

UNIVERSITATEA BABEȘ-BOLYAI
FACULTATEA DE MATEMATICĂ ȘI INFORMATICĂ

Dezvoltare de soft bazată pe instruire automată

Rezumat Extins

Doctorand: Zsuzsanna-Edit Marian
Conducător de doctorat: Prof. Dr. Gabriela Cibula

Septembrie 2014

Lista publicațiilor

Publicații indexate în ISI Web of Knowledge

Publicații indexate în ISI Science Citation Index Expanded

1. [CMC14b] Gabriela Czibula, **Zsuzsanna Marian** și Istvan-Gergely Czibula, Software defect prediction using relational association rule mining. *Information Sciences* publicat de *Elsevier*, Vol. 264, pp. 260-278, DOI: 10.1016/j.ins.2013.12.031, 2014. (**IF = 3.643**)
2. [CMC14a] Gabriela Czibula, **Zsuzsanna Marian** și Istvan-Gergely Czibula. Detecting software design defects using relational association rule mining. *Knowledge and Information Systems* publicat de *Springer*, DOI: 10.1007/s10115-013-0721-z, publicat on-line în Ianuarie, 2014. (**IF = 2.225**)
3. [MCC12] **Zsuzsanna Marian**, Gabriela Czibula și Istvan-Gergely Czibula. Using software metrics for automatic software design improvement. *SIC Journal, Studies in Informatics and Control*, România, Vol. 21, Number 3, pp. 249-258, 2012. (**IF = 0.578**)
4. [MCC14] **Zsuzsanna Marian** Gabriela Czibula și Istvan-Gergely Czibula, Software packages refactoring using a hierarchical clustering-based approach. *Fundamenta Informaticae*, În curs de recenzie, 2014. (**IF = 0.399**)

Publicații indexate în ISI Conference Proceedings Citation Index

5. [SMV09] Adrian Sterca, **Zsuzsanna Marian** și Alexandru Vancea. Distortion-based media-friendly congestion control. *Proceedings of the International Conference on Knowledge Engineering, Principles and Techniques*, Cluj-Napoca, România, pp. 265-267, 2009.
6. [ȚIM09] Radu Țurcaș, Oana Iova și **Zsuzsanna Marian**. The autonomous robotic tank (ART): an innovative lego mindstorm NXT battle vehicle. *Proceedings of the International Conference on Knowledge Engineering, Principles and Techniques*, Cluj-Napoca, România, pp. 95-98, 2009.

Articole publicate în jurnale internaționale și la conferințe internaționale

7. [Mar12a] **Zsuzsanna Marian**. Aggregated metrics guided software restructuring. *Proceedings of 8th IEEE International Conference on Intelligent Computer Communication and Processing*, Cluj-Napoca, România, pp. 259-266, 2012. (**Indexat IEEE**)

8. [Mar12c] **Zsuzsanna Marian**. A study on hierarchical clustering based software restructuring. *Studia Universitatis "Babes-Bolyai", Informatica*, România, Vol. LVII, Number 2, pp. 20-31, 2012. (**Indexat MathSciNet**)
9. [Mar13b] **Zsuzsanna Marian**. A study on Relational Association Rule Mining Based Software Design Defect Detection. *Studia Universitatis "Babes-Bolyai", Informatica*, România, Vol. LVIII, Number 1, pp. 42-57, 2013. (**Indexat MathSciNet**)
10. [Mar13a] **Zsuzsanna Marian**. On the software metrics influence in Relational Association Rule-based Software Defect Prediction. *Studia Universitatis "Babes-Bolyai", Informatica*, România, Vol. LVIII, Number 4, pp. 35-48, 2013. (**Indexat MathSciNet**)
11. [Mar14b] **Zsuzsanna Marian**. On evaluating the structure of software packages. *Studia Universitatis "Babes-Bolyai", Informatica*, România, Vol. LIX, Number 1, pp. 46-58, 2014. (**Indexat MathSciNet**)
12. [Mar14a] **Zsuzsanna Marian**. FAOS - A framework for analyzing object-oriented software systems. *Studia Universitatis "Babes-Bolyai", Informatica*, România, Under review, 2014. (**Indexat MathSciNet**)
13. [MS10] **Zsuzsanna Marian** și Christian Săcărea. Using contextual topology to discover similarities in modern music. *Proceedings of the IEEE International Conference on Automation Quality and Testing, Robotics*, Cluj-Napoca, România, vol. 3, pp. 1-6, 2010. (**Indexat IEEE**)
14. [TLM11] Doina Tatar, Mihaiela Lupea și **Zsuzsanna Marian**. Text summarization by formal concept analysis approach. *Studia Universitatis "Babes-Bolyai", Informatica*, România, Vol. LVI, Number 2, pp. 7-12, 2011. (**Indexat MathSciNet**)
15. [MCB11] **Zsuzsanna Marian**, Cosmin Coman și Attila Bartha. Learning to play the guessing game. *Studia Universitatis "Babes-Bolyai", Informatica*, România, Vol. LVI, Number 2, pp. 119-124, 2011. (**Indexat MathSciNet**)

Articole publicate la conferințe naționale

16. [Mar12b] **Zsuzsanna Marian**. Software metrics based refactoring: a case study. *Proceedings of the National Symposium ZAC 2012*, Cluj-Napoca, pp. 59-64, 2012.
17. [Mar10] **Zsuzsanna Marian**. Solving the subset sum problem with DNA computation, *Proceedings of the National Symposium ZAC 2010*, Cluj-Napoca, pp. 25-29, 2010.

În cazul publicațiilor cu mai mulți autori, am contribuit la: conceperea, proiectarea și execuția experimentelor; analiza datelor; dezvoltarea structurii și argumentelor din articole; scrierea manuscrisului; revizuirea articolelor și aprobarea variantei finale.

Cuvinte cheie

- Instruire Automată
- Clustering
- Metrici Soft
- Reguli de Asociere Relaționale
- Dezvoltare de Soft
- Restructurare Soft
- Detectarea Defectelor în Sisteme Informatice
- Detectarea Defectelor de Proiectare
- Predicția Defectelor
- Seturi de date Nasa
- Framework-uri Soft

Cuprins

Lista publicațiilor	1
Introducere	4
1 Instruire Automată în Ingineria Soft	7
1.1 Inteligența Computațională în Inginerie Soft	7
1.2 Remodularizare Soft	8
1.3 Detectarea Defectelor în Sisteme Informatice	9
1.4 Prezentare metodelor relevante din Instruirea Automată	10
2 Abordări noi pentru Remodularizare Soft	12
2.1 Remodularizare la nivel de pachete	12
2.2 Remodularizare la nivel de clasă	15
2.3 Concluzii și cercetări ulterioare	19
3 Abordări noi pentru Detectarea Defectelor în Sisteme Informatice	20
3.1 Model teoretic	20
3.2 Detectarea Defectelor de Proiectare	21
3.3 Predicția Defectelor	25
3.4 Concluzii și cercetări ulterioare	27
4 Un Framework pentru Analiza Sistemelor Informatice Orientate Obiect	29
4.1 <i>Framework</i> -ul FAOS	29
4.2 Comparație cu framework-uri similare	29
4.3 Concluzii și cercetări ulterioare	30
Concluzii	31
Bibliografie	32

Introducere

Această teză este rezultatul cercetării mele în domeniul dezvoltării de *modele de instruire automată pentru rezolvarea diferitor probleme legate de dezvoltare de soft*. Această cercetare a început în 2011, sub coordonarea lui Prof. Dr. Gabriela Czibula.

Direcția principală de cercetare asupra căreia m-am concentrat a fost aplicarea diferitor modele și algoritmi de instruire automată pentru probleme din ingineria soft. Problemele abordate în această teză sunt de o importanță majoră pentru întreținerea și evoluția sistemelor informatice. Soluții la aceste probleme ar ajuta programatorii în diferite probleme de dezvoltare a sistemelor informatice.

Instruirea Automată este o ramură a Inteligenței Artificiale care încearcă să găsească un răspuns la întrebarea “Cum putem construi sisteme informatice care se îmbunătățesc automat prin experiență, și care sunt legile fundamentale care conduc toate procesele de învățare?” [Mit06]. Prin urmare, domeniul de Instruire Automată este alcătuit dintr-o serie de abordări și algoritmi care pot învăța din date existente.

Pornind de la articolul publicat de Mark Harman și Brian Jones în anul 2001 s-a născut un nou domeniu de cercetare, numit *Search Based Software Engineering - SBSE* (Inginerie Soft Bazată pe Căutare) [HJ01]. Acest articol descrie că sunt multe probleme de inginerie soft care pot fi reformulate ca și probleme de căutare, iar apoi diferiți algoritmi de căutare pot fi aplicați asupra lor. Dezvoltarea acestui nou domeniu a dus la noi abordări inteligente pentru probleme de inginerie soft: există probleme de inginerie soft care pot fi reformulate ca probleme de *clustering* asupra cărora pot fi aplicați algoritmi de *clustering*. De asemenea, există probleme de inginerie soft care pot fi reformulate ca probleme de clasificare sau predicție asupra cărora pot fi aplicați diferiți clasificatori.

Dar chiar dacă metode inteligente *pot* fi aplicate asupra problemelor de inginerie soft, *de ce* ar trebui să fie aplicate? De ce ar trebui să dezvoltăm abordări pentru rezolvarea problemelor care pot fi rezolvate de dezvoltatori soft? Răspunsul este că aceste probleme nu tot timpul pot fi rezolvate de dezvoltatori soft și ei pot folosi ajutorul acestor abordări. Din cauza complexității și dimensiunii sistemelor informatice, sunt multe situații când dezvoltatorii unui sistem nu au o privire de ansamblu asupra sistemului, nu înțeleg toate legăturile între elementele sistemului și nu văd care sunt părțile care ar trebui îmbunătățite sau care ar trebui testate în mod intensiv. De aceea metode bazate pe căutare, pe inteligență computațională și pe instruire automată sunt binevenite: ele pot ajuta dezvoltatorii soft în munca lor.

Dintre diferitele probleme de dezvoltare soft asupra cărora au fost aplicate metode inteligente, am ales două direcții principale asupra cărora m-am concentrat în această teză. Prima direcție a fost dezvoltarea abordărilor pentru restructurarea sistemelor informatice orientate obiect folosind algoritmi de *clustering*. Am lucrat atât la restructurare la nivel de pachete, cât și la restructurare la nivel de clase. A doua direcție principală abordată în această teză este problema detectării defectelor în sisteme informatice. Această direcție poate fi împărțită în detectarea defectelor de proiectare (*design defect detection*) și predicția defectelor (*defect prediction*), și pentru fiecare am propus o nouă metodă bazată pe reguli de asociere relaționale. O contribuție originală spre dezvoltarea sistemelor informatice este prezentată de asemenea, în forma unui *framework*, numit *FAOS (Framework for Analyzing Object-oriented Software systems)*. Acest *framework* a fost dezvoltat pentru a sprijini evaluarea experimentală a ma-

jorității abordărilor prezentate în această teză.

Această teză este alcătuită din patru capitole, după cum urmează.

Primul capitol, **Instruire Automată în Inginerie Soft. Context General**, începe cu o prezentare scurtă a diferitor metode de inteligență computațională în ingineria soft. Aceasta este urmată de prezentarea celor două direcții de cercetare abordate în această teză: *Remodularizare Soft* și *Detectarea Defectelor în Sisteme Informaticе*. Pentru ambele direcții problemele abordate sunt detaliate împreună cu o scurtă prezentare a metodelor existente în literatura de specialitate. În cele din urmă este o scurtă descriere a celor două metode de instruire automată care au fost folosite în abordările noastre: reguli de asociere relaționale și *clustering*.

Contribuțiile originale sunt prezentate în Capitolele 2, 3 și 4, unde descriem modelele de instruire automată propuse pentru rezolvarea problemelor de inginerie soft ce fac studiul acestei teze.

Capitolul 2, **Abordări noi pentru Remodularizare Soft**, este original și prezintă abordările noastre în cadrul primei direcții de cercetare principală a acestei teze, și anume restructurarea sistemelor informatice folosind algoritmi de *clustering*. Prezentarea începe cu prima problemă din acest domeniu: remodularizarea la nivel de pachete. Prezentăm doi algoritmi, unul care identifică o împărțire bună a claselor în pachete și unul care, plecând de la o structură existentă de pachete, poate să identifice pachetul potrivit pentru o clasă care urmează să fie adăugată în sistem. Aceasta este urmată de evaluarea experimentală a acestor algoritmi, o analiză a abordării și o comparație cu alte metode existente în literatura de specialitate pentru restructurare la nivel de pachete. A doua problemă abordată în Capitolul 2, restructurare soft la nivel de clasă, este prezentată în continuare. Prezentăm trei algoritmi care pot identifica o structură de clasă îmbunătățită pentru un sistem informatic, împreună cu evaluarea lor experimentală. În primul rând se face o comparație a celor trei algoritmi urmând ca aceștia să fie comparați cu alți algoritmi prezentați în literatura de specialitate. Ultima secțiune prezintă concluziile din capitol și direcții pentru cercetări ulterioare.

Capitolul 3, **Abordări noi pentru Detectarea Defectelor în Sisteme Informaticе**, este original și prezintă abordările mele pentru a doua direcție principală de cercetare a acestei teze, detectarea defectelor în sisteme informatice. Ambele probleme abordate din această direcție folosesc același model teoretic, care este prezentat în prima secțiune. Aceasta este urmată de introducerea abordării noastre pentru identificarea claselor cu defecte de proiectare folosind reguli de asociere relaționale. După aceea este prezentată evaluarea experimentală a abordării și o comparație cu metode existente. În continuare se prezintă un studiu despre efectele modificării valorilor parametrilor asupra abordării propuse în teză. A doua parte a Capitolului 3 prezintă abordarea noastră pentru a doua problemă din această direcție, și anume, predicția defectelor în sisteme informatice. Descriem abordarea noastră originală, un clasificator care poate clasifica entitățile dintr-un sistem soft ca defecte sau non-defective folosind reguli de asociere relaționale. Urmează evaluarea experimentală și o comparație cu abordări existente. La sfârșit prezentăm concluziile capitolului.

Capitolul 4, **Un framework pentru Analiza Sistemelor Informaticе Orientate Obiect**, este original și prezintă *framework*-ul *FAOS (Framework for Analyzing Object-oriented Software Systems)*, un *framework* general dezvoltat în limbajul de programare Java pentru analiza sistemelor orientate obiect. Capitolul prezintă cele trei module principale ale *framework*-ului, unul pentru analiza codului Java compilat și extragerea unei liste de elemente (clase, metode, attribute) și relațiile între ele, unul pentru calcularea valorilor diferitor metrici soft, și unul pentru implementarea abordării noastre pentru restructurarea unui sistem informatic la nivel de pachete. O comparație scurtă cu alte *framework*-uri și unelte similare este prezentată după cele trei module. Concluziile și direcțiile de cercetare ulterioare se găsesc în ultima secțiune.

