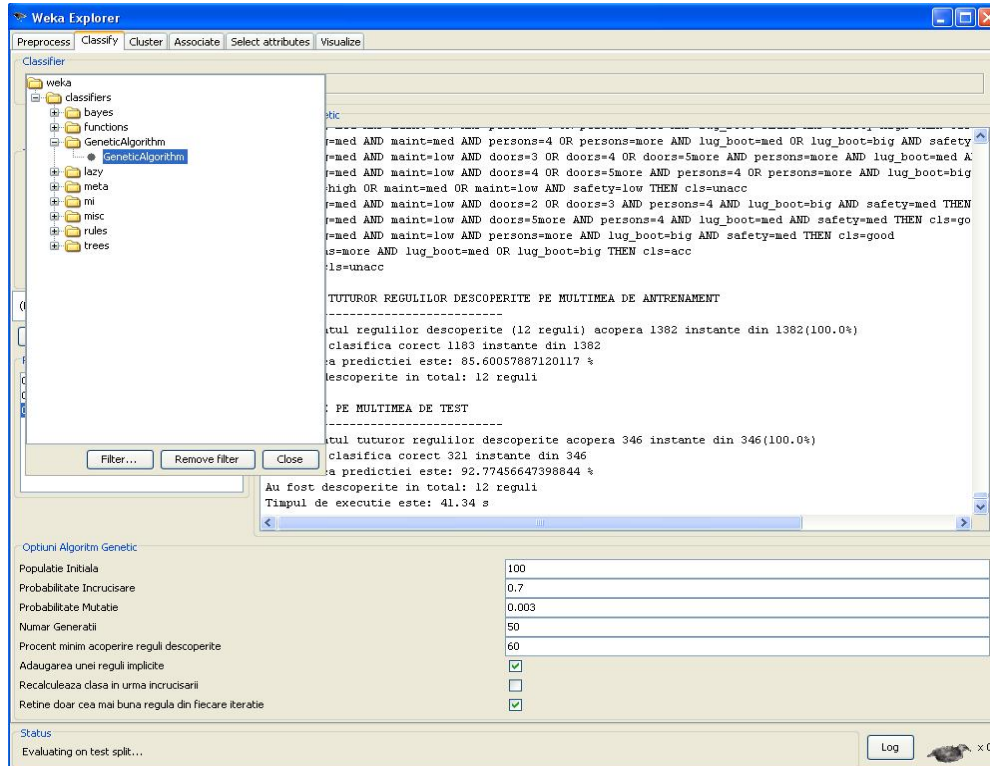


Fig. 4. 2 - Clasificatorul GeneticAlgorithm

Pentru a accesa instanțele din fișierul *ARFF* deschis s-a utilizat un obiect de tip *Instances*. Clasa *Instances* conține metode ce permit determinarea numărului de atribute, a valorilor atributelor, a numărului de valori pe care îl are atributul clasă, a valorilor acestuia, a numărului de instanțe, permit accesul la orice instanță, etc.

Pentru toate procesările algoritmului au fost create metode în clasa *ClassifierPanel*, iar algoritmul a fost implementat în metoda *run*.

4.6. Analiza coeficienților algoritmului și rezultate experimentale

Analiza comportamentului algoritmului s-a realizat în 3 situații: întâi s-au reținut în lista de reguli toate regulile descoperite de fiecare iterație a algoritmului (prin iterație se înțelege execuția algoritmului pentru numărul de generații stabilit). Al doilea experiment s-a realizat reținând în lista de reguli doar cea mai bună regulă descoperită de fiecare iterație. Ultimul experiment s-a realizat utilizând 3 funcții *fitness* și comparând rezultatele obținute cu ele. Comportamentul algoritmului a fost analizat utilizând în total 14 variante de parametri (8 variante de parametri pentru

primul experiment, 4 variante de parametri pentru cel de-al doilea experiment și 2 variante de parametri pentru ultimul experiment).

Algoritmul a fost executat pe 5 seturi de date reale obținute de pe UCI Machine Learning Repository. Aceste seturi de date sunt *Car*, *Zoo*, *Mushrooms*, *Breast Cancer*, *Heart Disease*. La fiecare rulare, 80 % din instanțe au fost folosite ca mulțime de antrenament și 20 % ca mulțime de test. Fiecare variantă de parametri a fost rulată de 5 ori pe fiecare set de date și tabelele de mai jos prezintă pentru toți parametrii valorile medii obținute la cele 5 rulări. Analiza făcută a condus la concluzii și recomandări cu privire la alegerea parametrilor de intrare, care sunt de asemenea prezentate în cadrul acestui subcapitol.

4.6.1. Utilizarea tuturor regulilor descoperite de fiecare iterație

În cazul utilizării tuturor regulilor descoperite de fiecare iterație se consideră că toate regulile descoperite după ce algoritmul genetic a fost rulat pentru numărul de generații stabilit sunt suficient de bune încât să fie considerate reguli de ieșire ale algoritmului și să fie adăugate listei regulilor descoperite. Configurațiile de parametri de intrare alese sunt prezentate în tabelul 4.3.

Varianta de parametri	1	2	3	4	5	6	7	8
PI	100	100	100	100	100	100	100	100
Pc	0.7	0.7	0.7	0.7	0.7	0.7	0.7	0.7
Pm	0.003	0.003	0.003	0.003	0.003	0.003	0.003	0.003
Număr generații	50	50	50	50	50	50	50	50
Procent minim acoperire reguli	0	0	0	60	60	60	60	80
Adăugare regulă implicită	nu	nu	da	nu	da	da	nu	da
Recalculare clasă	nu	da	nu	nu	nu	da	da	nu
Reține doar cea mai bună regulă din fiecare iterație	nu	nu	nu	nu	nu	nu	nu	nu

Tabel 4.3 - Parametrii de intrare

Rezultatele din tabelul 4.4 arată că acuratețea predicției pe mulțimea de antrenament poate lua valori cuprinse între 52.74 % și 100 % în funcție de modul de alegere a parametrilor. Varianta de parametri care a produs cele mai slabe rezultate în ce privește acuratețea pe setul de antrenare este varianta 3, iar cea care a produs cele mai bune rezultate este varianta 1.

Varianta de parametri	v1	v2	v3	v4	v5	v6	v7	v8
Car	91.10	89.52	73.22	82.04	98.14	99.61	88.57	97.58
Zoo	100.00	100.00	52.74	100.00	87.00	84.50	100.00	100.00
Mushrooms	100.00	100.00	57.04	100.00	93.12	99.40	100.00	94.55
Breast cancer	88.00	83.15	69.55	83.55	83.85	79.21	83.77	92.19
Heart disease	96.74	93.76	67.35	86.89	81.32	78.01	82.34	90.58

Tabel 4.4 - Acuratețea pe mulțimea de antrenament (%)

Varianta de parametri care a produs cele mai bune rezultate pe mulțimea de test este varianta 2, iar cea care a produs cele mai slabe rezultate este la fel ca și în cazul setului de antrenare varianta 3 (vezi tabelul 4.5).

Varianta de parametri	v1	v2	v3	v4	v5	v6	v7	v8
Car	94.44	94.62	88.43	90.89	97.80	99.88	93.83	100.00
Zoo	100.00	100.00	63.80	100.00	91.42	86.66	100.00	100.00
Mushrooms	100.00	100.00	88.51	100.00	92.65	99.74	100.00	93.47
Breast cancer	60.00	89.59	78.27	84.13	87.93	82.41	89.22	95.86
Heart disease	100.00	100.00	70.82	88.64	83.27	79.67	81.63	91.14

Tabel 4.5 - Acuratețea pe mulțimea de test (%)

Acoperirea regulilor descoperite arată capacitatea lor de a clasifica instanțe. Așa cum era de așteptat cea mai slabă acoperire o au regulile descoperite în variantele de parametri 1 și 2, adică după rularea celor 50 de generații ale algoritmului genetic. Cea mai bună acoperire o au variantele de parametri ce utilizează pe lângă regulile descoperite rulând algoritmul genetic și o regulă implicită. Nivelul de acoperire al regulilor descoperite poate să fie stabilit de utilizator la valoarea dorită. În cazul în care se impune execuția repetată a algoritmului genetic până la descoperirea unei liste de reguli care acoperă minim 60 % din setul de antrenare și nu este bifată opțiunea de determinare și adăugare a unei reguli implicite la lista de reguli (versiunile 4 și 7), acoperirea este cuprinsă între 60.28 și 88.00 de procente (vezi tabelul 4.6).

Varianta de parametri	v1	v2	v3	v4	v5	v6	v7	v8
Car	25.53	26.05	100.00	77.09	100.00	100.00	88.00	100.00
Zoo	7.25	7.75	100.00	63.00	100.00	100.00	64.00	100.00
Mushrooms	3.45	1.77	100.00	60.28	100.00	100.00	61.05	100.00
Breast cancer	3.86	30.70	100.00	64.39	100.00	100.00	60.97	100.00
Heart disease	9.83	11.74	100.00	67.11	100.00	100.00	71.16	100.00

Tabel 4.6 - Acoperire pe mulțimea de antrenament (%)

Pe mulțimea de test, cea mai slabă acoperire o au variantele 1 și 2 de parametri. Toate instanțele din mulțimea de test pot fi clasificate de variantele de parametri ce utilizează o regulă implicită. În variantele de parametri care nu utilizează o regulă implicită, dar în care algoritmul se execută repetat până la descoperirea unor reguli care acoperă minim 60 % din setul de antrenare, acoperirea pe mulțimea de test este cuprinsă între 23.76 % și 94.22 %. Procentul de acoperire al regulilor pe mulțimea de test este prezentat în tabelul 4.7.

Varianta de parametri	v1	v2	v3	v4	v5	v6	v7	v8
Car	60.35	59.54	100.00	88.61	100.00	100.00	94.22	100.00
Zoo	16.19	17.15	100.00	70.48	100.00	100.00	71.43	100.00
Mushrooms	6.71	6.53	100.00	33.60	100.00	100.00	23.76	100.00

Breast cancer	2.76	24.14	100.00	57.59	100.00	100.00	57.59	100.00
Heart disease	11.80	11.80	100.00	59.67	100.00	100.00	68.85	100.00

Tabel 4.7 - Acoperire pe mulțimea de test (%)

Timpul de execuție cel mai mic se obține în varianta 1 de parametri. Un timp de execuție mai ridicat au variantele de parametri care asigură acoperirea dorită de utilizator și care execută repetat algoritmul genetic. Variantele de parametri în care se recalculază clasa cromozomilor în urma încrucișării au un timp de execuție mai ridicat decât variantele similare care nu fac acest lucru. Tabelul 4.8 prezintă valorile medii obținute pentru timpul de execuție cu variantele de parametri alese.

Varianta de parametri	v1	v2	v3	v4	v5	v6	v7	v8
Car	5.83	16.61	5.90	37.60	34.58	144.46	112.52	39.07
Zoo	0.75	3.53	0.76	7.49	7.24	30.16	33.30	9.06
Mushrooms	117.90	151.09	73.18	2669.89	2981.55	9580.11	9833.57	4045.12
Breast cancer	1.37	3.23	1.42	21.89	17.38	21.11	18.92	27.06
Heart disease	1.92	6.88	1.96	7.59	7.59	19.81	18.30	6.48

Tabel 4.8 - Timpul de execuție (s)

Tabelul 4.9 prezintă numărul de reguli descoperite. Variantele de parametri care descoperă cele mai puține reguli sunt variantele 1 și 2, iar varianta care descoperă cele mai multe reguli, este varianta 8. Recalcularea clasei în urma încrucișării conduce la reguli care au o acuratețe mai bună dar au o acoperire mai slabă. În cazul în care se impune ca acoperirea să aibă o anumită valoare minimă, variantele de parametri în care clasa se recalculază vor descoperi mai multe reguli pentru a asigura aceeași acoperire ca și celelalte (comparație între numărul mediu de reguli obținut cu variantele 6 și 7 față de varianta 5).

Varianta de parametri	v1	v2	v3	v4	v5	v6	v7	v8
Car	2.20	2.40	4.20	18.20	15.40	25.00	15.40	26.80
Zoo	2.20	3.80	3.80	33.60	42.75	39.00	43.60	65.40
Mushrooms	3.80	2.40	4.00	180.33	186.67	311.20	346.00	493.40
Breast cancer	3.00	1.60	4.80	125.40	85.80	43.00	34.00	184.60
Heart disease	1.80	3.00	4.40	17.60	22.20	14.00	13.60	25.80

Tabel 4.9 - Numărul regulilor descoperite

Numărul de condiții din reguli este direct proporțional cu numărul regulilor descoperite și este prezentat în tabelul 4.10.

Varianta de parametri	v1	v2	v3	v4	v5	v6	v7	v8
Car	7.80	10.00	17.20	78.60	53.20	88.60	54.40	142.40
Zoo	35.00	57.80	41.00	511.80	626.50	557.20	640.40	929.20
Mushrooms	75.80	48.40	61.80	-	-	-	-	-
Breast cancer	27.20	10.20	33.20	892.40	710.00	334.60	392.60	1634.00
Heart disease	13.20	21.80	22.40	107.60	125.20	67.40	64.40	150.20

Tabel 4.10 - Numărul de condiții din regulile descoperite

Variantele de parametri 1 și 2 reprezintă ieșirea algoritmului genetic după rularea numărului de generații stabilit, fără nicio restricție suplimentară. După cum se observă în tabelele de mai sus, aceste variante descoperă într-un timp de execuție foarte bun, un număr mic de reguli, cu o acuratețe foarte bună. Acoperirea acestor reguli, atât pe setul de antrenare cât și pe cel de testare, este foarte slabă, ceea ce înseamnă că foarte multe instanțe nu pot fi clasificate doar cu aceste reguli, deci utilizarea acestor variante de parametri nu reprezintă o soluție.

Adăugarea unei reguli implicite la lista regulilor descoperite cu varianta 1 (varianta 3), rezolvă problema acoperirii instanțelor dar scade acuratețea predicției. Din punct de vedere al acurateții predicției varianta 3 este cea mai slabă. Dacă se analizează global toți parametrii se poate constata că sunt situații în care varianta 3 este varianta optimă. De exemplu în cazul setului de date *Mushrooms*, variantele 4, 5, 6, 7 și 8 produc rezultate ale acurateții predicției între 92.65 și 100 %, dar au un timp de execuție foarte mare (minim 2669 s) și descoperă minim 180 de reguli. Varianta 3 de parametri a descoperit în 73.18 secunde 4 reguli cu o acuratețe a predicției 88.51 % pe mulțimea de test. Deși acuratețea predicției este puțin mai slabă decât în cazul celorlalte variante, ceilalți parametri sunt mult mai buni și de aceea varianta 3 este de preferat.

Variantele 4, 5, 6 și 7 de parametri rulează în mod repetat algoritmul genetic până când se descoperă reguli ce acoperă minim 60 % din instanțele de antrenament. Această condiție face ca regulile descoperite în variantele 4 și 7 de parametri să acopere între 60.28 % și 88 % din setul de antrenare, și între 23.76 % și 94.22 din setul de testare. Variantele de parametri 5 și 6 acoperă toate instanțele din mulțimile de antrenament și test, pentru că fiind bifată opțiunea de adăugare a unei reguli implicite se analizează datele care au rămas neclasificate de reguli și se construiește o regulă care acoperă toate instanțele și le clasifică ca aparținând de clasa cea mai frecventă.

Acoperirea impusă de 60 % este o acoperire medie. Dacă regulile descoperite cu această restricție acoperă prea puțin din mulțimea de test, atunci poate mări restricția de 60 %, și / sau se poate adăuga o regulă implicită. Cu cât se impune ca regulile descoperite să acopere mai mult din setul de antrenare, cu atât timpul de execuție crește, la fel și numărul regulilor descoperite (comparație între v5 și v8).

În încercarea de a determina secvența optimă de parametri, dintre variantele care impun o acoperire de minim 60% (variantele 4, 5, 6 și 7), trebuie tratați împreună toți indicatorii. În ceea ce privește acuratețea predicției, rezultatele sunt bune și sunt comparabile în cazul variantelor menționate. Din punct de vedere al acoperirii variantele de parametri care sunt preferate sunt cele care acoperă toate instanțele, deci 5 și 6. Timpul de execuție este mult mai bun în cazul variantei 5, decât în cazul variantei 6, la fel și numărul regulilor descoperite este mult mai bun. Deci varianta optimă dintre cele menționate este varianta 5.

Varianta 8 de parametri este similară variantei 5 cu mențiunea că impune algoritmului să descopere și mai multe reguli decât versiunea 5, astfel încât regulile descoperite să acopere minim 80 % din instanțele de antrenare, spre deosebire de versiunea 5 care permite ca minim 60 % din instanțele de antrenare să fie acoperite de regulile descoperite. Ambele versiuni permit clasificarea tuturor instanțelor datorită regulii implicite care se adaugă listei regulilor descoperite. În ceea ce privește acuratețea predicției, ambele variante de parametri conduc la valori bune,

pe mulțimea de antrenament varianta 8 a obținut valori mai mari decât varianta 5 pe 4 din 5 seturi de date, iar pe mulțimea de test varianta 8 a obținut valori mai bune pentru acuratețea predicției pe toate seturile de date considerate. Timpul de execuție este mai mare în cazul utilizării variantei 8 decât în cazul variantei 5. Numărul regulilor descoperite și al condițiilor este mult mai mare în cazul variantei 8 decât în cazul variantei 5.

Scopul oricărui algoritm de clasificare este să găsească cât mai repede reguli, care să clasifice cât mai corect datele, care să acopere toate instanțele, care să fie cât mai puține, mai simple, mai ușor de înțeles și utilizat. De multe ori nu se pot asigura valori maxime simultan, pentru toți parametrii de ieșire și trebuie alese configurații de compromis pentru parametrii de intrare, ținând cont de datele care trebuiesc clasificate, astfel încât parametrii considerați mai importanți dintre parametrii de ieșire să aibă valori bune iar ceilalți să aibă valori satisfăcătoare. Dacă pe setul de date considerat, varianta de parametri 8 conduce la o acuratețe puțin mai bună decât varianta 5, dar la un timp de execuție și un număr de reguli mult mai mare, atunci varianta 5 este de preferat. Altfel este de preferat varianta 8.

4.6.2. Utilizarea celei mai bune reguli descoperite de fiecare iterație

În cazul în care se reține în lista de reguli doar cea mai bună regulă descoperită în urma rulării algoritmului genetic pentru 50 de generații, s-au considerat din nou variantele de parametri 4, 5, 6, 7 și 8 la care a fost modificată doar valoarea ultimului parametru (parametrul *Reține doar cea mai bună regulă din fiecare iterație* a primit valoarea da). Noile variante de parametri au fost numite 4', 5', 6', 7' și 8'. Tabelele de mai jos prezintă comparativ rezultatele obținute utilizând toate regulile descoperite respectiv doar cea mai bună regulă din fiecare rulare. Valorile cele mai mari obținute comparând versiuni similare de parametri (de ex. versiunea 4 cu versiunea 4' au fost subliniate).

Tabelul 4.11 prezintă acuratețea obținută pe mulțimea de antrenament. Toate variantele de parametri asigură o acuratețe a predicției între 76.22 și 100 %. Variantele de parametri ce rețin doar cea mai bună regulă conduc de 12 ori la rezultate mai bune decât cele care rețin toate regulile, iar variantele care rețin toate regulile, conduc de 6 ori la rezultate mai bune decât celelalte.

Varianta de parametri	v4	v4'	v5	v5'	v6	v6'	v7	v7'	v8	v8'
Car	82.04	90.97	98.14	91.28	99.61	98.00	88.57	85.44	97.58	95.67
Zoo	100.00	100.00	87.00	89.25	84.50	84.50	100.00	100.00	100.00	100.00
Mushrooms	100.00	100.00	93.12	99.12	99.40	94.40	100.00	100.00	94.55	99.91
Breast cancer	83.55	88.70	83.85	94.38	79.21	76.22	83.77	83.00	92.19	95.96
Heart disease	86.89	87.06	81.32	81.90	78.01	82.97	82.34	84.77	90.58	92.47

Tabel 4.11 - Acuratețea pe mulțimea de antrenament (%)

Rezultatele din tabelul 4.12 arată că variantele de parametri care folosesc toate regulile descoperite după 50 de generații au obținut de 10 ori rezultate mai bune decât celelalte, iar variantele care rețin doar cea mai bună regulă, au obținut

de 10 ori rezultate mai bune. Acuratețea pe mulțimea de test este cuprinsă între 79.67 și 100 % pe cele 5 seturi de date considerate.

