

Universitatea Babeş-Bolyai Cluj-Napoca
Facultatea de Matematică și Informatică

Metode evolutive pentru rezolvarea problemelor de optimizare

Candidat:

Flavia Zamfirache

Coordonator:

Prof. Dr. Militon Frențiu

Rezumat

2011

Abstract

Metodele evolutive sunt utilizate cu succes pentru a rezolva probleme de optimizare statică, a căror scop este de a găsi soluția optimă globală. În cazul problemelor de optimizare în medii dinamice și incerte scopul algoritmului este de a urmări modificările din mediu sau de a găsească un comportament robust. În cazul dinamic, algoritmi pot întâmpina probleme din cauza convergenței premature a metodelor evolutive. Diferite mecanisme sunt propuse pentru a rezolva problemele care apar în mediile dinamice și incerte, metode care permit menținerea diversității populației în timpul rulării, utilizarea arhivelor de memorie sau multi-populații.

Multe probleme din lumea reală sunt dinamice, proprietățile lor se schimbă/variază în timp. Această teză se bazează pe trei probleme reale: planificarea taskurilor în grid, gruparea datelor și extragerea de informații din date, analizând în același timp comportamentul lor în medii incerte și dinamice. Incertitudinea în problemele de optimizare poate fi cauzată de diferite surse: funcția obiectiv poate fi afectată de zgomot, funcția obiectiv poate fi aproximată sau schimbată de-a lungul timpului.

În această teză am propus și analizat comportamentul unor algoritmi inspirați din natură în medii incerte și dinamice. Am propus un mecanism simplu de perturbare a poziției unei particule pentru euristica Particle Swarm Optimization și introducerea de elemente aleatoare în populație în cazul euristicii Differential Evolution pentru a încetini convergența euristiciilor și am testat metodele propuse pe două probleme de benchmark. În cazul problemei de planificare a datelor pe grid, am propus și studiat comportamentul unui planificator bazat pe algoritmi evolutivi în cazul planificării online. De asemenea, pentru planificatorul evolutiv și pentru planificatorul bazat pe euristica Ant Colony Optimization am studiat robustețea și am analizat comportamentul într-un mediu dinamic. În cazul problemei de grupare a datelor, am studiat influența zgomotului și valorilor lipsă din date pentru euristica AntClust și am propus un

mecanism de calcul al densității care îmbunătățește rezultatele grupării în cazul datelor afectate de zgomot. În cazul problemei de extragere a informațiilor din date am propus un algoritm evolutiv multi-obiectiv care încorporează evaluarea regulilor de către experții umani și am analizat influența valorilor lipsă în date și impactul intervenției utilizatorului în procesul de căutare. De asemenea, am propus utilizarea unui clasificator bazat pe o selecție evolutivă de reguli generate folosind tehnica prototipurilor neimbricate ca alternativă la algoritmi de interschimbare din cadrul unui planificator.

Cuprins

1	Introducere	3
1.1	Cuvinte cheie	3
1.2	Motivație	3
1.3	Structura tezei	4
1.4	Contribuții	5
2	Probleme de optimizare în medii dinamice și incerte	7
2.1	Probleme de optimizare în medii incerte	7
2.2	Probleme de optimizare în medii dinamice	8
2.2.1	Probleme de optimizare continue	8
2.2.2	Probleme de optimizare combinatoriale	8
2.3	Gruparea datelor - problemă de optimizare	8
2.3.1	Abordări ale problemelor de grupare a datelor	9
2.3.2	Gruparea datelor afectate de zgomot	9
2.4	Găsirea de reguli - problemă de optimizare	9
2.5	Concluzii	10
3	Metaeuristici inspirate din natură utilizate în medii dinamice	11
3.1	Metaeuristici inspirate din natură	11
3.1.1	Algoritmi evolutivi	11
3.1.2	Ant Colony Optimization	12
3.1.3	Particle Swarm Optimization	12
3.2	Mecanisme de adaptare la medii dinamice	12
3.2.1	Algoritmi evolutivi	13
3.2.2	Ant Colony Optimization	13
3.2.3	Particle Swarm Optimization	13
3.3	Experimente numerice	14
3.3.1	Differential Evolution	14
3.3.2	Mecanisme de adaptare la mediul dinamic	14
3.3.3	Rezultate numerice	16
3.4	Concluzii	17
4	Planificarea taskurilor în medii dinamice distribuite	19
4.1	Problema planificării pe grid	19
4.2	Abordări într-un mediu static	20
4.2.1	Simple Population-based Scheduler	20
4.2.2	Ant Colony Optimization	21