Contribuțiile originale introduse în această teză se găsesc în Capitolele 2, 3 și 4 și sunt

următoarele:

- O abordare nouă bazată pe *clustering* pentru remodularizare soft la nivel de pachete - abordarea *CASP* (*Clustering Approach for Software Package Restructuring*) (Subsecțiunea 2.1.1) împreună cu un algoritm bazat pe *clustering* ierarhic, *HASP* (*Hierarchical Clustering Algorithms for Software Packages Restructuring*) care poate fi folosit în abordarea *CASP* (Subsecțiunea 2.1.2) [MCC14].
 - Definiția a șapte caracteristici, care pot fi agregate într-o singură valoare, numită *overallScore*, care poate fi folosită ca măsură de similaritate în algoritmul *HASP* (Subsecțiunea 2.1.2) și definiția unei măsuri noi, *CIP* (*Cohesion of Identified Packages*), care măsoară cât de aproape este o partiție de o altă partiție (Subsecțiunea 2.1.4). [MCC14].
 - Algoritmul *AssignClass*, care poate găsi pachetul potrivit pentru o clasă nou adăugată, folosind aceiași caracteristici și scor ca și algoritmul *HASP* (Subsecțiunea 2.1.3) [MCC14].
 - Evaluarea experimentală a algoritmilor *HASP* și *AssignClass* (Subsecțiunea 2.1.4), analiza rezultatelor (Subsecțiunea 2.1.5) și o comparație a măsurii *overallScore* cu alte măsuri existente (Subsecțiunea 2.1.6). [MCC14, Mar14b].
- Trei algoritmi noi pentru restructurarea unui sistem informatic la nivel de clasă: *ARI*, *HAC* (Subsecțiunea 2.2.1) și o abordare bazată pe metrice agregate (Subsecțiunea 2.2.2) [MCC12, Mar12c, Mar12b, Mar12a].
 - Evaluarea experimentală a celor trei algoritmi și o analiză comparativă a lor (Subsecțiunea 2.2.3).
- O metodă pentru detectarea defectelor de proiectare bazată pe reguli de asociere relaționale - *SDDRAR* (Subsecțiunea 3.2.2) [CMC14a].
 - Un algoritm nou, similar cu Apriori [AS94], pentru a extrage regulile de asociere relaționale de orice lungime dintr-un set de date - *DRAR* (Subsecțiunea 3.2.1) [CMC14a].
 - O evaluare experimentală detaliată a abordării *SDDRAR* (Subsecțiunea 3.2.3) urmată de o analiză a abordării din mai multe puncte de vedere (Subsecțiunea 3.2.4) [CMC14a].
 - Un studiu despre efectele modificării valorilor parametrilor asupra rezultatelor abordării *SDDRAR* (Subsecțiunea 3.2.5) [Mar13b].
- O abordare de clasificare nouă, care, folosind reguli de asociere relaționale, poate să clasifice entitățile dintr-un sistem informatic ca fiind defecte sau non-defective - *DPRAR* (Subsecțiunea 3.3.1) [CMC14b].
 - O evaluare experimentală amănunțită a abordării *DPRAR*, folosind 10 seturi de date *NASA* (Subsecțiunea 3.3.2) și o comparație a rezultatelor cu cele raportate în literatura de specialitate pentru aceleași seturi de date (Subsecțiunea 3.3.3) [CMC14b].
 - Studii asupra efectului eliminării unor caracteristici și a modificării formulei de calcul a scorului (Subsecțiunea 3.3.4) [Mar13a].
- Un *framework* general, *FAOS*, pentru analiza sistemelor informatice orientate obiect (Secțiunea 4.1) [Mar14a]. Acest *framework* a fost implementat pentru a sprijini evaluarea experimentală a abordărilor dezvoltate.

Capitolul 1

Instruire Automată în Inginerie Soft

Aplicarea metodelor inteligente asupra problemelor din inginerie soft este un domeniu de cercetare relativ nou, apărut în anul 2001 cu un articol publicat de Mark Harman și Bryan Jones, [HJ01]. De atunci acest domeniu a progresat mult. În [ZT05] autorii prezintă 44 probleme diferite din ingineria soft pentru care cel puțin un algoritm de instruire automată a fost aplicat. În mod similar, o analiză a tendințelor curente în *SBSE*, [HMZ09], a adunat peste 500 de publicații din domeniu.

1.1 Inteligența Computațională în Inginerie Soft

În anul 2001 Mark Harman și Bryan Jones au publicat un articol, [HJ01], care este considerat manifestul unui domeniu de cercetare nou, numit *Search-Based Software Engineering - SBSE (Inginerie Soft Bazată pe Căutare)*. Aceștia arată că tehnici de căutare (și anume Algoritmi Genetici, Căutare Tabu și Călirea Simulată) au fost aplicați în diferite domenii de inginerie și că ei ar fi potriviți și pentru inginerie soft.

În următoarea perioadă un număr mare de abordări au fost prezentate în literatura de specialitate, abordări care aplică diferiți algoritmi de căutare pentru probleme de inginerie soft. La 8 ani după apariția domeniului, în anul 2009, Harman et al. au publicat o analiză a curentelor, tehnicilor și aplicațiilor existente în *SBSE*, [HMZ09], luând în calcul peste 500 de publicații din domeniu.

Algoritmii de căutare nu sunt singurele metode inteligente aplicate pe probleme de inginerie soft, așa cum este prezentat în [Har12], unde autorul menționează trei domenii largi din Inteligența Artificială folosite de ingineri soft: căutare și optimizare computațională, metode probabilistice și *fuzzy*, învățare și predicție.

Pe lângă [Har12], o altă analiză recentă despre rolul Inteligenței Artificiale în ingineria soft [AAH12] prezintă că metode de inteligență computațională au fost aplicate pentru proiectarea arhitecturii sistemelor informatice, testare soft, estimarea și predicția costului și repararea automată a bug-urilor. În [DK05] Dick și Kandel prezintă exemple de cum pot fi aplicate metode din inteligența computațională și instruirea automată pentru asigurarea calității softului.

Și algoritmi de instruire automată au fost aplicați pentru diferite probleme din ingineria soft. În [ZT05] autorii prezintă 44 probleme diferite din ingineria soft pentru care cel puțin un algoritm de instruire automată a fost aplicat.

1.2 Remodularizare Soft

Un sistem informatic trebuie să se adapteze cerințelor care se modifică pe parcursul duratei de viață, altfel nu va mai fi folosit. Modificările legate de această adaptare sunt efectuate de cele mai multe ori după lansarea softului, la întreținere, și dacă sunt efectuate doar modificări care adaugă noi funcționalități sau corectează erorile existente, sistemul va deveni din ce în ce mai greu de întreținut. Pentru a preveni acest lucru, o altă activitate importantă este restructurarea codului sursă. Restructurarea este activitatea de a schimba (îmbunătăți) structura internă a sistemului soft, fără a afecta comportamentul exterior al sistemului.

Remodularizarea soft poate fi efectuată pe diferite nivele ale unui sistem informatic. Pentru un sistem orientat obiect exista remodularizare la nivel de metodă, clasă, pachet și arhitectură.

1.2.1 Remodularizare la nivel de pachete

Definiția și relevanța problemei. În zilele noastre sistemele informatice devin din ce în ce mai complexe, sunt alcătuite din mii de clase care sunt grupate în pachete. Fără o structură de pachete adecvată sistemul devine greu de întreținut. Astfel, remodularizarea soft la nivel de pachete este un proces important în întreținerea și evoluția sistemului. Cu cât e mai complex sistemul, cu atât sunt mai ridicate costurile întreținerii. Astfel, este destul de greu pentru dezvoltatorii soft să aleagă structura de pachete potrivită pentru sistem. Când numărul de clase este mare, găsirea pachetului potrivit pentru o clasă nouă nu este simplă, necesită o bună cunoaștere a sistemului.

Problema de restructurare a pachetelor se ivește din nevoi practice, astfel abordări inteligente pot fi de folos, ajutând dezvoltatorii în munca lor.

Abordări existente în literatură. Există câteva metode prezentate în literatura de specialitate pentru identificarea modului în care clasele ar trebui grupate în pachete. O astfel de metodă este cea prezentată în [AAM11], unde *clustering* este folosit pentru a găsi gruparea ideală a claselor. O altă metodă, bazată pe detectarea comunităților constrânse, este prezentată în [PJL13]. Chiar dacă nu poate restructura un sistem întreg, metoda prezentată de Bavota et al. în [BLMO10] poate să dividă un pachet cu coeziune mică în mai multe pachete mai coezive. Metode bazate pe algoritmi de căutare sunt prezentate în [MHH03, MMCG99].

O altă direcție de cercetare care trebuie menționată este definirea diferitor metrici, care măsoară calitatea pachetelor într-un sistem. Astfel de metrici sunt prezentate în [SKR08] și [DABH11].

1.2.2 Remodularizare la nivel de clasă

Definiția și relevanța problemei. Cele mai multe metode legate de remodularizare soft din literatura de specialitate identifică sau efectuează restructurare la nivel de clasă. Aceste abordări încearcă să identifice acele clase din sistemul informatic, care nu sunt bine proiectate, și să identifice acele restructurări care ar îmbunătăți proiectarea lor.

O carte excelentă și des menționată în legătură cu restructurare este *Refactoring: Improving the Design of Existing Code* de Martin Fowler [Fow99]. În această carte Fowler prezintă de ce este importantă restructurarea și cum se efectuează. De asemenea, el prezintă un catalog cu 72 de restructurări diferite. Chiar și cu un catalog atât de complet, există două întrebări: când și unde ar trebui efectuate restructurări? Pentru a răspunde la această întrebare, Fowler introduce noțiunea de *bad smell*, care este o structură în codul sursa care sugerează posibilitatea unei restructurări [Fow99]. În carte este și o listă cu 22 *bad smell*-uri.

Abordări existente în literatură. Identificarea oportunităților de restructurare în

sisteme informatice este un domeniu bine-studiat. Diferite metode sunt prezentate în literatura de specialitate care folosesc Teoria Jocurilor [BOL⁺10, HT07], *Concept Analysis* [ST98], Rețele Bayesiene [KVG09], *clustering* [RR11, FTCS09, AL99] sau diferite metode de căutare [MC11, OC08, GW11, SSB06, HT07]. Cu toate că aproape toate metodele enumerate mai sus folosesc și metrici soft, există abordări bazate doar pe metrici soft [ISIE12, SS07, HKI08, CLMM05, SSL01, TK03, DDN00].

1.3 Detectarea Defectelor în Sisteme Informatice

Detectarea defectelor în sisteme informatice este ramura ingineriei soft care încearcă să identifice sau să prezică defectele într-un sistem informatic. Este o activitate importantă, conform [Kan03], detectarea și corectarea defectelor pot consuma în jur de 75% din toate costurile ciclului de viață. În [Kan03] Kandt prezintă două tehnici complementare de detectarea defectelor: inspecția și testarea. Ambele necesită timp și efort considerabil, de cele mai multe ori greu de obținut, mai ales în cazul sistemelor mari. Această lipsă a timpului, combinată cu dimensiunea mare a sistemelor informatice este motivul pentru care sunt dezvoltate metode inteligente care să ajute dezvoltatorii prin sugerarea părților din sistem care necesită o testare mai intensivă.

1.3.1 Detectarea Defectelor de Proiectare

Definiția și relevanța problemei. Structura internă a entităților într-un sistem informatic se poate modifica de mai multe ori pe parcursul ciclului de viață și are un impact important asupra mentenabilității sistemului. Chiar dacă un sistem informatic are o proiectare bună la început, pe parcursul întreținerii proiectarea se degradează, ceea ce face întreținerea mai complicată și mai costisitoare pe viitor. De aceea monitorizarea continuă a proiectării și identificarea și corectarea defectelor de proiectare este esențială.

Detectarea defectelor de proiectare este înrudită cu o altă problemă de inginerie soft, și anume, restructurarea. În general, soluția pentru a elimina defectele de proiectare este restructurarea sistemului.

Ce sunt defectele de proiectare? Pentru a putea detecta defectele de proiectare, trebuie să definim ce este un defect de proiectare. Pentru a defini defectele de proiectare putem pleca de la un set de principii pentru o proiectare corectă, și căuta încălcarea acestor principii. De exemplu, Larman în [Lar04] prezintă nouă principii (sau șabloane), cunoscute sub numele de șabloane *GRASP* (*General Responsibility Assignment Software Patterns*).

Defectele de proiectare sunt adesea împărțite în două categorii: defecte de proiectare de nivel înalt și defecte de proiectare de nivel scăzut (numite și *bad smell*).

Abordări existente în literatură. Identificarea manuală a defectelor de proiectare durează foarte mult pentru sisteme mari, de aceea în literatura de specialitate diferite abordări sunt prezentate pentru a le găsi în mod automat. Asemenea abordări sunt *detection strategies* introduse de Marinescu [Mar02], *rule cards* introduse de Moha [MGL06, Moh06], *plugin*-ul pentru *Eclipse JDeodorant* [FTSC11] și abordarea bazată pe metrici soft a lui Munro [Mun05]. Aplicarea metodelor de căutare pentru detectarea defectelor de proiectare este prezentată în [KSBW11].

Fontana et al. prezintă în [FBZ12] o comparație între 4 instrumente care detectează *code smell*-uri, care sunt aplicate pe diferite sisteme informatice cu scopul de a detecta șase tipuri de *bad smell*. Una dintre cele mai importante concluzii ale articolului este că instrumente diferite, aplicate pe același sistem pentru a identifica același *bad smell*, produc rezultate diferite (singura excepție a fost la *bad smell God Class*).

1.3.2 Predicția Defectelor

Definiția și relevanța problemei. Predicția defectelor încearcă să determine acele părți ale unui sistem informatic unde se găsesc erori, astfel, testerii pot să petreacă mai mult timp testând acele componente care probabil conțin erori, și mai puțin timp testând componentele care probabil sunt fără erori. Multe abordări folosesc metrici soft, pentru a măsura calitatea softului pentru a prezice defectele. Astfel, predicția defectelor constă în clasificarea modulelor ca fiind defecte și non-defecte, folosind o clasificare bazată pe metrici [BMW02]. Cele mai multe abordări folosesc date de la alte proiecte sau versiuni anterioare, astfel reprezintă abordări de clasificare supervizată.

Seturile de date NASA. Multe abordări din literatura de specialitate folosesc pentru evaluarea experimentală seturile de date NASA. Aceste seturi de date pot fi găsite la [Nas], unde ele sunt accesibile tuturor, pentru a facilita crearea modelelor predictive.

Abordări existente în literatură. Una dintre primele metodele de predicție a defectelor este metoda CBA, prezentată în [LHM98], o metoda care folosește o extensie a regulilor de asociere. [LMW01] prezintă o extensie a acestei metode, CBA2. Un model hibrid, care combină reguli de asociere și regresia logistică este prezentată în [KMMiM08]. O altă metodă bazată pe reguli este EDER-SD care produce doar reguli care caracterizează modulele defecte [RRRAR12].

Pe lângă metode bazate pe reguli au fost aplicați și alți algoritmi din instruirea automată pentru predicția defectelor. O asemenea lucrare este [MGF07], în care autorii evaluează clasificatorii OneR, Naive Bayes și J48. Challagulla et al. evaluează în [CBYP05] niște modele de predicție pe patru seturi de date NASA. Haghighi et al. prezintă în [HDF12] o analiză comparativă a 37 de clasificatori și folosesc seturile NASA pentru experimente.

O altă direcție este folosirea metodelor de învățare semi-supervizate, bazate pe dezacord. Asemenea abordări sunt ROCUS [JLZ11], ACoForest [LZWZ12] și folosirea metodei *Random Forests* prezentată în [GMCS04].