Varianta de parametri	v4	v4'	v5	v5'	v6	v6'	v7	v7'	v8	v8'
Car	90.89	93.77	97.80	96.36	99.88	99.77	93.83	89.09	100.00	98.15
Zoo	100.00	100.00	91.42	94.28	86.66	84.76	100.00	100.00	100.00	100.00
Mushrooms	100.00	100.00	92.65	99.53	99.74	93.19	100.00	100.00	93.47	99.93
Breast cancer	84.13	93.58	87.93	95.51	82.41	78.96	89.22	86.36	95.86	97.58
Heart disease	88.64	84.68	83.27	81.63	79.67	83.93	81.63	83.63	91.14	91.47

Tabel 4.12. Acuratețea pe mulțimea de test (%)

În ceea ce privește acoperirea instanțelor din mulțimea de antrenament, de cinci ori s-a obținut o acoperire mai bună dacă s-au utilizat toate regulile și de 5 ori dacă s-a utilizat doar cea mai bună regulă. Acoperirea obținută pe datele de antrenare este prezentată în tabelul 4.13.

Varianta de parametri	v4	v4'	v5	v5'	v6	v6'	v7	v7'	v8	v8'
Car	77.09	87.06	100.00	100.00	100.00	100.00	88.00	77.39	100.00	100.00
Zoo	63.00	62.50	100.00	100.00	100.00	100.00	64.00	66.00	100.00	100.00
Mushrooms	60.28	62.86	100.00	100.00	100.00	100.00	61.05	62.45	100.00	100.00
Breast cancer	64.39	60.62	100.00	100.00	100.00	100.00	60.97	63.86	100.00	100.00
Heart disease	67.11	66.03	100.00	100.00	100.00	100.00	71.16	64.46	100.00	100.00

Tabel 4.13 - Acoperire pe mulțimea de antrenament (%)

Acoperirea pe mulțimea de test este mai bună dacă se reține cea mai bună regulă la fiecare rulare (de șase ori s-a obținut o acoperire mai bună dacă se reține cea mai bună regulă și de 4 ori s-a obținut o acoperire mai bună dacă se rețin toate regulile).

Varianta de parametri	v4	v4'	v5	v5'	v6	v6'	v7	v7'	v8	v8'
Car	88.61	91.73	100.00	100.00	100.00	100.00	94.22	83.93	100.00	100.00
Zoo	70.48	67.62	100.00	100.00	100.00	100.00	71.43	73.33	100.00	100.00
Mushrooms	33.60	41.24	100.00	100.00	100.00	100.00	23.76	37.54	100.00	100.00
Breast cancer	57.59	54.48	100.00	100.00	100.00	100.00	57.59	59.31	100.00	100.00
Heart disease	59.67	65.90	100.00	100.00	100.00	100.00	68.85	62.95	100.00	100.00

Tabel 4.14 - Acoperire pe mulțimea de test (%)

Timpul de execuție a fost în 10 situații mai mic atunci când s-a determinat și utilizat cea mai bună regulă din cele descoperite de algoritmul genetic și în 15 situații a fost mai mic când s-au utilizat toate regulile descoperite de algoritmul genetic. Valorile obținute pentru timpul de execuție sunt prezentate în tabelul 4.15.

Varianta de parametri	v4	v4'	v5	v5'	v6	v6'	v7	v7'	v8	v8'
Car	37.60	35.85	34.58	31.95	144.46	120.15	112.52	81.43	39.07	38.65

Zoo	7.49	8.19	7.24	7.62	30.16	31.87	33.30	33.97	9.06	33.97
Mushrooms	2669.89	3606.53	2981.55	4299.34	9580.11	7983.96	9833.57	11657.30	4045.12	5606.56
Breast cancer	21.89	28.41	17.38	40.52	21.11	19.07	18.92	18.03	27.06	38.72
Heart disease	7.59	5.60	7.59	6.34	19.81	24.32	18.30	19.60	6.48	9.68

Tabel 4.15 - Timpul de execuție

Numărul regulilor descoperite, la fel ca și numărul condițiilor este mai mic în cazul în care la fiecare rulare se adaugă în lista de reguli doar cea mai bună regulă descoperită (vezi tabele 4.16 și 4.17).

Varianta de parametri	v4	v4'	v5	v5'	v6	v6'	v7	v7'	v8	v8'
Car	18.20	9.60	15.40	9.20	25.00	10.40	15.40	6.20	26.80	10.40
Zoo	33.60	15.00	42.75	14.80	39.00	14.20	43.60	14.00	65.40	14.00
Mushrooms	180.33	80.60	186.67	99.60	311.20	87.40	346.00	127.75	493.40	156.00
Breast cancer	125.40	34.40	85.80	47.00	43.00	10.60	34.00	9.80	184.60	50.40
Heart disease	17.60	3.80	22.20	5.40	14.00	5.80	13.60	3.80	25.80	11.80

Tabel 4.16 - Numărul regulilor descoperite

Varianta de parametri	v4	v4'	v5	v5'	v6	v6'	v7	v7'	v8	v8'
Car	78.60	38.40	53.20	32.60	88.60	38.00	54.40	27.20	142.40	50.60
Zoo	511.80	240.40	626.50	214.40	557.20	199.00	640.40	215.00	929.20	215.00
Mushrooms	-	1700.80	-	2095.40	-	1827.20	-	2740.75	-	3311.5
Breast cancer	892.40	290.40	710.00	397.40	334.60	73.40	392.60	76.00	1634.00	421.4
Heart disease	107.60	21.60	125.20	26.20	67.40	29.00	64.40	20.80	150.20	68.2

Tabel 4.17. Numărul de condiții din regulile descoperite

	Număr valori maxime dacă se rețin toate regulile	Număr valori maxime dacă se reține doar cea mai bună regulă
Acuratețea pe mulțimea de antrenament	6	12
Acuratețea pe mulțimea de test	10	10
Acoperirea pe mulțimea de antrenament	5	5
Acoperirea pe mulțimea de test	4	6
Timpul de execuție	15	10
Numărul regulilor descoperite	0	25
Numărul de condiții din reguli	0	25

Tabel 4.18 - Numărul de valori maxime obținute când parametrul BR=true, respectiv când parametrul BR=false

Rezultatele prezentate în tabelul 4.18 arată comparativ în câte situații s-au obținut rezultate mai bune pe cele 5 seturi de date analizate dacă s-au considerat valide toate regulile descoperite după epuizarea celor 50 de generații, respectiv dacă s-a considerat validă doar cea mai bună regulă obținută după rularea celor 50 de generații. Rezultatele sunt comparabile în ceea ce privește acuratețea, acoperirea și

timpul de execuție. În schimb numărul regulilor descoperite și implicit și numărul condițiilor din reguli este mult mai mic în cazul în care se memorează doar cea mai bună regulă.

4.6.3. Analiza comportamentului algoritmului la utilizarea a diferite funcții fitness

Analiza comportamentului algoritmului în cazul utilizării a diferite funcții fitness s-a realizat utilizând parametrii care au dat cele mai bune rezultate bune în studiile precedente. Pe lângă funcția Q_1 utilizată deja în studiile precedente au mai fost testate încă două funcții fitness. Parametrii de intrare ai algoritmului au fost setați la valorile din tabelul 4.19. Rezultate obținute pentru parametrii de ieșire sunt prezentate în continuare.

Varianta de parametri	v8'	v9	v10
PI	100	100	100
Pc	0.7	0.7	0.7
Pm	0.003	0.003	0.003
Număr generații	50	50	50
Procent minim acoperire reguli	80	80	80
Adaugare regulă default	da	da	da
Recalculare clasă	nu	nu	nu
Reține doar cea mai bună regulă din fiecare rulare	da	da	da
Fitness	Q_1	Q_2	Q_3

Tabel 4.19 - Configurații de parametri

Q_1 este funcția fitness din relația 4.9, Q_2 este funcția fitness din relația 4.12 iar Q_3 este funcția fitness din relația 4.13.

$$Q_2 = PA \quad (4.12)$$

$$Q_3 = w_1 * PA + w_2 * CPH + w_3 * Senz \quad (4.13)$$

unde

- PA este acuratețea predicției, $PA = tp / (tp + fp)$
- CPH este comprehensibilitatea, $CPH = 1 - ANC / MNC$
- $ANC =$ numărul actual de condiții din antecedent
- $MNC =$ numărul maxim de condiții din antecedent
- $Senz$ este senzitivitatea, $Senz = tp / (tp + fn)$
- w_1, w_2 și w_3 sunt ponderi definite de utilizator, valorile alese pentru ele au fost $w_1 = 0.6, w_2 = 0.2, w_3 = 0.2$, punând astfel accentul mai mult pe acuratețea predicției și mai puțin pe comprehensibilitatea și senzitivitatea regulilor [ROB11b]

Tabelele de mai jos prezintă comparativ, acuratețea, timpul de execuție, numărul de reguli descoperite, numărul de condiții din reguli în cazul utilizării celor 3 funcții fitness propuse. Cele mai bune rezultate s-au obținut cu funcția fitness Q_1 ,

acesta este motivul pentru care a fost inclusă în algoritm. În ceea ce privește acuratețea predicției, atât pe mulțimea de antrenament cât și pe mulțimea de test, cu ajutorul funcției *fitness* Q_1 s-au obținut valori maxime pe 4 seturi de date, cu ajutorul funcției Q_2 s-a obținut o valoare maximă pe un set de date de test, iar cu ajutorul funcției Q_3 s-a obținut o valoare maximă pe un set de date de antrenare (vezi tabelele 4.17 și 4.18). Timpii de execuție cei mai buni s-au obținut cu funcția *fitness* Q_1 (valori minime pe 3 seturi de date, vezi tabelul 4.19). În ceea ce privește numărul de reguli descoperite și numărul de condiții din reguli, cele mai bune rezultate s-au obținut tot cu funcția *fitness* Q_1 (valori minime pentru numărul de reguli și numărul de condiții pe toate seturile de date). Cu ajutorul funcției *fitness* Q_3 s-au obținut cele mai bune rezultate pentru timpul de execuție pe două seturi de date (vezi tabele 4.20 și 4.21).

Parametrii	v8'	v9	v10
Car	95.67	88.98	90.78
Zoo	100.00	88.44	91.25
Mushrooms	99.91	90.23	
Breast cancer	95.96	91.01	91.07
Heart disease	92.47	89.75	93.08

Tabel 4.20- Acuratețea pe mulțimea de antrenament (%)

Parametrii	v8'	v9	v10
Car	98.15	93.49	94.60
Zoo	100.00	91.66	96.82
mushrooms	99.93	93.11	
Breast cancer	97.58	94.82	89.65
Heart disease	91.47	93.11	88.52

Tabel 4.21 - Acuratețea pe mulțimea de test (%)

Parametrii	v8'	v9	v10
Car	38.65	541.38	132.00
Zoo	33.97	11.52	7.17
Mushrooms	5606.56	6309.4	5501.32
Breast cancer	38.72	105.47	74.38
Heart disease	9.68	65.83	30.35

Tabel 4.22 - Timp execuție (s)

Parametrii	v8'	v9	v10
Car	10.40	207.75	58.67
Zoo	14.00	29.50	18.67
Mushrooms	156.00	289.4	254.7
Breast Cancer	50.40	145.25	106.33
Heart disease	11.80	71.00	36.33

Tabel 4.23 - Numărul regulilor descoperite

Parametrii	v8'	v9	v10
Car	50.60	1288.00	249.33
Zoo	215.00	442.00	264.33
Mushrooms	3311.5		
Breast cancer Cancer	421.4	1287.75	812.33
Heart disease	68.2	495.20	206.33

Tabel 4.24 - Numărul de condiții din reguli

4.6.4. Recomandări cu privire la alegerea parametrilor

Analiza realizată asupra parametrilor conduce la următoarele concluzii și recomandări cu privire la alegerea acestora:

- Variantele de parametri care utilizează doar cea mai bună regulă descoperită la fiecare iterație au performanțe comparabile cu cele care utilizează toate regulile descoperite, în ceea ce privește acuratețea și timpul de execuție dar net superioare în ceea ce privește numărul regulilor descoperite și numărului de condiții din reguli, de aceea se recomandă utilizarea acestei facilități
- Se recomandă impunerea execuției repetate a algoritmului până se descoperă reguli care acoperă între 60 și 80 % din instanțele de antrenare cumulat cu determinarea și adăugarea unei reguli implicite la lista de reguli, care să clasifice toate instanțele ce nu pot fi clasificate cu regulile descoperite prin rulare repetată a algoritmului genetic. În acest fel regulile descoperite vor putea clasifica toate instanțele. Procentul de acoperire minim impus depinde de acuratețea obținută, timpul de execuție, numărul de reguli descoperite etc. Dacă setul de date este mare și timpul de execuție este foarte mare se poate impune o acoperire mică, cu condiția ca acuratețea obținută să fie satisfăcătoare sau se poate lucra doar cu un eșantion reprezentativ al setului de date
- Recalcularea clasei în urma încrucișării poate conduce la reguli cu o acuratețe mai bună dar o acoperire mai slabă. Deci pentru a asigura aceeași acoperire va fi nevoie de mai multe reguli și de un timp de execuție mai mare dacă este bifată această opțiune. Se recomandă utilizarea pentru seturi de date care nu sunt foarte mari
- Pentru parametrii clasici ai unui algoritm genetic, cum ar fi numărul de indivizi din populație, probabilitățile de încrucișare și mutație, numărul de generații, se recomandă alegerea unor valori apropiate celor din studiile de mai sus, ținând cont de următoarele aspecte:
 - Rata încrucișării stabilește gradul de înnoire al populației. Cu cât rata încrucișării este mai mare cu atât șansa de a avea mai mulți indivizi noi în noua populație este mai mare. De obicei se alege între 0.6 și 1.
 - Cu cât probabilitatea mutației este mai mare cu atât riscul convergenței premature este redus, dar dacă rata mutației este prea mare algoritmul tinde să devină unul de căutare aleatoare. De obicei se aleg pentru mutație valori cuprinse între 0.001 și 0.1.

4.7. Comparație cu alți algoritmi

Întrucât *Weka* dispune de 117 algoritmi de clasificare și regresie, pentru a compara rezultatele obținute cu algoritmul genetic propus și cele obținute cu alți algoritmi s-a utilizat tot *Weka* și algoritmi *Naive Bayes*, *k-Nearest Neighbor* și *J48*. Algoritmii au fost rulați în condiții similare cu *AGR* propus, utilizând 80 % din instanțe la antrenament și 20 % din instanțe la testare. Comparația s-a realizat între parametrii de ieșire pe care îi oferă atât *AGR* cât și algoritmi aleși pentru comparație. *AGR* a fost rulat cu configurația considerată optimă (v8').

După cum se observă rezultatele obținute pentru acuratețe sunt mai bune în cazul algoritmului propus decât în cazul celorlalți algoritmi pe patru seturi de date

din 5. În cazul setului de date *breast cancer* se obțin rezultate mult mai bune cu algoritmul propus.

Algorimii	Acuratețea predicției			
	Naive Bayes	k-Nearest Neighbor	J48	AGR
Car	87.57	92.77	91.04	<u>98.15</u>
Zoo	90	100	90	<u>100.00</u>
Mushrooms	95.38	<u>100</u>	<u>100</u>	99.03
Breast Cancer	70.17	71.92	71.92	<u>97.58</u>
Heart disease	86.88	77.04	85.24	<u>91.47</u>

Tabel 4.25 - Acuratețea pe mulțimea de test

Weka furnizează pentru algoritmii implementați în ea doar timpul de execuție pentru construirea modelului nu și pentru testarea lui așa cum o face *AGR*. De aceea nu a putut fi realizată o comparație riguroasă la nivel de valori. În schimb se observă cu ușurință că timpul de execuție al lui *AGR* este mai mare decât al celorlalți algoritmi, acesta fiind dezavantajul algoritmului propus.

În ceea ce privește comparația între numărul de reguli și condiții descoperite de *AGR* și cel determinat de ceilalți algoritmi, aceasta se poate realiza doar între *J48* și *AGR*, întrucât modelele construite de *Naive Bayes* și *k-Nearest Neighbor* nu sunt compuse din reguli (vezi subcapitolele 1.3.2 și 1.3.3).

Algoritmul *J48* determină un set de reguli ce sunt reprezentate sub forma unui arbore. Numărul de reguli determinate de algoritm nu este furnizat direct, dar este egal cu numărul de frunze. Numărul de condiții nu este furnizat de *Weka*, el poate fi determinat prin numărare dar această operație poate consuma foarte mult timp dacă arborele este stufos.

Algoritmii	J48	AGR
Car	131	<u>10.40</u>
Zoo	17	<u>14.00</u>
Mushrooms	<u>25</u>	156.00
Breast Cancer	<u>4</u>	50.40
Heart disease	21	<u>11.80</u>

Tabel 4.26 - Număr reguli

Din punct de vedere al numărului de reguli descoperite există situații în care numărul de reguli descoperite cu *AGR* este mult mai mic decât cu *J48* și invers. Pe cele 5 seturi de date considerate, în 3 cazuri s-au obținut rezultate mai bune (mai puține reguli cu *AGR*) decât cu *J48*.

Comparațiile cu alți algoritmi arată că *AGR* este un algoritm performant, care poate determina reguli ce pot clasifica anumite seturi de date cu o acuratețe mai bună decât alți algoritmi, poate determina un set de reguli mai redus decât alți algoritmi, singura sa slăbiciune fiind timpul de execuție.

4.8. Concluzii

Algoritmii genetici sunt tehnici adaptive care pot fi folosite cu succes în rezolvarea unor probleme de căutare și optimizare complexe. O bază de date poate fi privită ca fiind un spațiu de căutare foarte mare și algoritmul genetic poate fi privit ca o strategie de căutare într-o bază de date. Ceea ce se caută sunt regulile de clasificare care vor constitui modelul clasificator. Capitolul de față prezintă structura generală a unui algoritm genetic, o scurtă sinteză cu privire la algoritmi genetici propuși pentru clasificarea datelor și propune un nou algoritm genetic în acest scop. Algoritmul genetic propus a fost implementat în *Java*, fiind adăugat instrumentului open source *Weka*. Mediul utilizat pentru dezvoltare a fost *Eclipse*. Testarea algoritmului și analiza comportamentului său la diferite valori ale coeficienților de intrare s-a realizat utilizând 5 seturi de date obținute din *UCI Machine Learning Repository*. În total au fost analizate 14 variante de parametri. Fiecare variantă de parametri a fost rulată de 5 ori pe fiecare set de date (la fiecare rulare au fost considerate 80 % din instanțele alese aleatoriu ca mulțime de antrenament și 20 % din instanțe ca mulțime de test), iar valorile medii obținute sunt prezentate în lucrare. Rezultatele obținute au fost foarte bune și în urma analizei s-au emis recomandări cu privire la alegerea parametrilor. Rezultatele obținute cu *AGR* au fost comparate cu cele obținute cu algoritmi clasici precum *Naive Bayes*, *k-Nearest Neighbor* și *J48*, iar rezultatele arată că algoritmul propus este competitiv și poate fi utilizat în studii de clasificare alături de acești algoritmi.

Contribuțiile aduse de acest capitol sunt următoarele:

- prezentarea succintă a structurii generale a unui algoritm genetic
- realizarea unei sinteze cu privire la algoritmi genetici utilizați la clasificarea datelor
- propunerea unui nou algoritm genetic pentru clasificarea datelor, care după cunoștințele autorului, diferă de ceilalți algoritmi genetici de clasificare prin următoarele:
 - permite recalcularea clasei în urma încrucișării
 - permite stabilirea unui procent minim de acoperire pentru regulile descoperite, rularea repetată a algoritmului până se atinge acest procent
 - permite adăugarea unei reguli implicite listei regulilor descoperite
 - permite adăugarea la lista de reguli descoperite doar a celei mai bune reguli obținută după rularea algoritmului pentru numărul de generații stabilit sau a tuturor regulilor descoperite
 - utilizează o funcție *fitness* diferită de cele utilizate de alți algoritmi
- îmbunătățirea tehnologiei de clasificare a datelor prin implementarea algoritmului genetic propus într-un instrument Open Source utilizat pe scară largă (*Weka*)
- analizarea comportamentului algoritmului la diferite valori ale parametrilor de intrare, și la utilizarea a diferite funcții *fitness* și emiterea unui set de recomandări cu privire la alegerea acestora.

