

UNIVERSITATEA BABEȘ-BOLYAI  
FACULTATEA DE MATEMATICĂ ȘI INFORMATICĂ

# Noi abordări în vederea extragerii cunoștințelor din date. Aplicații în domeniul educațional

Rezumatul tezei de doctorat

Student doctorand: Liana-Maria Crivei  
Coordonator științific: Prof. Dr. Czibula Gabriela

Septembrie 2020

**Cuvinte cheie:** Extragerea cunoștințelor din date de natură academică, predicția performanțelor studenților, reguli de asociere relaționale, ansambluri de arbori de decizie, rețele neuronale profunde, rețele cu autoorganizare.

# Cuprins

<b>Cuprinsul tezei</b>	<b>2</b>
<b>Lista publicațiilor</b>	<b>5</b>
<b>Introducere</b>	<b>7</b>
<b>1 Concepte introductive</b>	<b>12</b>
1.1 Extragerea de cunoștințe din date de natură academică (EDM)	12
1.2 Extragerea regulilor de asociere din date	12
1.3 Modelele de învățare automată folosite	14
1.3.1 Rețele neuronale artificiale	14
1.3.2 Păduri aleatoare de arbori	14
1.3.3 Rețele cu auto-organizare	14
<b>2 Noi abordări în extragerea cunoștințelor din date</b>	<b>16</b>
2.1 O abordare privind extragerea incrementală a regulilor de asociere relațională	16
2.2 O abordare privind extragerea concurentă a regulilor de asociere relațională	17
<b>3 Noi modele de învățare nesupervizată în extragerea cunoștințelor din date de natură academică</b>	<b>18</b>
3.1 Utilizarea modelelor de învățare nesupervizată pentru analiza performanței academice a studenților	18
3.2 Extragerea incrementală a regulilor de asociere relațională din seturi de date de natură academică	19
<b>4 Modele de învățare supervizată pentru predicția performanței studenților</b>	<b>21</b>
4.1 O analiză a metodelor de învățare supervizată pentru predicția performanței studenților în medii academice	21
4.2 <i>SPRAR</i> : Un model de clasificare bazat pe extragerea regulilor de asociere relațională aplicat pentru predicția performanței academice	22
4.3 Un studiu privind predicția performanței studenților folosind <i>SPRAR</i>	23
<b>Concluzii</b>	<b>25</b>
<b>Bibliografie</b>	<b>27</b>

# Cuprinsul tezei

<b>List of publications</b>	<b>8</b>
<b>Introduction</b>	<b>10</b>
<b>Acknowledgements</b>	<b>15</b>
<b>1 Background</b>	<b>16</b>
1.1 Educational data mining (EDM)	16
1.1.1 Predicting students' performance	17
1.1.1.1 Supervised learning approaches	17
1.1.1.2 Unsupervised learning approaches	18
1.1.2 Predicting instructors' performance	19
1.1.3 Other EDM tasks	19
1.2 Background on association rule mining	20
1.2.1 Frequent itemsets and association rules	20
1.2.2 Relational association rules (RARs)	21
1.2.2.1 Example	22
1.2.3 Literature review on incremental association rule mining	23
1.2.4 Literature review on concurrent association rule mining	25
1.3 Machine learning models used	26
1.3.1 Artificial neural networks	27
1.3.2 Random forests	27
1.3.3 Self-organizing maps	29
<b>2 New approaches in Data Mining</b>	<b>31</b>
2.1 A new incremental relational association rule mining approach	31
2.1.1 Methodology	33
2.1.1.1 Example of <i>IRARM</i> 's execution	36
2.1.2 Results and discussion	37
2.1.2.1 Experimental results	37
2.1.2.2 Comparison to related work	39
2.1.3 Conclusions and future work	40
2.2 A novel concurrent relational association rule mining approach	41
2.2.1 Examples of RAR mining	43
First example	43
Second example	44
2.2.2 Methodology	44
2.2.2.1 Theoretical considerations	45

2.2.2.2	<i>CRAR</i> algorithm . . . . .	45
2.2.3	Experimental evaluation . . . . .	51
2.2.3.1	Data sets . . . . .	52
2.2.3.2	Experimental setup . . . . .	52
2.2.3.3	Results . . . . .	53
2.2.4	Discussion . . . . .	54
2.2.4.1	Evaluation measures . . . . .	55
2.2.4.2	Analysis of <i>CRAR</i> method . . . . .	55
2.2.4.3	Comparison to related work . . . . .	59
2.2.5	Conclusions and future work . . . . .	60
<b>3</b>	<b>New unsupervised learning models in Educational Data mining</b>	<b>61</b>
3.1	Using unsupervised learning models for analyzing students' academic performance . . . . .	61
3.1.1	Methodology . . . . .	62
3.1.1.1	Data set . . . . .	62
3.1.1.2	Experiments and setup . . . . .	62
3.1.2	Results and discussion . . . . .	63
3.1.2.1	First experiment . . . . .	63
3.1.2.2	Second experiment . . . . .	64
3.1.3	Conclusions and future work . . . . .	65
3.2	Incremental relational association rule mining of educational data sets . . . . .	65
3.2.1	Computational experiments . . . . .	66
3.2.1.1	Case studies and data sets . . . . .	66
	First case study . . . . .	66
	Second case study . . . . .	67
3.2.1.2	Experimental results . . . . .	67
3.2.1.3	Comparison to related work . . . . .	70
3.2.2	Conclusions and future work . . . . .	71
<b>4</b>	<b>New supervised learning models for students' performance prediction</b>	<b>72</b>
4.1	An analysis of supervised learning methods for predicting students' performance in academic environments . . . . .	72
4.1.1	Methodology . . . . .	73
4.1.1.1	Data sets . . . . .	74
	Data analysis . . . . .	74
4.1.1.2	Performance measures . . . . .	76
	RMSE . . . . .	76
	F-score . . . . .	76
4.1.1.3	Experimental methodology . . . . .	77
4.1.2	Results and discussion . . . . .	77
4.1.2.1	Comparison to related work . . . . .	78
4.1.3	Conclusions and future work . . . . .	79
4.2	<i>SPRAR</i> : A relational association rule mining classification model applied for academic performance prediction . . . . .	80
4.2.1	Methodology . . . . .	81
4.2.1.1	Classification using <i>SPRAR</i> . . . . .	83

4.2.1.2	Evaluation measures . . . . .	83
4.2.2	Example of classification using <i>SPRAR</i> . . . . .	84
4.2.3	Results and discussion . . . . .	84
4.2.3.1	Data sets . . . . .	85
4.2.3.2	Results . . . . .	85
4.2.3.3	Discussion and comparison to related work . . . . .	87
4.2.4	Conclusions and further work . . . . .	89
4.3	A study on students' performance prediction using <i>SPRAR</i> . . . . .	89
4.3.1	Methodology . . . . .	90
4.3.1.1	First method proposed ( <i>score1</i> ) . . . . .	90
4.3.1.2	Second method proposed ( <i>score2</i> ) . . . . .	91
4.3.1.3	Third method proposed ( <i>score3</i> ) . . . . .	91
4.3.1.4	Example . . . . .	92
4.3.2	Experimental evaluation . . . . .	92
4.3.2.1	Data sets and experimental setup . . . . .	92
4.3.2.2	Results . . . . .	94
4.3.3	Discussion . . . . .	94
4.3.4	Conclusions and future work . . . . .	97
	<b>Conclusions</b>	<b>98</b>
	<b>Bibliography</b>	<b>100</b>

# Lista publicațiilor

Toate rangurile sunt listate conform clasificării revistelor <sup>1</sup> și conferințelor <sup>2</sup> din Informatică din 2014.

## Publicații în Web of Science - Science Citation Index Expanded

- [14] Czibula, G., Czibula, I.G., Miholca, D.L., **Crivei, M.L.**. A novel concurrent relational association rule mining approach. *Expert systems with Applications* 125(6), 2019, pp. 142–156. (indexată Web of Science, FI=3.768)

Rang **A**, 4 puncte.

## Publicații în Web of Science - Conference Proceedings Citation Index

- [25] Diana-Lucia Miholca, Gabriela Czibula, **Liana Maria Crivei**. *A new incremental relational association rules mining approach*. International Conference on Knowledge Based and Intelligent Information and Engineering Systems (KES), 2018, Procedia Computer Science, Volume 126, 2018, pp. 126–135. (indexată Web of Science)

Rang **B**, 4 puncte.

- [15] Gabriela Czibula, Andrei Mihai, **Liana Maria Crivei**. *SPRAR: A novel relational association rule mining classification model applied for academic performance prediction*. International Conference on Knowledge Based and Intelligent Information and Engineering Systems (KES), Procedia Computer Science, Volume 159, 2019, pp. 20–29. (indexată Web of Science)

Rang **B**, 4 puncte.

- [11] **Liana Maria Crivei**, Andrei Mihai, Gabriela Czibula. *A study on applying relational association rule mining based classification for predicting the academic performance of students*. The 12th International Conference on Knowledge Science, Engineering and Management (KSEM), LNAI 11775, 2019, pp. 287–300. (indexată Web of Science)

Rang **B**, 4 puncte.