1.4 Prezentare metodelor relevante din Instruirea Automată

1.4.1 Reguli de Asociere Relaționale

Extragerea regulilor de asociere din baze de date care conțin tranzații de articole a fost descrisă pentru prima dată de Agrawal et. al în anul 1993 [AIS93]. O regulă de asociere este o implicație de forma $X \Rightarrow Y$, unde X și Y sunt mulțimi disjuncte de articole și înseamnă că prezența lui X implică prezența lui Y .

Una dintre dezavantajele regulilor de asociere este că acestea consideră doar prezența sau absența unui articol, și ignoră orice relație care ar putea exista între elementele dintr-o mulțime de articole. Pentru a depăși acest dezavantaj Marcus et al. propune în [MML01] o extensie a regulilor de asociere, *regulile de asociere ordinale*, unde există relații ordinale între elementele dintr-o mulțime de articole.

Totuși, relațiile ordinale nu sunt tot timpul suficiente, pot exista alte relații, care nu sunt ordinale, între articole. De aceea, *reguli de asociere relaționale*, o extensie a regulilor de asociere ordinale, au fost introduse în [SCC06]. Regulile de asociere relaționale permit orice tip de relație între articolele dintr-o mulțime.

Formal, așa cum este prezentat în [SCC06], să considerăm că $R = \{r_1, r_2, \dots, r_n\}$ este o mulțime de instanțe. Fiecare instanță este caracterizată de o listă de m atribute, (a_1, \dots, a_m) . Valoarea atributului a_i în instanța r_j este notată cu $\Phi(r_j, a_i)$. Pentru fiecare atribut a_i există un domeniu D_i de unde provin valorile. Acest domeniu conține și valoarea vidă, notată cu ε . Între doi domenii D_i și D_j relații, cum ar fi *mai mic sau egal* (\leq), *egal* ($=$) și *mai mare sau egal* (\geq) pot fi definite. Mulțimea tuturor relațiilor posibile care pot fi definite între D_i

x D_j este notată cu M .

Definiție 1 [SCC06] O regulă de asociere relațională este o expresie $(a_{i_1}, a_{i_2}, a_{i_3}, \dots, a_{i_\ell}) \Rightarrow (a_{i_1} \mu_1 a_{i_2} \mu_2 a_{i_3} \dots \mu_{\ell-1} a_{i_\ell})$, unde $\{a_{i_1}, a_{i_2}, a_{i_3}, \dots, a_{i_\ell}\} \subseteq \mathcal{A} = \{a_1, \dots, a_m\}$, $a_{i_j} \neq a_{i_k}$, $j, k = 1..l$, $j \neq k$ și $\mu_i \in M$ este o relație pe $D_{i_j} \times D_{i_{j+1}}$, D_{i_j} este domeniul atributului a_{i_j} dacă attributele $a_{i_1}, a_{i_2}, a_{i_3}, \dots, a_{i_\ell}$ apar împreună (nu sunt vide) în $s\%$ dintre cele n instanțe (s fiind numit suportul regulii) și notăm $R' \subseteq R$ mulțimea de instanțe unde $a_{i_1}, a_{i_2}, \dots, a_{i_\ell}$ apar împreună și $\Phi(r_j, a_{i_k}) \mu_k \Phi(r_j, a_{i_{k+1}})$ este adevărat pentru fiecare $1 \leq k \leq l - 1$ pentru fiecare instanță r_j din R' ($c = |R'|/|R|$ fiind numit confidența regulii).

În afară de suport și confidență, o regulă poate fi caracterizată și de lungimea ei, care este numărul de atribute care apar în ea. Într-un set de date un număr mare de reguli de asociere relaționale pot fi găsite, dar în general suntem interesați doar de acele reguli ale căror confidență și suport este mai mare decât un prag ales de utilizator c_{min} și s_{min} . Aceste reguli sunt numite interesante.

1.4.2 Clustering

Clusteringul este considerat cea mai importantă metodă de învățare nesupervizată, a cărei scop principal este gruparea entităților dintr-un set de date în diferite grupe (numite *clusteri*) în așa fel încât entitățile într-un *cluster* să fie mai similare între ele decât cu entitățile din afara *clusterului*.

Fie $\mathcal{O} = \{O_1, O_2, \dots, O_n\}$ mulțimea obiectelor care trebuie să fie *clusterizați*. Pe parcursul procesului de clusterizare, un *cluster* C este o mulțime nevidă de obiecte din \mathcal{O} , care aparțin împreună pe baza unei funcții de distanță d . O *partiție* $P = (C_1, C_2, \dots, C_k)$ este o mulțime de *clusteri*, astfel încât fiecare obiect din \mathcal{O} apare într-un singur *cluster* din P .

Algoritmii de *clustering* pot fi împărțiți în două grupe mari: algoritmi *ierarhici* și *partiționali*, iar algoritmi ierarhici pot fi împărțiți în algoritmi *aglomerativi* și *divizivi*. În cazul algoritmilor aglomerativi, procesul începe cu o partiție în care fiecare entitate este plasată în propriul *cluster*, iar numărul *clusterilor* scade treptat, până se ajunge la o partiție în care toate entitățile sunt în același *cluster*.

Un pas important pe parcursul *clusterizării* ierarhice este calculul distanței între doi *clusteri*, notată cu $dist(C_i, C_j)$, pentru ca această distanță decide care *clusteri* vor fi uniți. Cele mai cunoscute metrice de distanță în literatura de specialitate sunt *single-link* - unde distanța minimă dintre două entități din C_i și C_j este considerată -, *complete-link* - unde distanța maximă este considerată - și *average-link* - unde media distanțelor se folosește.

Algoritmi de *clustering* partiționali [Han05] crează o singură partiție care conține K *clusteri*. Algoritmul pornește cu K centre inițiale pentru *clusteri* și fiecare entitate este adăugată în *clusterul* cu centrul cel mai apropiat. După aceea, centrele sunt recalculat și entitățile sunt adăugate din nou. Aceste pași se repetă până se ajunge la *clusteri* stabili.

Capitolul 2

Abordări noi pentru Remodularizare Soft

În acest capitol, care este original, abordăm problema de remodularizare soft. Remodularizare soft înseamnă modificarea structurii interne a unui sistem informatic, fără a afecta comportamentul extern, și este o activitate importantă, pentru că ajută la menținerea unei proiectări bune a sistemului. În acest capitol sunt prezentate abordările noastre pentru remodularizare la nivel de pachet și la nivel de clasă folosind *clustering*.

Abordările prezentate în acest capitol reprezintă lucrările noastre originale, publicate în [Mar14b], [MCC12], [Mar12c], [Mar12b], [Mar12a], sau în curs de recenzie: [MCC14].

Motivație. În zilele noastre sistemele informatice trebuie să se adapteze cerințelor utilizatorilor sau a mediului în care sunt folosite. Majoritatea modificărilor legate de această adaptare sunt făcute după lansarea softului, și în multe cazuri influențează negativ calitatea sistemului. Acest lucru face întreținerea mai complicată și costisitoare. Pentru a evita aceste costuri, întreținerea calității sistemului este un scop important, dar, din cauza dimensiunii sistemelor, e aproape imposibil pentru un dezvoltator să aibă o viziune de ansamblu asupra sistemului și să decidă modificările necesare. De aceea, diferite metode și unelte care găsesc ce modificări ar trebui făcute pentru îmbunătățirea calității sistemului sunt binevenite.

Am plecat de la ideea că un sistem informatic poate fi considerat un set de date, care conține înregistrări care pot fi elementele sistemului informatic. În acest context aplicarea algoritmilor de *clustering* pentru a grupa entitățile care aparțin împreună pare naturală. Pentru a ghida acest proces, am decis să folosim metrici soft, pentru că s-au dovedit a fi utile în alte abordări.

2.1 Remodularizare la nivel de pachete

În această secțiune se prezintă abordarea originală pentru restructurarea unui sistem informatic la nivel de pachete. Această abordare primește un sistem informatic și îl remodularizează la nivel de pachet folosind *clustering*-ul ierarhic pentru a obține pachete structurate mai bine. Considerând o structură de pachete existentă, metoda propusă în această secțiune poate fi folosită și pentru a sugera dezvoltatorului pachetul potrivit pentru o clasă nou adăgată.

2.1.1 Abordarea *CASP*

În această subsecțiune se introduce o abordare bazată pe *clustering* (*CASP* - *Clustering Approach for Software Packages Restructuring*) pentru remodularizare la nivel de pachete, alcătuită din doi pași:

- **Colectarea datelor** - Sistemul informatic este analizat pentru a extrage din el informația relevantă despre clasele, metodele, atributele și relațiile dintre ele.
- **Gruparea** - Clasele din sistemul informatic sunt grupate în pachete folosind informația extrasă la pasul anterior și un algoritm de *clustering* (*HASP* în abordarea noastră). Scopul acestui pas este obținerea unei împărțiri a sistemului în pachete.

2.1.2 Gruparea în pachete

În continuare se introduce un algoritm nou de *clustering* ierarhic aglomerativ (*HASP* - *Hierarchical Clustering Algorithm for Software Packages Restructuring*) care urmărește identificarea unei partiții a unui *framework* S , care corespunde unei structuri de pachete bune. În această partiție fiecare *cluster* reprezintă câte un pachet al sistemului informatic.

Au fost identificate șapte caracteristici relevante pentru problema împărțirii claselor în pachete. Aceste caracteristici măsoară coeziunea, cuplarea, potențialul de refolosire și cât de similare sunt elementele dintr-un pachet. Ele au fost agregate într-o singură valoare, numită *score*, în felul următor:

$$score(K_i \cup K_j, \mathcal{K}^*) = \frac{\sum_{i=1}^2 w_i \cdot F_i - w_3 \cdot F_3}{|K_i \cup K_j|^2 - 1} + \sum_{i=4}^7 w_i \cdot F_i \quad (2.1)$$

unde K_i și K_j reprezintă doi *clusteri* care voi fi uniți, \mathcal{K}^* reprezintă partiția fără K_i și K_j , F_i , $1 \leq i \leq 7$, sunt valorile celor șapte caracteristici și w_i ($0 \leq w_i \leq 1$) sunt valorile ponderilor asociate acestor caracteristici. Aceste ponderi au fost determinate folosind o căutare de tip *grid* și un sistem informatic pentru care se cunoaște o structură corectă a pachetelor. Pentru a ghida procesul de căutare de tip *grid* am introdus măsura *CIP* (*Cohesion of Identified Package*) care determină distanța dintre două partiții.

Algoritmul *HASP* este bazat pe ideea de *clustering* ierarhic aglomerativ. La fiecare pas, perechea de *clusteri* pentru care *score*-ul este cel mai mare este unită. Procesul de *clustering* se efectuează până când se ajunge la o partiție cu un singur *cluster*, și creează o serie de partiții cu un număr de *clusteri* care tot scade. Pentru a identifica “cea mai bună” partiție dintre toate cele generate, am introdus o nouă măsură, *overallScore*, care e calculată considerând valoarea pentru o variantă ușor modificată a măsurii *score*.

2.1.3 Atribuirea claselor în pachete

O altă problemă pe care dezvoltatorii soft o întâlnesc des este găsirea pachetului potrivit pentru o nouă clasă. Pentru a rezolva această problemă am introdus algoritmul *AssignClass*. Să considerăm că $\mathcal{K} = \{K_1, K_2, \dots, K_v\}$ este partiția care reprezintă structura actuală de pachete a sistemului. O clasă nouă, C , este adăugată în sistem. Pentru a găsi cel mai potrivit pachet pentru C , se calculează $score(K_i \cup C, \mathcal{K}^*)$ pentru fiecare pachet $K_i \in \mathcal{K}$. Pachetul pentru care *score* este maxim este cel în care C ar trebui adăugată.

2.1.4 Evaluarea experimentală

Pentru evaluarea experimentală a abordării noastre au fost folosite două *framework*-uri disponibile public, *Commons DbUtils* [DbU] și un *Reinforcement Learning* (*Învățare cu întărire*) *framework* [rlf]. Pentru ambele *framework*-uri au fost efectuate două experimente: am executat algoritmul *HASP* pentru a identifica o împărțire bună a claselor în pachete și am identificat pachetul potrivit pentru anumite clase folosind algoritmul *AssignClass*.

2.1.4.1 Commons DbUtils. Rezultate

Am aplicat algoritmul *HASP* pe *framework*-ul *DbUtils*; acesta a returnat o partiție cu patru *clusteri* (pachete) ca structura optimă, în locul structurii originale alcătuite din trei pachete. Analizând rezultatele algoritmului am concluzionat ca partiția sugerată de algoritmul *HASP* este mai bună decât cea originală.

La al doilea experiment am executat de mai multe ori algoritmul *AssignClass*, de fiecare dată eliminând câte o clasă din *DbUtils* și încercând să găsim pachetul potrivit pentru ele. Din 19 rulări ale algoritmului în 16 cazuri el a identificat pachetul corect, având o acuratețe de 84.2 %.

2.1.4.2 Reinforcement Learning. Rezultate

Am aplicat algoritmul *HASP* pentru *framework*-ul Reinforcement Learning; acesta a returnat o partiție cu 11 pachete, în locul partiției originale (despre care știm că este cea corectă) cu 10 pachete. Totuși, cele două partiții sunt foarte similare. Am calculat valoarea metricii *CIP* pentru a vedea cât de aproape este partiția generată de algoritmul *HASP* de cea originală, și a dat valoarea 0.95, care este destul de mare.

Pentru al doilea experiment am creat trei clase noi și am folosit algoritmul *AssignClass* pentru a sugera pachetul potrivit pentru ele. Toate cele trei clase au fost plasate în pachetul corect de algoritmul *AssignClass*.

2.1.5 Analiză și comparație cu metode existente

Am efectuat anumite experimente adiționale pe lângă cele prezentate mai sus pentru a investiga eficacitatea algoritmului *HASP*. Pentru aceste experimente am folosit pe lângă *framework*-ul *DbUtils* alte două *framework*-uri disponibile public *Email* [Ema] și *EL* [EL].

În această analiză am calculat valoarea diferitor metrici și măsuri din literatura de specialitate pentru structura originală a sistemelor și pentru structura sugerată de algoritmul *HASP*. Am calculat în total 29 de valori, și în 25 de cazuri partiția generată de algoritmul *HASP* avea o valoare mai bună sau egală cu partiția originală.

Dintre metodele raportate în literatura de specialitate, cea mai apropiată de metoda noastră este abordarea din [AAM11], care folosește *clustering* ierarhic precum o face și metoda noastră. Acest articol prezintă două abordări; ambele au fost aplicate pe *framework*-ul *DbUtils*. Prima abordare este bazată pe numărul de apeluri la constructorul claselor, dar aceasta nu este o direcție bună în cazul *framework*-urilor, care în general au multe interfețe și clase abstracte care nu sunt inițializate niciodată. Pentru *framework*-ul *DbUtils* această abordare nu sugerează nicio modificare.

A doua metodă prezentată în [AAM11] folosește o reprezentare multi-dimensională bazată pe dependențe create de apeluri de metode. Toate cele trei măsuri de distanță la *clustering* (*single-link*, *complete-link* și *average-link*) au fost încercate; de asemenea a fost introdus și un algoritm nou cu o complexitate computațională mai mică. Am încercat cele trei metode de *clustering* pe sistemul *DbUtils*, folosind ca și criteriu de oprire momentul când distanța între toate perechile de clusteri a fost 1 (cea mai mare distanță posibilă). După aceea am calculat valoarea măsurii *CIP* între rezultatele primite și structurile originale și sugerate de *HASP*. Valorile între 0.3 și 0.43 sugerează că rezultatele sunt departe de ambele.