4.2.3	Rezultate numerice	22
4.3	Analiza robusteții euristiciilor	23
4.4	Planificare Online	23
4.5	Adaptarea în mediul dinamic	24
4.5.1	Simple Population-based Scheduler	24
4.5.2	Ant Colony Optimization	25
4.5.3	Experimente în mediul dinamic	25
4.6	Concluzii	26
5	Algoritmi de grupare a datelor inspirați din natură	27
5.1	Coloniile de furnici și gruparea datelor	27
5.2	Algoritmul AntClust	27
5.3	Tratarea zgomotelor în tehnicile de grupare inspirate de comportamentul furnicilor	28
5.3.1	Adăugarea informației de densitate la algoritmul AntClust	29
5.4	AntClust aplicat în cazul datelor medicale	31
5.4.1	Măsuri de similaritate/disimilaritate potrivite în cazul datelor medicale	31
5.5	Concluzii	32
6	Abordări bazate pe algoritmi evolutivi pentru extragerea de reguli	33
6.1	Extragerea de reguli din date	33
6.1.1	Măsuri de evaluare	34
6.1.2	Valori lipsă în date	34
6.2	Un algoritm evolutiv pentru extragerea de reguli din date	35
6.3	Interacțiunea utilizatorului cu procesul de extragere de informații	36
6.4	O abordare evolutivă de reducere a regulilor extrase din date	37
6.5	Experimente numerice	37
6.5.1	Setul de date medicale	38
6.5.2	Setul de date utilizat pentru planificarea taskurilor	41
6.6	Concluzii	42
7	Concluzii și direcții viitoare	43
	Bibliografie	47

Capitolul 1

Introducere

1.1 Cuvinte cheie

Algoritmi evolutivi, Particle Swarm Optimization, Ant Colony Optimization, Planificare pe grid, Determinarea de reguli pentru date medicale, Determinarea de reguli pentru aplicarea algoritmilor de planificare

1.2 Motivație

Multe probleme întâlnite în lumea reală sunt supuse schimbărilor. Aceste schimbări pot fi determinate de simulări, erori de măsură sau aproximare, ceea ce duce la apariția zgomotelor în funcția obiectiv. De asemenea, variabilele care caracterizează problema sau condițiile din mediu pot fi perturbate sau schimbate de-a lungul timpului. De exemplu, în cazul unei probleme de planificare a unor aplicații, noi aplicații pot să apară altele pot să dispară, mașinile pot funcționa mai încet, materialele își pot schimba calitatea, blocaje de trafic pot să apară în cazul problemelor de rutare, etc. În aceste cazuri, când rezolvăm probleme de optimizare, va trebui să luăm în considerare faptul că se execută în medii incerte și dinamice. În acest caz, obiectivul problemei nu este de a găsi un optim global, ci de a urmări modificările în mediul dinamic sau de a găsi soluții robuste când mediul este incert.

Metodele evolutive au fost folosite cu succes pentru a rezolva probleme de optimizare. Ele sunt o alternativă la algoritmi determinați, care găsesc soluții suboptimale ale problemelor. Printre avantajele lor se numără robustețea raportată la dimensiunea problemei (o bună alternativă în cazul problemelor mari), la instanțele problemei și variabilele aleatoare. Calitatea

algoritmilor evolutivi depinde de balansul dintre proprietățile de exploatare și explorare; exploatarea permite îmbunătățirea soluției în timp ce explorarea încurajează căutarea în regiuni noi ale spațiului de căutare. Există diferite tehnici de căutare ce aparțin metodelor evolutive: algoritmi genetici, strategiile evolutive, Particle Swarm Optimization, Ant Colony Optimization, sistemele imune, etc. Metodele evolutive s-au dovedit a fi potrivite pentru a rezolva probleme de optimizare în medii incerte și dinamice, deoarece permit adăugarea de mecanisme ce duc la obținerea de rezultate bune fără a fi nevoie de restartarea algoritmului.

În această teză se analizează abilitatea unor metode evolutive de a rezolva diferite probleme de optimizare în medii incerte și dinamice. Această problemă este una de actualitate deoarece multe probleme care se regăsesc în practică sunt dinamice și diferiți factori pot să influențeze calitatea soluției, fezabilitatea soluției depinzând de criterii multiple.