- [12] **Liana Maria Crivei**, Gabriela Czibula, George Ciubotariu, Mariana Dindelegan. *Un-supervised learning based mining of academic data sets for students' performance analysis*. IEEE 14th International Symposium on Applied Computational Intelligence and Informatics, SACI 2020, pp. 11–16. (indexată Web of Science)

---

<sup>1</sup>[hfpop.ro/standarde/doctorat/2014-jurnale.pdf](http://hfpop.ro/standarde/doctorat/2014-jurnale.pdf)

<sup>2</sup><https://hfpop.ro/standarde/doctorat/2014-conferinte.pdf>

**Rang C, 1 punct.**

**Publicații în reviste și proceedings de conferințe indexate în baze de date internaționale**

- [13] **Liana Maria Crivei**, Vlad-Sebastian Ionescu, Gabriela Czibula. An analysis of supervised learning methods for predicting students' performance in academic environments. *ICIC Express Letters*, 13(3), 2019, pp. 181–189 (**indexată Scopus**)

**Rang C, 2 puncte.**

- [9] **Liana Maria Crivei**. Incremental relational association rule mining of educational data sets. *Studia Universitatis Babes-Bolyai Series Informatica*, 63(2), 2018, pp. 102–117. (**indexată Mathematical Reviews**)

**Rang D, 1 punct.**

- [8] George Ciubotariu, **Liana-Maria Crivei**. *Analysing the academic performance of students using unsupervised data mining*. *Studia Universitatis Babes-Bolyai Series Informatica*, Vol. 64(2), 2019, pp. 34–48. (**indexată Mathematical Reviews**)

**Rang D, 1 punct.**

- [10] **Liana-Maria Crivei**. *Using unsupervised learning models for analyzing students' academic performance*. PhD Coloquium at 20th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC 2018), 2018, pp. 1–4.

**Scorul publicațiilor: 21 puncte.**

# Introducere

Domeniul principal de cercetare al tezei de doctorat este extragerea de cunoștințe din date de natură academică. *Extragerea de cunoștințe din date de natura academică* (EDM) este un domeniu de cercetare în curs de dezvoltare care se concentrează pe analiza *mediilor educaționale* din perspectiva *extragerii cunoștințelor din date*. Scopul principal este de a înțelege mai bine fenomenele educaționale și de a descoperi, prin tehnici de extragere a datelor, șabloane ascunse relevante din datele educaționale. Tehnicile de *extragere a cunoștințelor din date* (DM) sunt aplicate în zilele noastre în diferite domenii, inclusiv medicină, bioinformatică, inginerie software, pentru a descoperi șabloane relevante în bazele de date mari, în special datorită potențialului lor de a descoperi informații ascunse din date.

Aplicarea tehnicilor de *învățare automată* în educație [1] a atras în mod continuu cercetătorii din domeniul *extragerii de cunoștințe din date de natura academică*. Scopul principal în acest domeniu este acela de a descoperi șabloane relevante în date care provin din diverse medii educaționale. Unul dintre scopurile *extragerii cunoștințelor din date educaționale* este de a oferi informații suplimentare asupra procesului de învățare a studenților și de a oferi astfel o mai bună înțelegere a fenomenelor educaționale. Extragerea tiparelor relevante din procesele educaționale ar fi de asemenea utilă pentru înțelegerea studenților și a modului lor de învățare, precum și pentru îmbunătățirea rezultatelor educaționale (de exemplu, rezultatelor învățării). Domeniul *extragerii cunoștințelor din date educaționale* a primit în ultima vreme o atenție considerabilă din partea comunității de cercetare, deoarece extragerea cunoștințelor ascunse din datele educaționale este de interes deosebit pentru instituțiile academice și este utilă și pentru îmbunătățirea metodologiilor lor de predare și a proceselor de învățare [29].

Extragerea *regulilor de asociere* (AR) reprezintă o tehnică importantă de analiză și de extragere a datelor [7] utilă în procesele de *învățare automată* pentru descoperirea unor tipare semnificative bazate pe reguli de asociere în seturile de date. *Regulile de asociere ordinale* (OAR) [6] au fost propuse ca o clasă particulară a *regulilor de asociere* care exprimă relații ordinale între atributele care caracterizează un set de date. *Regulile de asociere relațională* (RAR) [4, 3] au fost introduse ca o extensie a *regulilor de asociere ordinale*, fiind capabile să exprime diferite tipuri de relații non-ordinale între atributele datelor. Atât extragerea regulilor de asociere ordinală, cât și relațională, au fost aplicate cu succes pentru a rezolva probleme care variază de la sarcinile de extragere a datelor (de exemplu, detectarea defectelor de proiectare software [18] sau curățarea datelor [4]) până la clasificarea supervizată (de exemplu, predicția defectelor software [26]).

## Motivație

Datorită capacității lor de a exprima relații relevante în date, aplicarea regulilor de asociere relaționale (RAR) pe seturi de date academice ar putea oferi abordări noi și eficiente pentru rezolvarea diverselor sarcini EDM. Un exemplu concret în care RAR-urile



ar fi mult mai semnificative decât AR-urile clasice este legat de domeniul *extragerii de cunoștințe din date de natura academică* (EDM) care se referă la seturile de date din mediile educaționale. Într-un scenariu în care setul de date care trebuie extras conține informații privind performanța academică a studenților (de exemplu, fiecare student este caracterizat prin notele sale), RAR-urile ar putea să exprime relații relevante între notele studenților. Astfel, RAR-urile pot oferi informații suplimentare în ceea ce privește predicția performanței studenților.

Literatura EDM dezvăluie că DM este foarte utilă în domeniul educațional, în special atunci când se explorează mediul de învățare online [29]. În acest sens, datorită relevanței sale practice pentru extragerea de date *online*, problema *extragerii incrementale a regulilor de asociere relațională* (IRARM) în contextul unui set de date dinamic, la care se adaugă noi instanțe este atractivă pentru cercetătorii DM. Întrucât procesele de învățare în mediile educaționale sunt prin natura lor procese online, ideea investigării perspectivei RAR incrementale asupra EDM apare în mod natural.

Datorită NP-completitudinii sale, problema extragerii tuturor *regulilor de asociere relațională* interesante dintr-un set de date este dificilă din punct de vedere computațional. Pe măsură ce dimensionalitatea setului de date analizat crește, algoritmul clasic de *Descoperire a regulilor de asociere relațională* (DRAR) pentru extragerea RAR-urilor interesante eșuează să furnizeze setul de reguli într-un timp rezonabil. În scopul reducerii duratei generale de extragere, vom investiga o perspectivă de extragere a *regulilor de asociere relaționale concurente* folosind concurența pentru procesul de descoperire a RAR-urilor.

În cadrul cercetării EDM, dezvoltarea unor modele de învățare automată pentru analiza și predicția performanței studenților este de mare interes. Ca etapă naturală anterioară dezvoltării de noi modele de învățare automată pentru predicția performanței studenților, este investigată utilitatea aplicării unor *metode de învățare automată nesupervizată* în analiza datelor de performanță academică ale studenților. Ca pas ulterior, datorită capacității lor de a descoperi cunoștințe semnificative din date, RAR-urile vor fi explorate în contextul prezicerii performanței academice a studenților. Mai precis, sunt avute în vedere *modele de învățare automată* supervizată bazate pe reguli de asociere relațională.

Pentru a rezuma, cercetarea realizată în teza de doctorat va fi concentrată pe trei direcții principale: (1) contribuții conceptuale în domeniul extragerii RAR prin introducerea de noi abordări pentru extragerea *incrementală* și *concurentă* RAR; (2) cercetarea și evidențierea relevanței și eficienței extragerii *regulilor de asociere relațională* din seturi de date educaționale; și (3) dezvoltarea modelelor de *învățare automată* supervizate și nesupervizate pentru predicția performanței academice a studenților bazată pe RAR-uri și sisteme bazate pe reguli. Mai precis, obiectivele majore ale cercetării realizate în această teză sunt: utilizarea extragerii regulilor de asociere relațională în domeniul EDM; dezvoltarea de noi modele de învățare automată (bazate în principal pe extragerea RAR) pentru a prezice performanța academică a studenților; dezvoltarea unei versiuni concurente a algoritmului de extragere a RAR pentru a reduce timpul de calcul al procesului de extragere și investigarea unei versiuni incrementale a algoritmului de extragere RAR, care ar fi adecvată în scenariile de extragere online în care sunt adăugate noi instanțe la setul de date care urmează să fie extras.

## Structura tezei

Restul tezei este structurat după cum urmează.

Capitolul 1 prezintă principalele concepte privind EDM și extragerea regulilor de asociere,

împreună cu problemele pe care intenționăm să le abordăm din perspectivă computațională. Importanța EDM în contextul învățării automate este subliniată în secțiunea 1.1. Abordăm concepte de învățare supervizată și nesupervizată și prezentăm o imagine cuprinzătoare a literaturii asupra lucrărilor de cercetare existente în acest domeniu. De asemenea, includem o prezentare a literaturii despre predicția performanței instructorilor și alte studii EDM. Secțiunea 1.2 prezintă conceptele fundamentale de seturi frecvente și reguli de asociere clasice, precum și o descriere cuprinzătoare a regulilor de asociere relaționale, definiția acestora și un exemplu practic. În această secțiune este de asemenea prezentată o imagine de ansamblu a literaturii existente privind extragerea regulilor de asociere incrementale și concurente. Ultima parte a acestui capitol, secțiunea 1.3 introduce trei modele de învățare automată fundamentale: rețele neuronale artificiale, păduri aleatoare de arbori și rețele cu auto-organizare, care vor fi ulterior folosite pentru predicția performanței studenților.