2.1.6 O comparație a măsurii *overallScore* cu alte măsuri

În continuare este prezentat un studiu efectuat pentru a evalua cât de bine sunt structurate pachetele într-un sistem informatic, folosind cele două măsuri introduse mai sus: *overallScore*

și *CIP*. Prin experimentele efectuate, urmărim să accentuăm că cele două măsuri sunt bine corelate. Astfel, în locul măsurii *CIP*, care necesită cunoașterea unei structuri de pachete corecte, *overallScore* poate fi folosită pentru a evalua pachetele dintr-un sistem informatic.

Pentru experimente au fost folosite cele trei sisteme disponibile public utilizate și mai sus: *DbUtils*, *Email* și *EL*. Pentru fiecare sistem am considerat patru structuri de pachete diferite: originala, cea sugerată de algoritmul *HASP* și două “partiții greșite” create manual prin modificarea structurii originale a pachetelor.

Pentru experimente am calculat valoarea măsurilor *overallScore* și *CIP* pentru fiecare partiție. După aceea am calculat valoarea a două măsuri de corelație bazate pe ranguri: corelația Spearman [Spe04] și *Footrule*-ul lui Spearman [DG77]. Ambele măsuri de corelație sugerează că măsurile *overallScore* și *CIP* sunt bine corelate.

Pentru a compara aceste măsuri cu alte măsuri din literatura de specialitate, am calculat valoarea a șapte metrici introduse în [DABH11] pentru cele patru partiții pentru cele trei sisteme. După aceea am calculat *Footrule*-ul lui Spearman pentru aceste metrici și măsura *overallScore*. Comparația valorilor este prezentată în Figura 2.1.

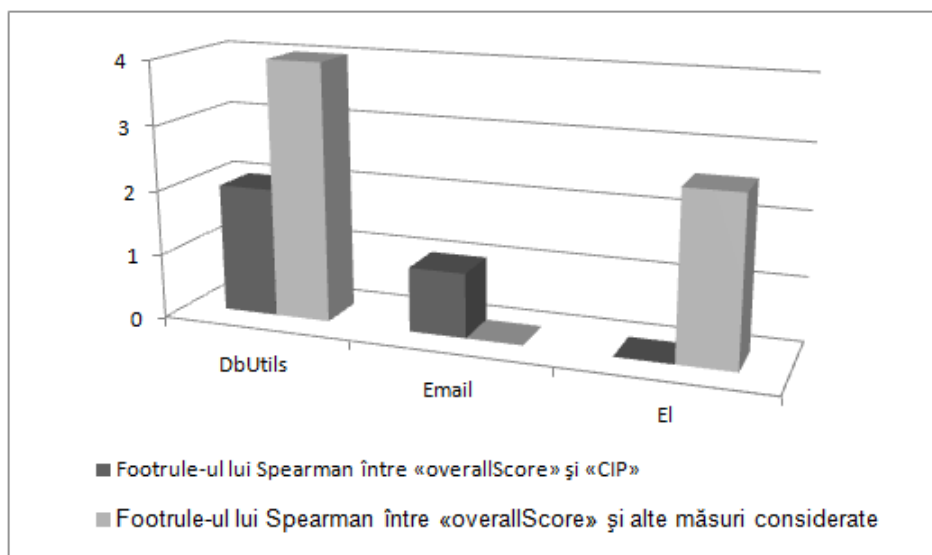


Figura 2.1: Comparația valorilor *Footrule*.

Considerând comparația valorilor *Footrule* prezentate în Figura 2.1 putem concluziona că măsura *overallScore* este mai potrivită pentru evaluarea unei partiții a sistemului informatic decât restul metricilor considerate. În plus, *overallScore* poate fi folosit în locul măsurii *CIP*, deoarece cele două sunt bine corelate.

2.2 Remodularizare la nivel de clasă

În această secțiune este prezentată o abordarea originală pentru remodularizarea la nivel de clasă a unui sistem informatic [MCC12, Mar12c, Mar12b, Mar12a]. Vor fi prezentate metode care sunt capabile să identifice automat o structură îmbunătățită a sistemului informatic folosind valoarea diferitor metrici soft. Această structura internă îmbunătățită poate fi realizată aplicând diferite restructurări pe sistemul original.

2.2.1 Restructurare bazată pe metrici soft

În următoarele subsecțiuni următoarele notații vor fi folosite: vom considera că un sistem informatic este o mulțime $S = \{s_1, s_2, \dots, s_n\}$ unde s_i , $1 \leq i \leq n$, se numește o entitate și

poate fi o clasă sau o metodă dintr-o clasă. Am dorit să folosim valoarea diferitor metrici pentru a reprezenta aceste entități, deci am avut nevoie de metrici care pot fi calculate atât pentru clase cât și pentru metode. Am ales următoarele 5 metrici: *Relevant Properties (RP)*, *Depth in Inheritance Tree (DIT)* [CK91], *Number of Children (NOC)* [CK91], *Fan-In (FI)* [HK81, Mai09], *Fan-Out (FO)* [HK81, Mai09]. Dintre acestea, valoarea lui *RP* este o mulțime, iar restul metricilor au ca valoare numere naturale.

2.2.1.1 Modelul vectorial și funcția de distanță

Ideea de bază a abordării noastre este să caracterizăm entitățile dintr-un sistem informatic printr-o listă de valori a metricilor prezentate mai sus. Astfel, fiecare entitate s_i ($1 \leq i \leq n$) din sistemul informatic va fi reprezentată ca un vector cu 5 componente: $(rp(s_i), dit(s_i), noc(s_i), fi(s_i), fo(s_i))$. Valorile metricilor au fost scalate în intervalul $[0,1]$ acolo unde a fost necesar. Pentru a aplica un algoritm de *clustering* pe această reprezentare am definit o funcție de distanță, care este prezentată pe Ecuația 2.2.

$$d(s_i, s_j) = \begin{cases} 0 & \text{dacă } i = j, \\ \sqrt{\frac{1}{5} \cdot \left(1 - \frac{|s_{i1} \cap s_{j1}|}{|s_{i1} \cup s_{j1}|} + \sum_{k=2}^5 (s_{ik} - s_{jk})^2 \right)} & \text{dacă } s_{i1} \cap s_{j1} \neq \emptyset, \\ \infty & \text{altfel} \end{cases} \quad (2.2)$$

2.2.1.2 Algoritmul ARI

Folosind metricile soft prezentate mai sus, modelul vectorial și funcția de distanță am definit algoritmul *ARI* (*Automatic Refactoring Identification*). Ideea principală a algoritmului este identificarea unei partiții $\mathcal{K} = \{K_1, K_2, \dots, K_m\}$ a sistemului informatic S care corespunde unei structuri îmbunătățite a lui S . Fiecare *cluster* K_i corespunde unei clase din această structură îmbunătățită.

Algoritmul *ARI* pornește cu o partiție vidă. După aceea pentru fiecare metodă din S clasa cea mai apropiată (considerând funcția de distanță din Ecuația 2.2) este căutată. Dacă se găsește o clasă cu o distanță mai mică decât ∞ metoda ori este plasată în *clusterul* clasei (daca clasa este deja într-un *cluster*) ori un *cluster* nou este creat unde cele două entități sunt plasate. Dacă distanța metodei față de toate clasele este ∞ , căutăm în *clusterii* fără clase metoda cea mai apropiată de această metodă. Dacă toate distanțele sunt ∞ vom crea un *cluster* nou și vom adăga metoda în acel *cluster*. Altfel o adăugăm în *clusterul* găsit. Când toate metodele au fost adăugate, calculăm distanța dintre toate perechile de clase. Dacă este vreo pereche cu distanță sub ∞ vom uni *clusterii* respectivi.

2.2.1.3 Algoritmul HAC

Folosind aceleași metrici, reprezentare și funcție de distanță, am introdus o extensie a algoritmului *ARI*, numit algoritmul *HAC* (*Hierarchical Agglomerative Clustering*). Este bazată pe *clustering* ierarhic aglomerativ și folosește o funcție euristică care exprimă că doi clusteri sunt uniți doar dacă distanța dintre ei este mai mică ca 1. Motivul pentru folosirea acestei funcții este că distanțe mai mari ca 1 apar numai pentru entități necoezive.

2.2.1.4 Restructurări identificate

Atât algoritmul *ARI* cât și algoritmul *HAC* creează o partiție a sistemului informatic S , care corespunde unei structuri interne bune. Comparând această partiție cu partiția originală, o listă de restrukturări poate fi identificată, care ar îmbunătăți structura sistemului

Metrică	Ex 1	Ex 2	Ex 3	Ex 4	Metrică	Ex 1	Ex 2	Ex 3	Ex 4
NOA	A	A	A	A	ICH	C	A	A	A
NOM	C	A	A	C	CA	A	A	A	A
CBO	A	A	A	A	INS	A	A	A	A
DAC	A	A	A	A	LCOM1	A	S	S	S
TCC	A	C	C	C	LCOM2	A	S	S	S
LCC	A	C	C	C	LCOM4	A	A	S	S
MPC	S	A	S	S	LCOM5	A	C	A	C
RFC	C	A	C	A	LD	A	A	C	A

Tabela 2.1: Direcția în care valoarea metricilor s-a modificat pentru cele două variante a exemplurilor.

informatic. Pentru ambii algoritmi aceste restructurări sunt: *Move Method*, *Inline Class* și *Extract Class*.

2.2.2 Restructurare bazată pe metrici soft agregate

Această subsecțiune prezintă un al treilea algoritm pentru restructurarea sistemelor informatice, care, similar cu algoritmi prezentați mai sus, folosește metrici soft, dar agreghează valorile lor într-o singură valoare.

Pentru această abordare am avut nevoie de metrici soft care pot fi calculate pentru o clasă. De asemenea, din moment ce aceste valori vor fi folosite să caracterizeze calitatea unei clase (și calitatea sistemului informatic prin clasele lui) intenția noastră a fost să alegem metrici care măsoară dimensiunea, coeziunea și cuplarea sistemului informatic. Studiind lista metricilor implementate de diferite instrumente am selectat următoarele 16 metrici: *Afferent Coupling (CA)* [BJWD99], *Coupling Between Objects (CBO)* [CK94], *Data Abstraction Coupling (DAC)* [LH93], *Information Flow Based Cohesion (ICH)* [LLWW95], *Instability (INS)*, *Tight Class Coupling (TCC)* [BK95], *Loose Class Coupling (LCC)* [BK95], *Lack of Cohesion in Methods 1 (LCOM1)* [CK94], *Lack of Cohesion in Methods 2 (LCOM2)* [HS96], *Lack of Cohesion in Methods 4 (LCOM4)* [HM95], *Lack of Cohesion in Methods 5 (LCOM5)*, *Locality of Data (LD)* [HM95], *Message Passing Coupling (MPC)* [LH93], *Number of Attributes (NOA)* [LH93], *Number of Methods (NOM)* [LH93], *Response for a Class (RFC)* [CK94].

2.2.2.1 Studiu de caz

Deoarece valori ideale pentru aceste metrici soft nu sunt cunoscute, am efectuat un studiu, pentru a analiza cum se modifică valoarea acestor metrici dacă îmbunătățim sistemul informatic. Am ales patru cazuri de studii, fiecare fiind un exemplu mic cu 2-3 clase, și în fiecare exemplu a fost cel puțin o metodă care era plasată într-o clasă diferită față de clasa unde ar fi trebuit să fie.

Am calculat valoarea metricilor prezentate în subsecțiunea anterioară pentru aceste exemple. După aceea, am “corectat” exemplele, mutând metodele în clasele unde ar trebui să fie, și am calculat valoarea metricilor din nou. Am fost interesate să vedem cum se modifică valoarea metricilor pentru cele două versiuni. Rezultatul studiului este prezentat pe Tabelul 2.1, unde eticheta *A* înseamnă că valoarea metricii a rămas aceeași, *C* înseamnă că a crescut, iar *S* înseamnă că a scăzut pentru varianta corectă.

Cel mai important lucru de observat pe Tabelul 2.1 este că pentru fiecare metrică direcția schimbării este consistentă: pentru fiecare exemplu valoarea unei metrici ori nu s-a modificat deloc, sau, dacă s-a modificat, tot timpul s-a modificat în aceeași direcție. Acest lucru este importantă, pentru că ne ajută la comparația a două versiuni a aceluiași sistem.

2.2.2.2 Algoritm de restructurare bazat pe metrici soft agregate

Ideea principală a algoritmului nostru este să caracterizăm fiecare clasă din sistemul informatic printr-o singură valoare agregată a acestor metrici soft. Pentru asta, prima dată am normalizat valorile metricilor. Am efectuat normalizarea în așa fel încât cu cât e mai mică valoarea unei metrici, cu atât e mai bună proiectarea sistemului.

După normalizarea valorilor am introdus o măsură, $M(C)$, definită pentru o clasă C ca suma valorilor normalizate pentru cele 16 metrici soft. Cu cât e mai mică valoarea acestei metrici, cu atât e mai bună structura internă a clasei C . Deoarece un sistem informatic S este alcătuit din mai multe clase, putem defini o măsură agregată, $AM(S)$, ca media valorilor $M(C)$ pentru toate clasele din sistem.

Folosind măsura $AM(S)$ și *clustering* ierarhic aglomerativ, abordarea noastră prima dată plasează fiecare metodă în propriul *cluster* (care reprezintă o clasă cu o singură metodă), și va uni *clusteri* până când ajunge la o partiție unde toate metodele sunt în același *cluster*. Din seria partițiilor generate se returnează ca soluție partiția cu cea mai mică valoare pentru măsura $AM(S)$.

2.2.3 Evaluare experimentală

Algoritmul *ARI*. Pentru a evalua performanța algoritmului *ARI* am folosit două studii de caz. Primul este un exemplu simplu, cu două clase, unde o clasă conține o metodă care ar trebui să fie în cealaltă clasă. Pentru acest exemplu algoritmul *ARI* sugerează mutarea metodei în cealaltă clasă.

Al doilea studiu de caz folosit pentru evaluarea algoritmului *ARI* este sistemul informatic *JHotDraw* versiunea 5.1 disponibil public [Gam]. Am ales sistemul *JHotDraw* pentru că este un exemplu bine cunoscut pentru folosirea șabloanelor de proiectare și pentru că este un proiect complex, alcătuit din 173 clase, 1375 metode și 475 atribute.

Algoritmul *ARI* a identificat 20 de restructurări *Move Method* în sistemul *JHotDraw*. După o analiză a rezultatelor am ajuns la concluzia că 11 sunt justificate și 9 sunt greșite.

Algoritmul *HAC*. Am folosit sistemul informatic *JHotDraw* ca studiu de caz și la algoritmul *HAC*. Rezultatul algoritmului a conținut doar trei *clusteri* (reprezentând clase) care erau diferiți de structura existentă a sistemului. Acești trei *clusteri* reprezintă restructurări posibile de tip *Extract Class*, adică clase noi ar trebui create și metodele selecționate ar trebui mutate în aceste clase noi. După o analiză a rezultatelor am ajuns la concluzia că o restructurare este justificată, iar celelalte două sunt parțial justificate, adică o parte din metodele sugerate ar trebui mutate în clase noi.