1.3 Structura tezei

Teza este structurată în următoarele capitole:

Capitolul 2 realizează o prezentare generală a conceptelor dezbătute în această teză: probleme de optimizare, medii incerte și dinamice. Acest capitol conține de asemenea o motivare a faptului că problemele de grupare și găsimă de reguli pot fi formulate ca probleme de optimizare.

Capitolul 3 conține o prezentare generală a euristiciilor folosite în această teză: algoritmi evolutivi, Ant Colony Optimization și Particle Swarm Optimization. De asemenea, sunt descrise câteva mecanisme ce permit adaptarea euristiciilor prezentate anterior în medii dinamice. Acest capitol conține și un exemplu de aplicare a două euristici: Differential Evolution și Particle Swarm Optimization într-un mediu dinamic pentru două probleme de optimizare.

Capitolul 4 dezbate problema unui planificator pe grid în medii statice și dinamice. Capitolul conține o trecere în revistă a metodelor evolutive aplicate problemei de planificare în grid, formalizarea problemei și descrierea particularităților a două euristici inspirate din natură: algoritmi evolutivi și Ant Colony Optimization, necesare pentru a mapa euristicele pe problema de planificare. Rezultatele numerice prezentate în acest capitol se referă la robustețea și comportamentul celor două euristici în medii statice și dinamice. De asemenea, este analizat comportamentul euristiciilor și în cazul online.

Capitolul 5 descrie o problemă de grupare a datelor afectate de zgomot și valori lipsă, aceasta fiind abordată folosind un algoritm inspirat din comportamentul furnicilor AntClust. În Secțiunea 5.1 se regăsește o trecere în revistă a algoritmilor de grupare inspirați din compor-

tamentul furnicilor. De asemenea, este prezentat algoritmul **AntClust**, câteva recomandări în setarea parametrilor și variante ale acestuia. Experimentele numerice prezintă comportamentul algoritmul **AntClust** în cazul: datelor afectate de zgomot, în acest caz propunându-se un mecanism de îmbunătățire a performanței; și în cazul datelor medicale incomplete.

Capitolul 6 prezintă o abordare evolutivă interactivă, multiobiectiv pentru problema găsirii de reguli în datele medicale. Capitolul conține o prezentare a problemei de găsim a regulilor, măsuri utilizate pentru evaluarea regulilor, câteva metode de corectare a valorilor lipsă din date și un review al abordărilor evolutive utilizate în găsim a reguli. Sunt prezentate de asemenea rezultate numerice pentru găsim a regulilor de asociere pentru date de test (UCI Machine Learning Repository) și date reale, colectate de la un spital de Obstetrică-Ginecologie pe parcursul unui an. S-a analizat influența diferitelor tehnici de tratare a valorilor lipsă din date asupra regulilor. Acest capitol prezintă un clasificator hibrid bazat pe determinarea regulilor folosind prototipuri generalizate neimbricate și o selecție a atributelor și prototipurilor care utilizează un algoritm evolutiv. Experimentele numerice s-au bazat pe determinarea acurateții algoritmului față de alte tehnici de clasificare și interpretarea rezultatelor pentru a determina caracteristicile taskurilor care determină selecția unui anumit algoritm de planificare sau a altuia.

Capitolul 7 prezintă concluziile și câteva direcții pentru studiile viitoare.

1.4 Contribuții

Această teză tratează problema aplicării metodelor evolutive pentru rezolvarea problemelor de optimizare. Teza se bazează pe aplicarea următoarelor metode evolutive: algoritmi evolutivi, ant colony optimization and particle swarm optimization în medii incerte. Deoarece incertitudinea în probleme poate apărea din diferite motive am analizat următoarele situații: influența zgomotului în procesul de grupare a datelor, robustețea în cazul algoritmilor de planificare, comportamentul unor planificatoare evolutive într-un mediu dinamic și descoperirea de reguli de asociere într-un mediu interactiv.