Capitolul 2 prezintă cercetarea originală asupra extragerii incrementale a regulilor de asociere relațională în contextul seturilor de date dinamice, precum și extragerea concurentă a RAR. În secțiunea 2.1 introducem o nouă metodă incrementală numită *extragerea incrementală a regulilor de asociere relațională (IRARM)*. Prin *IRARM* intenționăm să ajustăm progresiv RAR-urile interesante descoperite într-un set de date original care este ulterior extins cu noi instanțe primite. Obiectivul nostru este obținerea unui timp de calcul redus în comparație cu cel obținut prin intermediul metodei *DRAR* aplicate de la zero pe setul de date extins. Mai mult, este prezentată metodologia utilizată pentru *IRARM*, cu algoritmul de calcul și o defalcare a celor două etape numite *filtering* și *extending*. Sunt efectuate experimentele pe două studii de caz pentru a evalua și evidenția eficiența timpului de extragere *IRARM*. Secțiunea 2.2 introduce o nouă abordare de extragere concurentă numită *CRAR extragerea concurentă a regulilor de asociere relațională*. Noua metodă folosește concurența pentru descoperirea RAR-urilor și scade considerabil timpul de extragere. Performanța *CRAR* din punct de vedere al timpului de calcul este evaluată și comparată cu algoritmul secvențial *DRAR*. O prezentare exhaustivă a algoritmului *CRAR* este inclusă împreună cu funcțiile pentru generarea candidaților (*candidate generation*), partiționare (*partitioning*) și generarea regulilor (*rule generation*) pe mai multe fire de execuție. Evaluările experimentale sunt efectuate pe nouă seturi de date disponibile public din diverse domenii: inginerie software, extragere de date educaționale și bioinformatică. Implementarea *CRAR* pe seturile de date menționate anterior demonstrează că timpul de rulare al algoritmului *CRAR* este semnificativ mai mic în comparație cu metoda *DRAR* aplicată de la zero.

Capitolul 3 prezintă contribuțiile și rezultatele originale privind dezvoltarea modelelor de învățare nesupervizate în contextul EDM. Secțiunea 3.1 explorează două modele de învățare automată nesupervizată, în special *rețele cu auto-organizare* și *extragerea regulilor de asociere relațională* pentru a evalua performanța studenților în medii academice. Analiza se realizează pe un set de date academice reale în care atributele fiecărei instanțe reprezintă performanța studenților la un anumit curs academic, de-a lungul unui semestru. Studiul nostru dezvăluie potențialul celor două tehnici de învățare nesupervizată de a identifica șabloane relevante și relații în cadrul datelor academice care sunt valoroase pentru predicția realizărilor academice ale studenților. În secțiunea 3.2 extragerea RAR clasică și metoda *IRARM* anterior introduse în secțiunea 2.1 sunt testate și evaluate în contextul datelor educaționale. Experimentele sunt efectuate având în vedere două studii de caz și un total de patru seturi de date academice. Studiul subliniază importanța și eficacitatea extragerii RAR (atât incrementale, cât și clasice) pentru descoperirea cunoștințelor și șabloanelor semnificative din datele educaționale.

Capitolul 4 descrie contribuțiile noastre în dezvoltarea modelelor de învățare supervizată

pentru predicția performanței studenților. Secțiunea 4.1 analizează metode supervizate utile pentru predicția performanței studenților în medii academice. Două modele de regresie, *păduri aleatoare de arbori* (RF) și *rețele neuronale artificiale* (ANN), sunt cercetate din punct de vedere al performanței predicției. Experimentele sunt efectuate pe trei seturi de date academice reale care conțin notele obținute de studenți în timpul semestrului la diferite cursuri la nivel de licență. Performanța predicției este măsurată folosind *rădăcina pătrată a erorii medii pătratice* (RMSE) și măsura *F - score*. Este prezentată, de asemenea, o discuție cu privire la rezultatele obținute, precum și o comparație cu alte lucrări din domeniu. Secțiunea 4.2 introduce un nou model de clasificare, *SPRAR* (*predicția performanței studenților folosind reguli de asociere relațională*) pentru predicția rezultatului final al unui student la o anumită disciplină academică (pe baza notelor primite în timpul semestrului) folosind *reguli de asociere relațională* (RAR). Experimentele sunt efectuate utilizând trei seturi reale de date academice. Performanța clasificatorului *SPRAR* pe studiile de caz considerate este comparată cu lucrări similare existente, fiind superioară rezultatelor obținute de abordări propuse în literatură. Un studiu comparativ asupra modelului de clasificare *SPRAR* introdus anterior este efectuat în continuare în secțiunea 4.3. În această secțiune sunt introduse trei noi scoruri de clasificare, care sunt utilizate în etapa de clasificare a *SPRAR*. Performanța lor este analizată pe aceleași seturi de date academice reale utilizate în cercetarea inițială din secțiunea 4.2 și este comparată cu rezultatele similare existente în literatura.

## Contribuții originale

Teza are trei contribuții majore, care sunt cuprinse în Capitolele 2, 3 și 4 după cum urmează:

1. Contribuții conceptuale în domeniul DM, în special în direcția exragerii *incrementale* și *concurente* a RAR, introduse în capitolul 2 și publicate în lucrările originale [25] și [14].
  - O metodă de extragere incrementală, numită *IRARM* *extragerea incrementală a regulilor de asociere relațională*, este introdusă în secțiunea 2.1 pentru a extrage eficient toate RAR-urile interesante dintr-un set de date care a fost mărit (adică s-au adăugat noi entități) [25]. *IRARM* adaptează setul tuturor RAR-urilor interesante descoperite anterior în setul de date inițiale pentru a calcula noul set al tuturor RAR-urilor care sunt interesante în setul de date extins. Algoritmul *IRARM* este evaluat experimental pe două seturi de date disponibile public. Experimentele efectuate subliniază faptul că *IRARM* este capabil să returneze setul complet de RAR-uri interesante mai eficient decât în cazul aplicării algoritmului de extragere de la zero pe setul de date extins. Lucrarea originală a fost publicată în [25].
  - În secțiunea 2.2 este introdusă o nouă abordare denumită *CRAR* care presupune extragerea concurentă de *reguli de asociere relaționale*. Eficacitatea extragerii este validată empiric pe nouă seturi de date disponibile public. Reducerea timpului de extragere atunci când se utilizează *CRAR* comparativ cu algoritmul secvențial *DRAR* subliniază faptul că acesta poate fi aplicat cu succes în diferite scenarii practice de extragere a datelor. Lucrarea originală a fost publicată în [14].
2. Noi modele de învățare nesupervizată în domeniul EDM, mai precis pentru *predicția*

*performanței studenților* (SPP) sunt introduse în capitolul 3 și au fost publicate în lucrările originale [9], [8] și [10].

- Modele de *învățare nesupervizată* pentru analiza datelor de performanța academică a studenților sunt propuse în secțiunea 3.1. Experimentele efectuate utilizând un set de date academice reale au evidențiat potențialul modelelor de *învățare nesupervizată* pentru descoperirea cunoștințelor semnificative în cadrul datelor educaționale. Lucrările originale a fost publicate în [10] și [8].
  - Am investigat și am evidențiat relevanța și eficacitatea *extragerii regulilor de asociere relațională* din seturile de date educaționale în secțiunea 3.2. *Extragerea incrementală a regulilor de asociere relațională* folosind abordarea *IRARM* introdusă anterior în Capitolul 2 este aplicată în contextul EDM, subliniind relevanța și eficacitatea ei. Lucrarea originală a fost publicată în [9].
3. Noi *modele de învățare automată* supervizată pentru SPP, construite pe sisteme bazate pe reguli și RAR-uri sunt dezvoltate în Capitolul 4 și au fost publicate în lucrările originale [13], [15] și [11].
- În secțiunea 4.1 realizăm un studiu asupra aplicării a două modele de regresie, *păduri aleatoare de arbori* (RF) și *rețele neuronale artificiale* (ANN), pentru predicția performanței academice a studenților. Pentru fiecare regresor, două modele de calcul sunt introduse pentru predicția notei finale la examen pentru un student, bazată pe notele sale primite în timpul unui semestru academic. Lucrarea originală a fost publicată în [13].
  - În secțiunea 4.2 introducem un nou model de clasificare, *SPRAR* (*predicția performanței studenților folosind reguli de asociere relațională*) pentru predicția, pe baza notelor primite în timpul semestrului, rezultatului final al unui student la o anumită disciplină academică folosind *regulile de asociere relațională* (RAR). Experimentele sunt efectuate pe trei seturi de date academice reale colectate de la Universitatea Babeș-Bolyai din România. Performanța clasificatorului *SPRAR* pe studiile de caz considerate este comparată cu cea a clasificatorilor existenți pentru predicția performanței studenților, fiind superioară acestora. Lucrarea originală a fost publicată în [15].
  - Un studiu comparativ asupra modelului de clasificare introdus anterior *SPRAR* este efectuat în continuare în secțiunea 4.3. Trei noi scoruri de clasificare sunt introduse în această secțiune și sunt utilizate în etapa de clasificare a *SPRAR*. Performanța lor este analizată pe aceleași seturi reale de date academice utilizate în lucrarea originală [15] și comparată cu rezultatele similare existente în literatura de specialitate. Lucrarea originală a fost publicată în [11].