Abordare bazată pe metrici soft agregate. Evaluarea abordării bazate pe metrici soft agregate a fost efectuată până acum doar pe exemple mici (cele patru studii de caz prezentate în Subsecțiunea 2.2.2), pentru că anumiți pași din algoritm trebuie să fie executați manual. Motivul este că pentru a avea rezultate valide, trebuie să ne asigurăm la fiecare pas că nu am schimbat comportamentul sistemului.

Am aplicat algoritmul pentru trei exemple din Subsecțiunea 2.2.2 și în toate cele 3 cazuri a identificat partiția corectă pentru studiul de caz.

2.2.3.1 Analiză comparativă

Analiză comparativă a celor trei algoritmi. Algoritmii *ARI* și *HAC* sunt destul de similari, folosesc același model vectorial pentru a reprezenta entitățile din sistemul informatic și aceleași metrici soft. Cu toate acestea, algoritmul *HAC* folosește *clustering* ierarhic aglomerativ, o metodă de învățare nesupervizată.

Ambii algoritmi au identificat majoritatea sistemului *JHotDraw* ca fiind corect proiectat, așa cum ne-am așteptat, deoarece *JHotDraw* este considerat un exemplu de proiectare corectă. La restructurările sugerate, algoritmul *HAC* a avut doar 5 greșeli, în timp ce algoritmul *ARI* a avut 9. De asemenea, algoritmul *HAC* a identificat 3 restructurări *Extract Class*, care nu au fost identificate de algoritmul *ARI*, astfel demonstrând că folosirea metodelor de învățare nesupervizate este folositoare, pentru că pot identifica șabloane ascunse în date.

Algoritmul bazat pe metrice soft agregate are un studiu de caz comun cu algoritmul *ARI*, acel exemplu artificial simplu care a fost primul caz de studiu pentru *ARI*. Pentru acest exemplu ambii algoritmi au identificat restructurarea corectă.

Analiză comparativă cu metode existente. Seng et al. aplică în [SSB06] o căutare multiobiectivă ponderată unde metrice soft sunt combinate într-o singură funcție obiectivă. Această abordare produce parțial aceleași rezultate ca și algoritmul *ARI*, dar *ARI* are mai puține sugestii greșite, este deterministic și are un timp de rulare mai mic.

În [CS06] este prezentată o abordare care folosește *clustering* pentru a îmbunătăți structura de clase într-un sistem informatic. Această abordare folosește un algoritm de *clustering* partițional, *kRED* (*k-means for REfactorings Determination*). Spre deosebire de *kRED*, algoritmul *ARI* poate să identifice și restructurări de tip *Extract Class*. Comparat cu [RR11] și [FTCS09] algoritmi *HAC* și *ARI* pot identifica trei tipuri de restructurări (nu numai *Extract Class*).

Abordarea cea mai apropiată de abordarea noastră este algoritmul *HARS*, [CS07], care folosește doar metrica *Relevant Properties* pentru a calcula distanța dintre două entități, și care consideră și atributele ca fiind entități. În [CS07] sistemul informatic *JHotDraw* este folosit ca și caz de studiu, și găsește doar una dintre cele trei restructurări *Extract Class* pe care algoritmul *HAC* le găsește (dar găsește două restructurări *Move Attribute* pe care nici *ARI* nici *HAC* nu le pot găsi).

2.3 Concluzii și cercetări ulterioare

Acest capitol a prezentat diferite metode pentru restructurarea automată a unui sistem informatic pentru a îmbunătăți calitatea lui. Am început cu prezentarea abordării *CASP* și a algoritmului *HASP*, care este o metodă de restructurare la nivel de pachete. Am evaluat experimental abordarea noastră folosind două studii de caz, și rezultatele arată că algoritmul *HASP* poate găsi o structură bună a pachetelor.

Am prezentat de asemenea trei algoritmi care folosesc valoarea diferitor metrice soft și *clustering* ierarhic pentru restructurarea la nivel de clase a unui sistem informatic. Evaluarea experimentală a acestor metode arată că ele pot identifica o partiție a sistemului care are o proiectare mai bună.

În viitor dorim să extindem evaluarea experimentală pentru toate abordările și pe alte sisteme informatice disponibile public. De asemenea dorim să îmbunătățim abordările noastre eliminând deficiențele posibile.

Capitolul 3

Abordări noi pentru Detectarea Defectelor în Sisteme Informatice

În acest capitol, care este original, abordăm problema de detectare a defectelor în sisteme informatice, o problemă importantă în inginerie soft, deoarece, din cauza dimensiunii și complexității mari a sistemelor informatice, găsirea manuală a acestor defecte este din ce în ce mai complicată. Prin urmare, abordările inteligente, care îi ajută dezvoltatorii soft în identificarea părților defective care necesită atenție, sunt binevenite. În literatura de specialitate se pot găsi abordări bazate pe instruire automată pentru ambele direcții, *detectarea defectelor de proiectare* și *predicția defectelor*. În acest capitol prezentăm o abordare nouă, și anume, aplicarea regulilor de asociere relaționale pentru ambele direcții.

Cele două abordări prezentate în acest capitol sunt lucrări originale, publicate în [CMC14a], [Mar13b], [CMC14b] și [Mar13a].

Motivație. Toate sistemele informatice trec prin multe schimbări pe parcursul ciclului de viață, altfel devin învechite și nu mai sunt folosite. Aceste modificări, executate în general pe parcursul întreținerii, pot adăuga funcționalități noi, corecta defectele existente, sau pot fi făcute pentru adaptarea sistemului la mediul în care este folosit. Ele pot degrada structura sistemului, ceea ce, la rândul ei, va face modificările următoare și mai grele. Pentru a evita acest ciclu vicios, identificarea entităților cu defecte este foarte importantă.

Un sistem informatic poate fi reprezentat ca un set de date unde elementele sunt clasele, folosind o reprezentare multi-dimensională bazată pe metrici soft. Din acest set de date informația importantă poate fi extrasă. Diferite tipuri de relații pot fi definite între valorile metricilor soft și reguli de asociere relaționale pot fi extrase. Aceste reguli pot furniza șabloane interesante, care pot indica clasele cu defecte de proiectare, sau ele pot descrie entitățile defective și non-defective din sistemul informatic.

3.1 Model teoretic

În această secțiune prezentăm modelul teoretic care este folosit atât la abordarea de detectare a defectelor de proiectare, cât și la abordarea de predicția defectelor.

Ideea principală a modelului este reprezentarea entităților dintr-un sistem informatic ca un vector multidimensional, elementele cărui sunt valorile diferitor metrici soft calculate pentru entitatea respectivă. Considerăm că un sistem informatic este o mulțime de componente (entități) $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$. Deoarece dorim să identificăm entitățile rău-structurate, vom considera o mulțime de metrici soft care caracterizează structura internă a entităților. Prin urmare, avem un set de caracteristici alcătuit din metrici soft $\mathcal{SM} = \{sm_1, sm_2, \dots, sm_k\}$ și astfel, fiecare entitate $s_i \in \mathcal{S}$ din sistemul informatic poate fi reprezentată ca un vector k -dimensional, compus din valoarea metricilor soft din \mathcal{SM} .

3.2 Detectarea Defectelor de Proiectare

Această secțiune prezintă abordarea noastră, numită *SDDRAR* (*Software Design Defect detection using Relational Association Rules*), pentru detectarea defectelor de proiectare într-un sistem informatic [CMC14a].

3.2.1 Algoritmul *DRAR*

Regulile de asociere relaționale sunt o extensie a regulilor de asociere ordinale (așa cum a fost prezentat și în Secțiunea 1.4.1), care permit definirea relațiilor mai generale între atribute, nu doar relații ordinale. Un algoritm similar cu algoritmul Apriori, intitulat *DOAR* (*Discovery of Ordinal Association Rules*) a fost introdus în [CSTM06]. Acest algoritm găsește într-un mod eficient toate regulile de asociere ordinale care există într-un set de date.

Algoritmul *DOAR* a fost extins în [CMC14a] spre algoritmul *DRAR* (*Discovery of Relational Association Rules*) pentru a găsi reguli de asociere relaționale interesante de orice lungime într-un set de date.

3.2.2 Abordarea *SDDRAR*

Considerăm o mulțime de sisteme informatice bine-proiectate, notată cu \mathcal{S}_{good} . Folosind modelul teoretic din Secțiunea 3.1, considerăm că $\mathcal{DS} = \{ds_1, ds_2, \dots, ds_n\}$ este setul de date compus din reprezentare k -dimensională a entităților din sistemele informatice din \mathcal{S}_{good} . Setul de caracteristici folosit pentru acest set de date este o mulțime a metricilor soft $\mathcal{SM} = \{sm_1, sm_2, \dots, sm_k\}$ relevante pentru procesul de detectare a defectelor de proiectare.

Pentru a detecta clasele rău-structurate folosind reguli de asociere relaționale, vom executa următorii pași

1. Colectarea și preprocesarea datelor
2. Construirea modelului *SDDRAR*
3. Testare

3.2.2.1 Colectarea și preprocesarea datelor

În cursul acestui pas, o analiză statistică este efectuată pe setul de date \mathcal{DS} pentru a găsi o submulțime a caracteristicilor, care sunt relevante pentru această problema. Pentru a determina dependențele dintre caracteristici folosim coeficientul de corelație Pearson [Tuf11]. Pentru fiecare metrică soft, sm , calculăm *corelația ei medie absolută* (notată cu $avg(sm)$) ca media corelațiilor metricii cu restul metricilor. Calculăm media, m , și deviația standard, $stdev$, a acestor valori, și eliminăm acele metrici soft pentru care $|avg(sm) - m| > stdev$.

3.2.2.2 Construirea modelului *SDDRAR*

Primul pas pentru a construi modelul *SDDRAR* este definirea relațiilor posibile dintre metricile soft. Aceste relații sunt necesare pentru procesul de extragere a regulilor de asociere relaționale. După ce relațiile au fost definite, extragem toate regulile de asociere relaționale, de orice lungime din setul de date \mathcal{DS} , având suport minimum s_{min} și confidență minimă c_{min} . Mulțimea regulilor este notată cu \mathcal{RAR} . Aceste reguli vor fi folosite pentru a identifica entitățile rău-proiectate din sistemul informatic.

3.2.2.3 Testare

În cursul testării un nou sistem informatic, \mathcal{S}_{new} , trebuie analizat pentru a determina entitățile care au o proiectare greșită. Prima dată \mathcal{S}_{new} este analizată și reprezentarea k -dimensională a entităților din \mathcal{S}_{new} este calculată. După aceea, următorii trei pași sunt executați:

1. **Pasul de calcul al erorilor.** Pentru fiecare entitate e_i din sistemul analizat, *numărul de erori*, $err(e_i)$, este calculat ca numărul de reguli de asociere relaționale din \mathcal{RAR} care nu sunt verificate de reprezentarea k -dimensională. După aceea, *procentul de erori*, pe_{e_i} , este calculată ca numărul de erori împărțite la numărul de reguli.
2. **Pasul de detectare.** După determinarea *numărului de erori* pentru fiecare entitate, raportăm ca entități cu posibile erori de proiectare acele entități care încalcă un număr destul de mare de reguli.
3. **Analiza defectelor.** Pentru entitățile raportate la pasul anterior, calculăm pentru fiecare metrică soft *numărul de erori*, ca numărul de reguli de asociere relaționale binare din \mathcal{RAR} care nu sunt respectate de entitatea respectivă. Aceste numere sunt analizate mai departe pentru a determina cauza defectului.

3.2.3 Evaluare experimentală

Pentru evaluarea experimentală a abordării *SDDRAR*, am ales ca entități clasele dintr-un sistem informatic. După cum a fost precizat în Secțiunea 3.1, pentru a folosi modelul vectorial avem nevoie de o mulțime de metrici soft relevante pentru a caracteriza proiectarea unei clase. Am ales să folosim aceeași 16 metrici care au fost folosite și pentru abordarea noastră bazată pe metrici agregate (Secțiunea 2.2.2).

3.2.3.1 Date de antrenament

Pentru pasul de antrenament este nevoie de o mulțime de sisteme informatice bine-proiectate, \mathcal{S}_{good} , de unde reprezentarea k -dimensională a claselor este calculată. Pentru implementarea curentă a abordării noastre am ales un singur sistem informatic bine-proiectat, sistemul disponibil public *JHotDraw*, versiunea 5.1 [Gam].

Pentru primul pas al abordării *SDDRAR* prima dată am analizat sistemul informatic *JHotDraw* și am extras din el setul de date \mathcal{DS} , alcătuit din reprezentarea 16-dimensională a claselor. Am scalat valoarea metricilor să fie în intervalul $[0,1]$, și am efectuat analiza statistică care a dus la eliminarea a patru metrici soft: *CA*, *NOM*, *RFC* și *TCC*.

Al doilea pas al abordării *SDDRAR* constă în extragerea din setul de date \mathcal{DS} a regulilor de asociere relaționale interesante, de orice lungime, având suport minim $s_{min} = 0.9$ și confidență minimă $c_{min} = 0.85$.

3.2.3.2 Studii de caz

Al treilea pas al abordării *SDDRAR* este *Testarea*, unde modelul *SDDRAR*, construit la pasul anterior, este folosit pentru a detecta entitățile rău-proiectate într-un sistem informatic. Pentru acest pas am considerat șase studii de caz, dintre care două au fost exemple simple, iar celelalte patru au fost sisteme informatice disponibile public pe pagina *SourceForge*. Rezultatele pentru cele șase studii de caz sunt prezentate în continuare:

- Primul studiu de caz simplu preluat din [SSL01] - abordarea *SDDRAR* a identificat clasa cu defect de proiectare.

- Al doilea studiu de caz simplu preluat din [Fow99] pagini 22-26 - abordarea *SDDRAR* a identificat clasa cu defect de proiectare.
- Sistemul informatic disponibil public *FTP4J* [FTP13] versiuni 1.5, 1.5.1, 1.6 și 1.6.1 - abordarea *SDDRAR* a identificat clasa *FTPClient* pentru toate versiunile. Analiza noastră a arătat că această clasă are un defect de proiectare, este *God Class*, deci abordarea *SDDRAR* a identificat-o corect.
- Sistemul informatic disponibil public *ISO8583* [ISO13] versiuni 1.5.2, 1.5.3 și 1.5.4 - abordarea *SDDRAR* a identificat clasa *MessageFactory* pentru toate versiunile. Analiza noastră a arătat că această clasă are un defect de proiectare, este *God Class*, deci abordarea *SDDRAR* a identificat-o corect.
- Sistemul informatic disponibil public *Profiler4J-Agent* [Pro13] versiuni 1.0-alpha5, 1.0-alpha6, 1.0-alpha7 și 1.0-beta1 - pentru versiunea 1.0-alpha5 clasa *Server* este raportată; pentru versiunea 1.0-alpha6 clasa *MemoryInfo* este raportată; pentru ultimele două versiuni clasa *Config* este raportată. Analiza noastră a arătat că *Server* este o clasa foarte cuplată, iar clasele *MemoryInfo* și *Config* sunt *Data Class*, deci au fost identificate corect de abordarea *SDDRAR*.
- Sistemul informatic disponibil public *WinRun4J* [Win13] versiuni 0.4.0, 0.4.1, 0.4.2, 0.4.3 și 0.4.4. - abordarea *SDDRAR* a identificat clasa *NativeBinder* pentru toate versiunile. Analiza noastră a arătat că această clasă are anumite probleme de proiectare, dar nu are un singur defect serios, deci am considerat că nu se poate decide clar dacă identificarea ei este corectă sau nu.