Direcții de dezvoltare ulterioară sunt realizarea unei facilități care să-i permită utilizatorului atât să aleagă funcția *fitness* pe care dorește să o utilizeze în analiză, dintr-o listă de funcții *fitness* propuse, dar și să-și poată construi funcția *fitness* pe care dorește să o utilizeze putând combina termenii *tp*, *fp*, *tn*, *fn*

împreună cu ponderi și cu operatori de bază precum adunare, scădere, înmulțire împărțire, modulo, logaritm, etc. O altă direcție de dezvoltare este implementarea mai multor metode de selecție și oferirea posibilității de a alege metoda de selecție dorită într-un studiu.

5. Îmbunătățirea procesului de realizare a predicțiilor

5.1. Noțiuni introductive

Weka permite construirea cu ușurință a unui model clasificator pe baza datelor de antrenament. Pentru a construi acest model clasificator se poate folosi unul din multitudinea de algoritmi de clasificare pe care *Weka* îi oferă. Testarea clasificatorului construit se poate face fie pe setul de antrenare, fie pe un set de test obținut prin împărțirea setului inițial de date în set de antrenare și set de testare, fie pe un set de date bine precizat. Informații cu privire la rezultatele testării pot fi vizualizate în secțiunea *Classifier output*. Acestea sunt informații legate de rulare, modelul clasificator construit, rezultatele testării modelului, acuratețea detaliată pe clasă, matricea de confuzie, etc.

Dacă în urma analizării informațiilor legate de testarea modelului se constată că acesta a obținut o performanță bună, el poate fi utilizat pentru a realiza predicții.

Procesul de realizare a predicțiilor din *Weka* este foarte ne-sugestiv, se pot realiza predicții doar în mod indirect re-testând modelul pe setul pe care dorim să-l prezicem. Întrebarea "How do I make predictions with a trained model?" de pe site-ul lui *Weka*, rubrica întrebări frecvente întărește afirmația făcută.

Capitolul prezintă modul în care se pot realiza predicții cu *Weka*, soluțiile care au fost găsite pentru a îmbunătăți acest proces și îmbunătățirile care au fost aduse interfeței de clasificare a datelor. Pentru a îmbunătăți procesul de realizare a predicțiilor au fost realizate modificări în secțiunea *Classify a Explorer*-ului și acestea constau în transformarea acestei interfețe într-o interfață dinamică în funcție de setul de date analizat, interfață cu care se pot realiza predicții într-un mod foarte simplu și intuitiv. Alte modificări realizate au transformat această interfață într-o interfață prietenoasă, ușor accesibilă chiar și începătorilor în *Weka*.

5.2. Realizare de predicții cu *Weka*

Utilizarea unui model construit pentru a efectua predicții în *Weka* se poate realiza doar indirect, folosind opțiunea de re-testare a modelului pe un set de date bine precizat. (nu există nici o opțiune clară de realizare de predicții cu un model construit).

Realizarea predicțiilor necesită trecerea prin pașii următori:

- Se creează un nou fișier *ARFF* în care se pun instanțele care se doresc a fi prezise

- Întrucât este necesară completarea unei valori în dreptul clasei fiecărei instanțe din fișierul *ARFF*, se poate completa o clasă presupusă din cele existente sau se poate pune ? (vezi figura 5.1)

```

@relation breast-cancer;
@attribute age {'10-19','20-29','30-39','40-49','50-59','60-69','70-79','80-89','90-99'};
@attribute menopause {'140','40','premeno'};
@attribute tumor-size {'0-4','5-9','10-14','15-19','20-24','25-29','30-34','35-39','40-44','45-49','50-54','55-59'};
@attribute inv-nodes {'0-2','3-5','6-8','9-11','12-14','15-17','18-20','21-23','24-26','27-29','30-32','33-35','36-39'};
@attribute node-caps {'yes','no'};
@attribute deg-malig {'1','2','3'};
@attribute breast {'left','right'};
@attribute breast-quad {'left_up','left_low','right_up','right_low','central'};
@attribute irradiat {'yes','no'};
@attribute 'Class' {'no-recurrence-events','recurrence-events'};
@data
'40-49','premeno','15-19','0-2','yes','3','right','left_up','no','?'

```

Fig. 5. 1 – Fișierul *ARFF* cu instanțe de prezis

- Se bifează opțiunea *Output prediction* (vezi figura 5.1)

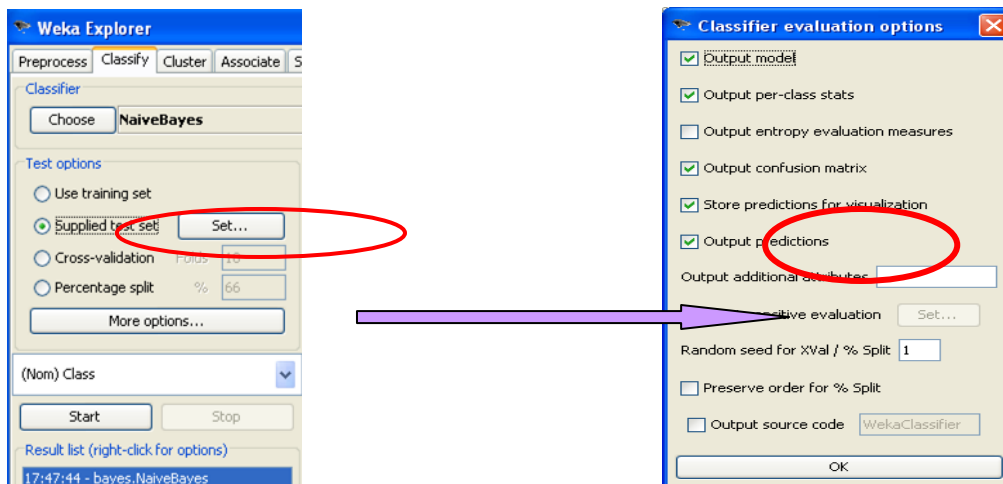


Fig.5.2 - Configurări necesare pentru realizarea de predicții

- Se încarcă fișierul *ARFF* care conține instanțele pentru care vrem să facem predicții cu ajutorul opțiunii *Supplied test set*, din panel-ul *Test options* (vezi figura 5.2)
- Se face click dreapta pe modelul rulat și se alege opțiunea de reevaluare a modelului pe setul curent de testare (vezi figura 5.3)

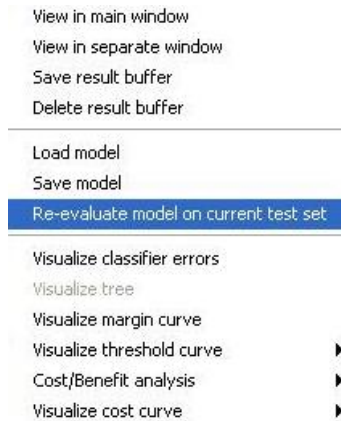


Fig. 5.3 - Reevaluarea modelului pe setul curent de testare

- modelul este de fapt retestat pe setul pe care noi vrem să-l prezicem, dar deoarece a fost bifată opțiunea *Output Prediction* în zona *Classifier output* se vor putea vizualiza valorile prezise (vezi figura 5.4)

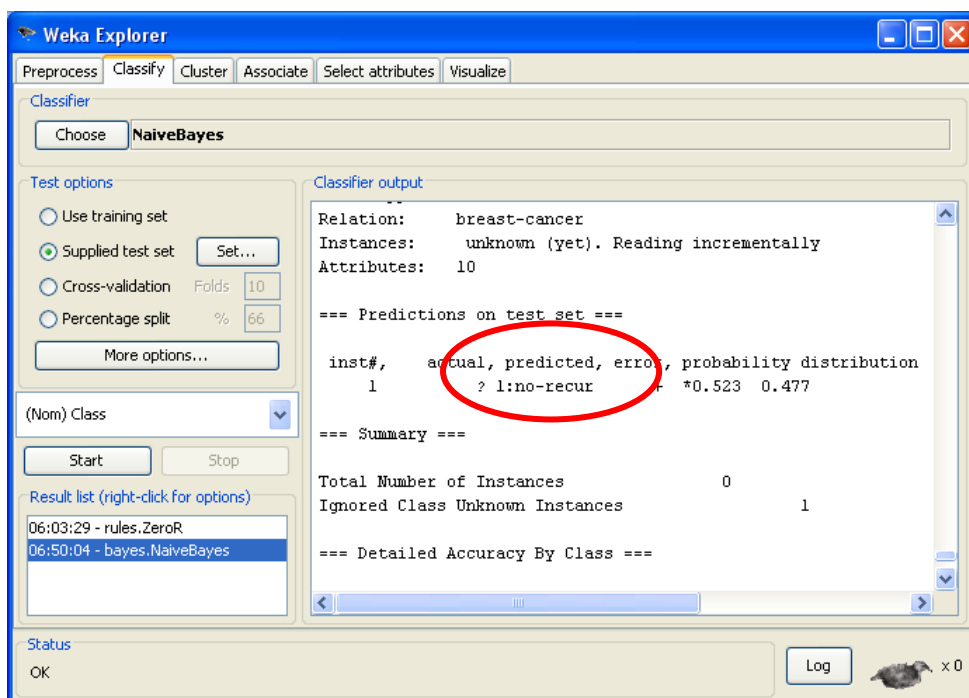


Fig. 5.4 - Vizualizarea valorilor prezise

Procesul prezentat este nesugestiv, predicțiile se realizează într-un mod indirect, printr-un artificiu, de fapt se retestează modelul pe setul pe care dorim să-l prezicem și de aceea întrebarea "Cum se pot face predicții cu un clasificator?" se găsește în rândul întrebărilor frecvent adresate de utilizatorii de *Weka*. Procesul de

realizare a predicțiilor prezentat mai sus a fost mult simplificat și transformat în unul rapid și intuitiv, datorită faptului ca *Weka* este un instrument Open Source. Îmbunătățirile aduse sunt prezentate în continuare.

5.3. Îmbunătățirea procesului de realizare a predicțiilor

Pentru a îmbunătăți acest proces, interfața de clasificare a datelor a fost modificată adăugând un nou panel, numit *Prediction of an instance*. Acest panel are un conținut dinamic în funcție de setul de date analizat. Panelul conține pentru fiecare set de date un număr de label-uri egal cu numărul de attribute din setul de date, iar pentru fiecare atribut numeric se va afișa o casetă de text care permite introducerea valorii numerice, pentru fiecare atribut nominal în interfață va exista un control *JComboBox* cu valorile nominale posibile, din care utilizatorul va putea selecta valoarea nominală dorită. Utilizatorul introduce în interfața dinamică valorile instanței pe care dorește să o prezică, apasă comanda *Predict* și clasa prezisă pentru acea instanță este afișată imediat lângă butonul *Predict* (vezi figura 5.5) [ROB10a].

O altă modificare care a fost adusă secțiunii de clasificare a datelor a urmărit transformarea acesteia într-o secțiune prietenoasă. Pentru aceasta, panelul *Classifier output* care afișează multe informații care nu sunt accesibile decât utilizatorilor avansați de *Weka*, a fost înlocuit cu panelul *Test Results*, care afișează sub forma unui grafic 3d acuratețea clasificatorului construit. Accesul la informațiile afișate în secțiunea *Classifier Output* nu a fost eliminat, dar utilizatorul care dorește să vizualizeze aceste informații trebuie să apese un nou buton numit *Advanced information* care a fost introdus în interfață (vezi figurile 5.5 și 5.6).

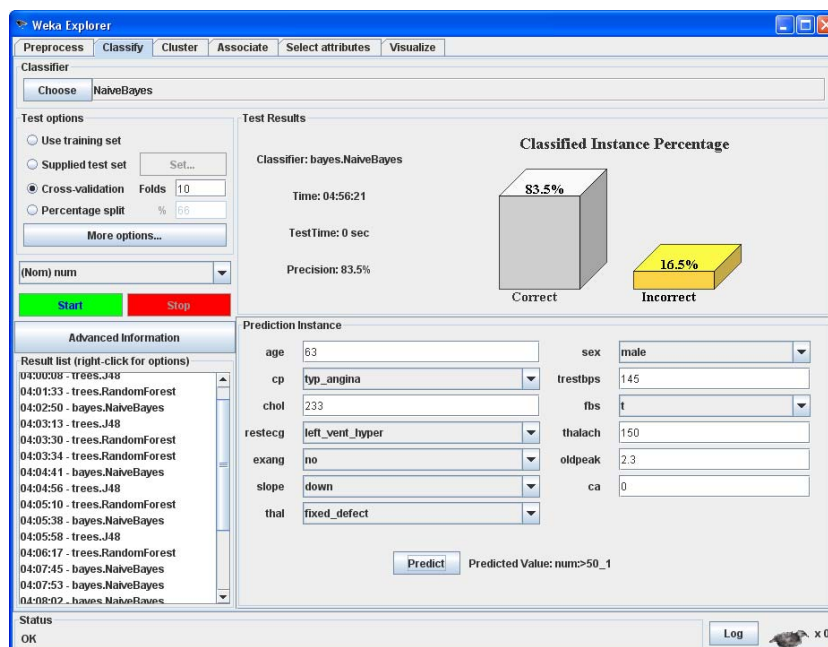


Fig. 5.5 - *Weka* cu interfață dinamică de realizare a predicțiilor

În continuare sunt prezentate o parte din informațiile afișate în secțiunea *Classifier Output* cu scopul de a argumenta modificarea făcută.

Panel-ul *Classifier Output* din cadrul *GUI*-ului *Classify* afișează în mod implicit informații legate de rulare (numărul de instanțe, numărul de atribute, modul de testare al modelului), modelul clasificator, rezultatul testării (procentul de instanțe clasificate corect, procentul de instanțe clasificate greșit, *Kappa statistic*, *Mean absolute error*, *Root mean squared error*, *Relative absolute error*, *Root relative squared error*), acuratețea detaliată pe clasă (*True Positive Rate*, *False Positive Rate*, *Precision*, *Recall*, *F-Measure*, *Class*) și matricea de confuzie [TAN05].

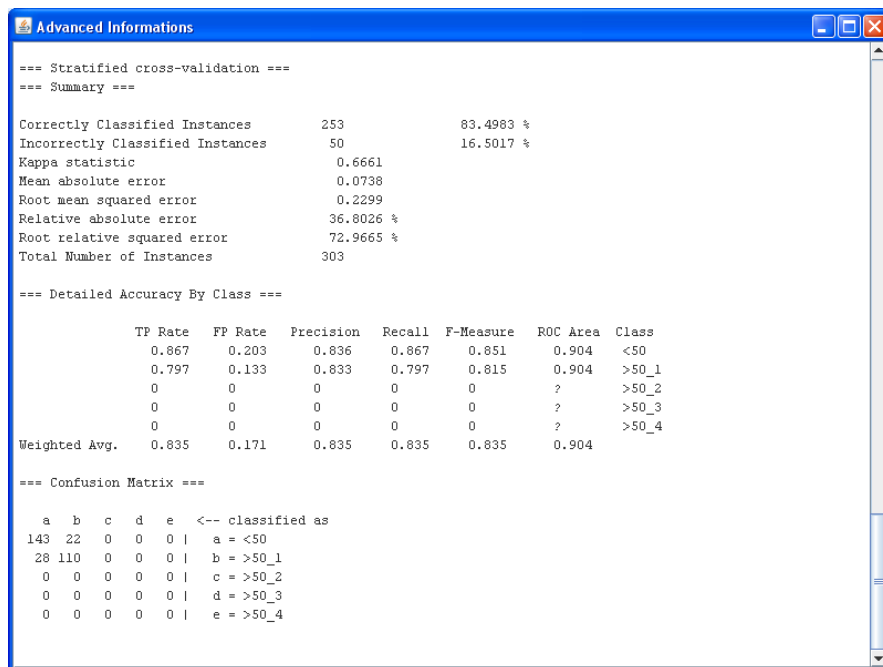
Ieșirile selectate în mod implicit pentru vizualizare sunt modelul, statisticile pe clasă, matricea de confuzie și stocarea predicțiilor pentru vizualizare (vezi figura 5.2). Utilizatorul are posibilitatea să selecteze output-urile care le dorește pentru vizualizare cu ajutorul comenzii *More Option*. Pe lângă ieșirile menționate mai pot fi afișate măsurile de evaluare ale entropiei, predicțiile, atributele adiționale, evaluarea cost-senzitivă, codul sursă, etc.

În cazul în care nu se selectează nici un output (sunt de-bifate toate opțiunile din interfața *Classifier evaluation options* (vezi fig. 5.2)) setul minimal de informații afișate în secțiunea *Classifier Output* cuprinde informații legate de rulare și rezultatul testării exprimat atât prin procentul de instanțe clasificate corect respectiv greșit cât și prin măsuri precum *Kappa statistic*, *Mean absolute error*, *Root mean squared error*, *Relative absolute error*, *Root relative squared error* (vezi figura 5.6).

Semnificația măsurilor menționate este următoarea:

- *Kappa statistic*- măsoară acordul predicției cu clasa adevărată, valoarea 1.0 înseamnă acord complet [BOU08]
- *Mean absolute error* este o măsură utilizată pentru a analiza cât de aproape au fost predicțiile de valorile reale. Se calculează împărțind suma valorilor absolute a erorilor la numărul de predicții
- *Root mean squared error* - este o modalitate de a cuantifica diferența dintre valorile reale și de cele estimate. Eroarea este cantitatea prin care valoarea estimată diferă de cea reală [LEH98]. Mărimea se calculează extrăgând radical din suma pătratelor erorilor raportate la numărul de predicții [WIT05].
- *Relative absolute error* – suma valorilor absolute a erorilor, raportată la suma valorilor absolute a diferențelor față de medie) [WIT05]
- *Root relative squared error*. – Se calculează împărțind suma pătratelor erorilor la suma pătratelor diferențelor față de medie. [WIT05]

Majoritatea măsurilor afișate în mod implicit, dar și a celor afișate în cazul în care nu se selectează nici un output, sunt de folos doar utilizatorilor experimentați care le cunosc însemnătatea și știu cum să le interpreteze. În rândul întrebărilor frecvente de pe site-ul *Weka* se află foarte multe întrebări legate de aceste măsuri, modul lor de calcul și însemnătatea lor. În consecință autorul consideră că este mai bine ca accesul la aceste informații să se realizeze ca urmare a unei solicitări și în mod implicit să fie prezentate informațiile care pot fi interpretate chiar și de utilizatori neexperimentați. În locul afișării măsurilor de mai sus a fost afișat un grafic 3d care redă acuratețea predicției, care poate fi înțeles și interpretat chiar și de începători. Cei interesați să vizualizeze măsurile menționate, trebuie să apese butonul nou introdus *Advanced Information* [ROB10c].



```

Advanced Informations

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      253      83.4983 %
Incorrectly Classified Instances    50      16.5017 %
Kappa statistic                    0.6661
Mean absolute error                 0.0738
Root mean squared error            0.2299
Relative absolute error             36.8026 %
Root relative squared error        72.9665 %
Total Number of Instances          303

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area  Class
          0.867   0.203   0.836     0.867   0.851     0.904    <50
          0.797   0.133   0.833     0.797   0.815     0.904    >50_1
          0         0         0         0         0         ?        >50_2
          0         0         0         0         0         ?        >50_3
          0         0         0         0         0         ?        >50_4
Weighted Avg.   0.835   0.171   0.835     0.835   0.835     0.904

=== Confusion Matrix ===

  a  b  c  d  e  <-- classified as
143 22  0  0  0 | a = <50
 28 110 0  0  0 | b = >50_1
  0  0  0  0  0 | c = >50_2
  0  0  0  0  0 | d = >50_3
  0  0  0  0  0 | e = >50_4

```

Fig.5.6 - Informațiile din secțiunea Classifier Output

Modificările realizate simplifică mult acțiunile pe care un utilizator de *Weka* trebuie să le facă pentru a realiza predicții cu un model construit și testat. Totodată procesul a devenit unul clar, simplu și intuitiv. Interfața a fost transformată în una simplă, prietenoasă care încurajează utilizarea lui *Weka*.