# Capitolul 1

## Concepte introductive

În acest capitol prezentăm conceptele fundamentale cunoscute legate de domeniul *extragerii de cunoștințe din date de natură academică* (EDM), precum și principalele concepte legate de modelele de calcul pe care le vizăm în cercetarea noastră.

### 1.1 Extragerea de cunoștințe din date de natură academică (EDM)

Tehnicile *extragerii de cunoștințe din date* (DM) sunt des aplicate în prezent în diverse domenii, inclusiv medicină, inginerie software, bioinformatică, pentru a descoperi șabloane relevante în bazele de date mari, în special datorită potențialului acestora de a descoperi informații ascunse în seturile de date.

Aplicarea tehnicilor DM în educație [1] a atras cercetători atât din DM cât și din cercetarea educațională și astfel a apărut o nouă disciplină de cercetare interdisciplinară cunoscută sub numele de *extragerea de cunoștințe din date de natură academică* (EDM). Principalul accent în EDM este de a dezvolta metode de extragere a cunoștințelor din date care provin din diverse sisteme de informații și medii educaționale. Prin extragerea de cunoștințe din seturile de date educaționale, scopul EDM este de a înțelege mai bine procesul de învățare al studenților și, astfel, de a oferi informații suplimentare asupra fenomenelor legate de educație.

EDM este un domeniu de cercetare în dezvoltare în care conceptul care stă la baza lui este acela de a aduce perspectiva *extragerii de cunoștințe din date* în *medii educaționale*. Scopul principal este de a înțelege mai bine fenomenele legate de educație prin extragerea tiparelor ascunse semnificative din datele educaționale folosind tehnici de extragere a datelor.

Extragerea șabloanelor relevante din procesele educaționale ar fi de asemenea utilă pentru înțelegerea studenților și a modului în care învață, precum și îmbunătățirea rezultatelor educaționale (de exemplu, rezultatele învățării). EDM a primit în ultimul timp o atenție considerabilă din partea comunității de cercetare, deoarece extragerea cunoștințelor ascunse din datele educaționale este de interes deosebit pentru instituțiile academice și este de asemenea utilă pentru îmbunătățirea metodologiilor de predare și a proceselor de învățare [29].

### 1.2 Extragerea regulilor de asociere din date

În această secțiune prezentăm conceptele fundamentale legate de extragerea *regulilor de asociere relațională* (RAR) [3]. Extragerea *regulilor de asociere* (AR) reprezintă o importantă metodă de analiză a datelor și tehnică de extragere a cunoștințelor [7] aplicată în

diferite sarcini legate de *învățarea automată* pentru descoperirea modelelor relevante bazate pe reguli în seturile de date. *Regulile de asociere ordinală* (OAR) [6] au fost propuse ca o clasă particulară de AR care exprimă relații ordinale între atributele care caracterizează un set de date. *Regulile de asociere relațională* (RAR) [5, 3] au fost introduse ca o extensie a OAR-urilor capabile să exprime diferite tipuri de relațiile neordinale între atributele datelor.

*Regulile de asociere* (AR) clasice [33] nu iau în considerare relațiile care pot apărea între valorile atributelor, ci numai co-ocurența lor. Spre deosebire de acestea, *regulile de asociere ordinală* (OAR) [24] exprimă relații ordinale care apar frecvent în date. Dar pot exista și relații informative între valorile atributelor care nu sunt ordinale.

Noțiunea de *regulă de asociere relațională* (RAR) este definită în următoarele paragrafe.

Considerăm o mulțime  $\mathcal{D} = \{d_1, d_2, \dots, d_n\}$  de *instanțe* sau *înregistrări*. Fie  $\Omega = (a_1, \dots, a_m)$  un șir de  $m$  atribute care caracterizează fiecare instanță din setul de date  $\mathcal{D}$ . Fiecare atribut  $a_i$  ia valori într-un domeniu nevid și non-fuzzy  $\Delta_i$ , care conține și o valoare *nulă* (*vidă*). Notăm cu  $\Psi(d_j, a_i)$  valoarea atributului  $a_i$  pentru o instanță  $d_j$ .

Notăm cu  $\mathcal{T}$  mulțimea tuturor relațiilor posibile (care nu sunt neapărat ordinale) care pot fi definite între două domenii  $\Delta_i$  și  $\Delta_j$ .

**Definiție 1.** O *regulă de asociere relațională* [3] este o expresie  $(a_{i_1}, a_{i_2}, a_{i_3}, \dots, a_{i_h}) \Rightarrow (a_{i_1}\tau_1 a_{i_2}\tau_2 a_{i_3} \dots \tau_{h-1} a_{i_h})$ , unde  $\{a_{i_1}, a_{i_2}, a_{i_3}, \dots, a_{i_h}\} \subseteq \Omega$ ,  $a_{i_k} \neq a_{i_p}$ ,  $k, p = 1, \dots, h$ ,  $k \neq p$  și  $\tau_k \in \mathcal{T}$  este o relație pe  $\Delta_{i_k} \times \Delta_{i_{k+1}}$ ,  $\Delta_{i_k}$  fiind domeniul atributului  $a_{i_k}$ .

- Dacă  $a_{i_1}, a_{i_2}, \dots, a_{i_h}$  nu lipsesc din  $m$  instanțe din setul de date atunci numim  $s = \frac{m}{n}$  *suportul* regulii.
- Dacă notăm cu  $D' \subseteq \mathcal{D}$  setul de instanțe unde  $a_{i_1}, a_{i_2}, a_{i_3}, \dots, a_{i_h}$  nu lipsesc și toate relațiile

$$\Psi(d_j, a_{i_1})\tau_1\Psi(d_j, a_{i_2}), \Psi(d_j, a_{i_2})\tau_2\Psi(d_j, a_{i_3}), \dots, \Psi(d_j, a_{i_{h-1}})\tau_{h-1}\Psi(d_j, a_{i_h})$$

au loc pentru fiecare instanță  $d$  din  $D'$  atunci numim  $c = \frac{|D'|}{n}$  *confidența* regulii.

RAR-urile *interesante* au fost definite în [3] ca acele reguli care au atât suportul cât și confidența mai mari sau egale cu un prag minim specificat. Pentru extragerea RAR-urilor interesante, un algoritm de tip Apriori, numit *DRAR* (descoperirea regulilor de asociere relațională) a fost propus în [16] ca o extensie a algoritmului *DOAR* introdus în [6] pentru descoperirea OAR-urilor.

Algoritmul *DRAR* constă în generarea iterativă a RAR-urilor, începând cu regulile de lungime 2, apoi continuând cu regulile de lungime 3, ș.a.m.d. Setul de reguli de lungime 2 este filtrat pentru a determina regulile interesante (adică acele reguli care verifică suportul minim și cerințele legate de confidența). La un pas dat, pentru determinarea RAR-urilor de o anumită lungime  $l$ , începem de la setul de RAR-uri interesante de lungime  $l - 1$  generate la pasul anterior. Acest set este folosit pentru a genera prin alăturare noi posibile RAR-uri interesante, numite reguli candidat, care vor fi filtrate pentru a păstra doar regulile care sunt interesante. După ce se generează RAR-uri interesante de lungime  $l$ , procesul iterativ continuă cu generarea regulilor de lungime  $l + 1$ . Procesul se oprește atunci când nu mai sunt descoperite noi RAR-uri interesante.

### 1.3 Modelele de învățare automată folosite

Această secțiune prezintă modelele de învățare supervizată și nesupervizată utilizate în teză pentru predicția performanței academice a studenților: *rețele neuronale artificiale*, *păduri aleatoare de arbori* și *rețele cu auto-organizare*.

#### 1.3.1 Rețele neuronale artificiale

*Rețelele neuronale artificiale* sunt utilizate pe scară largă ca modele de învățare supervizată pentru diverse aplicații cum ar fi recunoașterea tiparelor din date, recunoașterea vorbirii [31], predicția [23], identificarea și controlul sistemelor. În mod similar sistemelor neuronale biologice, rețelele neuronale artificiale [28] constau dintr-un set de unități de calcul dens interconectate, numite *neuroni*.

O *rețea neuronală artificială* (ANN) [35] este un sistem adaptiv care învață o corespondență (o funcție de intrare/ieșire) din date, prin reglarea autonomă a parametrilor sistemului în timpul *fazei de antrenament*. Parametrii rețelei obținuți după finalizarea antrenamentului sunt folosiți în continuare pentru a rezolva problema curentă (*faza de testare*).

Într-un scenariu de învățare supervizată, o instanță de intrare este prezentată rețelei neuronale artificiale împreună cu ținta corespunzătoare de ieșire [27]. Aceste exemple de intrare-ieșire sunt adesea furnizate de un supervisor extern. O eroare este reprezentată prin diferența dintre ieșirea dorită și ieșirea sistemului. Această informație de eroare este propagată înapoi în rețea și parametrii sistemului sunt ajustați prin *regula de învățare*. Rețeaua neuronală artificială este construită folosind o procedură iterativă al cărei obiectiv este de a reduce la minimum eroarea obținută, care este denumită în mod obișnuit *regulă de învățare*. Acest proces se repetă până când se obține o performanță acceptabilă. Algoritmul pentru construirea modelului ANN se numește algoritmul de *propagare înapoi* (*backpropagation*).