3.2.4 Analiză și comparație cu abordări existente

În această subsecțiune prezentăm o analiză din mai multe puncte de vedere a abordării *SDDRAR* și o comparație cu alte abordări similare din literatura de specialitate.

Precizia de detectare. Am considerat că pentru primele 5 studii de caz precizia este 1 (clasa identificată are defecte de proiectare), iar pentru ultimul am considerat o precizie de 0.5. Astfel, *precizia de detectare* medie este 0.917 pentru cele șase studii de caz.

Măsura de media erorilor. În afara determinării entităților cu posibile defecte de proiectare, abordarea *SDDRAR* a fost folosită pentru un studiu în modul următor. Pentru fiecare versiune a sistemelor informatice disponibile public am calculat media numărului de erori a claselor. Această valoare scade (sau rămâne la fel) pentru fiecare versiune nouă a sistemului. Acest lucru confirmă că valori mai mici pentru erori corespund unor proiectări mai bune, pentru că ne așteptăm ca versiunile mai mari să aibă o proiectare mai bună.

Robustețe. Procesul de detectare a entităților rău-structurate depinde de entitățile din setul de date *DS*. Aceste entități sunt considerate bine-proiectate, dar trebuie să studiem ce se întâmplă dacă introducem zgomot în această set de date. Am adăugat niște entități rău-proiectate în *DS* și am aplicat metoda *SDDRAR* la studii de caz prezentate mai sus. Așa cum ne-am așteptat, metoda *SDDRAR* este robustă, prezența unei cantități mici de zgomot nu influențează rezultatele ei.

Mulțimea relațiilor. Ultima perspectivă constă în analizarea mulțimii de relații \mathcal{R} considerate pentru extragerea regulilor de asociere relaționale. Am efectuat niște experimente cu combinații diferite a relațiilor posibile și am ajuns la concluzia că și aceste combinații pot identifica defecte de proiectare în seturi de date.

3.2.4.1 Comparație cu abordări existente

Pentru a compara abordarea noastră cu alte abordări din literatura de specialitate, am ales două instrumente disponibile public: *JDeodorant* un *plugin Eclipse* [JDe13] și *iPlasma*

	Ftp4J	ISO8583	WinRun4J
<i>SDDRAR</i>	FtpClient	MessageFactory	NativeBinder
JDeodorant	FtpClient (36) NVTASCIWriter (0) NVTASCIReader (0) FTPProxyConnector (2) FTPCommunication - - Channel (0)	ISOMessage (28) MessageFactory (38)	Launcher (28) FileAssociation (23) RegistryKey (0)
iPlasma	FTPClient - Brain Class (36) FTPFile - Data Class (19)	MessageFactory - God Class (38)	Closure - God Class (27) FileVerb - Data Class (26) NativeBinder - God Class (31)

Tabela 3.1: Clasele raportate de JDeodorant, iPlasma și metoda *SDDRAR*.

un instrument prezentat în [Mar02] și disponibil la [IP113]. Clasele cu defecte de proiectare raportate de cele două instrumente și abordarea noastră pentru prima versiune a trei studii de caz sunt prezentate în Tabelul 3.1. Pentru *JDeodorant* am considerat doar defectele *God Class*. Pentru *iPlasma* am specificat după fiecare clasă numele defectului. După fiecare clasă, între paranteze am adăugat numărul de erori raportate de abordarea *SDDRAR* pentru clasa respectivă.

În Tabelul 3.1 se poate observa că nici măcar cele două instrumente nu sunt de acord în legătură cu care clase au defecte de proiectare și care nu au. Ca și majoritatea abordărilor, am decis să analizăm manual clasele raportate, pentru a decide care instrument să fie folosit pentru comparație. În cazul instrumentului *iPlasma* suntem de acord cu clasele raportate. Instrumentul *JDeodorant* pare că încearcă să găsească restructurări *Extract Class*, nu *God Class*.

Comparând rezultatele raportate de instrumentul *iPlasma* cu rezultatele metodei *SDDRAR*, putem vedea că clasele raportate de *SDDRAR* sunt raportate și de *iPlasma*, iar pentru clasele raportate doar de *iPlasma* metoda noastră raportează un număr mare de erori (chiar dacă nu suficient de mare pentru a raporta clasa ca fiind defectivă). Folosind rezultatele raportate de *iPlasma* am calculat *recall*-ul pentru abordarea *SDDRAR*, care este 0.85.

3.2.5 Studiu despre variația parametrilor

Abordarea *SDDRAR* depinde de o serie de parametri, valoarea cărora influențează rezultatele raportate, astfel am efectuat un studiu pentru a investiga următoarele aspecte:

- Folosirea valorilor originale sau valorilor normalizate pentru metrici soft.
- Folosirea regulilor de asociere relaționale binare sau de orice lungime.
- Folosirea doar regulilor de asociere relaționale maximale, sau a tuturor regulilor extrase.
- Efectul modificării parametrului τ și a confidenței minime pe acuratețea detectării.

Ca rezultat al studiului am ajuns la concluzia că este mai bine să folosim valori normalizate pentru metrici soft, să extragem reguli de asociere relationale de orice lungime, și să folosim doar reguli maximale. Pentru ultimul aspect, dacă valoarea lui τ este scăzută, mai multe clase sunt raportate ca fiind defective, și majoritatea lor într-adevăr a avut mai multe sau

mai puține defecte. Scăzând valoarea minimă pentru confidență mai multe clase au fost raportate, dintre care multe erau *Data Class*.

3.3 Predicția Defectelor

Această secțiune prezintă abordarea noastră, numită *DPRAR - Defect Prediction using Relational Association Rules*, de învățare supervizată pentru predicția defectelor într-un sistem informatic [CMC14b]. Această abordare folosește modelul teoretic prezentat în Secțiunea 3.1 și algoritmul *DRAR* pentru extragerea regulilor de asociere relaționale.

În abordarea noastră, problema de *predicția defectelor* este considerată o problemă de clasificare binară supervizată, deoarece fiecare entitate trebuie clasificată *defectivă*, notată cu “+” sau *non-defectivă* notată cu “-”. Similar altor abordări, modelul nostru este construit folosind date de antrenament, adică un set de date care conține entități despre care se știe dacă sunt defecte sau nu.

3.3.1 Classificatorul *DPRAR*

Deoarece predicția defectelor este o problemă de clasificare supervizată, vom avea un pas de antrenament și un pas de testare. Pentru antrenament considerăm două seturi de date: DS_+ care conține entități k -dimensionale defecte și DS_- care conține entități k -dimensionale non-defective. Pentru a clasifica o entitate ca fiind defectivă sau non-defectivă, următorii pași sunt efectuați:

1. Preprocesarea datelor
2. Antrenament/Construirea clasificatorului *DPRAR*
3. Testare/Clasificare

3.3.1.1 Preprocesarea datelor

În cadrul acestui pas, datele de antrenament sunt scalate ca să aibă valori în intervalul $[0,1]$ și o analiză statistică este efectuată pe seturile de antrenament DS_+ și DS_- pentru a determina o submulțime a caracteristicilor care sunt corelate cu variabila de prezis. Pentru a determina dependențele dintre caracteristici și variabila țintă, se folosește coeficientul de corelație Spearman [Spe04].

Pentru a decide care caracteristici ar trebui eliminate, pentru fiecare caracteristică (metrică soft) $sm_i \in \mathcal{SM}$ calculăm corelația Spearman, $(cor(sm_i, target))$, între valorile caracteristicii și a variabilei țintă. Vom nota cu m media și cu $stdev$ deviația standard acestor valori de corelație. O caracteristică sm_i este eliminată dacă $abs(cor(sm_i, target)) < m - stdev$.

3.3.1.2 Antrenament/Construirea clasificatorului *DPRAR*

Prima dată definim mulțimea relațiilor dintre caracteristici care va fi folosită în procesul de extragere a regulilor de asociere relaționale. După aceea efectuăm extragerea de reguli de asociere relaționale separat pe cele două seturi de date. Astfel vom avea două seturi de reguli de asociere relaționale: RAR_+ și RAR_- . La fiecare regulă r asociem o valoare, numită *ratio*(r), calculată prin împărțirea confidenței regulii la suportul ei.

3.3.1.3 Testare/Clasificare

La pasul de clasificare, după ce antrenamentul s-a terminat și clasificatorul *DPRAR* a fost construit, o entitate nouă e trebuie clasificată. Pentru clasificare am calculat două scoruri

Nume	Entități defective	Entități non-defective
CM1	42 (12.84 %)	285 (87.16%)
KC1	314 (26.54%)	869 (73.46%)
KC3	36 (18.56 %)	158 (81.44 %)
PC1	61 (8.65 %)	644 (91.35%)
JM1	1672 (21.49 %)	6110 (78.51 %)
MC2	44 (35.2 %)	81 (64.8%)
MW1	27 (10.67%)	226 (89.33 %)
PC2	16 (2%)	729 (98%)
PC3	134 (12.4%)	943 (87.6 %)
PC4	177 (13.8%)	1110 (86.2 %)

Tabela 3.2: Seturile de date NASA folosite pentru evaluarea experimentală.

Set de date	c_{min}^+	c_{min}^-	Lungime	Acc	Pd	Spec	Prec	AUC
CM1	0.927	0.94	orice	0.8716	0.929	0.8632	0.5	0.896
KC1	0.8	0.822	2	0.823	0.818	0.825	0.628	0.822
KC3	0.885	0.96	2	0.83	0.889	0.8165	0.5246	0.85225
PC1	0.95	0.995	2	0.956	0.885	0.963	0.692	0.924
JM1	0.95	0.995	orice	0.96	0.842	0.992	0.967	0.917
MC2	0.96	0.99	orice	0.896	0.773	0.9632	0.919	0.868
MW1	0.97	0.975	orice	0.941	0.889	0.947	0.667	0.918
PC2	0.95	0.995	orice	0.984	0.938	0.985	0.577	0.962
PC3	0.95	0.995	2	0.967	0.85	0.983	0.877	0.917
PC4	0.95	0.995	2	0.961	0.814	0.985	0.894	0.899

Tabela 3.3: Rezultatele obținute pentru clasificatorul *DPRAR* pentru seturile de date considerate.

$score_+(e)$ (similaritatea lui e cu clasa defectivă) și $score_-(e)$ (similaritatea lui e cu clasa non-defectivă). Pentru a calcula $score_+(e)$ considerăm *ratio*-ul mediu al regulilor din RAR_+ care sunt verificate de e , dar și *ratio*-ul mediu al regulilor din RAR_- care nu sunt verificate de e . $Score_-(e)$ este calculat în mod similar.

Dacă $score_+ > score_-$ atunci e va fi clasificată defectivă, altfel va fi clasificată non-defectivă.

3.3.2 Evaluare Experimentală

Pentru evaluarea experimentală a abordării noastre am folosit zece seturi de date NASA [Nas]. Numele și caracteristicile acestor seturi de date sunt prezentate în Tabelul 3.2.

Pentru evaluarea performanței clasificatorului *DPRAR*, o validare încrucișată de tip *leave-one-out* a fost folosită, și următoarele măsuri de performanță au fost calculate: acuratețea (Acc), probabilitatea de detecție (Pd), *specificity* (Spec), precizie (Prec), și *AUC*. Tabelul 3.3 prezintă cele mai bune rezultate obținute de clasificatorul *DPRAR* pentru seturile de date considerate.

3.3.3 Analiză și comparație cu abordări existente

Am comparat rezultatele clasificatorului *DPRAR* cu rezultate raportate în literatura de specialitate pentru alte abordări pentru care evaluarea experimentală a fost făcută folosind aceleași seturi de date. Am ales metoda *CBA2* [BDVB11], clasificatorul *1R* [CBYP05], *Bagging* [HDF12] și *EDER – SD* [RRRAR12].

Nu toate abordările au fost testate pe toate cele 10 seturi de date pe care le-am folosit, și nu toate abordările raportează toate măsurile de performanță pe care le-am folosit. Considerând toate studiile de caz și toate măsurile de performanță, *DPRAR* a avut performanțe mai bune în 45 de cazuri, similară o dată și performanțe mai slabe în 23 de cazuri dintr-un total de 69 de comparații. În plus, măsura *AUC* (considerată în literatură una dintre cele mai bune măsuri de evaluare pentru clasificatori) pentru clasificatorul *DPRAR* depășește media valorilor *AUC* raportate de celelalte metode la toate studiile de caz. Acest lucru indică eficiența bună a clasificatorului *DPRAR*.

3.3.4 Un studiu despre clasificatorul *DPRAR*

Am efectuat un studiu asupra clasificatorului *DPRAR* pentru a analiza cum influențează eliminarea diferitor caracteristici rezultatele raportate, și care sunt efectele folosirii unei alte formule pentru a calcula acele scoruri pe baza cărora o entitate nouă este clasificată.

Primul pas al clasificatorului *DPRAR* este preprocesarea datelor, când diferite metrice soft pot fi eliminate. Pentru metoda *DPRAR* am ales să eliminăm acele caracteristici pentru care corelația Spearman cu valoarea țintă este mai mică decât diferența dintre media corelațiilor, m , și deviația standard a corelațiilor, $stdev$. În studiul nostru am comparat această metodă de eliminare a caracteristicilor cu trei alte metode (nicio caracteristică nu este eliminată, sunt eliminate caracteristicile ale căror corelație este mai mare decât $m + stdev$, sunt eliminate și caracteristicile ale căror corelație este mai mică decât $m - stdev$ și caracteristicile ale căror corelație este mai mare decât $m + stdev$). Pentru evaluarea experimentală am folosit setul de date PC3.

Măsurile de evaluare folosite au sugerat că cea mai bună opțiune este eliminarea caracteristicilor care au corelația sub $m - stdev$ sau peste $m + stdev$, iar cea mai slabă variantă este să nu fie folosită nicio metodă de eliminare.

Am efectuat încă un studiu, în legătură cu calculul scorului, folosind ca studiu de caz setul de date PC3, și folosind ca metodă de eliminare a caracteristicilor varianta cu cele mai bune rezultate din primul studiu. De data asta am folosit numărul regulilor de asociere relaționale (și nu *ratio*-ul lor) verificate și neverificate pentru a calcula scorul pe baza căruia se face clasificarea. Măsurile de performanță au sugerat că această variantă pentru calculul scorului este mai bună.

Am comparat rezultatele acestor două variante de a calcula scorul și cu metode din literatura de specialitate: *CBA2*, *ROCUS*, *Random Forests* și *Dynamic AdaBoost.NC*. Clasificatorul *DPRAR* cu calculul scorului bazat pe numărul de reguli a avut cele mai bune rezultate, urmat de *DPRAR* cu calculul scorului bazat pe *ratio*-ul regulilor.

3.4 Concluzii și cercetări ulterioare

Acest capitol a prezentat abordările noastre bazate pe extragerea regulilor de asociere relaționale pentru detectarea defectelor în sisteme informatice. Acest domeniu are două direcții principale, *detectarea defectelor de proiectare* și *predicția defectelor*, ambele fiind acoperite în cadrul cercetării prezentate în acest capitol.

În prima parte a capitolului am prezentat abordarea *SDDRAR*, care poate să identifice entitățile cu defecte de proiectare într-un sistem informatic. Abordarea a fost evaluată folosind șase studii de caz și analiza manuală a rezultatelor a arătat că *SDDRAR* e capabil să

identifice entitățile cu defecte de proiectare. Am comparat abordarea noastră cu alte două abordări din literatura de specialitate, și rezultatele au fost similare cu rezultatele raportate de una dintre instrumentele folosite. Aceste rezultate arată potențialul abordării propuse.