Aspecte legate de implementare

Weka realizează diverse prelucrări asupra unor seturi de date. În *Weka* setul de date este implementat de clasa *weka.core.Instances* [HAL09]. Fiecare instanță este constituită dintr-un număr de atribute care pot fi nominale, numerice sau șiruri de caractere. Reprezentarea externă a unei clase *Instances* este un fișier ARFF. Algoritmii de clasificare din *Weka* sunt derivați din clasa abstractă *weka.classifiers.Classifier* care conține la rândul ei metode pentru a genera un model de clasificare, pentru a evalua un model pe un set de date sau pentru a genera o distribuție de probabilitate [WIT05]. Preprocesarea datelor este un pas important pentru algoritmi de învățare automată. Un suport folositor pentru etapa de preprocesare este disponibil în pachetul *weka.filters* care este constituit din clase cu ajutorul cărora se pot face modificări asupra setului de date.

Arhitectura aplicației *Weka* este cea din fig 1. *Weka* este alcătuită dintr-o multitudine de pachete. Pachetele sunt organizate ierarhic începând de la pachetul principal numit *weka*. Orice pachet poate să conțină la rândul lui alte pachete, fișiere cu cod sursă Java sau pe acestea amândouă în același timp. Fișierele în care au fost realizate modificări sunt: *GuiChoser.java*, *ClasifierPanel.java*, *Explorer.java*, *PreprocessPanel.java*. Fișierele *SimpleBarChart.java* și *SpringUtilities.java* au fost adăugate.

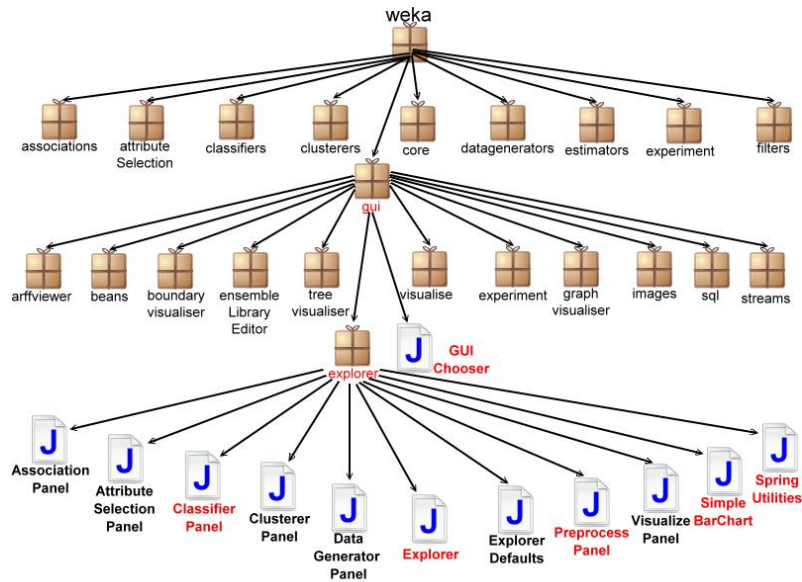


Fig. 5.7 – Arhitectura lui Weka

Testarea

Testarea aplicației *Weka* extinsă s-a făcut construind clasificatori și realizând predicții pe 3 seturi de date medicale obținute de pe *UCI Machine Learning Repository*. Acestea sunt *Ljubljana Breast Cancer*, *Heart Disease*, *Diabetes*. Informații cu privire la seturile de date menționate se găsesc în introducere. Algoritmii cu care au fost construite modelele clasificatoare au fost *Naive Bayes* [WU07], *Random Forest* și *J48*. Întrucât modelele care s-au obținut cu acești algoritmi au obținut rezultate bune, ele au fost utilizate pentru a realiza predicții. Predicțiile au fost realizate cu ajutorul noii interfețe dinamice adăugate în *Weka*.

Tabelul 5.1 arată rezultatele care s-au obținut pentru acuratețea predicției folosind tehnica de testare cross validation. Figura 5.5 ilustrează rezultatele obținute cu algoritmul *NaiveBayes* pe setul *Cleveland Heart Disease*. Figura 5.8 ilustrează rezultatele obținute pe setul de date *Ljubljana breast cancer* de algoritmul *J48* și rezultatele obținute de algoritmul *Random Forest* pe setul de date *Diabetes*. Aplicația îmbunătățită a funcționat corect.

Test data	Naive Bayes	J48	Random Forest
Ljubljana breast cancer	71.68	75.52	69.23
Diabetes	76.3	73.83	73.83
Cleveland heart disease	83.5	77.56	81.52

Tabel 5.1 - Acuratețea predicției

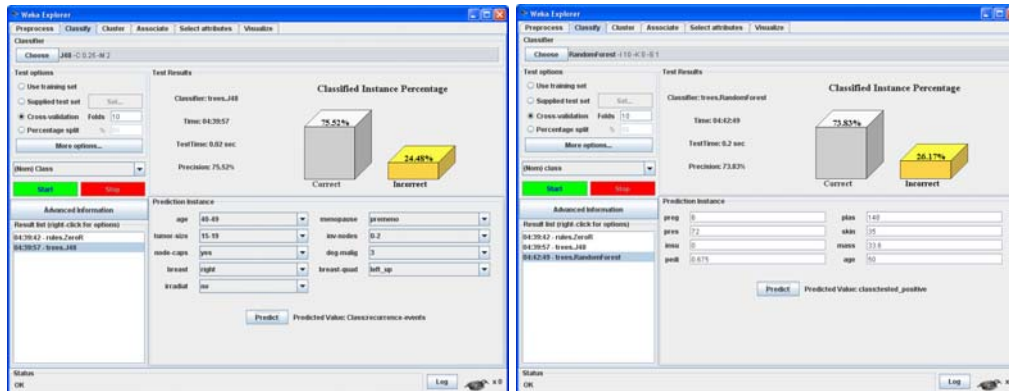


Fig. 5.8 - Predicții pe baza seturilor de date Ljubljana breast cancer (stânga) și diabetes (dreapta)

5.4. Concluzii

În cadrul acestui capitol a fost tratată problema realizării de predicții cu ajutorul unui model construit și testat de *Weka*. A fost prezentat modul în care clasificatorii pot fi utilizați pentru a realiza predicții. A fost concepută și implementată o interfață dinamică care conține controale în acord cu setul de date analizat și care ușurează și clarifică mult procesul de realizare a predicțiilor. Interfața de clasificare a fost transformată într-o interfață mult mai prietenoasă, fiind adăugat în ea un grafic 3d care arată acuratețea modelului și afișând indicatorii "complicați" din *Weka* într-o nouă fereastră *Advanced information*. Testarea modificărilor făcute s-a realizat construind modelele clasificatoare pe trei seturi de date medicale, obținute de pe *UCI Machine Learning Repository* și realizând predicții cu ele și cu ajutorul interfeței dinamice adăugate. Rezultatele testării au fost foarte bune.

Contribuțiile aduse în cadrul acestui capitol sunt:

- Prezentarea detaliată a modului în care *Weka* poate fi utilizat pentru a realiza predicții
- Concepția unei interfețe dinamice în funcție de setul de date analizat care permite realizarea cu ușurință a predicțiilor
- Prezentarea sintetică a mărimilor de ieșire furnizate de *Weka* în urma testării modelului
- Concepția unor soluții de transformare a interfeței de clasificare a datelor într-o interfață mult mai simplă, inteligibilă și prietenoasă, prin vizualizarea grafică a celei mai importante măsuri (acuratețea predicției) obținute în urma testării modelului și prin afișarea informațiilor din secțiunea *Classifier output* într-o nouă fereastră *Advanced Information*.
- Modificările menționate mai sus, realizate în unul din cele mai populare instrumente de data mining din lume, contribuie la îmbunătățirea tehnologiei de clasificare a datelor.

6. Analiza și clasificarea datelor legate de nașteri

6.1. Noțiuni introductive

Capitolul prezintă un studiu realizat pe datele legate de nașterile care au avut loc în anul 2010, la *Clinica de Obstetrică - Ginecologie Bega, Timișoara*. Setul de date disponibil pentru analiză cuprinde date despre 2326 de nașteri. Pentru fiecare naștere s-au reținut în total 19 atribute care conțineau date despre mamă, date despre nou-născut și date legate de intervenții medicale. Capitolul prezintă setul inițial de date, felul în care acestea au fost preprocesate pentru a putea fi și ulterior analizate cu *Weka*, analiza statistică realizată pe fiecare atribut, modelele clasificatoare construite și utilizarea acestora pentru a realiza predicții.

Pentru realizarea analizei prezentate în cadrul acestui capitol au fost folosite instrumentele dezvoltate și prezentate în capitolele 2, 3, 4 și 5 și noțiunile teoretice prezentate în toate capitolele anterioare.

Scopul analizei acestor date a fost de a extrage noi cunoștințe și informații folositoare din ele. Analiza s-a realizat în două etape: într-o primă etapă s-a realizat o analiză statistică la nivel de atribut, iar în a doua etapă s-a urmărit construirea unor modele clasificatoare din datele disponibile.

Analiza statistică realizată la nivel de atribut a permis determinarea unor informații statistice valoroase cum ar fi procentul din nașterile naturale în care este necesară realizarea unei tăieturi în zona pelvisului, care este procentul de femei care nasc natural, respectiv prin cezariană în zilele noastre, care este numărul mediu de ore de travaliu, etc.

Modelele clasificatoare permit realizarea de predicții și pot fi utilizate ca instrumente de suport ale actului medical. Modelul de clasificare construit permite estimarea cu o acuratețe de aprox 85 %, a indicelui apgar al nou născutului, în funcție de parametrii cunoscuți despre mamă, despre făt și în funcție de intervențiile medicale care se realizează asupra mamei și / sau fătului pentru a ajuta nașterea. Indicele apgar reprezintă o apreciere a funcțiilor vitale și a capacității de adaptare a fătului la condițiile din mediul extrauterin și este de dorit ca valoarea acestui indice să fie maximă pentru fiecare nou născut.

Valoarea medicală a modelului de clasificare construit este următoarea: medicii pot analiza înainte de naștere care este intervalul în care se estimează că indicele apgar al viitorului nou născut va lua valori completând la intrarea modelului datele mamei, ale fătului, și date legate de intervențiile care intenționează se le realizeze asupra mamei (de genul cezariană, sau naștere naturală cu epiziotomie, etc) pentru a ajuta nașterea. Modelul va furniza la ieșire intervalul în care se estimează că indicele apgar al nou născutului se va situa (1 din 5 intervale). Dacă se constată că acest interval este prea mic se poate verifica dacă se fac alt fel de

intervenții (cum ar fi: se renunță la nașterea naturală în favoarea cezarienei), sau dacă se ajută la îmbunătățirea parametrilor mamei sau a fătului (de pildă se face un masaj care să ajute la ajungerea fătului în poziția normală pentru nașterea naturală) ce valoare va lua indicele apgar după aceste intervenții. Modelul construit nu este unul decizional, este un model consultativ pentru medici, care să-i ajute în luarea deciziilor în așa fel încât fătul să aibă o trecere cât mai lină și mai bună în mediul extrauterin, contorizată printr-un indice apgar cât mai mare.

6.2. Setul de date

Datele pe care a fost realizată analiza provin de la *Clinica de Obstetrică - Ginecologie Bega, Timișoara*. Setul de date inițial cuprinde date despre 2326 de nașteri, realizate la acest spital în anul 2010. Datele despre nașteri au fost stocate într-un sheet al unui fișier Microsoft Excel.

Pentru fiecare naștere au fost disponibile următoarele informații:

- Un număr de ordine
- Luna nașterii
- Numele mamei
- Vârsta mamei
- Localitatea domiciliu
- Mediul - urban sau rural
- Gesta - a câta sarcină este
- Para - al câtelea copil este
- Numărul de săptămâni gestație
- Prezența - Prezența este așezarea fătului la ieșirea din uter, mai precis partea din corpul copilului care va ieși prima. Valorile posibile pentru acest câmp sunt: cefalică, pelviană, facială și transversală. Prezența **cefalică** este cea normală, fătul are coloana vertebrală paralelă cu cea a mamei și capul în jos cu bărbia în piept. Prezența **pelviană** înseamnă că fătul iese cu picioarele sau fundul înainte. Prezența **facială** este când fătul privește înainte, fața lui va ieși prima. În cazul în care coloana fătului nu este paralelă cu coloana mamei (fătul este poziționat oblic în burtă) prezența este **transversală**.
- Indicele apgar - Imediat după naștere, chiar în primele 60 de secunde după expulzie, în sala de travaliu, se face o apreciere a stării de sănătate a nou-născutului, prin aprecierea funcțiilor vitale și a capacității de adaptare la condițiile din mediul extrauterin. Astfel, simultan cu acordarea primelor îngrijiri, medicul neonatolog va nota starea clinică și comportamentul nou-născutului, cuantificând funcțiile vitale cu ajutorul scorului sau indicelui Apgar. Indicele apgar are valori cuprinse de la 0 la 10.
- Sexul copilului
- Greutatea
- Dacă a fost născut prin naștere naturală sau cezariană
- Videx - coloană care indică dacă fătului i-a fost aplicată pe cap sau nu o căciuliță de metal care să ajute nașterea lui naturală
- Motivul pentru care a fost indicată nașterea prin cezariană
- Epiziotomia - coloană care indică dacă s-a făcut sau nu tăietura în perineu pentru a ajuta nașterea naturală (0 - nu s-a făcut, 1 - s-a făcut)

- Numărul de ore de travaliu
- EMP = extracție manuală de placentă, valoarea 0 indică că nu s-a făcut o astfel de intervenție, iar valoarea 1 că s-a făcut

6.3. Preprocesarea datelor

Preprocesarea datelor a fost realizată printr-o serie de acțiuni:

- Au fost eliminate coloanele neimportante pentru studiu, precum numărul de ordine, numele mamei, localitatea de domiciliu, recomandarea pentru cezariană care era șir de caractere. Coloana luna nașterii a fost păstrată pentru analiza statistică, dar a fost eliminată din rândul coloanelor folosite la construirea modelului
- Coloana număr săptămâni de gestație conținea de regulă numărul exact al săptămânii de gestație (de exemplu săptămâna 40) dar la un număr de 178 de persoane nu se cunoștea săptămâna cu exactitate și valorile completate erau de genul 37/38, 38/39, etc. Acestea valori au fost înlocuite cu valori medii (de exemplu săptămâna 37.5).
- Au fost eliminate instanțele care conțineau valori lipsă, sau valori introduse greșit. Valori lipsă sau introduse greșit au fost în coloanele:
 - luna nașterii avea valoarea F pentru o persoană
 - vârsta pentru o persoană avea valoarea 0 și pentru 10 persoane nu a fost completată
 - pentru 6 persoane coloana număr săptămâni gestație avea fie valoarea 0 fie era necompletată
 - coloana mediul urban – rural nu a fost completată pentru două persoane
 - indicele apgar – nu a fost completat pentru 3 nou-născuți și pentru 4 nou-născuți în loc să fie completată nota în coloana indice apgar a fost completat motivul pentru care acesta nu a putut fi acordat și anume, datorită faptului că doi dintre aceștia au fost născuți acasă iar alții 2 pe ambulanță
 - coloana ore travaliu avea valori lipsă pentru 4 persoane
- Au fost corectate greșelile pentru care s-a putut face acest lucru. De exemplu:
 - Greutatea consemnată la naștere pentru doi copii a fost de 34000 g, respectiv 34450 g. Am presupus că a fost adăugat din greșeală un zero la finalul greutății reale și am eliminat această valoare.
 - Coloana ore travaliu avea date numerice, cuprinse între 0 și 30, dar la două persoane era scris 20 min. Coloana trebuie să fie uniformă, să conțină fie numărul de ore de travaliu fie numărul de minute de travaliu, prin urmare am transformat totul în ore și am înlocuit cele două valori 20 min cu 0.33 ore
- Au fost date coloanelor nume sugestive. De exemplu în fișierul primit litera G a fost folosită pentru a denumi atât coloana ce indică numărul de sarcini (gesta) cât și cea care indică greutatea nou născutului. Noile denumiri alese pentru cele două coloane au fost Gesta și Greutate.

Prelucrarea fișierului Excel a fost realizată cu ajutorul opțiunii *Filter*. A fost pus filtru pe fiecare coloană și verificată fiecare coloană în parte ca să conțină doar date valide.

În urma eliminării instanțelor care conțineau date lipsă sau greșite a rămas un număr de 2277 instanțe valide (din 2324).

Următoarea etapă a preprocesării a constat în construirea fișierului *ARFF*, cu datele valide rămase și încărcarea acestuia în *Weka*. Pentru această etapă a fost folosită aplicația *Arff Convertor* [ROB10b], propusă în capitolul 2. Figura 6.1 permite vizualizarea datelor încărcate din fișierul xls în *Arff Convertor*.

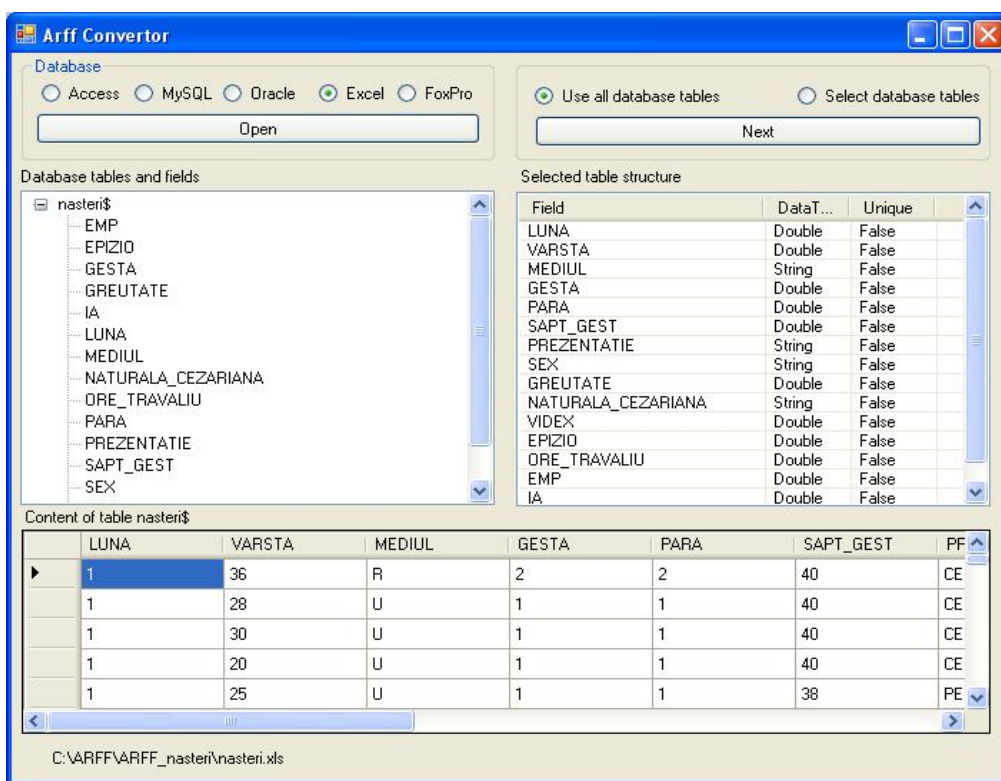


Fig. 6.1 - *Arff convertor* cu datele despre nașteri

Arff convertor a analizat valorile din fiecare câmp și a propus ca acestea să fie ori numerice, ori categoriale (vezi figura 6.2), utilizatorul fiind cel care alege cu ajutorul unor controale de tip ComboBox tipul acestor atribute în fișierul *ARFF*.

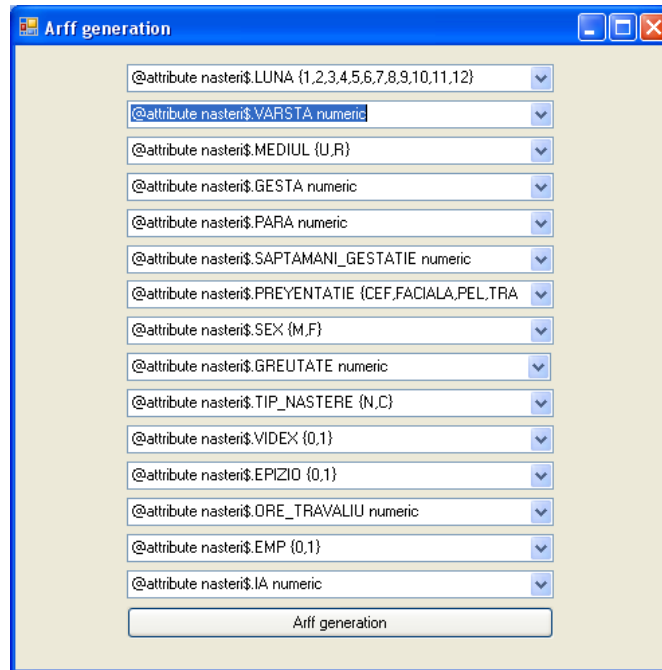


Fig. 6.2. - Alegerea tipului atributelor din fișierul ARFF

Figura 6.3 arată fișierul ARFF generat cu ajutorul lui Arff Convertor.