#### 1.3.2 Păduri aleatoare de arbori

*Pădurile aleatoare de arbori* (RF) [2] sunt o metodă de învățare bazată pe ansambluri constând din agregarea mai multor predictorii de tip arbore de decizie. În timpul procesului de construcție al fiecărui arbore individual, se consideră pentru analiză numai o submulțime aleatoare de caracteristici și o submulțime aleatoare de exemple de instruire. În acest fel, se evită supraantrenarea și se obține o stabilitate mai bună pentru generalizare. Eroarea de generalizare a pădurilor converge la limită, deoarece numărul arborilor din pădure devine mare, datorită legii numerelor mari. Această eroare depinde de puterea arborilor individuali din pădure și de corelația dintre ei. Folosind o selecție aleatoare de funcții pentru a împărți fiecare nod se produc rate de eroare care se compară favorabil cu algoritmul Adaboost [34], dar sunt mai robuste în ceea ce privește zgomotul. Fiind construite pe arbori de decizie, pădurile aleatoare de arbori pot fi utilizate în probleme de clasificare și regresie.

#### 1.3.3 Rețele cu auto-organizare

*Rețelele cu auto-organizare* (SOM) sunt modele de învățare nesupervizată, de asemenea cunoscute în literatura de specialitate sub numele de *hărți topografice Kohonen*, au fost introduse de Teuvo Kohonen și sunt utilizate pe scară largă ca instrumente pentru vizualizarea datelor de dimensiuni mari. SOM-urile sunt legate de *rețelele neuronale artificiale* (ANN) din literatură și de *învățarea competitivă*. În *învățarea competitivă*, neuronii de ieșire concurează pentru a fi activați. O *rețea cu auto-organizare* [32] este instruită folosind un algoritm de

învățare nesupervizată (algoritmul lui Kohonen) pentru a realiza o corespondență, folosind o asociere neliniară, între spațiul de intrare continuu format din instanțe multidimensionale și un spațiu de ieșire discret (de obicei bidimensional) numit *hartă* [19]. *Corespondența care păstrează topologia* este principala caracteristică oferită de un SOM. Aceasta înseamnă că instanțele de intrare care sunt aproape una de cealaltă în spațiul de intrare vor fi, de asemenea, aproape una de cealaltă pe hartă (spațiul de ieșire). Astfel, un SOM este capabil să furnizeze grupe de instanțe de date similare [22].

Rețelele cu auto-organizare sunt o clasă de rețele neuronale artificiale care utilizează o funcție de vecinătate pentru a menține caracteristicile topologice ale spațiului de intrare. SOM-urile folosesc învățarea nesupervizată pentru antrenament și ele sunt unul dintre cele mai inovatoare tipuri de rețele neuronale artificiale din literatura de specialitate actuală. Algoritmul SOM folosește nodurile (neuronii) din rețea în manieră nesupervizată și competitivă, neuronii concurând între ei pentru a răspunde la un subset de date primite. Rețeaua cu auto-organizare păstrează topologia spațiului de intrare, care este pus în corepondență cu unul de dimensiune mică, de obicei un spațiu bidimensional, denumit hartă. Reducerea dimensionalității prin păstrarea topologiei înseamnă că punctele de date de intrare vecine sau similare vor fi transformate în unități apropiate în spațiul de ieșire [21]. Prin urmare, SOM-urile pot fi privite ca un instrument de vizualizare a datelor de dimensiune mare, precum și ca un instrument de grupare.



## Capitolul 2

# Noi abordări în extragerea cunoștințelor din date

Acest capitol introduce contribuțiile tezei noastre de doctorat în domeniul *extragerii de cunoștințe din date*, mai precis în dezvoltarea de noi abordări privind extragerea *incrementală* și *concurrentă* a regulilor de asociere relațională.

### 2.1 O abordare privind extragerea incrementală a regulilor de asociere relațională

În primul rând, o nouă metodă incrementală numită *IRARM* (*extragerea incrementală a regulilor de asociere relațională*) este introdusă și evaluată în termeni de eficiență a timpului de execuție. *IRARM* își propune să identifice RAR-urile interesante dintr-un set de date de natură dinamică, ceea ce înseamnă că noi instanțe sunt adăugate treptat la setul de date inițial. Caracteristic pentru *IRARM* este că adaptează RAR-urile relevante descoperite anterior fără a necesita o rescansare a întregului set de date actualizat, spre deosebire de metoda *DRAR*. Eficacitatea *IRARM* este evaluată prin experimente efectuate pe două seturi de date disponibile public, iar rezultatele evidențiază capacitatea *IRARM* de a îmbunătăți semnificativ timpul de extragere comparativ cu *DRAR*.

Pe viitor, vor fi luate în considerare mai multe seturi de date, pentru a extinde și mai mult evaluarea *IRARM*. De asemenea, ne propunem să investigăm condițiile exacte în care este mai eficient să se aplice metoda incrementală *IRARM* decât să se ruleze de la zero algoritmul offline *DRAR*. Mai mult, ne propunem să utilizăm *IRARM* în sarcini de extragere a datelor dinamice concrete, cum ar fi *predicția defectelor software* incrementală și extragerea incrementală de cunoștințe din date educaționale (*extragerea cunoștințelor din date educaționale* [30]). În plus, intenționăm să investigăm o abordare incrementală adaptivă pentru extragerea *regulilor de asociere relațională*, în care atât funcții noi, cât și obiecte noi sunt adăugate la setul de date extrase. Pentru îmbunătățirea eficienței procesului de extragere, va fi de asemenea avută în vedere o abordare distribuită.

## 2.2 O abordare privind extragerea concurentă a regulilor de asociere relațională

A doua secțiune a acestui capitol introduce o nouă abordare privind extragerea concurentă a RAR, numită *CRAR* (*extragerea concurentă a regulilor de asociere relațională*). Programarea concurentă este frecvent utilizată pentru a îmbunătăți performanța, exploatarea resurselor și, de asemenea, pentru a dezvolta sisteme tolerante la erori și scalabile. Algoritmul *CRAR* a fost dezvoltat pe aceste principii cu scopul de a furniza un timp de extragere mai bun în procesul de descoperire RAR. Efectuăm o evaluare experimentală detaliată a *CRAR*, realizând experimente pe nouă seturi de date disponibile public din diverse domenii: extragerea datelor educaționale, predicția defectelor software, bioinformatică. *CRAR* oferă o valoare semnificativ îmbunătățită a timpului de execuție, prin reducerea timpului de extragere în medie cu 52,3 % în comparație cu algoritmul clasic *DRAR*. Abordările prezentate în capitolul curent sunt lucrări de cercetare originale publicate în [25] și [14].

Pe viitor, ne propunem să investigăm optimizări ale *CRAR*, cum ar fi creșterea volumului de calcule efectuate în paralel, în special calculul RAR-urilor binare, care ar crește *eficiența scalării și accelerarea paralelizării* algoritmilor. Pentru a evidenția mai bine utilitatea *CRAR*, intenționăm să utilizăm extragerea regulilor de asociere relațională pentru probleme concrete de clasificare supervizată din *inginerie software bazată pe căutare* (de ex. *predicția defectelor software* [27]) și *extragerea de cunoștințe din date educaționale* [9] (cum ar fi *predicția performanței studenților*). Cercetări suplimentare vor fi, de asemenea, axate pe extinderea abordării introduse de *CRAR* în această lucrare pentru extragerea *regulilor de asociere relațională graduală* [17] în loc de RAR-urile clasice (non-fuzzy sau negraduale).

## Capitolul 3

# Noi modele de învățare nesupervizată în extragerea cunoștințelor din date de natură academică

În acest capitol prezentăm contribuțiile noastre la dezvoltarea modelelor de învățare nesupervizată pentru analiza datelor de performanță academică a studenților. Secțiunea 3.1 evaluează rețele cu auto-organizare și modele de extragere a regulilor de asociere relațională în ceea ce privește potențialul lor de a detecta tipare relevante din seturi de date academice ale studenților, mai precis performanțele studenților reflectate în notele obținute pe parcursul unui semestru. Secțiunea 3.2 investighează algoritmul *IRARM* (*extragerea incrementală a regulilor de asociere relațională*) în contextul datelor educaționale. Experimentele sunt efectuate pe patru seturi de date și evidențiază potențialul și relevanța extragerii incrementale a regulilor de asociere relațională în descoperirea cunoștințelor semnificative din datele academice. Abordările prezentate în capitolul curent sunt lucrări de cercetare originale publicate în [9], [8] și [10].

### 3.1 Utilizarea modelelor de învățare nesupervizată pentru analiza performanței academice a studenților

În această secțiune investigăm utilitatea a două tehnici de *învățare automată nesupervizată* (*rețele cu auto-organizare și extragerea regulilor de asociere relațională din cunoștințe*) în analiza performanței academice a studenților, cu scopul mai larg de a dezvolta noi modele de învățare automată pentru predicția performanței studenților. Rezultatele experimentale obținute prin aplicarea modelelor de învățare nesupervizată menționate mai sus pe un set de date reale colectate de la Universitatea Babeș-Bolyai subliniază eficacitatea lor în extragerea relațiilor și regulilor relevante din datele educaționale care poate fi utilă pentru a prezice performanța academică a studenților.

Abordarea din această secțiune a fost introdusă în lucrările noastre originale [10] și [8].

RAR-urile nu au fost utilizate până în prezent pentru analiza seturilor de date academice. În conformitate cu cunoștințele noastre, un studiu similar nu a fost efectuat în literatura de specialitate.