În partea a doua a capitolului am prezentat abordarea *DPRAR*, un model de clasificare binară, bazată pe reguli de asociere relaționale pentru predicția defectelor. Pentru evaluarea experimentală am folosit zece seturi de date *NASA*. Comparația rezultatelor cu alte rezultate raportate în literatura de specialitate arată că abordarea *DPRAR* este comparabilă sau mai bună decât metodele existente.

Pentru cercetări ulterioare vrem să extindem pe alte seturi de date sau sisteme informatice evaluarea experimentală a acestor metode. Deasemenea, vrem să extindem abordările noastre spre reguli de asociere relaționale *fuzzy*.

Capitolul 4

Un Framework pentru Analiza Sistemelor Informatice Orientate Obiect

În capitolele anterioare am prezentat tehnici bazate pe inteligență computațională pentru a rezolva diferite probleme de o importanță ridicată pe parcursul întreținerii și evoluției sistemelor informatice. Pentru dezvoltarea acestor tehnici am folosit *framework*-ul *FAOS* (*Framework for Analyzing Object-oriented Software systems*) care a fost proiectat ca să fie suficient de general pentru a oferi un suport pentru analiza unui sistem informatic orientat obiect și extragerea informațiilor relevante din el, precum și pentru a oferi implementări pentru mai multe metrice soft care pot fi folosite pentru a măsura calitatea sistemului informatic. În acest capitol prezentăm acest *framework* mai detaliat.

4.1 *Framework*-ul *FAOS*

Framework-ul *FAOS* este alcătuit din trei module, *Analyzer*, *Metrics* și *PackageRestructuring*.

Modulul *Analyzer*. Primul modul din *FAOS* este *Analyzer*, implementarea căruia poate să analizeze un sistem informatic realizat în limbajul de programare Java, și necesită ori fișierele compilate *.class* ori o arhivă *jar* pentru această analiză. Am definit un model de date care este alcătuit din clase care reprezintă atribute, metode, clase din sistemul analizat. Scopul principal al modulului *Analyzer* este să efectueze analiza sistemului informatic și să extragă o listă de clase împreună cu metodele și atributele lor, și relațiile între ele. Pentru analiza codului compilat am folosit *framework*-ul de manipulare a *bytecode*-ului *ASM* [ASM13].

Modulul *Metrics*. Al doilea modul din *FAOS* este *Metrics*, un modul care conține implementarea diferitor metrice soft. Modulul este împărțit în două pachete principale, pachetul *classLevel* conține implementarea a 17 metrice soft, calculate pentru o clasă, iar pachetul *packageLevel* conține implementarea a 20 de metrice și măsuri soft calculate pentru un pachet.

Modulul *PackageRestructuring*. Al treilea modul din *FAOS* se numește *packageRestructuring* și abordarea de restructurare la nivel de pachete prezentată în Secțiunea 2.1 este implementată aici, de aceea este un modul mai puțin abstract ca primele două.

4.2 Comparație cu *framework*-uri similare

În inginerie soft bazată pe inteligență computațională cercetătorii adeseori implementează abordarea lor sub forma unor instrumente și *framework*-uri, care pot fi folosite atât de alți

cercetători pentru comparația rezultatelor, cât și de dezvoltatori soft. Cu toate că nu există un alt *framework* cu exact aceeași funcționalități ca *FAOS*, există anumite instrumente care sunt similare cu diferite părți din el, de exemplu *JDeodorant* - un *plugin Eclipse* pentru a identifica *bad smell*-uri în sisteme informatice -, *iPlasma* - un mediu proiectat pentru analiza calității sistemelor orientate obiect - și instrumentul *Bunch* care încearcă să găsească o partiție bună a unui sistem informatic folosind algoritmi de căutare. Primii doi au propria reprezentare internă (similară cu modelul nostru de date), iar *Bunch* folosește *Module Dependency Graph*-ul sistemului de restructurat. Astfel *Bunch* este independent de orice limbaj de programare, *JDeodorant* analizează sisteme realizate în limbajul Java și *iPlasma* poate să analizeze sisteme realizate în Java sau C++. Avantajul lor față de *framework*-ul *FAOS* este că sunt instrumente complete, gata de a fi folosite. Pe de altă parte *FAOS* este mai flexibil, poate fi extins ușor pentru a efectua alte tipuri de analize, nu numai cele implementate momentan.

În concluzie, *JDeodorant* și *iPlasma* sunt potrivite pentru situații când cineva vrea să analizeze codul sursă, să caute metrici soft cu valori anormale sau *bad smell*-uri. *Framework*-ul *FAOS* este mai potrivit pentru situații când cineva vrea să analizeze un sistem și să facă ceva cu lista elementelor extrase. Astfel, avantajul principal al *framework*-ului *FAOS* este că poate fi extins ușor: noi metrici soft pot fi adăgate, lista entităților returnate de modulul *Analyzer* poate fi folosită pentru analize ulterioare. Așa s-a întâmplat și la modulul *PackageRestructure* unde această listă este folosită pentru a găsi o împărțire bună a claselor în pachete.

4.3 Concluzii și cercetări ulterioare

În acest capitol am prezentat *framework*-ul *FAOS*, un *framework* pentru analiza sistemelor informatice orientate obiect. *FAOS* a fost dezvoltat pentru a sprijini evaluarea experimentală a majorității abordărilor bazate pe instruire automată prezentate în capitolele anterioare. Astfel, modulul *Analyzer* a fost folosit pentru toate abordările cu excepția predicției de defecte, unde seturile de date *NASA* sunt folosite. Modulul *Metrics* a fost folosit pentru restructurare la nivel de pachete și clase dar și la abordarea de detectare a defectelor de proiectare. În cele din urmă, modulul *PackageRestructuring* conține implementarea abordării noastre pentru restructurarea la nivel de pachete a unui sistem informatic.

Ca cercetare ulterioară dorim să facem acest *framework* disponibil public, ca alte persoane să poată să-l folosească și să-l extindă în funcție de context. Pentru a facilita folosirea lui, ne gândim să dezvoltăm un *plugin Eclipse* bazat pe *FAOS*. Deasemenea dorim să extindem partea de analiză la alte limbaje de programare, deoarece dacă reprezentarea internă a entităților este construită, restul *framework*-ului poate fi folosit fără probleme.

Concluzii

În această teză am prezentat lucrările noastre originale pentru aplicarea metodelor și algoritmilor de instruire automată pentru rezolvarea diferitor probleme din inginerie soft. Asemenea abordări inteligente sunt necesare pentru că în zilele noastre sistemele informatice sunt destul de complexe, alcătuite din multe componente și relații între aceste componente și abordările care sugerează care sunt acele părți care necesită atenție pot fi foarte utile dezvoltatorilor soft. Prin urmare, putem spune că asemenea abordări pot fi considerate instrumente importante și sunt benefice dezvoltatorilor soft.

Am ales două probleme din ingineria soft care au reprezentat principalele direcții de cercetare în această teză. În cadrul fiecărei direcții am lucrat la două probleme separate. Prima direcție principală este folosirea algoritmilor de *clustering* pentru modularizare soft atât la nivel de pachete, cât și la nivel de clasă. A doua direcție principală a fost folosirea regulilor de asociere relaționale pentru detectarea defectelor în sisteme informatice. În cadrul acestei direcții am lucrat la problema de detectare a defectelor de proiectare și predicția defectelor în sisteme informatice. În cele din urmă am prezentat o contribuție spre dezvoltarea sistemelor informatice, *framework*-ul *FAOS*, dezvoltat pentru analiza sistemelor informatice orientate obiect.

Evaluarea experimentală a abordărilor dezvoltate pentru prima direcție de cercetare au arătat că tehnici de *clustering*, combinate cu o măsură bună a similarității/disimilarității dintre entități poate să identifice o împărțire bună a claselor în pachete, sau o împărțire bună a metodelor și atributelor în clase.

În a doua direcție principală de cercetare am folosit reguli de asociere relaționale pentru detectarea defectelor în sisteme informatice. Prima problemă în cadrul acestei direcții a fost identificarea claselor cu defecte de proiectare. Experimentele efectuate arată că acele clase care au fost identificate de abordarea noastră au avut defecte de proiectare. În ceea ce privește predicția defectelor, o problemă de clasificare, abordarea noastră a fost testată pe 10 seturi de date *NASA* [Nas]. Am efectuat o comparație cu alte abordări folosind multe seturi de date și măsuri de performanță, care a arătat că performanța clasificatorului nostru este comparabilă sau mai bună decât metodele existente, demonstrând potențialul abordării noastre.

Framework-ul *FAOS* prezentat în această teză a fost dezvoltat pentru a sprijini evaluarea experimentală a abordărilor propuse în această teză. Cu toate acestea, a fost proiectat să fie general și extensibil ca să poată fi folosit și de alții. *Framework*-ul oferă suport pentru analiza unui sistem informatic, pe parcursul căreia o reprezentare internă a sistemului analizat este construită. Aceasta poate fi folosită ulterior pentru diferite probleme, problemele implementate momentan permit calculul valorilor diferitor metrici soft și restructurarea la nivel de pachete a unui sistem informatic.

În ceea ce privește cercetările ulterioare, intenționăm să lucrăm la neajunsurile abordărilor noastre pentru a le îmbunătăți. De asemenea dorim să extindem evaluarea experimentală, să folosim alte sisteme informatice sau seturi de date. Vom considera aplicarea variantelor *fuzzy* ale abordărilor propuse, acolo unde este posibil. În cele din urmă vrem să dezvoltăm abordări noi și pentru alte probleme din ingineria soft.

Bibliografie

- [AAH12] H. H. Ammar, W. Abdelmoez, and M. S. Hamdi. Software engineering using artificial intelligence techniques: Current state and open problems. In *Proceedings of First Taibah University International Conference on Computing and Information Technology*, 2012.
- [AAM11] A. Alkhalid, M. Alshayeb, and S.A. Mahmoud. Software refactoring at the package level using clustering techniques. *IET Software*, 5(3):274–286, 2011.
- [AIS93] Rakesh Agrawal, Tomasz Imielinski, and Arun Swami. Mining association rules between sets of items in large databases. In *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data*, pages 207–216, 1993.
- [AL99] Nicolas Anquetil and Timothy Lethbridge. Experiments with clustering as a software remodularization method. In *Proceedings of 6th Working Conference on Reverse Engineering*, pages 235–255, Atlanta, USA, October 1999.
- [AS94] R. Agrawal and R. Srikant. Fast algorithm for mining association rules. In *Proceeding of the 20th VLDB Conference*, pages 487–499, 1994.
- [ASM13] ObjectWeb: Open Source Middleware, 2013. <http://asm.objectweb.org/>.
- [BDVB11] Ma Baojun, Karel Dejaeger, Jan Vanthienen, and Bart Baesens. Software defect prediction based on association rule classification. Technical report, Katholieke Universiteit Leuven, February 2011.
- [BJWD99] L.C. Briand and and al. J. W. Daly. A unified framework for coupling measurement in object-oriented systems. 25(1):91–121, 1999.
- [BK95] J. M. Bieman and B. K. Kang. Cohesion and reuse in an object-oriented system. *ACM SIGSOFT Software Engineering Notes*, 20(SI):259–262, 1995.
- [BLMO10] Gabriele Bavota, Andrea De Lucia, Andrian Marcus, and Rocco Oliveto. Software remodularization based on structural and semantic metrics. In *17th Working Conference on Reverse Engineering*, pages 195–204, 2010.
- [BMW02] Lionel C. Briand, Walcelio L. Melo, and J. Wust. Assessing the applicability of fault-proneness models across object-oriented software projects. *IEEE Trans. Softw. Eng.*, 28(7):706–720, July 2002.
- [BOL⁺10] Gabriele Bavota, Rocco Oliveto, Andrea De Lucia, Giuliano Antoniol, and Yann-Gaël Guéhéneuc. Playing with refactoring: Identifying extract class opportunities through game theory. In *Software Maintenance, 2010 IEEE International Conference*, pages 1–5, 2010.
- [CBYP05] Venkata U. B. Challagulla, Farokh B. Bastani, I-Ling Yen, and Raymond A. Paul. Empirical assessment of machine learning based software defect prediction techniques. In *Proceedings of the 10th IEEE International Workshop on Object-Oriented Real-Time Dependable Systems, WORDS '05*, pages 263–270, Washington, DC, USA, 2005. IEEE Computer Society.
- [CK91] Shyam R. Chidamber and Chris F. Kemerer. Towards a metrics suite for object-oriented design. In *Conference Proceedings on Object Oriented Programming Systems, Languages, and Applications*, pages 197–211, 1991.
- [CK94] Shyam R. Chidamber and Chris F. Kemerer. A metrics suite for object oriented design. *IEEE Transactions on Software Engineering*, 20(6):476–493, June 1994.

- [CLMM05] Yania Crespo, Carlos López, Esperanza Manso, and Raúl Marticorena. Language independent metric support towards refactoring inference. In *Proceedings of the 9th Workshop on QAOOSE*, pages 18–29, 2005.
- [CMC14a] Gabriela Czibula, Zsuzsanna Marian, and Istvan Gergely Czibula. Detecting software design defects using relational association rule mining. *Knowledge and Information Systems*, 2014. DOI: 10.1007/s10115-013-0721-z (available online at <http://link.springer.com/article/10.1007/s10115-013-0721-z>).
- [CMC14b] Gabriela Czibula, Zsuzsanna Marian, and Istvan Gergely Czibula. Software defect prediction using relational association rule mining. *Information Sciences*, 264:260–278, 2014.
- [CS06] I.G. Czibula and G. Serban. Improving systems design using a clustering approach. *International Journal of Computer Science and Network Security*, 6(12):40–49, 2006.
- [CS07] Istvan Gergely Czibula and Gabriela Serban. Hierarchical clustering for software systems restructuring. *INFOCOMP Journal of Computer Science*, 6(4):43–51, 2007.
- [CSTM06] Alina Campan, Gabriela Serban, Traian Marius Truta, and Andrian Marcus. An algorithm for the discovery of arbitrary length ordinal association rules. In *Proceedings of the 2006 International Conference on Data Mining (DMIN)*, pages 107–113, 2006.
- [DABH11] Stéphane Ducasse, Nicolas Anquetil, Usman Bhatti, and Andre Cavalcante Hora. Software metrics for package remodularization. Technical report, Institut National de Recherche en Informatique et en Automatique, 2011.
- [DbU] Commons dbutils. <http://commons.apache.org/proper/commons-dbutils/index.html>.
- [DDN00] Serge Demeyer, Stephance Ducasse, and Oscar Nierstrasz. Finding refactorings via change metrics. In *Proceedings of the ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages and Applications*, pages 166–177, 2000.
- [DG77] Persi Diaconis and R. L. Graham. Spearman’s footrule as a measure of disarray. *Journal of the Royal Statistical Society*, 39(2):262–268, 1977.
- [DK05] Scott Dick and Abraham Kandel. *Computational Intelligence in Software Quality Assurance*. Series in Machine Perception and Artificial Intelligence. World Scientific Publishing, 2005.
- [EL] El. <http://commons.apache.org/proper/commons-el/>.
- [Ema] Email. <http://commons.apache.org/proper/commons-email/>.
- [FBZ12] Francesca Arcelli Fontana, Pietro Braione, and Marco Zanoni. Automatic detection of bad smells in code: An experimental assessment. *Journal of Object Technology*, 11:5:1–38, 2012.
- [Fow] Martin Fowler. Refactoring malapropism. <http://martinfowler.com/bliki/RefactoringMalapropism.html>.
- [Fow99] Martin Fowler. *Refactoring: Improving the Design of Existing Code*. Addison-Wesley, Boston, 1999.
- [FTCS09] Marios Fokaefs, Nikolaos Tsantalis, Alexander Chatzigeorgiu, and Jorg Sander. Decomposing object-oriented class modules using an agglomerative clustering technique. In *Proceedings of International Conference on Software Maintenance*, pages 93–101, Edmonton, Canada, 2009.
- [FTP13] Ftp4j, 2013. <http://sourceforge.net/projects/ftp4j/>.
- [FTSC11] Marios Fokaefs, Nikolaos Tsantalis, Eleni Stroulia, and Alexander Chatzigeorgiou. Ideodorant: Identification and application of extract class refactorings. In *Proceedings of the 33rd International Conference on Software Engineering, ICSE*, pages 1037–1049, 2011.
- [Gam] E. Gamma. JHotDraw Project. <http://sourceforge.net/projects/jhotdraw>.
- [GMCS04] Lan Guo, Yan Ma, Bojan Cukic, and Harshinder Singh. Robust prediction of fault-proneness by random forests. In *ISSRE*, pages 417–428, 2004.