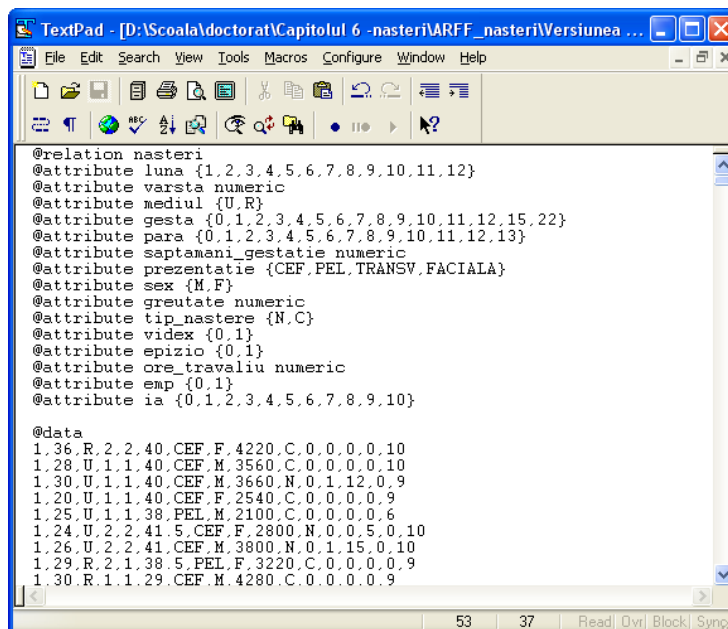


Fig. 6.3 - Fișierul ARFF generat

Preprocesarea datelor a continuat cu ajutorul lui *Weka*, dar abia după ce s-a încercat fără succes construirea unor modele clasificatoare rulând algoritmi clasici din *Weka* precum Naive Bayes [ROB11d], J48, k-Nearest Neighbour pe setul de date din figura 6.3. Atributele numerice au fost discreditate fiind împărțite în 5 intervale egale.

6.4. Analiza statistică

În continuare datele preprocesate au fost încărcate în *Weka* și cu ajutorul acestuia a fost realizată o analiză statistică la nivelul fiecărui atribut (vezi figura 6.4). Numărul de instanțe este de 2277.

După cum se observă în figura 6.5, în luna mai s-au înregistrat cele mai multe nașteri, (245 de nașteri) ceea ce înseamnă că luna septembrie este luna optimă pentru a plănuși un copil.

Vârsta femeilor gravide este cuprinsă între 12 și 45 de ani. Cele mai multe femei rămân însărcinate în jurul vârstei de 28 de ani. 70 de femei minore au rămas însărcinate și dintre acestea 17 au vârste cuprinse între 12 și 15 ani. La polul opus 48 de femei cu vârste mai mari de 40 de ani au rămas însărcinate.

Din punct de vedere al mediului din care provin 1493 din femeile care au născut în 2010 la Spitalul Bega provin din mediul urban și 784 provin din mediul rural.

În ceea ce privește indicatorul gesta (numărul de sarcini avute până în prezent) se observă că majoritatea femeilor au avut o sarcină (932 femei) sau două (646 femei), dar există și situații extreme de femei care au avut peste 10 sarcini (24 de femei).

Indicatorul Para arată că majoritatea covârșitoare a femeilor sunt la prima (1313 femei) sau a doua naștere (646 femei). Ca de obicei situațiile extreme nu lipsesc, avem 8 femei care au născut peste 10 copii.

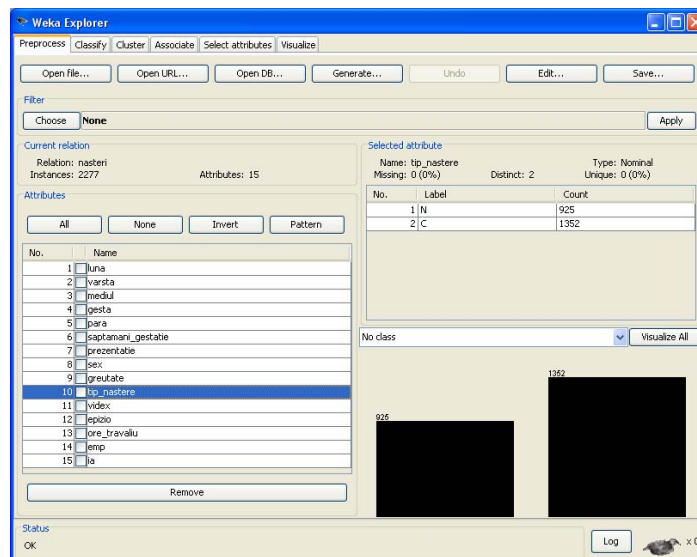


Fig. 6.4 - Analiză statistică la nivelul fiecărui atribut cu *Weka*

Numărul de săptămâni gestație este o valoare numerică cuprinsă între 19 și 44.5, iar valoarea medie este de 38.442 cu o deviație standard de 2.438 %.

Prezența normală, cefalică este cea mai răspândită (2115 cazuri), urmată de cea pelviană 156 de cazuri, prezența transversală cu 6 cazuri, iar cea mai rară este prezența facială (un singur caz).

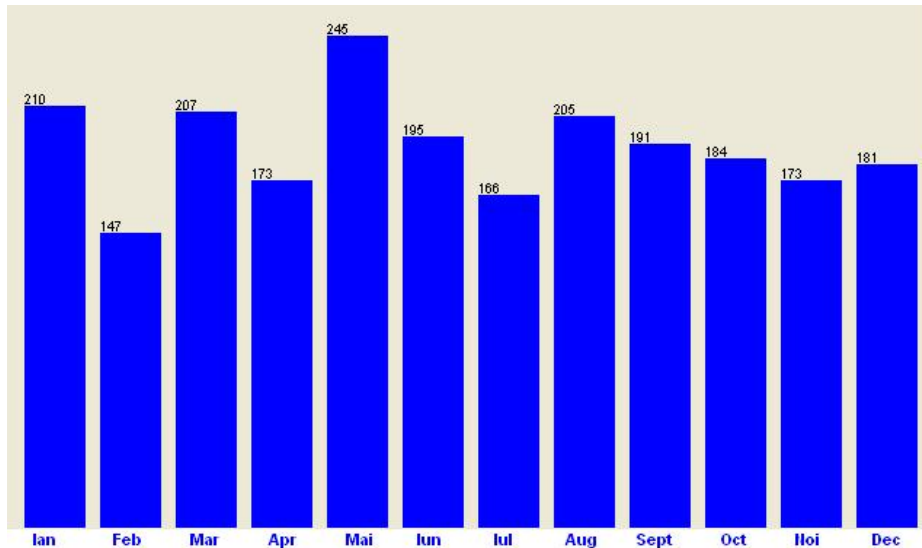


Fig. 6.5 - Numărul de nașteri pe lunile anului

În ceea ce privește sexul nou-născuților se observă un rezultat surprinzător, au fost născuți mai mulți băieți (1188) decât fete (1089) la Spitalul Bega în 2010.

Greutatea medie a nou născuților este de 3182 de grame, cu o deviație standard de 589 de grame, dar au fost născuți 2 bebeluși care aveau peste 5000 de grame (vezi figura 6.6).

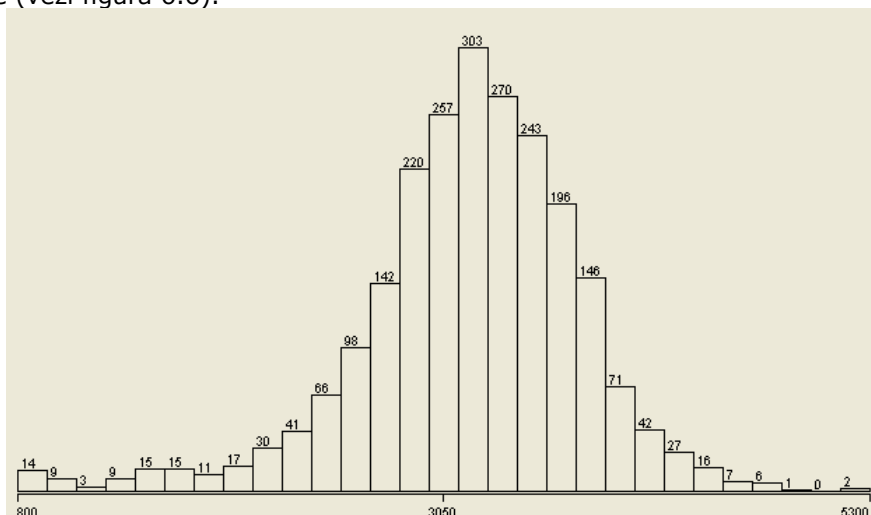


Fig. 6.6 - Greutatea nou născuților

Analizând numărul de femei care nasc normal respectiv prin cezariană în zilele noastre, am putea spune ca nașterea naturală este tot mai puțin recomandată sau dorită de femei întrucât în 2010 au născut natural 925 de femei și prin cezariană 1352 de femei.

Căciuliță de metal (videx) s-a pus doar pe capul a 20 de bebeluși din cei 925 care au fost născuți natural.

Epiziotomia indică dacă a fost necesară realizarea unei tăieturi în zona perineului la femeile care au născut natural. După cum se observă în 665 de cazuri din 925, (aprox. 70 %) a fost necesară această intervenție.

Numărul de ore de travaliu este cuprins între 0 și 30 de ore. De obicei 0 ore de travaliu sunt la nașterile prin cezariană.

Extracție manuală de placentă a fost făcută doar la 7 femei care au născut natural.

Nota primită la naștere este de obicei 10 (990 de cazuri) sau 9 (811 de cazuri), dar după cum se vede în figura 6.7 toate notele cuprinse între 0 și 10 au fost acordate, inclusiv cele foarte slabe.

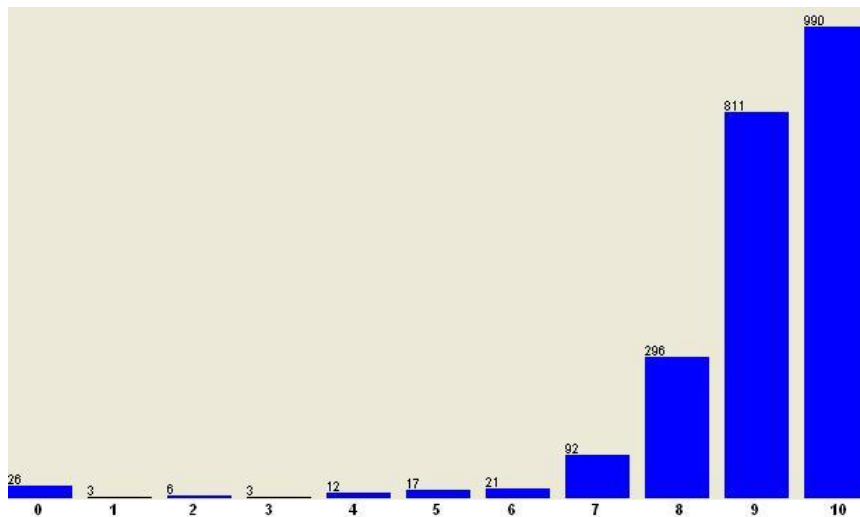


Fig. 6.7 - Indicele apgar

6.5. Clasificarea

În continuare s-a urmărit construirea unor modele clasificatoare, cu ajutorul algoritmilor propuși dar și cu ajutorul unor algoritmi clasici, implementați în *Weka*, care să clasifice cu o acuratețe cât mai mare datele propuse. Până la determinarea soluției optime au fost necesare multiple reveniri la etapa de preprocesare. Luna nașterii a fost păstrată în fișierul *arff* inițial pentru analiza statistică dar ulterior a fost eliminată, cu ajutorul lui *Weka*, fiind considerată nerelevantă pentru studiu. Setul de date considerat întâi pentru analiză, are 2277 de instanțe și atributele de tipurile de mai jos. Au fost construiți clasificatori cu ajutorul algoritmilor clasici *Naive Bayes*, *J48*, *k-Neares-Neghbour* din *Weka*. Acuratețea acestor clasificatori este extrem de slabă și este prezentată în tabelul 5.1. Întrucât au fost considerate atribute numerice, rularea algoritmului *Ant-Miner*, *Ant-r-Miner*, *AGR* nu a fost

posibilă. Acuratețea predicției prezentată în tabelele de mai jos este o acuratețe medie obținută în urma rulării fiecărui algoritm de 5 ori pe fiecare set de date. La fiecare rulare s-au utilizat 80 % din instanțe pentru antrenare și 20 % din instanțe pentru testare.

Atributele din fișierul *nasteri_1.arff*:

```
@attribute varsta numeric
@attribute mediul {U,R}
@attribute gesta {0,1,2,3,4,5,6,7,8,9,10,11,12,15,22}
@attribute para {0,1,2,3,4,5,6,7,8,9,10,11,12,13}
@attribute saptamani_gestatie numeric
@attribute prezentatie {CEF,PEL,TRANSV,FACIALA}
@attribute sex {M,F}
@attribute greutate numeric
@attribute tip_nastere {N,C}
@attribute videx {0,1}
@attribute epizio {0,1}
@attribute ore_travaliu numeric
@attribute emp {0,1}
@attribute ia {0,1,2,3,4,5,6,7,8,9,10}
```

Algoritm	Acuratete
Naive Bayes	46.81
J48	39.78
k-Nearest Neighbour	40.21

Tabel 5.1 – acuratețea predicției pe setul de date *nasteri_1.arff*

Întrucât clasificatorii obținuți au fost extrem de slabi, s-a revenit în etapa de preprocesare și s-au discreditat atributele numerice *varsta*, *saptamani_gestatie*, *greutate*, *ore_travaliu* cu ajutorul lui *Weka*. Discretizarea s-a făcut împărțind fiecare din atributele numerice menționate în 5 intervale egale (fiecare atribut numeric a devenit atribut nominal cu 5 valori posibile), după care s-a reîncercat construirea unui model clasificator. Algoritmii clasici din *Weka* au fost rulați din nou, dar datorita faptului că toate atributele obținute sunt acum nominale au putut fi rulați și algoritmii *Ant-Miner*, *Ant-r-Miner* și *AGR*. Rezultatele obținute sunt prezentate în tabelul 5.2.

Atributele din fișierul *nasteri_2.arff*:

```
@attribute varsta {"(-inf-18.6]","(18.6-25.2]","(25.2-31.8]","(31.8-38.4] ",
"(38.4-inf)"}
@attribute mediul {U,R}
@attribute gesta {0,1,2,3,4,5,6,7,8,9,10,11,12,15,22}
@attribute para {0,1,2,3,4,5,6,7,8,9,10,11,12,13}
@attribute saptamani_gestatie {"(-inf-24.1]","(24.1-29.2]","(29.2-33.4]","
"(34.3-39.4]","(39.4-inf)"}
@attribute prezentatie {CEF,PEL,TRANSV,FACIALA}
@attribute sex {M,F}
@attribute greutate {"(-inf-1700]","(1700-2600]","(2600-3500]","
```

```

"\(3500-4400]\",""\(4400-inf)\")
@attribute tip_nastere {N,C}
@attribute videx {0,1}
@attribute epizio {0,1}
@attribute ore_travaliu {"\(-inf-6]\",""\(6-12]\",""\(12-18]\",""\(18-24]\",""\(24-inf)\")
@attribute emp {0,1}
@attribute ia {0,1,2,3,4,5,6,7,8,9,10}

```

Algoritm	Acuratețe
Naive Bayes	47.91
J48	45.93
k-Nearest Neighbour	43.73
AGR	54.6
Ant-Miner	43.3
Ant-r-Miner	45.98

Tabel 5.2 - Acuratețea predicției pe setul de date nasteri_2.arff

Chiar dacă algoritmul AGR propus a obținut cel mai bun rezultat, și acesta este neacceptabil de slab, prin urmare încercările de prelucrare a datelor și de rulare a algoritmilor pe datele prelucrate au continuat. S-a revenit la etapa de preprocesare, atributele gesta și para au fost transformate în atribute numerice și apoi discreditate. Intervalele obținute după discreditare pot fi vizualizate mai jos. Rezultatele obținute rulând algoritmi pe aceste date discreditate pot fi vizualizate în tabelul 5.3

Atributele din fișierul *nasteri_3.arff*:

```

@attribute varsta {"\(-inf-18.6]\",""\(18.6-25.2]\",""\(25.2-31.8]\",""\(31.8-38.4]\",""\(38.4-inf)\")
@attribute mediul {U,R}
@attribute gesta {"\(-inf-4.4]\",""\(4.4-8.8]\",""\(8.8-13.2]\",""\(13.2-17.6]\",""\(17.6-inf)\")
@attribute para {"\(-inf-2.6]\",""\(2.6-5.2]\",""\(5.2-7.8]\",""\(7.8-10.4]\",""\(10.4-inf)\")
@attribute saptamani_gestatie {"\(-inf-24.1]\",""\(24.1-29.2]\",""\(29.2-34.3]\",""\(34.3-39.4]\",""\(39.4-inf)\")
@attribute prezentatie {CEF,PEL,TRANSV,FACIALA}
@attribute sex {M,F}
@attribute greutate {"\(-inf-1700]\",""\(1700-2600]\",""\(2600-3500]\",""\(3500-4400]\",""\(4400-inf)\")
@attribute tip_nastere {N,C}
@attribute videx {0,1}
@attribute epizio {0,1}
@attribute ore_travaliu {"\(-inf-6]\",""\(6-12]\",""\(12-18]\",""\(18-24]\",""\(24-inf)\")
@attribute emp {0,1}
@attribute ia {0,1,2,3,4,5,6,7,8,9,10}

```

Algoritm	Acuratețe
Naive Bayes	48.35
J48	44.61
k-Nearest Neighbour	42.85

AGR	67.1
Ant-Miner	43.21
Ant-r-Miner	46.34

Tabel 5.3 - Acuratețea predicției pe setul de date *nasteri_3.arff*

Rezultatele obținute arată că în urma reducerii împrăștierii valorilor, pentru atributele *gesta* și *para* modelele construite au o acuratețe mai bună, dar nesatisfăcătoare în continuare. În următoarea încercare s-a discretizat atributul *clasă*, indicele *apgar* și atributele *gesta* și *para* au fost lăsate cu valori nominale la fel ca în fișierul *nașteri_2.arff*. Rezultatele obținute sunt prezentate în tabelul 5.4.

Atributele din fișierul *nasteri_4.arff*:

```
@attribute varsta {"\(-inf-18.6]\","(18.6-25.2]\","(25.2-31.8]\","(31.8-38.4]\","(38.4-inf)\"}
@attribute mediul {U,R}
@attribute gesta {0,1,2,3,4,5,6,7,8,9,10,11,12,15,22}
@attribute para {0,1,2,3,4,5,6,7,8,9,10,11,12,13}
@attribute saptamani_gestatie {"\(-inf-24.1]\","(24.1-29.2]\","(29.2-34.3]\","(34.3-39.4]\","(39.4-inf)\"}
@attribute prezentatie {CEF,PEL,TRANSV,FACIALA}
@attribute sex {M,F}
@attribute greutate {"\(-inf-1700]\","(1700-2600]\","(2600-3500]\","(3500-4400]\","(4400-inf)\"}
@attribute tip_nastere {N,C}
@attribute videx {0,1}
@attribute epizio {0,1}
@attribute ore_travaliu {"\(-inf-6]\","(6-12]\","(12-18]\","(18-24]\","(24-inf)\"}
@attribute emp {0,1}
@attribute ia {"\(-inf-2]\","(2-4]\","(4-6]\","(6-8]\","(8-inf)\"}
```

Algoritm	Acuratețe
Naive Bayes	78.78
J48	78.43
k-Nearest Neighbour	75.49
AGR	84.86
Ant-Miner	78.88
Ant-r-Miner	79.01

Tabel 5.4 - Acuratețea predicției pe setul de date *nasteri_4.arff*

Rezultatele din tabelul 5.4 arată că problema principală pentru care modelele construite pe seturile de date *nasteri_1.arff*, *nasteri_2.arff* și *nasteri_3.arff*, au avut rezultate foarte slabe a fost numărul mare de valori al atributului *clasa*. Odată ce numărul de valori al acestui atribut a fost înjumătățit prin transformarea valorilor în intervale, toți algoritmi au obținut rezultate substanțial mai bune.