Contribuțiile principale ale acestei teze legate de aplicarea algoritmilor inspirației din natură pentru probleme de optimizare în medii incerte sunt:

- pentru a încetini convergența euristicii Differential Evolution într-un mediu dinamic am propus utilizarea de elemente aleatoare (Capitolul 3 pagina 15);

- pentru a încetini convergența euristicii Particle Swarm Optimization într-un mediu dinamic am propus utilizarea unei perturbații care afectează poziția particulei (Capitolul 3 pagina 15);
- s-a propus un algoritm evolutiv de planificare bazat pe populații și s-a analizat comportamentul lui în cazul planificării online (Capitolul 4 pagina 20) și influența strategiilor de perturbare a planificării (Capitolul 4 pagina 22) și structura populației inițiale asupra calității soluției.
- pentru un algoritm evolutiv și un algoritm bazat pe modelul coloniilor de furnici am analizat comportamentul într-un mediu static și dinamic (Capitolul 4 pagina 22) și robustețea euristiciilor (Capitolul 4 pagina 23);
- am analizat abilitatea euristicii **AntClust** de a grupa date afectate de zgomot, un parametru bazat pe densitate a fost propus pentru a îmbunătăți rezultatele grupării (Capitolul 5 pagina 29);
- am analizat abilitatea euristicii **AntClust** de a grupa date medicale cu valori lipsă (Capitolul 5 pagina 31);
- s-a propus un algoritm evolutiv multiobiectiv interactiv de găsim de reguli în date medicale, care pe lângă selecția automată implică și utilizatorul în selecția datelor (Capitolul 6 pagina 38). Particularitatea abordării noastre constă în introducerea unei distanțe de crowding între reguli, care se comportă ca un criteriu de menținere a diversității, astfel încât mulțimea Pareto este stimulată să aleagă elemente din zone mai puțin dense pentru a le adăuga în arhivă (Capitolul 6 pagina 36);
- s-a analizat influența câtorva mecanisme de gestionare a datelor cu valori lipsă în cazul găsimii de reguli (Capitolul 6 pagina 40);
- s-a propus un mecanism de selecție evolutiv pentru clasificarea regulilor generate pe baza euristicii NNGE și s-a aplicat pentru selectarea automată a algoritmului de planificare pentru un set de taskuri (Capitolul 6 pagina 37).

Capitolul 2

Probleme de optimizare în medii dinamice și incerte

Obiectivul unei probleme într-un mediu dinamic este de a urmări modificările optimului când au loc modificări în mediu și de a găsi soluții robuste. Obiectivul unei probleme într-un mediu incert este de a găsi cea mai bună soluție globală care va fi în continuare cea mai bună soluție și atunci a când mediul se schimbă.

2.1 Probleme de optimizare în medii incerte

Există diferite motive care pot induce incertitudine în mediu:

- *Zgomotul*. Funcția de evaluare este afectată de zgomot;
- *Perturbările*. Variabilele problemei pot fi perturbate sau se pot schimba după ce s-a găsit soluția optimă;
- *Aproximarea funcției obiectiv*. Evaluarea funcției obiectiv este costisitoare sau nu există o expresie analitică a acesteia; în aceste cazuri, funcția este aproximată pe baza unor date generate prin experimente sau simulări.
- *Funcția obiectiv se modifică în timp*. Funcția obiectiv se poate determina în orice moment de timp, dar depinde de acesta.

2.2 Probleme de optimizare în medii dinamice

Datorită convergenței premature a algoritmilor care rezolvă probleme de optimizare, identificarea optinelor care se schimbă este o problemă pentru algoritmi. Astfel, algoritmi trebuie să găsească un balans între exploatarea și explorarea spațiului de căutare. Dacă proprietatea de exploatare este prea mult încurajată, există posibilitatea blocării algoritmului într-un punct de minim local și transformarea căutării într-o căutare locală. Dacă proprietatea de explorare este prea mult încurajată, există posibilitatea de a transforma căutarea într-una aleatoare. Particularitățile unei probleme dinamice sunt determinate de: reprezentarea problemei, caracteristicile mediului și măsurile de calitate utilizate pentru a evalua performanța algoritmului.

2.2.1 Probleme de optimizare continue

Problemele de optimizare continuă se ocupă de funcții definite pe domenii continue, care trebuie optimizate (maximizate sau minimizate). În acest caz, optimul poate urma o traiectorie *continuă* sau *discontinuu* în spațiul de căutare.

2.2.2 Probleme de optimizare combinatoriale

Problemele de optimizare combinatorială (alocarea resurselor, planificare, etc) se ocupă de o alocare eficientă a unui set finit de resurse pentru a realiza un obiectiv.