Această secțiune a examinat două modele de învățare nesupervizată, *rețele cu auto-organizare* și *extragerea regulilor de asociere relațională din cunoștințe*, în contextul analizei seturilor de date referitoare la performanța studenților. Experimentele efectuate pe un set de date real colectate de la Universitatea Babeș-Bolyai, România, au evidențiat potențialul instrumentelor de extragere a cunoștințelor bazate pe învățare nesupervizată pentru a detecta șabloane relevante cu privire la performanța academică a studenților. Studiul curent reprezintă punctul de plecare al cercetărilor noastre ulterioare în vederea obținerii unei mai bune înțelegeri a proceselor de învățare ale studenților și a dezvoltării unor noi metode de învățare automată pentru precizarea performanțelor academice ale studenților.

Cercetări viitoare vor fi realizate pentru a extinde evaluarea experimentală la alte seturi de date academice și pentru a interpreta RAR-urile interesante extrase. Pentru creșterea performanței procesului de învățare nesupervizat, vom continua cercetarea metodelor de detectare a anomaliilor și a instanțelor aberante în date. În plus, vom analiza o fază post-procesare pentru filtrarea setului de RAR-uri extrase pentru eliminarea regulilor care se suprapun în mai multe clase. În plus, un model de clasificare bazat pe RAR va fi avut în vedere pentru a prezice performanța academică a studenților. De asemenea, vom cerceta alte tehnici de învățare automată bazate pe *clasificare* și *regresie* pentru predicția performanței studenților (de exemplu păduri aleatoare de arbori).

### 3.2 Extragerea incrementală a regulilor de asociere relațională din seturi de date de natură academică

*Extragerea incrementală a regulilor de asociere relațională (IRARM)* a fost introdusă anterior ca o metodă eficientă de *extragere online a cunoștințelor din date* pentru extragerea dinamică a *regulilor de asociere relațională* (RAR-uri) interesante dintr-un set de date dinamic (și anume un set de date care este extins cu noi instanțe). Studiul realizat în această secțiune a fost introdus în lucrarea originală [9] și are ca scop să sublinieze eficacitatea RAR și *IRARM* ca metode de extragere în contexte de *extragere a cunoștințelor din date de natură academică*. Experimentele efectuate pe diverse seturi de date academice evidențiază potențialul utilizării *regulilor de asociere relațională* pentru descoperirea cunoștințelor relevante din date academice.

În această secțiune abordăm problema *extragerii incrementale a regulilor de asociere relațională (IRARM)* în contextul EDM. Procesul de extragere incrementală a RAR-urilor este adecvat în special pentru scenariile *DM online*, unde setul de date care urmează să fie extras este dinamic și astfel este extins continuu cu fluxuri de date sosite în timp real. În astfel de situații, abordarea *IRARM* urmărește adaptarea progresivă a RAR-urilor interesante identificate într-un set de date, atunci când este extins cu noi instanțe. Întrucât procesele de învățare din mediile educaționale sunt prin natură procese online, ideea investigării perspectivei *IRARM* în EDM apare în mod natural. Literatura EDM dezvăluie și ea că *DM* este foarte utilă în domeniul educațional, în special atunci când se explorează mediul de învățare online [29].

Contribuția secțiunii este rezumată după cum urmează. În primul rând, subliniem relevanța extragerii RAR-urilor în domeniul *extragerii cunoștințelor din date de natură academică* (EDM) cu scopul de a descoperi tipare semnificative în cadrul seturilor de date educaționale. În al doilea rând, extindem evaluarea experimentală a *extragerii incrementale a regulilor de asociere relațională (IRARM)* [25] pentru mai multe studii de caz EDM. Eficacitatea

*IRARM* este subliniată prin reducerea timpului de extragere obținut atunci când se utilizează *IRARM* comparativ cu extragerea RAR-urilor de la zero când setul de date este extins cu noi instanțe. Studiul realizat în această lucrare este inedit în literatura de specialitate EDM, întrucât nici abordarea prin extragerea incrementală a RAR-urilor și nici cea prin extragerea clasică a acestora nu au fost aplicate pe seturi de date academice până în prezent.

În această secțiune am investigat aplicarea extragerii clasice și incrementale a RAR-urilor pentru descoperirea cunoștințelor în seturi de date din medii educaționale, cu scopul de a descoperi șabloane semnificative din seturi de date educaționale. Relevanța descoperirii RAR-urilor în seturile de date academice a fost accentuată în contextul procesului de învățare a studenților, oferind informații suplimentare asupra fenomenelor legate de educație. În plus, eficacitatea extragerii *incrementale* a RAR-urilor în scenariile EDM online a fost evidențiată prin mai multe studii de caz.

Cercetări viitoare vor fi efectuate pentru a extinde evaluarea experimentală a *IRARM* la alte sarcini EDM, pentru a-i testa în continuare performanța. Va fi cercetată extragerea *incrementală adaptativă* a RAR-urilor pentru seturi de date academice, în cazul în care la setul de date sunt adăugate atât instanțe noi, cât și funcții noi. În plus, intenționăm să aplicăm RAR, RAR gradual [17] și algoritmul de extragere *IRARM* în scenarii EDM de învățare supervizată, cum ar fi prezicerea performanței academice a studenților.

## Capitolul 4

# Modele de învățare supervizată pentru predicția performanței studenților

În acest capitol prezentăm contribuțiile noastre în dezvoltarea modelelor de învățare supervizată pentru predicția performanței studenților.

Efectuăm în [13] o analiză a metodelor de învățare supervizată pentru predicția performanței studenților în medii academice. Experimentele sunt efectuate pe trei seturi de date academice reale și investighează eficacitatea a două modele de regresie *păduri aleatoare de arbori* (RF) și *rețele neuronale artificiale* (ANN) în predicția performanței academice a studenților. În lucrarea de cercetare [15] introducem un nou model de clasificare numit *SPRAR* (*predicția performanței studenților folosind reguli de asociere relațională*). *SPRAR* este utilizat pentru a prezice rezultatele finale obținute de studenți la un anumit curs academic, pe baza notelor lor primite în timpul semestrului. Evaluarea este realizată pe trei seturi de date academice reale colectate de la Universitatea Babeș-Bolyai din România. Se compară clasificatorul *SPRAR* din punct de vedere al performanței cu clasificatorii supervizați din lucrări existente. Experimentele efectuate evidențiază că *SPRAR* depășește predictorii performanței studenților propuși anterior într-o mare majoritate a cazurilor. Lucrarea originală fost publicată în [15].

Abordările prezentate în capitolul curent sunt lucrări de cercetare originale publicate în [13], [15] și [11].

### 4.1 O analiză a metodelor de învățare supervizată pentru predicția performanței studenților în medii academice

Această secțiune analizează performanța modelelor de clasificare și regresie supervizate pentru predicția performanței academice a studenților. Investigăm eficiența diferitelor modele de învățare automată supervizată și ne propunem să prezicem nota finală a studenților la o anumită disciplină academică pe baza performanței lor din timpul semestrului. Prezentarea din această secțiune se bazează pe lucrarea originală [13].

În această secțiune realizăm un studiu privind aplicarea a două modele de regresie, *păduri aleatoare de arbori* (RF) și *rețele neuronale artificiale* (ANN), pentru predicția performanței academice a studenților. Pentru fiecare regresor, sunt prezentate câte două modele de calcul (unul de regresie și unul de clasificare) pentru a prezice nota studentului la examenul final pe baza notelor sale primite în timpul unui semestru academic. Rezultatele obținute pe trei

seturi de date reale colectate de la Universitatea Babeș-Bolyai, România dezvăluie că regresorii supervizați sunt utili pentru identificarea mijloacelor de creștere a calității proceselor educaționale. În conformitate cu cunoștințele noastre, un studiu similar nu a fost realizat până acum în literatura EDM. Clasificatorii RF au fost aplicați în literatură, dar în alte sarcini decât cea avută în vedere în această lucrare.

Pentru a rezuma, scopul studiului realizat în această secțiune este de a răspunde la următoarele întrebări de cercetare:

- RQ1** Care este potențialul modelelor de regresie supervizată pentru a prezice nota de la examenul final al studenților pe baza notelor primite în timpul semestrului?
- RQ2** În ce măsură selectarea caracteristicilor îmbunătățește acuratețea prezicerii performanței studenților?
- RQ3** Cum se compară modelele de învățare supervizată utilizate în această lucrare cu alte lucrări conexe din literatura de specialitate în termeni de predicție a performanței?

Studiul din această secțiune a fost realizat cu scopul de a evidenția eficacitatea modelelor de regresie supervizată în prezicerea performanței academice a studenților.

Experimentele efectuate pe trei seturi de date conținând date academice reale colectate de la o universitate din România au subliniat că, în general, RF sunt cel mai bun model de regresie pentru a prezice nota finală a studenților pe baza notelor primite în timpul unui semestru academic. Pentru experimentul efectuat am observat că regresorul SVM furnizează rezultate ușor mai bune decât RF. Studiul a relevat, de asemenea, dificultatea sarcinii de predicție a performanței studenților și importanța creșterii numărului de caracteristici utilizate în procesul de învățare.