O ultima încercare de a obține rezultate mai bune a constat în revenirea în etapa de preprocesare și reducere numărul de valori al atributelor *gesta* și *para* din setul de date *nasteri_4.arff*, prin discretizare.

Atributele din fișierul *nasteri_5.arff*:

```
@attribute varsta {"\(-inf-18.6]\"", \"(18.6-25.2]\"", \"(25.2-31.8]\"",
 \"(31.8-38.4]\"", \"(38.4-inf)\""}
@attribute mediul {U,R}
@attribute gesta {"\(-inf-4.4]\"", \"(4.4-8.8]\"", \"(8.8-13.2]\"", \"(13.2-17.6]\"",
 \"(17.6-inf)\""}
@attribute para {"\(-inf-2.6]\"", \"(2.6-5.2]\"", \"(5.2-7.8]\"", \"(7.8-10.4]\"",
 \"(10.4-inf)\""}
@attribute saptamani_gestatie {"\(-inf-24.1]\"", \"(24.1-29.2]\"",
 \"(29.2-34.3]\"", \"(34.3-39.4]\"", \"(39.4-inf)\""}
@attribute prezentatie {CEF,PEL,TRANSV,FACIALA}
@attribute sex {M,F}
@attribute greutate {"\(-inf-1700]\"", \"(1700-2600]\"", \"(2600-3500]\"",
 \"(3500-4400]\"", \"(4400-inf)\""}
@attribute tip_nastere {N,C}
@attribute videx {0,1}
@attribute epizio {0,1}
@attribute ore_travaliu {"\(-inf-6]\"", \"(6-12]\"", \"(12-18]\"", \"(18-24]\"", \"(24-inf)\""}
@attribute emp {0,1}
@attribute ia {"\(-inf-2]\"", \"(2-4]\"", \"(4-6]\"", \"(6-8]\"", \"(8-inf)\""}

```

Algoritm	Acuratete
Naive Bayes	77.14
J48	78.46
k-Nearest Neighbour	77.14
AGR	82.67
Ant-Miner	78.7
Ant-r-Miner	79.53

Tabel 5.5 - Acuratețea predicției pe setul de date *nasteri_5.arff*

Cel mai bun rezultat l-a obținut modelul creat cu algoritmul AGR (acuratețe a predicției de 84.86 %), pe setul de date 4.arff.

6.6. Predicții

Setul de reguli determinat poate fi utilizat în domeniul medical pentru a prezice cu o acuratețe de aprox 85 % care va fi intervalul în care se va situa indicele apgar al nou-născutului, în funcție de datele mamei, datele estimative ale copilului (greutatea) și în funcție de intervențiile care se fac mamei.

Pentru realizarea predicțiilor s-a utilizat și se recomandă utilizarea pe viitor a versiunii lui *Weka*, ce dispune de o interfață dinamică [ROB10a] cu care se pot realiza predicții cu ușurință (vezi figura 6.8 și capitolul 5) [ROB10c].

O altă utilitate a setului de reguli descoperit de algoritmul AGR a fost utilizarea lui ca set de date de intrare pentru un sistem care le transpune în sintaxă ARDEN [HEE05]. Pentru ca setul de reguli descoperit să poată fi utilizat ca set de date de intrare pentru sistemul menționat a fost necesară punerea acestuia în

format XML. Pentru a genera acest fișier XML a fost dezvoltată o aplicație în C care primește la intrare un fișier text cu setul de reguli descoperit de AGR și creează un fișier XML cu structura din figura 6.9.

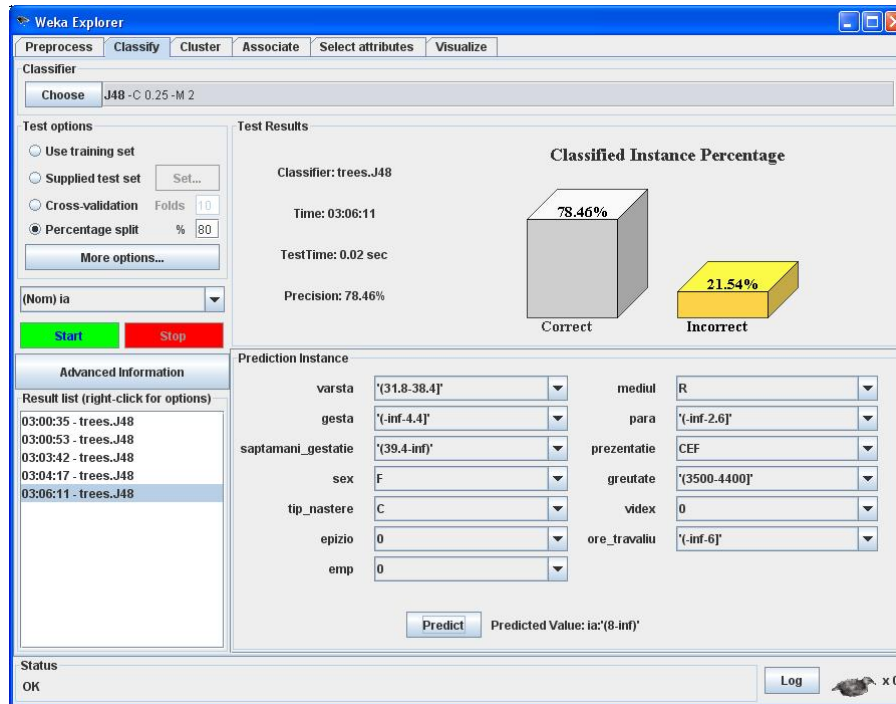


Fig. 6.8 – Realizarea de predicții

O regulă obținută de algoritmul AGR are structura de mai jos.

IF atribut1=valoare1 or atribut1=valoare2 and atribut2=valoare1 and...THEN clasa=calsa2

Prelucrările pe care aplicația dezvoltată le face sunt următoarele:

- deschide fișierul de intrare
- scrie în fișierul de ieșire versiunea XML și codificarea
- scrie în fișierul de ieșire tag-ul rădăcină <nașteri>
- cât timp nu s-a parcurs tot fișierul de intrare:
- citește un grup de litere (grupurile de litere sunt separate prin spații)
- dacă grupul de litere este egal cu "IF" atunci scrie tag-ul
<regula id=contor>
- dacă grupul de litere conține caracterul "=" și în regula curentă NU s-a ajuns la THEN, atunci acesta se împarte în două grupuri cu ajutorul delimitatorului "=", primul grup

reprezintă atributul iar al doilea valoarea lui: <atribut>valoare</atribut>

- dacă grupul de litere curent este egal cu AND sau cu OR se adaugă în fișierul XML tag-ul <operator >valoarea_operatorului</operator >

- dacă grupul de litere este egal cu THEN, se reține acest lucru

- dacă grupul de litere conține caracterul "=" și în regula curentă s-a ajuns la THEN, atunci înseamnă că s-a ajuns la clasă, se împarte aceasta în două grupuri cu ajutorul delimitatorului "=", primul grup reprezintă clasa, iar al doilea valoarea ei și se adaugă în fișierul XML tag-ul <atribut>valoare</atribut> și se închide tag-ul </regula>

- dacă s-a parcurs tot fișierul de intrare, atunci închide tag-ul rădăcină </nasteri>

```

<?xml version="1.0" encoding="utf-8"?>
<nasteri>
  <regula id="1">
    <mediul>U</mediul>
    <operator>AND</operator>
    <gesta>'(-inf-4.4]}'</gesta>
    <operator>AND</operator>
    <para>'(-inf-2.6]}'</para>
    <operator>AND</operator>
    <saptamani_gestatie>'(34.3-39.4]}'</saptamani_gestatie>
    <operator>AND</operator>
    <prezentatie>CEF</prezentatie>
    <operator>AND</operator>
    <greutate>'(2600-3500]}'</greutate>
    <operator>AND</operator>
    <tip_nastere>C</tip_nastere>
    <operator>AND</operator>
    <videx>0</videx>
    <operator>AND</operator>
    <epizio>0</epizio>
    <operator>AND</operator>
    <ore_travaliu>'(-inf-6]}'</ore_travaliu>
    <operator>AND</operator>
    <emp>0</emp>
    <ia>'(8-inf)'</ia>
  </regula>
  <regula id="2">
    <varsta>'(31.8-38.4]}'</varsta>
    <operator>AND</operator>
  </regula>
</nasteri>

```

Fig. 6.9 - Fișierului XML ce conține regulile descoperite

6.7. Concluzii

În cadrul acestui capitol s-a analizat un set de date privind nașterile realizate la *Clinica de Obstetrică - Ginecologie Bega, Timișoara* în anul 2010. Setul de date de la care a pornit analiza conține date despre 2326 de nașteri. Datele obținute au fost preprocesate deoarece conțineau valori lipsa, valori introduse greșit, valori neimportante pentru analiză și apoi aduse în format *arff* ca să poată fi analizate cu ajutorul lui *Weka*. O primă etapă a analizei datelor a constat în realizarea unei analize statistice la nivel de atribut. Această analiză a condus la determinarea unor informații interesante legate de nașterile din 2010, precum: în 70 % din cazuri femeilor care au născut natural a fost nevoie de ajutorul nașterii naturale printr-o tăietură în zona perineului, 60 % din femei au născut prin cezariană, majoritatea femeilor au născut în jurul vârstei de 28 de ani, s-au născut mai mulți băieți decât fete, etc. Următoarea etapă a constat în construirea cu ajutorul lui *Weka* a unui model clasificator, care să permită determinarea indicelui apgar pe baza unor variabile predictive precum gesta, para, epizionomie, etc. cu ajutorul unuia din algoritmi propuși sau a celor existenți în *Weka*. Cele mai bune rezultate s-au obținut cu ajutorul algoritmului AGR (vezi capitolul 4). Modelul clasificator construit permite determinarea cu o precizie de aprox 85 % a intervalului în care se va situa indicele apgar (unul din 5 intervale) al nou-născutului, în funcție de parametrii mamei, ai fătului și ai intervențiilor medicale care se vor face asupra mamei respectiv asupra nou-născutului.

Indicele apgar reprezintă o apreciere a funcțiilor vitale și a capacității de adaptare a fătului la condițiile din mediul extrauterin. Este de dorit ca valoarea acestui indicator să fie cât mai mare pentru fiecare nou-născut. Modelul construit poate fi utilizat înainte de naștere, pentru a realiza predicții ale intervalului în care se va plasa acest indice (unul din cele 5 intervale) în funcție de datele care se cunosc despre mamă, copil și intervențiile medicale care urmează să se realizeze pentru a ajuta nașterea. Dacă rezultatele predicției dau un indice apgar mic este necesar să se realizeze mai multe predicții cu modelul construit, la fiecare predicție modificând parametrii de intrare care pot fi modificați în practică de medici, pentru ca acest indice apgar să se maximizeze (de exemplu se poate modifica tipul de naștere sau se poate ajuta fătul să ajungă în poziția normală, etc). Ceea ce se caută este modul în care pot fi modificate intrările astfel încât la ieșire să se obțină o valoare multumitoare. Un exemplu al utilizării modelului este: se realizează predicții și se analizează în ce interval ar primi nou născutul nota (indicele apgar) dacă fătul s-ar naște natural respectiv prin cezariană. Dacă se estimează că ar primi notă mai mare dacă s-ar naște natural, atunci poate că mamele și / sau medicii se vor gândi de două ori înainte de a face cezariană. Situația poate fi valabilă și invers.

O altă utilitate a modelului dezvoltat, pe lângă utilizarea lui pentru a realiza predicții cu privire la intervalul în care se va situa indicele apgar, este utilizarea setului de reguli, transpuse în format XML, ca intrare într-un sistem care le transpune în Arden Syntax [HEE05]. Pentru aceasta a fost dezvoltată o aplicație care primește la intrare un set de reguli în formatul *if...then...* și le transformă în format XML.

Dezvoltarea acestui capitol a fost realizată utilizând cunoștințele și îmbunătățirile tehnologice aduse și prezentate în toate capitolele anterioare. Punctual, contribuțiile aduse de autor în cadrul acestui capitol sunt:

- Preprocesarea în Microsoft Excel a datelor legate de nașterile din 2010 prin analiza fiecărei coloane, eliminarea instanțelor cu valori lipsă sau valori introduse greșit, corectarea valorilor introduse greșit, eliminarea atributelor neimportante pentru analiză, crearea unor noi attribute, etc
- Transformarea datelor din format *xls* în formatul *arff*, specific *Weka* cu ajutorul aplicației *Arff Convertor*, concepută, implementată și prezentată de autor în cadrul capitolului 2
- Preprocesarea datelor cu ajutorul lui *Weka*, prin discreditarea atributelor numerice
- Analiza statistică la nivel de atribut, care a condus la obținerea unor informații noi și interesante privind nașterile
- Dezvoltarea unor modele de clasificare cu algoritmi *Naive Bayes*, *J48*, *k-Nearest-Neighbour* din *Weka*, algoritmi care au fost prezentați în introducere
- Dezvoltarea unor modele de clasificare cu algoritmi *Ant-Miner* și *Ant-r-Miner*, prezentați în capitolul 3
- Dezvoltarea unui model de clasificare care poate estima cu o precizie de 85 % unul din cele 5 intervale în care indicele apgar al nou-născutului va lua valoare. Modelul a fost construit cu ajutorul algoritmului *AGR*, care a fost conceput, implementat și prezentat de autor în cadrul capitolului 4
- Realizarea de predicții – modelul construit a fost utilizat pentru a realiza predicții cu ajutorul noii interfețe dinamice a lui *Weka*, concepută și prezentată de autor în cadrul capitolului 5
- Setul de reguli descoperit a fost transpus în format XML cu ajutorul unei aplicații concepute și implementate de autor, pentru a servi ca intrare unui sistem care le transpune în Arden Syntax [HEE05].

7. Concluzii

Clasificarea este una dintre cele mai populare tehnici de data mining. Clasificarea permite construirea unui model clasificator pe baza datelor de antrenare, testarea acestuia pe datele de test și apoi utilizarea lui pentru a realiza predicții. Construirea modelului clasificator se poate face cu ajutorul unor algoritmi specializați precum *Naive Bayes*, *k-Nearest Neighbor*, *C4.5* etc. Modelul construit poate fi utilizat pentru a prezice valoarea nominală a unui atribut clasă.

Algoritmii de clasificare au un domeniu de aplicabilitate foarte larg. Ei pot fi utilizați cu succes pentru clasificarea amprentelor, bolilor de inimă, a riscului de reapariție a unor evenimente specifice cancerului de sân, emoțiilor de pe fața umană, imaginilor de la sateliți, etc. Prin intermediul algoritmilor de clasificare o universitate poate prezice cu o precizie între 75 și 80 % care studenți vor absolvi și care nu. Universitatea poate folosi această informație pentru a-și concentra atenția asupra acelor studenți care sunt în pericol de a nu absolvi [DEK09].

În continuare sunt prezentate aspectele tratate pe parcursul tezei și contribuțiile aduse la nivel de capitol.

În *capitolul 1* a fost realizată o sinteză amplă, folosind 47 de surse bibliografice, a celor mai importante aspecte legate de domeniul Data Mining. Au fost prezentate definiții și terminologia folosită în domeniu, etapele procesului de data mining, domenii în care tehnicile de data mining se aplică cu succes, aspecte legate de clasificarea datelor și de algoritmii folosiți pentru a construi modele clasificatoare, aspecte legate de instrumentele software utilizate în studii de data mining și de seturile de date utilizate pentru testarea și validarea algoritmilor. Structura tezei a fost de asemenea prezentată.

Contribuțiile aduse în acest capitol sunt următoarele:

- A fost conturat domeniul, subliniind importanța lui, au fost prezentați succint termenii existenți, un număr mare de definiții existente în domeniu, a fost introdusă o nouă definiție, au fost prezentate original etapele procesului de data mining, pe baza experienței practice a autorului
- A fost realizată o sinteză amplă asupra aplicării tehnicilor de data mining în domeniile e-commerce, educație, medicină, combaterea terorismului. De asemenea a fost prezentată o sinteză a sistemelor de *data mining* utilizate de agențiile federale din SUA cu scopul de sublinia puterea și importanța tehnologiei *data mining*, precum și nivelul de vârf la care s-a ajuns cu utilizarea ei
- Au fost prezentate clar, succint aspecte generale legate de clasificarea datelor, precum și principiile de bază privind funcționarea a trei algoritmi de clasificare de referință (*Naive Bayes*, *k-Nearest Neighbor*, *C4.5*) care au fost utilizați în studiile realizate în capitolele 4, 5 și 6
- Au fost prezentate pe scurt aspecte legate de două dintre cele mai populare instrumente, la nivel mondial, utilizate pentru a realiza studii de data mining (*Weka* și *R*). Au fost analizate avantajele și dezavantajele pe care le are

fiecare instrument și a fost realizată o comparație între cele două instrumente

- Au fost prezentate aspecte legate de cele 10 seturi de date donate, reale, obținute de pe *UCI Machine Learning Repository*, care au fost utilizate pe parcursul tezei, în capitolele 3, 4 și 5 pentru a testa și valida algoritmi propuși și îmbunătățirile tehnologice aduse.

În *capitolul 2* au fost prezentate aspecte legate de preprocesarea datelor, aceasta fiind o etapă foarte importantă a procesului de data mining, etapa cea mai mare consumatoare de timp. Contribuțiile care au fost aduse de capitolul 2 sunt următoarele:

- Realizarea unei sinteze ample și riguroase cu privire la principalele tehnici de preprocesare a datelor, bazată pe studiul a 30 de referințe. Au fost tratate următoarele tehnici de preprocesare, agregarea, eșantionarea, reducerea dimensionalității (prin analiza componentelor principale, descompunerea valorilor singulare, analiză factorială, încapsularea locală lineară, scalare multidimensională), selectarea de submulțimi de caracteristici, crearea caracteristicilor, discretizarea și binarizarea, transformarea variabilelor, etc.
- Îmbunătățirea tehnologiei de preprocesare a datelor prin concepția și dezvoltarea unui instrument software care permite conversia datelor din 5 formate de baze de date în formatul de date folosit de *Weka* (arff). Instrumentul propus dispune de interfețe dinamice, în acord cu datele care se prelucrează, iar utilizarea lui este simplă și nu necesită realizarea de configurări în sistemul de operare, în *Weka* sau cunoștințe de SQL
- Realizarea unei analize a datelor studenților de la facultatea de Automatică și Calculatoare din anul universitar 2008-2009 cu ajutorul lui *Arff Converter* și *Weka*, analiză care oferă o imagine de ansamblu asupra provenienței studenților, pregătirii lor, etc.

În *capitolul 3* sunt tratate aspecte legate de clasificarea datelor cu ajutorul algoritmului *Ant Colony Optimization*. Din anul 2001 când a fost propusă utilizarea lui pentru descoperirea regulilor de clasificare, au fost realizate multe studii cu ajutorul acestui algoritm și problema alegerii unei secvențe de parametri optime a fost lăsată ca o direcție de cercetare de foarte mulți cercetători. În cadrul acestui capitol a fost studiat codul algoritmului și prezentată în detaliu funcționarea lui, au fost analizate foarte multe configurații de parametri și s-a emis o listă de recomandări privind alegerea acestora. De asemenea problema îmbunătățirii algoritmului a fost atinsă de mulți cercetători. Au fost prezentate sintetic îmbunătățirile pe care o parte din aceștia le-au adus și au fost propuse și operate noi îmbunătățiri. Acest lucru a fost posibil datorită faptului că aplicația care implementează algoritmul se numește *Ant Miner* și este Open Source.

Contribuțiile aduse în cadrul acestui capitol sunt următoarele:

- A fost realizată o prezentare detaliată a algoritmului, bazată atât pe documentația existentă cât și pe studiul codului algoritmului.
- Au fost studiați parametri algoritmului rulând câte 64 de configurații pe fiecare din cele 10 seturi de date alese și analizând rezultatele. Analiza rezultatelor s-a făcut atât vizual, comparativ cât și matematic construind și analizând modele de regresie și matrice de corelare

- Analiza făcută a condus către o listă de recomandări privind modul de alegere al parametrilor, utile pentru toți utilizatorii acestui algoritm
- A fost realizată o sinteză a îmbunătățirilor aduse algoritmului din anul 2001 și până în prezent
- Au fost aduse îmbunătățiri algoritmului care vizează modul de calcul a cantități inițiale de feromon a termenilor, modul de actualizare a feromonului, modul de evaluare a calității regulilor descoperite și modul de selecție a termenilor. Modificările propuse au fost implementate în aplicația open source *Ant Miner*. Aplicația care conține algoritmul îmbunătățit a fost numită *Ant-r-Miner*. S-a demonstrat că îmbunătățirile aduse conduc la o mai bună performanță utilizând date reale.