2.3 Gruparea datelor - problemă de optimizare

Gruparea datelor este o tehnică a cărei scop este de a identifica grupuri de date similare, fără a avea informații despre aceste grupuri, prin măsurarea similarității (sau disimilarității) dintre date. Există o multitudine de metode de grupare care diferă prin: (i) informațiile necesare despre problemă (ex. numărul de clustere); (ii) modul de identificare și construire a grupurilor; (iii) modul în care tratează datele (numerice, categorice sau ambele); (iv) măsuri de similaritate utilizate. Referitor la modul în care algoritmi de grupare interpretează datele, există două abordări principale: (i) grupurile sunt mulțimi compacte, distante de date (această abordare poate fi văzută ca o *problemă de optimizare globală*); (ii) grupurile corespund unor regiuni dense de date separate prin regiuni mai puțin dense de date (această abordare poate fi văzută ca o

problemă de optimizare locală). Algoritmii de grupare pot fi de asemenea clasificați în funcție de numărul de criterii de grupare ca tehnici cu un singur obiectiv sau cu obiective multiple.

2.3.1 Abordări ale problemelor de grupare a datelor

Problemele de grupare a datelor pot fi clasificate în:

- Probleme de optimizare globală
- Probleme de optimizare locală
- Probleme de optimizare multi-obiectiv

În [ZZN⁺07] am analizat influența diferitelor criterii de calitate asupra grupării documentelor. Analiza comparativă se bazează pe un algoritm de grupare inspirat de algoritmul Differential Evolution care permite atât estimarea numărului de clustere cât și a reprezentanților. Rezultatele ilustrează particularitățile diferitelor criterii de grupare și abilitatea abordării propuse de a identifica atât numărul de clustere cât și reprezentanții clusterelor, sugerând că cel mai bun comportament este obținut prin combinarea măsurilor de conectivitate și separabilitate.

2.3.2 Gruparea datelor afectate de zgomot

Zgomotul în date poate fi introdus de erori aleatoare (codare necunoscută, valori care nu se încadrează în domeniul de definiție, date inconsistente) sau variații în măsurarea variabilelor. Pentru înlăturarea zgomotului, se pot folosi diferite metode precum: verificarea consistenței, informații despre domeniu, metode statistice, etc. În cazul metodelor de grupare, pentru eliminarea zgomotelor se pot folosi măsuri bazate pe densitatea datelor. Printre algoritmii de grupare care se bazează pe clasificare putem aminti: DBSCAN, DENCLUE și UNC (Unsupervised Niche Clustering). În [ZZ05b] am propus un mecanism bazat pe densitate care permite îmbunătățirea performanțelor algoritmului `AntClust` (capitolul 5).

2.4 Găsirea de reguli - problemă de optimizare

Extragerea de informații din date (data mining) este un proces automat sau semi-automat de găsire de informații necunoscute, potențial utile și de șabloane din bazele de date. Una din

problemele care se regăsește în extragerea de informații din date este găsirea de reguli de asociere. O regulă de asociere este o regulă de forma $(A \Rightarrow C)$ unde A (antecedent) și C (consecvent) sunt mulțimi disjuncte. Cum regulile sunt combinații de atribute, regulile de asociere pot fi văzute ca probleme de optimizare combinatorială. În acest caz, scopul problemei de optimizare este de a optimiza diferite măsuri de calitate ale regulilor (ex. suport, încredere). În capitolul 6 este prezentat în detaliu un algoritm de extragere de reguli de asociere multi-obiectiv.

2.5 Concluzii

În acest capitol am introdus contextul în care această teză este dezvoltată: probleme de optimizare, medii incerte, gruparea datelor și extragerea de reguli din date. Am început cu o prezentare a proprietăților care pot induce incertitudine în mediu: evaluarea funcției obiectiv poate conține zgomote, variabilele problemei pot fi perturbate, funcția obiectiv este aproximată sau funcția obiectiv se schimbă de-a lungul timpului. Au fost detaliate apoi proprietățile unui mediu dinamic, măsuri de evaluare, tipuri de probleme și benchmark utilizate pentru testare. Am prezentat apoi problemele de grupare a datelor și de extragere de informații din date ca probleme de optimizare.