Cercetările viitoare vor urmări adăugarea mai multor caracteristici la sarcinile noastre de învățare, precum și aplicarea unor tehnici de preprocesare pentru a detecta instanțele aberante din date. De asemenea, ne propunem să investigăm utilizarea altor modele de învățare pentru predicția performanțelor finale ale studenților, cum ar fi un clasificator bazat pe *extragerea regulilor de asociere relațională*.

## 4.2 *SPRAR*: Un model de clasificare bazat pe extragerea regulilor de asociere relațională aplicat pentru predicția performanței academice

Această secțiune abordează problema predicției performanței academice a studenților, o problemă intens investigată în literatura privind extragerea cunoștințelor din date de natură academică (EDM). Pentru o mai bună înțelegere a fenomenelor legate de educație, există un interes continuu în aplicarea metodelor de învățare *supervizată* și *nesupervizată* pentru obținerea unor informații suplimentare despre procesul de învățare al studenților. Problema prezicerii faptului că un student va promova sau nu la o anumită disciplină academică bazată pe notele studenților din timpul semestrului este una dificilă, și depinde foarte mult de diverse condiții, cum ar fi disciplina (cursul), numărul de examene pe semestru, profesorii și exigențele lor. Problema are o relevanță practică majoră în mediile educaționale, deoarece ar putea oferi feedback relevant studenților care ar putea să nu promoveze la un anumit curs academic. Având un astfel de sfat în timpul semestrului, studenții vor avea posibilitatea de a preveni un posibil eșec academic. Propunem un nou model de clasificare, *SPRAR* (*predicția*

performanței studenților folosind reguli de asociere relațională) pentru a prezice rezultatul final al unui student la o anumită disciplină academică folosind reguli de asociere relațională (RAR). Abordarea *SPRAR* a fost introdusă în lucrarea originală [15]. Performanța clasificatorului *SPRAR* în studiile de caz considerate este comparată cu cea din lucrările existente, fiind superioară predictorilor de performanță a studenților propuși anterior. Trebuie să notăm generalitatea clasificatorului *SPRAR*, care nu este specific sarcinii de predicție a performanței studenților. *SPRAR* poate fi adaptat cu ușurință la alte probleme de clasificare.

Modelul de clasificare introdus este unul binar (adică există doar două clase care pot fi prezise: *va promova* sau *nu va promova*), dar modelul propus este unul general, și poate fi extins pentru o problema de clasificare cu clase multiple (de exemplu, predicția notei finale a studentului). Dintr-o perspectivă de învățare supervizată, problema predicției finalizării cu succes a unui curs în timpul unui semestru academic este o problemă dificilă, în special datorită naturii dezechilibrate a datelor de antrenament (de exemplu numărul de studenți care au promovat examenul este în general mult mai mare decât numărul de studenți care *nu au promovat* examenul). Experimente efectuate pe trei seturi de date reale de la Universitatea Babeș-Bolyai, România evidențiază eficacitatea clasificării folosind *SPRAR*. Literatura de specialitate cu privire la predicția performanței studenților relevă faptul că utilizarea RAR-urilor în prezicerea performanței academice a studenților reprezintă o abordare inedită.

În concluzie, scopul cercetării efectuate este propunerea unui model de clasificare *SPRAR* bazat pe extragerea RAR pentru a prezice performanța academică a studenților, ca dovadă a conceptului. În acest scop, dovada conceptului ia în considerare doar trei seturi de date de dimensiuni medii pentru a evidenția faptul că *SPRAR* este potrivit pentru problema abordată. Dacă acest lucru este valabil, studiul aplicării *SPRAR* pentru predicția performanței studenților poate fi extins în continuare la o scară mai mare.

Cercetările viitoare vor urmări extinderea evaluării experimentale a *SPRAR* la alte studii de caz. În plus, ne propunem să generalizăm clasificatorul binar *SPRAR* la un clasificator cu mai multe clase, astfel încât să putem prezice nota finală a studenților la un anumit curs academic. În ceea ce privește procesul de descoperire a regulilor de asociere relațională, intenționăm să extindem modelul nostru, luând în considerare reguli de *asociere relațională graduală* [17]. Vor fi studiate în continuare măsuri alternative pentru definirea probabilităților folosite în procesul de clasificare a unei instanțe, precum și ideea folosirii unui ansamblu de clasificatori *SPRAR* pentru îmbunătățirea performanței predictive. *SPRAR* poate fi extins și pentru gestionarea relațiilor  $n$ -are dintre valorile atributelor.

### 4.3 Un studiu privind predicția performanței studenților folosind *SPRAR*

Această secțiune analizează modelul de clasificare *SPRAR* (predicția performanței studenților folosind reguli de asociere relațională) introdus în [15], Secțiunea 4.2, pentru a prezice rezultatele academice ale studenților utilizând reguli de asociere relațională. Trei noi scoruri de clasificare sunt introduse în acest capitol și sunt utilizate în etapa de clasificare a *SPRAR*. Performanța lor este analizată folosind trei seturi de date academice reale de la Universitatea Babeș-Bolyai, România și comparate cu rezultatele similare existente în literatura de specialitate.

Materialul prezentat în această secțiune este o lucrare originală publicată în [11]. În această lucrare am extins modelul de clasificare *SPRAR* prin introducerea a trei scoruri



alternative care pot fi utilizate în etapa de clasificare a *SPRAR* pentru a decide dacă un student *va promova* sau *nu va promova* la un anumit curs. Au fost utilizate trei studii de caz reale pentru evaluarea performanței *SPRAR* având în vedere scorurile de clasificare recent propuse. Cu scopul de a determina empiric care este scorul de clasificare cel mai potrivit, scorurile propuse sunt analizate comparativ și evaluate în raport cu varianta inițială introdusă în [15]. Studiul efectuat în această secțiune este nou în literatura EDM.

Această secțiune a extins modelul de clasificare *SPRAR* (predicția performanței studenților folosind reguli de asociere relațională) [15] pentru a prezice realizările studenților la un anumit curs bazat pe extragerea *regulilor de asociere relațională* (RAR) [3]. Apartenența unei instanțe de interogare la o clasă de trecere sau eșec este determinată folosind seturi de date academice obținute de la Universitatea Babeș-Bolyai, România. Trei noi scoruri de clasificare sunt introduse și utilizate în etapa de clasificare a *SPRAR* pentru a decide dacă un student *va promova* sau *nu va promova*. Performanța *SPRAR* este evaluată folosind mai multe măsuri de evaluare: *sensibilitate*, *specificitate* și *aria de sub curba ROC*. Cel mai bun scor de predicție este identificat și comparat cu rezultatele existente din literatura actuală.

Prin mai multe experimente efectuate pe date academice reale, am obținut o dovadă empirică conform căreia regulile de asociere relațională sunt relevante pentru a face diferența între performanța academică a studenților. De asemenea, am observat că metodologia de clasificare (identificarea celui mai bun scor pentru a discrimina între performanța studenților) depinde de seturile de date și prin urmare am putea investiga în continuare posibilitatea de a o învăța automat, folosind metode supervizate. Ca viitoare direcții de cercetare ne propunem, de asemenea, să extindem experimentele noastre pe alte seturi de date, pentru a adăuga atribute și relații suplimentare în procesul de extragere RAR pentru îmbunătățirea performanței predictive a *SPRAR*.

# Concluzii

Teza actuală a subliniat principalele contribuții ale noastre la dezvoltarea și aplicarea modelelor de extragere a cunoștințelor din date pentru rezolvarea unor probleme legate de domeniul educațional. Extragerea cunoștințelor din date de natură academică (EDM) este un domeniu de cercetare în curs de dezvoltare care se concentrează pe dezvoltarea unor modele pentru identificarea și extragerea informațiilor relevante din datele academice. Cercetarea în EDM își propune să descopere șabloane în procesul de învățare a studenților și să îmbunătățească realizările lor academice. Succesul și performanța academică a studenților pot fi prezise prin tehnici de extragere a cunoștințelor din date [20] folosind învățarea *supervizată* sau *nesupervizată*.

*Regulile de asociere relațională* (RAR-uri) [5, 3] pot exprima diferite tipuri de relații neordinale dintre atributele datelor. Atât extragerea regulilor de asociere ordinale, cât și a celor relaționale au fost aplicate cu succes pentru a soluționa probleme care variază de la sarcinile de extragere a cunoștințelor din date (de exemplu, detectarea defectelor de proiectare software [18] sau curățarea datelor [5]) la clasificarea supervizată (de exemplu, predicția defectelor software [26]).

Am introdus o nouă metodă de extragere incrementală a regulilor de asociere relațională, numită *IRARM*, care extrage eficient toate RAR-urile interesante dintr-un set de date dinamic. Deoarece procesul de extragere a RAR-urilor este costisitor din punct de vedere computațional, ne-am concentrat de asemenea pe dezvoltarea unei versiuni concurente a algoritmului de extragere RAR și anume *CRAR*, cu obiectivul principal de reducere a timpul de calcul al procesului de extragere.

Modelele de învățare nesupervizată și supervizată au fost, de asemenea, studiate în teză în contextul EDM și anume: extragerea incrementală RAR și rețelele cu auto-organizare pentru analiza datelor de performanță ale studenților. Experimentele efectuate pe mai multe seturi de date educaționale reale au evidențiat potențialul metodelor de învățare nesupervizată în descoperirea unor tipare semnificative în datele academice ale studenților. Ca dovadă a conceptului, am introdus un nou model de clasificare *SPRAR* [15] bazat pe descoperirea *regulilor de asociere relațională* pentru a prezice finalizarea cu succes a unui curs academic, în funcție de notele primite de studenți în timpul semestrului academic. Au fost introduse și utilizate trei scoruri de clasificare în etapa de clasificare a *SPRAR* pentru a decide dacă un student *va promova* sau *nu va promova* examenul [11].