În *capitolul 4* este propus un nou algoritm genetic pentru clasificarea datelor. Algoritmul propus a fost denumit *AGR*. Contribuțiile aduse de acest capitol sunt următoarele:

- Prezentarea succintă a structurii generale a unui algoritm genetic
- Realizarea unei sinteze cu privire la algoritmi genetici utilizați la clasificarea datelor
- Propunerea unui nou algoritm genetic pentru clasificarea datelor care după cunoștințele autorului diferă de ceilalți algoritmi genetici de clasificare prin următoarele:
 - permite recalcularea clasei în urma încrucișării
 - permite stabilirea unui procent minim de acoperire pentru regulile descoperite, rularea repetată a algoritmului până se atinge acest procent
 - permite adăugarea unei reguli implicite listei regulilor descoperite
 - permite adăugarea la lista de reguli descoperite doar a celei mai bune reguli obținută după rularea algoritmului pentru numărul de generații stabilit sau a tuturor regulilor descoperite
 - utilizează o funcție *fitness* diferită de cele utilizate de alți algoritmi
- Îmbunătățirea tehnologiei de clasificare a datelor prin implementarea lui *AGR* într-un instrument Open Source utilizat pe scară largă (*Weka*)
- Analizarea comportamentului algoritmului la diferite valori ale parametrilor de intrare, și la utilizarea a diferite funcții *fitness* și emiterea unui set de recomandări cu privire la alegerea acestora

În cadrul *capitolului 5* a fost tratată problema realizării de predicții cu ajutorul unui model construit și testat de *Weka*. Contribuțiile aduse în cadrul acestui capitol sunt:

- Prezentarea detaliată a modului în care *Weka* poate fi utilizat pentru a realiza predicții
- Conceperea unei interfețe dinamice care conține controale în acord cu setul de date analizat și care ușurează și clarifică mult procesul de realizare a predicțiilor
- Prezentarea sintetică a mărimilor de ieșire furnizate de *Weka* în urma testării modelului
- Concepția unor soluții de transformare a interfeței de clasificare a datelor într-o interfață mult mai simplă, inteligibilă și prietenoasă, prin vizualizarea sub formă de grafic 3d a celei mai importante măsuri (acuratețea predicției) obținute în urma testării modelului și prin afișarea informațiilor din secțiunea

Classifier output (a indicatorilor "complicații") într-o nouă fereastră *Advanced Information*.

- Modificările menționate mai sus, realizate în unul din cele mai populare instrumente de data mining din lume, contribuie la îmbunătățirea tehnologiei de clasificare a datelor. Modificările au fost testate pe 3 seturi de date medicale și rezultatele au fost foarte bune.

În *capitolul 6*, a fost utilizată tehnologia prezentată și îmbunătățirile tehnologice propuse în capitolele precedente pentru a analiza și clasifica datele despre cele 2326 nașteri care au avut loc la *Clinica de Obstetrică - Ginecologie Bega, Timișoara* în anul 2010. Punctual, contribuțiile aduse de autor în cadrul acestui capitol sunt:

- Preprocesarea în Microsoft Excel a datelor legate de nașterile din 2010 prin analiza fiecărei coloane, eliminarea instanțelor cu valori lipsă sau valori introduse greșit, corectarea valorilor introduse greșit, eliminarea atributelor neimportante pentru analiză, crearea unor noi atribute, etc
- Transformarea datelor din format *xls* în formatul *arff*, specific *Weka*, cu ajutorul aplicației *Arff Convertor*, concepută, implementată și prezentată de autor în cadrul capitolului 2
- Preprocesarea datelor cu ajutorul lui *Weka*, prin discreditarea atributelor numerice
- Analiza statistică la nivel de atribut, care a condus la obținerea unor informații noi și interesante privind nașterile
- Dezvoltarea unor modele de clasificare cu algoritmi *Naive Bayes*, *J48*, *k-Nearest-Neighbour* din *Weka*, algoritmi care au fost prezentați în introducere
- Dezvoltarea unor modele de clasificare cu algoritmi *Ant-Miner* și *Ant-r-Miner*, prezentați în capitolul 3
- Dezvoltarea unui model de clasificare care poate estima cu o precizie de aproximativ 85 % unul din cele 5 intervale în care indicele apgar al nou-născutului va lua valoare. Modelul a fost construit cu ajutorul algoritmului *AGR*, care a fost conceput, implementat și prezentat de autor în cadrul capitolului 4
- Realizarea de predicții – utilizarea interfeței dinamice a lui *Weka*, concepută și prezentată de autor în cadrul capitolului 5 pentru realizarea predicțiilor
- Setul de reguli descoperit a fost transpus în format XML cu ajutorul unei aplicații concepute și implementate de autor, pentru a servi ca intrare unui sistem care le transpune în *Arden Syntax* [HEE05]
- Valoarea medicală a modelului clasificator este următoarea: Indicele apgar reprezintă o apreciere a funcțiilor vitale și a capacității de adaptare a fătului la condițiile din mediul extrauterin. Este de dorit ca valoarea acestui indicator să fie cât mai mare pentru fiecare nou-născut. Modelul construit poate fi utilizat înainte de naștere, pentru a realiza predicții a intervalului în care se va plasa acest indice (unul din cele 5 intervale) în funcție de datele care se cunosc despre mamă, copil și intervențiile medicale care urmează să se realizeze pentru a ajuta nașterea. Dacă rezultatele predicției dau un indice apgar mic este necesar să se realizeze mai multe predicții cu modelul construit, la fiecare predicție modificând parametrii de intrare care pot fi modificați în practică de medici, pentru ca acest indice apgar să se maximizeze (de exemplu se poate modifica tipul de naștere, sau se poate ajuta fătul să ajungă în poziția normală, etc). Ceea ce se caută este modul

În care pot fi modificate intrările astfel încât la ieșire să se obțină o valoare mulțumitoare. Un exemplu al utilizării modelului este: se realizează predicții și se analizează în ce interval ar primi nou născutul nota (indicele apgar) dacă fătul s-ar naște natural respectiv prin cezariană. Dacă se estimează că ar primi notă mai mare dacă s-ar naște natural, atunci poate că mamele și / sau medicii se vor gândi de două ori înainte de a face cezariană. Situația poate fi valabilă și invers.

Direcții de dezvoltare

O direcție de dezvoltare ulterioară este continuarea studiului din capitolul 3 și extinderea aplicației *Ant-r-Miner* astfel încât utilizatorul să poată alege dintr-o listă de funcții, funcția pe care dorește să o folosească pentru a evalua calitatea regulilor, respectiv să-și poată construi în mod dinamic funcția de calitate, cu ajutorul termenilor *TP*, *FP*, *TN*, *FN* și a unor operații elementare gen adunare, scădere, înmulțire, împărțire. Analiza influenței pe care o au valorile parametrilor de intrare al algoritmului *Ant Colony Optimization* asupra numărului de reguli și asupra numărului de condiții este de asemenea o direcție de dezvoltare ulterioară.

Continuarea studiului din capitolul 4 este de asemenea o direcție de dezvoltare. La fel ca și în cazul lui *Ant-r-Miner* se dorește dezvoltarea unei facilități care să-i permită utilizatorului să aleagă sau să introducă funcția *fitness* pe care dorește să o utilizeze în analizele pe care le realizează cu *AGR*. O altă direcție de dezvoltare este implementarea mai multor metode de selecție și oferirea posibilității de a alege metoda de selecție dorită într-un anumit studiu cu *AGR*.

Autorul va aplica tehnologia propusă în domeniul medical, pe cât mai multe seturi de date medicale.

Un alt obiectiv care va fi urmărit este conceperea unei modalități de utilizare a algoritmului albină în problema clasificării datelor [PHA07].

Bibliografie

- [AIC10] Aici M., Inan C. and Avci M., *A hybrid classification method of k nearest neighbor, Bayesian methods and genetic algorithm*, Expert Systems with Applications, Vol. 37, 2010, 5061-5067
- [ALK10] Alkan A., "Cahit Arf", Available from: <http://sourceforge.net/projects/cahitarf/>, Accessed: 18-03-2010
- [ANS01] Ansari S., Kohavi R., Mason L. and Zheng Z. *Integrating E-Commerce and Data Mining: Architecture and Challenges*. Proceedings IEEE International Conference on Data Mining, ISBN:0-7695-1119-8, pp. 27-34, 2001
- [ATA10] Atallah M.J. and Blanton M., *Algorithms and Theory of Computation Handbook*, CRC Press, Taylor & Francis Group, LLC, ISBN 978-1-58488-822-2, 2010
- [ATI09] Atilgan Y. and Dogan F. *Data Mining on Distributed Medical Databases: Recent Trends and Future Directions*. IT Revolutions, Volume 11, ISBN978-3-642-03977-5, pp. 216-224, 2009
- [BEN06] Bendakir N. and Aimeur E. *Using Association Rules for Course Recommendation*, Proceedings of the Twenty-First National Conference on Artificial Intelligence, ISBN 978-1-57735-281-5, 2006
- [BER04] Berry M. and Linoff G., *Data Mining Techniques For Marketing, Sales, and Customer Relationship Management*, Wiley Publishing, ISBN 0-471-47064-3, 2004
- [BER09] Bertini E. and Lalanne D., *Surveying the complementary role of automatic data analysis and visualization in knowledge discovery*, Proceedings of the ACM SIGKDD Workshop on Visual Analytics and Knowledge Discovery: Integrating Automated Analysis with Interactive Exploration, 2009
- [BHA08] Bhavani Thuraisingham, *Data Mining for Counter-Terrorism*, <http://www.utdallas.edu/~bxt043000/key-papers/%5B18%5DMining-Terrorism-NGDM04.pdf> , 2008
- [BIN11] Bing L., *Web Data Mining Exploring Hyperlinks, Contents, and Usage Data*, Springer, Second Edition, ISBN 978-3-642-19459-7, 2011
- [BOU08] Bouckaert R., Frank E., Hall M., Kirkby R., Reutemann P.,

- Seewald A., Scuse D., *Weka Manual for Version 3-6-0*, University of Waikato, Hamilton, New Zealand, 2008
- [DEH06] Dehuri S., Ghosh A. and Mall R., *Genetic Algorithms for Multi-Criterion Classification and Clustering in Data Mining*, International Journal of Computing & Information Sciences, Vol 6, No 3, 2006, pp. 143-154
- [DEK09] Dekker G.W., Pechenizkiy M. and Vleeshouwers J.M., *Predicting Students Drop Out: A Case Study*, Proceedings of 2nd International Conference On Educational Data Mining, Cordoba, Spain, July 1-3, ISBN: 978-84-613-2308-1, pp 41-50, 2009
- [DEZ09] Deza E. and Deza M.M., *Encyclopedia of Distances*, Springer, ISBN 978-3-642-00233-5, 2009
- [DIN09] Ding R., Dong H., Feng X. and Yin G., *A Hybrid Particle Swarm Genetic Algorithm for Classification*, Proceedings of the 2009 International Joint Conference on Computational Sciences and Optimization, vol.2, 2009, pp. 301
- [DOR04] Dorigo M. and Stützle T., *Ant Colony Optimization*, Massachusetts Institute of Technology, ISBN 0-262-04219-3, 2004
- [DOR96] Dorigo, M. and Maniezzo, V., *The ant system: optimization by a colony of cooperating agents*, IEEE Transactions on Systems, Man, and Cybernetics, 26(1), pp. 1-13, 1996
- [DOU95] Dougherty J., Kohavi R. and Sahami M., *Supervised and Unsupervised Discretization of Continuous Features*, Proceedings of the Twelfth International Conference on Machine Learning (1995), pp. 194-202.
- [EZE05] Ezeife C.I. and Yi Lu (2005). *Mining Web Log Sequential Patterns with Position Coded Pre-Order Linked WAP-Tree*, Data Mining and Knowledge Discovery, Volume 10, Number 1 / January, ISSN1384-5810, pp 5-38, 2005
- [FAL95] Faloutsos C. and Lin K., *FastMap: A Fast Algorithm for Indexing, Data-Mining and Visualization of Traditional and Multimedia Datasets*, in Proceedings of 1995 ACM SIGMOD, SIGMOD RECORD (June 1995), vol.24, no.2, pp 163-174
- [FIL08] FILIP Ioan, ROBU Raul, SZEIDERT Iosif and ROBU Andreea, *Considerations Regarding an Unusual Cause of Database Concurrency Exception*, Annals of DAAAM for 2008 & Proceedings of the 19th International DAAAM Symposium, pag. 493-494, Trnava, Slovakia, 2008
- [FIL09] FILIP Ioan, VAŞAR Cristian and ROBU Raul, *Considerations about*

- an Oracle Database Multi-Master Replication*, Proc. 5th International Symposium on Applied Computational Intelligence and Informatics, SACI, Timisoara, Romania, pag. 147-152, 2009
- [FRA10] Frank, A. and Asuncion, A. (2010). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.
- [FRE02] Freitas A.A., *A survey of evolutionary algorithms for data mining and knowledge discovery*, Advances in Evolutionary Computation, Springer-Verlag, 2002, pp. 819-845
- [GAN07] Gang K., Yi P., Yong S. and Zhengxin C. *Privacy-Preserving Data Mining of Medical Data Using Data Separation-Based Techniques*. Data Science Journal, Vol. 6, pp.429-434, 2007
- [GAO04] United States General Accounting Office (GAO), *DATA MINING - Federal Efforts Cover a Wide Range of Uses*, Report to the Ranking Minority Member, Subcommittee on Financial Management, the Budget, and International Security, Committee on Governmental Affairs, U.S. Senate
http://www.gao.gov/docsearch/locate?searched=1&o=0&order_by=rel&search_type=publications&keyword=data+mining&Submit=Search
- [GOR06] Gorunescu F. *Data Mining, Concepte, Modele și Tehnici*, Editura Albastra, Cluj-Napoca, ISBN 973-650-169-8, 2006
- [HAL09] Hall M., Frank E., Holmes G., Pfahringer B., Reutemann P. and Witten I.H., *The Weka Data Mining Software: An Update*, ACM SIGKDD Explorations Newsletter, Volume 11 , Issue 1, 2009, pp. 10-18
- [HEE05] Hee K.P., Jae P.K., Jin W.C., Jae J.H. and Seung M.H., *Design of real-time arden syntax based decision support system with minimum workload and building cost*, Enterprise networking and Computing in Healthcare Industry, 2005. *HEALTHCOM 2005. Proceedings of 7th International Workshop on* , vol., no., pp. 430-433, 23-25 June 2005
- [HOL09] Holban S. *Data Mining and Machine Learning – Notițe de curs*, 2009
- [HOR09] Hornik K., Buchta C. and Zeileis A., *Open-Source Machine Learning: R Meets Weka*, Computational Statistics, ISSN 0943-4062, pp 225-232, 2009
- [HU09] Hu X., Zhang X., Lu C., Park E. K. and Zhou X., *Exploiting Wikipedia as External Knowledge for Document Clustering*, Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, Paris, France, ISBN:978-

- [JOL02] 1-60558-495-9, pp 389-396, 2009
Jolliffe I.T., *Principal Component Analysis*, Springer-Verlag, Second Edition, October 2002
- [KOT06] Kotsiantis S. B., Kanellopoulos D. and Pintelas P.E., *Data Preprocessing for Supervised Learning*, International Journal of Computer Science, pp.111-117, 2006
- [KOT09] Kotsiantis S., *Educational Data Mining: A case study for predicting drop-out Prone Students*, International Journal of Knowledge Engineering and Soft Data Paradigms, Inderscience Publishers, ISSN:1755-3210, pp 101-111, 2009
- [LAL11] Lalbakhsh, P., Fasaeei, M.S.K. and Fesharaki M.N., *Focusing on rule quality and pheromone evaporation to improve ACO rule mining*, Proceedings of the IEEE Symposium on Computers & Informatics (ISCI), ISBN: 978-1-61284-689-7, pp. 108 - 112, 2011
- [LAR06] Larose D.T., *Data Mining Methods and Model*, Published by John Wiley & Sons, Inc., Hoboken, New Jersey, ISBN-13 978-0-471-66656-1, 2006
- [LEH98] Lehmann E. L. and Casella G., *Theory of Point Estimation (2nd ed.)*, New York: Springer, ISBN 0-387-98502-6, 1998
- [LES06] Leskovec J., *Dimensionality reduction PCA, SVD, MDS, ICA, and friends, Machine Learning recitation*, April 27 2006
- [LIU02a] Liu B., Abbass H. A. and McKay B., *Density-based heuristic for rule discovery with ant-miner*, Proceedings of the 6th Australasia-Japan joint workshop on intelligent and evolutionary systems (AJWIS2002), Canberra, Australia, pp. 180-184, 2002
- [LIU02b] Liu H., Hussain F., Tan C. L. and Dash M., *Discretization: An Enabling Technique*, Data Mining and Knowledge Discovery, Springer Netherlands, Vol. 6, pp. 393-423, 2002
- [LIU03a] Liu, B., Abbass H. A. and McKay B. (2003). *Classification rule discovery with ant colony optimization*, Proceedings IEEE/WIC international conference on intelligent agent technology (pp. 83-88). IEEE Comput. Soc.: Los Alamitos.
- [LIU03b] Liu H., Yu L., and Motoda H., *Feature Extraction, Selection, and Construction*, The handbook of Data Mining, Edited by Nong Ye, Lawrence Erlbaum Associates, Inc., pp. 409-425, 2003
- [MAR05] Marcelon A. and Yacef K. *Educational Data Mining: a Case Study. Artificial Intelligence in Education: supporting learning through Intelligent and Socially Informed Technology*, IOS Press, 2005

- [MAR11] D. Martens, B. Baesens, and T. Fawcett, *Editorial survey: swarm intelligence for data mining*, Machine Learning, Volume 82, Number 1, pp. 1-42, 2011
- [MAT08] Matthews D.C., *Improved Lower Limits for Pheromone Trails in Ant Colony Optimization*, Proceedings of the 10th international conference on Parallel Problem Solving from Nature: PPSN X, Springer-Verlag Berlin, Heidelberg, pp. 508 – 517, 2008
- [MAU00] Maulik U. and Bandyopadhyay S., *Genetic algorithm-based clustering technique*, Pattern Recognition 33, 2000, pp. 1455-1465
- [MOB04] Mobasher B., *Web usage mining and personalization*. Practical handbook of internet computing (ed.) M P Singh (CRC Press), 2004
- [MOH09] Mohammed N., Fung B., Hung P. and Lee C., *Anonymizing Healthcare Data: A Case Study on the Blood Transfusion Service*, Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, Paris, France, ISBN:978-1-60558-495-9, pp 1285-1294, 2009
- [MOO09] Moore S., D'Addario D., Kurinskas J., and Weiss G. M., *Are Decision Trees Always Greener on the Open (Source) Side of the Fence?* In Proceedings of the 2009 International Conference on Data Mining, CSREA Press, pp185-188, 2009
- [MUK09] Mukhopadhyay D.M., Balitanas M.O., Farkhod A., Jeon S.H., and Bhattacharyya D., *Genetic Algorithm: A Tutorial Review*, International Journal of Grid and Distributed Computing, Vol.2, No.3, 2009, pp. 25-32
- [NEJ08] Nejad N. Z., Bakhtiary A.H. and Analoui M., *Classification Using Unstructured Rules and Ant Colony Optimization*, Proceedings of the International MultiConference of Engineers and Computer Scientists, Vol I, ISBN: 978-988-98671-8-8, Hong Kong, pp. 506-510, 2008
- [NIU02] Niu L., Yan X.W., Zhang C.Q. and Zhang S.C., *Product hierarchy-based customer profiles for electronic commerce recommendation*. In Int. Conf. on Machine Learning and Cybernetics, pp 1075–1080, 2002
- [PAR01] Parpinelli R.S., Lopes H.S. and Freitas A.A., *An ant colony algorithm for classification rule discovery*, Abbas, H.A., Sarker, R.A., Newton, C.S. (eds.), Data Mining: A Heuristic Approach., Hershey: Idea Group Publishing, v. 1, ISBN 1-9307-0825-4, pp. 190-208, 2001