Capitolul 3

Metaeuristici inspirate din natură utilizate în medii dinamice

Problemele de optimizare în medii incerte și dinamice sunt complexe și dificile și de multe ori algoritmi clasici bazați pe programarea dinamică sau abordările matematice reușesc să rezolve doar instanțe mici ale problemelor. Din acest motiv, în ultimi ani, metaeuristici ca Ant Colony Optimization, Evolutionary Computation, Simulated Annealing, Tabu Search și altele, s-au dovedit a fi alternative fezabile la abordările clasice.

3.1 Metaeuristici inspirate din natură

În această teză am folosit trei metaeuristici inspirate din natură: algoritmi evolutivi, Ant Colony Optimization și Particle Swarm Optimization. Am ales aceste metaeuristici datorită abilității lor de a se adapta la modificările din mediu, fapt care le face potrivite pentru experimentele numerice descrise în următoarele capitole.

3.1.1 Algoritmi evolutivi

Algoritmii evolutivi (EA) sunt tehnici de căutare euristice care se bazează pe ideea de evoluție naturală (selecție, supraviețuirea celui mai bun individ). Procesul de învățare colectivă a populației, auto-adaptarea și robustețea sunt câteva din avantajele algoritmilor evolutivi asupra tehnicilor globale de optimizare. Cei mai cunoscuți algoritmi evolutivi sunt: *algoritmii genetici* (Holland 1975), *programarea evolutivă* (Fogel 1966), *strategiile evolutive* (Rechenberg 1973, Schwefel 1981), *programarea genetică* (Koza 1999). Aplicații ale algoritmilor evolutivi

se pot găsi în probleme ca: gruparea datelor, planificare, proiectare, control, simulare și identificare, etc.

3.1.2 Ant Colony Optimization

Ant Colony Optimization (ACO) este o metaeuristică inspirată din comportamentul coloniilor de furnici care sunt capabile să găsească cel mai scurt drum de la mușuroi la o sursă de mâncare. Acest comportament este determinat de interacțiunea colectivă care duce la găsirea celui mai scurt drum. Euristică ACO a fost propusă de Dorigo și a fost larg folosită pentru a rezolva probleme de optimizare (ex. problema comis voiajorului, rutarea în graf, colorarea grafurilor, etc) și a devenit rapid o euristică populară.

3.1.3 Particle Swarm Optimization

Particle Swarm Optimization (PSO) este o tehnică de optimizare bazată pe populații propusă de Eberhart și Kennedy în 1995, inspirată de comportamentul social al stolurilor de păsări și al bancurilor de pești. PSO este o tehnică bazată pe populații în care indivizii sunt numiți particule și sunt grupați în stoluri. Fiecare particulă din swarm reprezintă o soluție candidat a problemei de optimizat, care "zboară" în spațiul de căutare urmărind particula optimă de la momentul curent. PSO este aplicat în diferite domenii ca: extragerea de reguli din rețele fuzzy, recunoașterea de imagini, optimizarea de rețele electrice, etc.

3.2 Mecanisme de adaptare la medii dinamice

Problemele de optimizare dinamice sunt probleme ale căror constrângeri, obiective sau reprezentare se modifică în timp de-a lungul procesului de optimizare. În cazul problemelor în medii statice, scopul algoritmului de optimizare este de a găsi cât mai repede optimul, în schimb, în cazul dinamic, scopul este de a detecta schimbarea și de a urmări optimul de-a lungul timpului. Euristicile inspirate din natură întâmpină probleme legate de convergența prematură în medii dinamice care nu le permit descoperirea de noi optime, deoarece populația pierde din diversitate în procesul iterativ. Cea mai simplă soluție de a rezolva convergența prematură este de a reporni algoritmul atunci când se observă o schimbare în mediu, dar această abordare are dezavantajul că pierde toate informațiile stocate de populație în momentul repornirii. Astfel, s-au propus diferite tehnici care să rezolve această problemă.

3.2.1 Algoritmi evolutivi

Există diferite tehnici de păstrare a proprietăților populației, de a reacționa la schimbările din mediu fără a fi necesară repornirea alogoritmului [JB05]: menținerea diversității populației, memorii sau multi-populații.

Există două direcții principale în cazul menținerii diversității în populație: una reactivă și una proactivă. Prima abordare constă în apelarea unui mecanism de inducere a diversității în populație când se observă o schimbare în mediu (ex. hypermutația). Abordarea proactivă constă în păstrarea diversității de-a lungul întregului proces de căutare. Acest lucru împiedică convergența prin diminuarea presiunii selecției (tehnici de crowding sau sharing) sau prin stimularea directă a diversității (prin imigranți aleatori în fiecare generație).