- [GWI11] Isaac Griffith, Scott Wahl, and Clemente Izurieta. Truerefactor: An automated refactoring tool to improve legacy system and application comprehensibility. In *Proceedings of the ISCA 24th International Conference on Computer Applications in Industry and Engineering*, pages 316–321, 2011.
- [Han05] Jiawei Han. *Data Mining: Concepts and Techniques*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2005.
- [Har12] Mark Harman. The role of artificial intelligence in software engineering. In *1st International Workshop on Realizing Artificial Intelligence Synergies in Software Engineering*, 2012.
- [HDF12] A.A. Shahrjooi Haghighi, M. Abbasi Dezfuli, and S.M. Fakhrahmad. Applying mining schemes to software fault prediction: A proposed approach aimed at test cost reduction. In *Proceedings of the World Congress on Engineering 2012 Vol I, WCE 2012*, pages 1–5, Washington, DC, USA, 2012. IEEE Computer Society.
- [HJ01] Mark Harman and Bryan F. Jones. Search-based software engineering. *Information and Software Technology*, 43(14):833–839, 2001.
- [HK81] S. Henry and D. Kafura. Software structure metrics based on information flow. *IEEE Transactions on Software Engineering*, 7(5):510–518, September 1981.
- [HKI08] Yoshiki Higo, Shinji Kusumoto, and Katsuro Inoue. A metric-based approach to identifying refactoring opportunities for merging code clones in a java software system. *Journal of Software Maintenance and Evolution: Research and Practice*, 20:435 – 461, 2008.
- [HM95] M. Hitz and B. Montazeri. Measuring coupling and cohesion in object-oriented systems. In *Proceedings of International Symposium on Applied Corporate Computing*, Monterrey, Mexico, October 1995.
- [HMZ09] Mark Harman, S. Afshin Mansouri, and Yuanyuan Zhang. Search based software engineering: A comprehensive analysis and review of trend techniques and applications. Technical report, Department of Computer Science, King’s College, London, 2009.
- [HS96] B. Henderson-Sellers. *Object-Oriented Metrics Measures of Complexity*. Prentice-Hall, 1996.
- [HT07] Mark Harman and Laurence Tratt. Pareto optimal search based refactoring at the design level. In *Conference on Genetic and Evolutionary Computation*, pages 1106–1113, 2007.
- [IP113] iplasma, 2013. <http://loose.upt.ro/reengineering/research/iplasma>.
- [ISIE12] Safwat M. Ibrahim, Sameh A. Salem, Manal A. Ismail, and Mohamed Eladawy. Identification of nominated classes for software refactoring using object-oriented cohesion metrics. *International Journal of Computer Science Issues*, 9(2):68–76, 2012.
- [ISO13] Iso8583, 2013. <http://sourceforge.net/projects/j8583/>.
- [JDe13] Jdeodorant, 2013. <http://www.jdeodorant.com/>.
- [JLZ11] Yuan Jiang, Ming Li, and Zhi-Hua Zhou. Software defect detection with rocus. *Journal of Computer Science and Technology*, 26(2):328–342, 2011.
- [Kan03] Ronald Kirk Kandt. A software defect detection methodology, 2003.
- [KMMiM08] Yasutaka Kamei, Akito Monden, Shuji Morisaki, and Ken ichi Matsumoto. A hybrid faulty module prediction using association rule mining and logistic regression analysis. In *Proceedings of the International Symposium on Empirical Software Engineering and Measurements (ESEM)*, pages 279–281, 2008.
- [KSBW11] Marouane Kessentini, Houari Sahraoui, Mounir Boukadoum, and Manuel Wimmer. Search-based design defects detection by example. In *Proceedings of the 14th International Conference on Fundamental Approaches to Software Engineering*, pages 401–415, 2011.
- [KVG09] Foutse Khomh, Stéphane Vaucher, Yann-Gaël Guéhéneuc, and Houari Sahraoui. A bayesian approach for the detection of code and design smells. In *Proceedings of the 9th International Conference on Quality Software*, pages 305–314, 2009.

- [Lar04] Craig Larman. *Applying UML and Patterns: An Introduction to Object-Oriented Analysis and Design and Iterative Development*. Addison Wesley Professional, third edition, 2004.
- [LH93] Wei Li and Sallie Henry. Object oriented metrics which predict maintainability. *Journal of Systems and Software*, 23(2):111–122, 1993.
- [LHM98] Bing Liu, Wynne Hsu, and Yiming Ma. Integrating classification and association rule mining. In *Proceedings of the 4th International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 80–86, 1998.
- [LLWW95] Y. S. Lee, B. S. Liang, S. F. Wu, and F. J. Wang. Measuring the coupling and cohesion of an object-oriented program based on information flow. In *Proceedings of International Conference on Software Quality*, Maribor, Slovenia, 1995.
- [LMW01] Bing Liu, Yiming Ma, and Ching-Kian Wong. *Data Mining for Scientific and Engineering Applications*, chapter Classification Using Association Rules: Weaknesses and Enhancements. Kluwer Academic, 2001.
- [LZWZ12] Ming Li, Hongyu Zhang, Rongxin Wu, and Zhi-Hua Zhou. Sample-based software defect prediction with active and semi-supervised learning. *Automated Software Engineering*, 19(2):201–230, 2012.
- [Mai09] Sayyed Garba Maisikeli. *Aspect Mining Using Self-Organizing Maps With Method Level Dynamic Software Metrics as Input Vectors*. PhD thesis, Graduate School of Computer and Information Sciences Nova Southeastern University, 2009.
- [Mar02] Radu Marinescu. *Measurement and Quality in Object-Oriented Design*. PhD thesis, Politehnica University Timisoara, Faculty of Automatics and Computer Science, 2002.
- [Mar10] Zsuzsanna Marian. Solving the subset sum problem with dna computation. In *Proceedings of the National Symposium ZAC*, pages 25–29, 2010.
- [Mar12a] Zsuzsanna Marian. Aggregated metrics guided software refactoring. In *Proceedings of the 8th IEEE International Conference on Intelligent Computer Communication and Processing*, pages 259–266, 2012.
- [Mar12b] Zsuzsanna Marian. Software metrics based refactoring: a case study. In *Proceedings of the National Symposium ZAC*, pages 59–64, 2012.
- [Mar12c] Zsuzsanna Marian. A study on hierarchical clustering based software restructuring. *Studia Universitatis “Babes-Bolyai” Informatica*, LVII(2):20–31, 2012.
- [Mar13a] Zsuzsanna Marian. On the software metrics influence in relational association rule-based software defect prediction. *Studia Universitatis “Babes-Bolyai” Informatica*, LVIII(4):35–48, 2013.
- [Mar13b] Zsuzsanna Marian. A study on association rule mining based software defect detection. *Studia Universitatis “Babes-Bolyai” Informatica*, LVIII(1):42–57, 2013.
- [Mar14a] Zsuzsanna Marian. Faos - a framework for analyzing object-oriented software systems. *Studia Universitatis “Babes-Bolyai” Informatica*, 2014. Under review.
- [Mar14b] Zsuzsanna Marian. On evaluating the structure of software packages. *Studia Universitatis “Babes-Bolyai” Informatica*, LIX(1):46–58, 2014.
- [MC11] Iman Hemati Moghadam and Mel Ó. Cinnéide. Code-imp: A tool for automated search-based refactoring. In *Proceeding of the 4th Workshop on Refactoring Tools*, pages 41–44, Honolulu, USA, 2011.
- [MCB11] Zsuzsanna Marian, Cosmin Coman, and Attila Bartha. Learning to play the guessing game. *Studia Universitatis “Babes-Bolyai” Informatica*, LVI(2):119–124, 2011.
- [MCC12] Zsuzsanna Marian, Gabriela Czibula, and Istvan-Gergely Czibula. Using software metrics for automatic software design improvement. *SIC Journal, Studies in Informatics and Control*, 21(3):249–258, 2012.

- [MCC14] Zsuzsanna Marian, Gabriela Czibula, and Istvan Gergely Czibula. Software packages refactoring using a hierarchical clustering-based approach. *Fundamenta Informaticae*, 2014. Under review.
- [MGF07] Tim Menzies, Jeremy Greenwald, and Art Frank. Data mining static code attributes to learn defect predictors. *IEEE Transactions on Software Engineering*, 33(1):2–13, 2007.
- [MGL06] Naouel Moha, Yann-Gaël Guéhéneuc, and Pierre Leduc. Automatic generation of detection algorithms for design defects. In *Proceedings of the 21st IEEE/ACM International Conference on Automated Software Engineering*, pages 297–300, 2006.
- [MHH03] Kiarash Mahdavi, Mark Harman, and Robert Mark Hierons. A multiple hill climbing approach to software module clustering. In *Proceedings of the International Conference on Software Maintenance*, pages 315–324, 2003.
- [Mit06] Tom M. Mitchell. The discipline of machine learning. Working paper, 2006.
- [MMCG99] S. Mancoridis, B. S. Mitchell, Y. Chen, and E. R. Gansner. Bunch: A clustering tool for the recovery and maintenance of software system structures. In *In Proceedings of the IEEE International Conference on Software Maintenance*, pages 50–59, 1999.
- [MML01] Andrian Marcus, Jonathan I. Maletic, and K. I. Lin. Ordinal association rules for error identification in data sets. In *Proceedings of the 10th International Conference on Information and Knowledge Management*, pages 589–591, 2001.
- [Moh06] Naouel Moha. Detection and correction of design defects in object-oriented architectures. In *Doctoral Symposium, 20th edition of the European Conference on Object-Oriented Programming*, 2006.
- [MS10] Zsuzsanna Marian and Christian Săcărea. Using contextual topology to discover similarities in modern music. In *Proceedings of the IEEE International Conference on Automation Quality and Testing, Robotics*, volume 3, pages 1–6, 2010.
- [Mun05] Matthew James Munro. Product metrics for automatic identification of “bad smell” design problems in java source code. In *Proceedings of the 11th IEEE International Software Metrics Symposium*, 2005.
- [Nas] Nasa software defect datasets.
- [OC08] Mark O’Keeffe and Mel Ó. Cinnéide. Search-based refactoring for software maintenance. *The Journal of Systems and Software*, 81:502–516, 2008.
- [PJJ13] Wei-Feng Pan, Bo Jiang, and Bing Li. Refactoring software packages via community detection in complex software networks. *International Journal of Automation and Computing*, 10(2):157–166, 2013.
- [Pro13] Profiler4j, 2013. <http://sourceforge.net/projects/profiler4j/>.
- [rlf] Reinforcement learning framework. <http://www.cs.ubbcluj.ro/~gabis/rl>.
- [RR11] Akepogu Ananda Rao and Kalam Narendar Reddy. Identifying clusters of concepts in a low cohesive class for extract class refactoring using metrics supplemented agglomerative clustering technique. *International Journal of Computer Science Issues*, 8(2):185–194, 2011.
- [RRRAR12] D. Rodríguez, R. Ruiz, J. C. Riquelme, and J. S. Aguilar-Ruiz. Searching for rules to detect defective modules: A subgroup discovery approach. *Information Sciences*, 191:14–30, May 2012.
- [SCC06] Gabriela Serban, Alina Câmpăan, and Istvan Gergely Czibula. A programming interface for finding relational association rules. *International Journal of Computers, Communications & Control*, I(S.):439–444, June 2006.
- [SKR08] Santonu Sarkar, Avinash C. Kak, and Girish Maskeri Rama. Metrics for measuring the quality of modularization of large-scale object-oriented software. *IEEE Transactions on Software Engineering*, 34(5):700–720, 2008.

- [SMV09] Adrian Sterca, Zsuzsanna Marian, and Alexandru Vancea. Distortion-based media-friendly congestion control. In *Proceedings of the International Conference on Knowledge Engineering, Principles and Techniques*, pages 265–267, 2009.
- [Spe04] C. Spearman. The proof and measurement of association between two things. *Amer. J. Psychol.***15**, pages 72–101, 1904.
- [SS07] Konstantinos Stroggylos and Diomidis Spinellis. Refactoring - does it improve software quality? In *Proceedings of the 5th International Workshop on Software Quality*, 2007.
- [SSB06] Olaf Seng, Johannes Stammel, and David Burkhart. Search-based determination of refactorings for improving the class structure of object-oriented systems. In *Conference on Genetic and Evolutionary Computation*, pages 1909 – 1916, 2006.
- [SSL01] Frank Simon, Frank Steinbrückner, and Claus Lewerentz. Metrics based refactoring. In *Proceedings of the 5th European Conference on Software Maintenance and Reengineering*, pages 30–38, 2001.
- [ST98] Gregor Snelting and Frank Tip. Reengineering class hierarchies using concept analysis. In *Proceedings of the 6th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 99–110, 1998.
- [ȚIM09] Radu Țurcaș, Oana Iova, and Zsuzsanna Marian. The autonomous robotic tank (art): an innovative lego mindstorm nxt battle vehicle. In *Proceedings of the International Conference on Knowledge Engineering, Principles and Techniques*, pages 95–98, 2009.
- [TK03] Ladan Tahvildari and Kostas Kontogiannis. A metric-based approach to enhance the design quality through meta-pattern transformations. In *Proceedings of the European Conference on Software Maintenance and Reengineering*, pages 183–192, 2003.
- [TLM11] Doina Tatar, Mihaiela Lupea, and Zsuzsanna Marian. Text summarization by formal concept analysis approach. *Studia Universitatis “Babes-Bolyai” Informatica*, LVI(2):7–12, 2011.
- [Tuf11] Stéphane Tuffry. *Data Mining and Statistics for Decision Making*. John Wiley and Sons, 2011.
- [Win13] Winrun4j, 2013. <http://sourceforge.net/projects/winrun4j/>.
- [ZT05] Du Zhang and Jeffrey J. P. Tsai. *Machine Learning Applications in Software Engineering*. Series on Software Engineering and Knowledge Engineering. World Scientific Publishing, 2005.