În ceea ce privește contribuțiile noastre conceptuale în domeniul DM, ne propunem să investigăm condițiile exacte în care este mai eficient să se aplice metoda incrementală *IRARM* decât să se ruleze algoritmul offline *DRAR* de la zero. În plus, intenționăm să investigăm o abordare incrementală adaptativă a extragerii *regulilor de asociere relațională*, în care sunt adăugate atât atribute noi, cât și obiecte noi la setul de date extrase. Preconizăm optimizări ulterioare ale metodei extragerii concurente RAR (*CRAR*), cum ar fi: creșterea volumului de rulare paralelă a *CRAR*, în special calculul RAR-urilor binare, care ar crește *eficiența*

*scalării și accelerarea paralelă.*

În ceea ce privește clasificatorul *SPRAR* pentru *predicția performanței studenților*, cercetări viitoare vor fi efectuate pentru a extinde evaluarea experimentală a *SPRAR* la alte studii de caz. În plus, ne propunem să generalizăm clasificatorul binar *SPRAR* către un clasificator cu mai multe clase, astfel încât să prezică notele studenților la examenul final la un anumit curs academic. În ceea ce privește procesul de descoperire a regulilor de asociere relațională, intenționăm să extindem modelul nostru având în vedere *reguli de asociere relațională graduală* [17]. *SPRAR* poate fi extins și pentru gestionarea relațiilor *n*-are dintre valorile atributelor.

# Bibliografie

- [1] Alejandro Bogarín, Rebeca Cerezo, and Cristóbal Romero. A survey on educational process mining. *Wiley Interdisc. Rev.: Data Mining and Knowledge Discovery*, 8(1), 2018.
- [2] Leo Breiman. Random forests. *Mach. Learn.*, 45(1):5–32, October 2001.
- [3] Gabriela Șerban, Alina Câmpan, and Istvan Gergely Czibula. A programming interface for finding relational association rules. *International Journal of Computers, Communications & Control*, I(S.):439–444, June 2006.
- [4] Alina Câmpan, Gabriela Șerban, and Andrian Marcus. Relational association rules and error detection. *Studia Universitatis Babes-Bolyai Informatica*, LI(1):31–36, 2006.
- [5] Alina Câmpan, Gabriela Șerban, and Andrian Marcus. Relational association rules and error detection. *Studia Universitatis Babes-Bolyai Informatica*, LI(1):31–36, 2006.
- [6] Alina Campan, Gabriela Șerban, Traian Marius Truta, and Andrian Marcus. An algorithm for the discovery of arbitrary length ordinal association rules. In *DMIN*, pages 107–113, 2006.
- [7] H. Y. Chang, J. C. Lin, M. L. Cheng, and S. C. Huang. A novel incremental data mining algorithm based on fp-growth for big data. *2016 International Conference on Networking and Network Applications (NaNA)*, pages 375–378, 2016.
- [8] George Ciubotariu and Liana Maria Crivei. Analysing the academic performance of students using unsupervised data mining. *Studia Universitatis Babes-Bolyai Series Informatica*, pages 1–14, 2019.
- [9] Liana Maria Crivei. Incremental relational association rule mining of educational data sets. *Studia Universitatis Babes-Bolyai Series Informatica*, 63(2):102–117, 2018.
- [10] Liana Maria Crivei. Using unsupervised learning models for analyzing students’ academic performance. In *PhD Colloquium at SYNASC’18, Symbolic and Numeric Algorithms for Scientific Computing*, pages 1–4, 2018.
- [11] Liana-Maria Crivei, Mihai Andrei, and Czibula Gabriela. A study on applying relational association rule mining based classification for predicting the academic performance of students. In *KSEM 2019 : The 12th International Conference on Knowledge Science, Engineering and Management, LNAI 11775*, pages 287–300, 2019.
- [12] Liana Maria Crivei, Gabriela Czibula, George Ciubotariu, and Mariana Dindelegan. Unsupervised learning based mining of academic data sets for students’ performance

- analysis. In *IEEE 14th International Symposium on Applied Computational Intelligence and Informatics (SACI 2019)*, pages 1–6. IEEE Hungary Section, 2019.
- [13] L.M. Crivei, V.-S. Ionescu, and G. Czibula. An analysis of supervised learning methods for predicting students' performance in academic environments. *ICIC Express Letters*, 13(3):181–190, 2019.
- [14] G. Czibula, I.G. Czibula, D.-L. Miholca, and L.M. Crivei. A novel concurrent relational association rule mining approach. *Expert Systems with Applications*, 125(6):142–156, 2019.
- [15] G. Czibula, A. Mihai, and L.M. Crivei. A novel relational association rule mining classification model applied for academic performance prediction. volume 159, pages 20–29, 2019.
- [16] Gabriela Czibula, Maria-Iuliana Bocicor, and Istvan Gergely Czibula. Promoter sequences prediction using relational association rule mining. *Evolutionary Bioinformatics*, 8:181–196, 04 2012.
- [17] Gabriela Czibula, Istvan Gergely Czibula, and Diana-Lucia Miholca. Enhancing relational association rules with gradualness. *International Journal of Innovative Computing, Communication and Control*, 13:289–305, 2017.
- [18] Gabriela Czibula, Zsuzsanna Marian, and István Gergely Czibula. Detecting software design defects using relational association rule mining. *Knowledge and Information Systems*, 2014.
- [19] N. Elfelly, J.-Y. Dieulot, and P. Borne. A neural approach of multimodel representation of complex processes. *International Journal of Computers, Communications & Control*, III(2):149–160, 2008.
- [20] Richard A Huebner. A survey of educational data-mining research. *Research in Higher Education Journal*, 19:1–13, 2013.
- [21] Andreas Khler, Matthias Ohrnberger, and Frank Scherbaum. Unsupervised feature selection and general pattern discovery using self-organizing maps for gaining insights into the nature of seismic wavefields. *Computers Geosciences*, 35(9):1757 – 1767, 2009.
- [22] J. Lampinen and E. Oja. Clustering properties of hierarchical self-organizing maps. *Journal of Mathematical Imaging and Vision*, 2(3):261–272, 1992.
- [23] Kristopher R. Linstrom and A. John Boye. A neural network prediction model for a psychiatric application. *International Conference on Computational Intelligence and Multimedia Applications*, 0:36–40, 2005.
- [24] Andrian Marcus, Jonathan I. Maletic, and King-Ip Lin. Ordinal association rules for error identification in data sets. In *Proceedings of the tenth international conference on Information and knowledge management, CIKM '01*, pages 589–591, New York, NY, USA, 2001. ACM.
- [25] Diana-Lucia Miholca, Gabriela Czibula, and Liana Maria Crivei. A new incremental relational association rules mining approach. In *22nd International Conference on*

- Knowledge-Based and Intelligent Information & Engineering Systems*, volume 126 of *KES2018*, pages 19–28. Procedia Computer Science, 2018.
- [26] Diana-Lucia Miholca, Gabriela Czibula, and István Gergely Czibula. A novel approach for software defect prediction through hybridizing gradual relational association rules with artificial neural networks. *Inf. Sci.*, 441:152–170, 2018.
- [27] Diana-Lucia Miholca, Gabriela Czibula, and Istvan Gergely Czibula. A novel approach for software defect prediction through hybridizing gradual relational association rules with artificial neural networks. *Information Sciences*, 441:152 – 170, 2018.
- [28] Tom M. Mitchell. *Machine Learning*. McGraw-Hill, New York, 1997.
- [29] Siti Khadijah Mohamad and Zaidatun Tasir. Educational data mining: A review. *Procedia - Social and Behavioral Sciences*, 97:320 – 324, 2013. The 9th International Conference on Cognitive Science.
- [30] Siti Khadijah Mohamad and Zaidatun Tasir. Educational data mining: A review. *Procedia - Social and Behavioral Sciences*, 97(Supplement C):320 – 324, 2013. The 9th International Conference on Cognitive Science.
- [31] Mark D. Skowronski and John G. Harris. Automatic speech recognition using a predictive echo state network classifier. *Neural Networks*, 20(3):414–423, April 2007.
- [32] Panu Somervuo and Teuvo Kohonen. Self-organizing maps and learning vector quantization for feature sequences. *Neural Processing Letters*, 10:151–159, 1999.
- [33] Anping Song, Xuehai Ding, Mingbo Li, Wei Cao, and Ke Pu. A novel binary bat algorithm for association rules mining. *ICIC Express Letters*, 9(9):2387–2394, September 2015.
- [34] Ruihu Wang. Adaboost for feature selection, classification and its relation with svm, a review. *Physics Procedia*, 25:800 – 807, 2012. International Conference on Solid State Devices and Materials Science, April 1-2, 2012, Macao.
- [35] Tiantian Yang, Ata Akabri Asanjan, Mohammad Faridzad, Negin Hayatbini, Xiaogang Gao, and Soroosh Sorooshian. An enhanced artificial neural network with a shuffled complex evolutionary global optimization with principal component analysis. *Information Sciences*, 418-419(Supplement C):302 – 316, 2017.