- [PAR02] Parpinelli R.S., Lopes H.S. and Freitas, A.A., *Data mining with an ant colony optimization algorithm*, IEEE Transactions on Evolutionary Computation, special issue on Ant Colony Algorithms, v. 6, n. 4, pp. 321-332, August, 2002
- [PHA07] Pham D.T., Otri S., Afify A., Mahmuddin M. and Al-Jabbouli H., *Data Clustering Using the Bees Algorithm*, Manufacturing Systems. Retrieved from <http://bees-algorithm.com/modules/2/6.pdf> , 2007
- [PLA10] Platt J.C., *FastMap, MetricMap, and Landmark MDS are all Nystrom Algorithms*, Microsoft Research, Available at: <http://research.microsoft.com/pubs/69185/nystrom2.pdf>, Accessed: 09-05-2010
- [PON07] Ponniah P., *Data Modeling Fundamentals, A Practical Guide for IT Professionals*, Wiley Interscience, A John Wiley and sons, Inc., Publication, ISBN-10: 0-471-79049-4, 2007
- [PRO09] PROȘTEAN G., ROBU R., PROȘTEAN O. and TAROATĂ A., *Considerations regarding the integration of the cost centers based on SAP management accounting*, Proc. 6th International Conference on Management of Technological Changes, vol.1, Alexandroupolis, Greece, pag. 699-702, 2009
- [QUI86] Quinlan, J.R., *Induction of Decision Trees*, Machine Learning 1, pp 81-106, 1986
- [QUI93] Quinlan J. R., *C4.5: Programs for Machine Learning* , Morgan Kaufmann Publishers, Inc., pp.235-240, 1993
- [QUI98] Quinlan. R., *C5.0: An Informal Tutorial*, <http://www.rulequest.com/see5-unix.html>, 1998, Accessed: 05-05-2012
- [RAG05] Raghavan Srinivasa N.R. (2005) "Data mining in e-commerce: A survey". Sadhana Vol. 30, Parts 2 & 3, pp. 275-289, 2005
- [ROB10a] Robu, R., Hora, C. and Stoicu - Tivadar, V. (2010). Improving the Classify User Interface in *Weka Explorer* (2010). 0171-0173, Annals of DAAAM for 2010 & Proceedings of the 21st International DAAAM Symposium, ISBN 978-3-901509-73-5, ISSN 1726-9679, pp 0086, Editor B. Katalinic, Published by DAAAM International, Vienna, Austria 2010
- [ROB10b] Robu R. and Stoicu-Tivadar V., *Arff Convertor Tool for WEKA Data Mining Software*, Proceedings of the IEEE International Joint Conferences on Computational Cybernetics and Technical Informatics, ICC-CNTI, pag. 247-251, Timisoara, Romania,

- 2010
- [ROB10c] Robu R. and Stoicu-Tivadar V., *Evaluation and Improvement of User Interfaces for Medical Data Mining with Open Source Programs*, Proceedings of the 31st National Conference on Medical Informatics, ISBN 978-606-8054-14-8, pp. 284-289, Arad, 18-20 November, 2010
- [ROB11a] ROBU R., FEIER F., STOICU-TIVADAR V., ILIE C., ENĂTESCU I., *The analysis of the new-borns' cry using NEONAT and data mining techniques*, Proceedings of the 15th IEEE International Conference on Intelligent Engineering Systems (INES), Pages 235 - 238, Poprad, Slovakia, 2011
- [ROB11b] ROBU R. and HOLBAN S., *A genetic algorithm for classification*, Proceedings of the International Conference on Computers and Computing, (ICCC), Pages 52-56 Lanzarote, Spania, 2011
- [ROB11c] ROBU R. and STOICU-TIVADAR V., *Regional data mining system*, Studies in Health Technology and Informatics, Volume 165, (STC), Pages 117-122, Lasko, Slovenia, 2011
- [ROB11d] ROBU R., TOLBUŞ A. and FILIP I., *DMCP - A data mining tool for classification and prediction*, Proceedings of the 7th International Conference on Management of Technological Changes (MTC 2011), vol.1, Pages 605-608, Alexandroupolis, Greece, 2011
- [ROM08] Romero C., Ventura S. and García E., *Data mining in course management systems: Moodle case study and tutorial*, Computers & Education, Volume 51, Issue 1, pp. 368-384, 2008, Elsevier Science
- [SAI11] Said N., Hammami M. and Ghedira K., *MuO-AntMiner: a new ant colony algorithm for the multi-objective classification problem*, Proceedings of the 2011 international conference on Computational science and its applications - Volume Part II, ISBN: 978-3-642-21886-6, 2011
- [SAL09] Salle P., Bringay S., Teisseire M., Chakkour F., Roche M., Rassoul R.A., Verdier J.M., Devau G., *GeneMining: Identification, Visualization, and Interpretation of Brain Ageing Signatures*, Proceedings of XXII International Conference of the European Federation for Medical Informatics, 2009
- [SAL11] Salama K., Abdelbar A. and Freitas A., *Multiple pheromone types and other extensions to the Ant-Miner classification rule discovery algorithm*, Swarm Intelligence (11 June 2011), pp. 1-34
- [SAR00] Sarukkai R.R., *Link prediction and path analysis using markov chains*. Proceedings of the 9th International World Wide Web Conference, Amsterdam, 2000
- [SAU00] Saul L.K. and Roweis S.T., *Nonlinear dimensionality reduction by*

- [SCH06] *locally linear embedding*, Science, Vol. 290, pp. 2323–2326, 2000.
Schloeffel P., Beale T., Hayworth G., Heard S. and Leslie H., *The relationship between CEN 13606, HL7 and openEHR*, HIC 2006 and HINZ 2006: Proceedings, Sydney, Australia, 2006, pp.24-28
- [SEI08] Seifert, Jeffery W. *Data Mining and Homeland Security: An Overview*. April 3, 2008. Congressional Research Reports for the People
http://assets.opencrs.com/rpts/RL31798_20080403.pdf ,
Accessed:17-05-2008
- [SKI03] Skillcorn D.B., *Cluster within Clusters: SVD and counterterrorism*, Proceedings of the Third SIAM International Conference on Data Mining, San Francisco, ISBN 0-89871-545-8, 2003
- [SMA06] Smaldon J. and Freitas, A. (2006). *A new version of the Ant-Miner algorithm discovering unordered rule sets*, Proceedings genetic and evolutionary computation conference (GECCO-2006) (pp. 43–50). San Francisco: Morgan Kaufmann
- [SON08] Song J., *A new dissimilarity measure in time-dependent experiments*, Journal of the Korean Statistical Society, Volume 37, Issue 2, pp. 145-153, 2008
- [STE02] Steyvers M, *Multidimensional scaling*, Encyclopedia of Cognitive Science, London, United Kingdom, Nature Publishing Group, 2002
- [TAN05] Tan P.-N., Steinbach M. and Kumar V., *Introduction to Data Mining*, Editura Addison Wesley, US, ISBN 0-321-32136-7,2005
- [TWO99] Two Crows Corporation, *Introduction to Data Mining and Knowledge Discovery, Third Edition*, ISBN: 1-892095-02-5, 1999
- [VAS09] VASAR Cristian, PROSTEAN Octavian, FILIP Ioan, ROBU Raul and POPESCU Dan, *Markov Models for Wireless Sensor Network Reliability*, Proc. IEEE 5th International Conference on Intelligent Computer Communication and Processing, (ICCP), Cluj-Napoca, Romania, pag.323-328, 2009
- [VER05] Verleysen M. and François D., *The Curse of Dimensionality in Data Mining and Time Series Prediction*, Proceedings of the 8th international conference on Artificial Neural Networks: computational Intelligence and Bioinspired Systems, pp. 758-770, ISBN: 978-3-540-26208-4, 2005
- [VIV10a] Vivekanandan P. and Nedunchezian D. R., *A New Incremental Genetic Algorithm Based Classification Model to Mine Data With Concept Drift*, Journal of Theoretical and Applied Information Technology, 2010, pp. 36-42
- [VIV10b] Vivekanandan P. and Nedunchezian D. R., *A Fast Genetic*

- Algorithm for Mining Classification Rules in Large Datasets*, International Journal on Soft Computing (IJSC), Vol.1, No.1, 2010, pp. 10-20
- [WAN04] Wang Z. and Feng B., *Classification rule mining with an improved ant colony algorithm*, Proceedings of the 17th Annual Australian Conference on Artificial Intelligence, 2004
- [WIT05] Witten I.H. and Frank E., *Data Mining Practical Machine Learning Tools and Techniques, Second Edition*, Elsevier Inc., ISBN-13:978-0-12-088407-0, 2005
- [WU07] Wu X., Kumar V., Quinlan R., Ghosh J., Yang Q., Motoda H., McLachlan G., Ng A., Liu B., Yu P., Zhou Z., Steinbach M., Hand D. Steinberg D. *Top 10 algorithms in data mining*, Knowledge and Information Systems, Volume 14, pp. 1-37, ISSN:0219-1377, 2007
- [www1] *Data mining*, <http://wordnetweb.princeton.edu/perl/webwn?s=data%20mining>, Accessed: 02-06-2012
- [www2] *Data mining*, <http://emeld.org/school/glossary.html#data>, Accessed: 28-06-2012
- [www3] *Web Mining*, <http://brodholt.com/WebMining.htm>, Accessed: 03-05-2011
- [www4] *Attribute-Relation File Format*, <http://www.cs.waikato.ac.nz/~ml/weka/arff.html>, Accessed: 22-06-2012
- [www5] *Machime Learning with Weka and R*, http://videlectures.net/bootcamp07_belanche_mldm, Accessed: 03-04-2011
- [www6] *Connecting Weka to databases*, Available from: <http://weka.wikispaces.com/Windows+Databases>, Accessed: 18-03-2010
- [www7] *DataFit*, <http://www.oakdaleengr.com/>, Accessed: 28-06-2012
- [YAN06] Yang Y. F., Lohmann P. and Heipke C., *Genetic Algorithms for the Unsupervised Classification of Satellite Images*, IntArchPhRS. Band XXXVI/3. Bonn, 2006, S. 179-184
- [ZHA03] Zhang Y.Q., Shteynberg M., Prasad S.K. and Sunderraman R. *Granular fuzzy web intelligence techniques for profitable data mining*. In 12th IEEE Int. Conf. on Fuzzy Systems, FUZZ '03 (New York:IEEE Comput. Soc.) pp 1462-1464, 2003

- [ZHA05] Zhang C., Yang Q. and Liu B., *Guest Editors' Introduction: Special Section on Intelligent Data Preparation*, IEEE Transactions On Knowledge And Data Engineering, Vol. 17, No. 9, pp. 1163-1165, 2005
- [ZLH07] Zaharie D., Lungeanu D. and Holban S., "Feature ranking based on weights **estimated** by multiobjective optimization", *Proceedings of IADIS First European Conference on Data Mining*, pp.124-128, 2007

Lista figurilor

Fig. 1.1 - Disciplinele care stau la baza Data Mining.....	8
Fig. 1.2 - Etapele procesului de Data Mining.....	10
Fig. 1.3 - Structura unui arbore de decizie construit cu C4.5.....	27
Fig. 1.4 - Structura unui fișier *.arff.....	29
Fig. 1.5 - Interfața principală a lui Weka.....	30
Fig. 1.6 - Aplicația Explorer.....	30
Fig. 1.7 - Model de clasificare.....	31
Fig. 1.8 - Arbore de clasificare.....	32
Fig. 1.9 - Interfața R – Exemplu de histogramă.....	34
Fig. 1.10 - Exemple de grafice trasate cu R.....	34
Fig. 2.1 - Preprocesarea datelor.....	42
Fig. 2.2 - Puncte în spațiul bidimensional.....	46
Fig. 2.3 - Proiectarea unui punct pe o dreaptă.....	54
Fig. 2.4 - Fluxul procesului de selectare a unui subset de caracteristici.....	58
Fig. 2.5 - Preprocesarea datelor în Weka.....	67
Fig. 2.6 - Funcționalitățile lui ARFF Convertor.....	68
Fig. 2.7 - Fereastra principală a aplicației.....	71
Fig. 2.8 - Selectarea tabelelor de interes.....	72
Fig. 2.9 - Selectarea câmpurilor din tabele diferite pentru ARFF.....	73
Fig. 2.10 - Stabilirea tipului atributelor.....	74
Fig. 2.11 - Specificarea legăturilor dintre câmpuri.....	74
Fig. 2.12 - Fișierul ARFF generat.....	75
Fig. 2.13 - Numărul de studenți care stau în cămin.....	76
Fig. 2.14 - Domiciliul studenților de la Facultatea de Automatică și Calculatoare.....	77
Fig. 2.15 - Nota la bacalaureat.....	77
Fig. 2.16 - Media de admitere.....	78
Fig. 2.17 - Numărul de credite la sfârșitul semestrului 1.....	79
Fig. 2.18 - Media anuală.....	79
Fig. 3.3 - Acuratețea predicției pentru varianta 1 de parametri.....	92
Fig. 3.4 - Timpul de execuție pentru varianta 1 de parametri.....	93
Fig. 3.3 - Ant Miner 1.2.1.....	108
Fig. 3.4 - Aplicația propusă, Ant-r-Miner.....	111
Fig. 4.3 - Weka cu algoritm genetic.....	123
Fig. 4.4 - Clasificatorul GeneticAlgorithm.....	124
Fig. 5.2 - Fișierul ARFF cu instanțele de prezis.....	139
Fig. 5.2 - Configurări necesare pentru realizarea de predicții.....	139
Fig. 5.3 - Reevaluarea modelului pe setul curent de testare.....	140
Fig. 5.4 - Vizualizarea valorilor prezise.....	140
Fig. 5.5 - Weka cu interfață dinamică de realizare a predicțiilor.....	141
Fig. 5.6 - Informațiile din secțiunea Classifier Output.....	143
Fig. 5.7 - Arhitectura lui Weka.....	144
Fig. 5.8 - Predicții pe baza seturilor de date Ljubljana breast cancer (stânga) și diabetes (dreapta).....	145

Fig. 6.1 - Arff convertor cu datele despre nașteri.....	149
Fig. 6.2. - Alegerea tipului atributelor din fișierul ARFF.....	150
Fig. 6.3 - Fișierul ARFF generat.....	150
Fig. 6.4 - Analiză statistică la nivelul fiecărui atribut cu Weka.....	151
Fig. 6.5 - Numărul de nașteri pe lunile anului.....	152
Fig. 6.6 - Greutatea nou născuților.....	152
Fig. 6.7 - Indicele apgar.....	153
Fig. 6.8 - Realizarea de predicții.....	158
Fig. 6.9 - Fișierului XML ce conține regulile descoperite.....	159

Lista tabelelor

Tabel 1.1 - O parte din departamentele și agențiile analizate.....	20
Tabel 1.2 - Scopul și numărul eforturilor de data mining.....	20
Tabel 1.3 - Proiecte de data mining în stadiul operațional.....	24
Tabel 1.4 - Avantajele și dezavantajele lui Weka.....	33
Tabel 1.5 - Exemple de comenzi în R.....	33
Tabel 1.6 - Avantajele și dezavantajele lui R.....	35
Tabel 1.7 - Comparatie între Weka și R.....	35
Tabel 1.8 - Bazele de date utilizate pentru analiza comportamentului algoritmilor.....	37
Tabel 2.1 - Binarizare calificative.....	63
Tabel 2.2 - Binarizare în probleme de asociere.....	63
Tabel 3.1 - Parametrii aleși pentru algoritm în diverse studii.....	89
Tabel 3.2 - Valorile parametrilor algoritmului.....	90
Tabel 3.3 - Acuratețea predicției pentru varianta 1 de parametri.....	91
Tabel 3.4 - Valorile maxime și minime ale acurateței predicției și numărul de furnici cu care s-au obținut acestea, pentru varianta 1 de parametri.....	91
Tabel 3.5 - Numerele de furnici și la câte valori maxime și minime au condus în ceea ce privește acuratețea predicției în varianta 1 de parametri.....	92
Tabel 3.6 - Timpul de execuție în varianta 1 de parametri.....	93
Tabel 3.7 - Valorile parametrilor algoritmului.....	93
Tabel 3.8 - Acuratețea predicției și timpul de execuție pentru variantele 2, 3 și 4 de parametri.....	94
Tabel 3.9 - Valorile parametrilor algoritmului în variantele 4 și 5 de parametri	94
Tabel 3.10 - Comparatie între valorile maxime obținute cu variantele 4 și 5 de parametri.....	95
Tabel 3.11 - Variantele 6 și 7 de parametri.....	95
Tabel 3.12 - Acuratețea predicției și timpul de execuție pentru variantele 5 și 6 de parametri.....	95
Tabel 3.13 - Variantele 7, 8, 9 și 10 de parametri.....	96
Tabel 3.14 - Acuratețea predicției și timpul de execuție pentru variantele 8, 9, 10 și 11 de parametri.....	96
Tabel 3.15 - Variantele 12, 13, 14 și 15 de parametri.....	97
Tabel 3.16 - Acuratețea predicției și timpul de execuție pentru variantele 12, 13, 14 și 15 de parametri.....	97
Tabel 3.17 - Variantele 16, 17, 18 și 19 de parametri.....	97
Tabel 3.18 - Acuratețea predicției și timpul de execuție pentru variantele 16, 17, 18 și 19 de parametri.....	98
Tabel 3.19 - Valorile parametrilor de intrare ai algoritmului.....	99
Tabel 3.20 - Modelele de regresie determinate și performanțele lor.....	100
Tabel 3.21 - Coeficienții parametrilor din modelele de regresie determinate.	101
Tabel 3.22 - Matricea de corelare aferentă setului de date car.....	102
Tabel 3.23 - Coeficienții de corelare dintre variabila dependentă și variabilele independente.....	102
Tabel 3.24 - Acuratețea predicției pentru diferite valori ale cantitatii initiale de feromon.....	109

Tabel 3.25 - Acuratețea predicției cu Ant-Miner și Ant-r-Miner.....	111
Tabel 4.1 - Codificarea cromozomilor.....	120
Tabel 4.2 - Încrucișarea cromozomilor.....	121
Tabel 4.3 - Parametrii de intrare.....	125
Tabel 4.4 - Acuratețea pe mulțimea de antrenament (%).....	125
Tabel 4.5 - Acuratețea pe mulțimea de test (%).....	126
Tabel 4.6 - Acoperire pe mulțimea de antrenament (%).....	126
Tabel 4.7 - Acoperire pe mulțimea de test (%).....	127
Tabel 4.8 - Timpul de execuție (s).....	127
Tabel 4.9 - Numărul regulilor descoperite.....	127
Tabel 4.10 - Numărul de condiții din regulile descoperite.....	128
Tabel 4.11 - Acuratețea pe mulțimea de antrenament (%).....	129
Tabel 4.12 - Acuratețea pe mulțimea de test (%).....	130
Tabel 4.13 - Acoperire pe mulțimea de antrenament (%).....	130
Tabel 4.14 - Acoperire pe mulțimea de test (%).....	130
Tabel 4.15 - Timpul de execuție.....	131
Tabel 4.16 - Numărul regulilor descoperite.....	131
Tabel 4.17 - Numărul de condiții din regulile descoperite.....	131
Tabel 4.18 - Numărul de valori maxime obținute când parametrul BR=true, respectiv când parametrul BR=false.....	131
Tabel 4.19 - Configurații de parametri.....	132
Tabel 4.20 - Acuratețea pe mulțimea de antrenament (%).....	133
Tabel 4.21 - Acuratețea pe mulțimea de test (%).....	133
Tabel 4.22 - Timp execuție (s).....	133
Tabel 4.23 - Numărul regulilor descoperite.....	133
Tabel 4.24 - Numărul de condiții din reguli.....	133
Tabel 4.25 - Acuratețea pe mulțimea de test.....	135
Tabel 4.26 - Număr reguli.....	135
Tabel 5.1 - Acuratețea predicției.....	144
Tabel 5.1 - Acuratețea predicției pe setul de date nasteri_1.arff.....	154
Tabel 5.2 - Acuratețea predicției pe setul de date nasteri_2.arff.....	155
Tabel 5.3 - Acuratețea predicției pe setul de date nasteri_3.arff.....	156
Tabel 5.4 - Acuratețea predicției pe setul de date nasteri_4.arff.....	156
Tabel 5.5 - Acuratețea predicției pe setul de date nasteri_5.arff.....	157