În cazul abordărilor bazate pe memorii, algoritmi evolutivi conțin o memorie care permite refolosirea informațiilor dobândite de generațiile anterioare, aceasta fiind extrem de utilă dacă optimul se reîntoarce în locațiile anterioare. Abordările bazate pe memorii se pot împărți în: *explicite* și *implicite*. În cazul memoriilor explicite, trebuie definită o strategie de stocare și scoatere de informații (ex. readăugarea în populație a indivizilor care s-au comportat bine în situații similare). În cazul memoriilor implicite, EA conține o reprezentare redundantă a populației (ex. diploidy).

Un alt mecanism de păstrare a diversității este bazat pe populații multiple. Dacă fiecare populație conține informații despre diferite zone promițătoare ale spațiului de căutare procesul de căutare poate urmări ușor optimul. Această abordare poate fi reactivă dacă populația este reinițializată când se detectează o schimbare în mediu sau proactivă când populațiile sunt mereu încurajate să caute în diferite spații ale mediului.

3.2.2 Ant Colony Optimization

În cazul problemelor dinamice, s-au propus diferite strategii pentru a îmbunătăți performanțele algoritmului ACO: multi-populații, strategii de modificare a matricii de feromon sau noi proprietăți ale furnicilor (ex. sensibilitate).

3.2.3 Particle Swarm Optimization

S-au propus multe strategii pentru aplicarea algoritmului PSO în medii dinamice deoarece PSO întâmpină probleme legate de neactualizarea memoriei și pierderea diversității. Problemele de

neactualizare a memoriei apar datorită faptului că informațiile stocate în elemente nu mai sunt valide și pot ghida în mod greșit căutarea. În acest caz, se consideră de obicei că se știe când se realizează o modificare în mediu iar elementele sunt reevaluate. Pentru a detecta modificările din mediu se alege câteva particule aleatoare, *santinele*, care sunt reevaluate. Pe lângă particulele aleatoare, s-a propus folosirea celei mai bune și celei de a doua cea mai bună particulă pentru a fi considerate santinele. Pentru evitarea convergenței populației, s-au propus diferite strategii: introducerea diversității după detectarea unei schimbări, menținerea diversității de-a lungul căutării sau multipopulații.

3.3 Experimente numerice

Existența așa multor mecanisme de menținere a diversității în populații, ridică problema găsirii mecanismului potrivit pentru o anumită problemă. În afară de a fi efectiv, un mecanism ar trebui să fie cât mai simplu, pentru a nu crește costul computațional al algoritmului. În același timp același mecanism ar putea avea un impact diferit asupra diferitelor abordări bazate pe populații. Astfel în [ZZ06] am analizat niște mecanisme simple de păstrare a diversității aplicate la două metode de optimizare bazate pe populații: Differential Evolution (DE) [SP95] și Particle Swarm Optimization (PSO) [KE95].

3.3.1 Differential Evolution

Differential evolution (DE) a fost introdusă în [SP95] și este o metodă de optimizare bazată pe populații care utilizează un operator de recombinare aleator și o selecție bazată pe competiția dintre copil și părinte.

3.3.2 Mecanisme de adaptare la mediul dinamic

Pentru a împiedica convergența prematură a celor doi algoritmi analizați, DE și PSO în medii dinamice, am analizat impactul unor mecanisme de menținere a convergenței. În cazul DE, am analizat o schemă de adaptare a parametrilor, introducerea de elemente aleatoare și multipopulații; pentru PSO au fost testate de asemenea o simplă schemă de perturbare și multipopulații.

Mecanisme de adaptare pentru DE

În [ZZ06], ca măsură de diversitate am folosit varianța. Pentru a menține diversitatea algoritmului la același nivel pe tot parcursul procesului, la fiecare generație, t , este calculat raportul $c_j = \gamma \cdot Var(x^i)/Var(z^i)$. Un alt mecanism testat în [ZZ06] este o abordare cu multipopulații ce se caracterizează prin faptul că pentru fiecare populație se aplică algoritmul adaptiv DE, periodic un proces de migrare aleatoare este început: fiecare element al fiecărei populații, cu o anumită probabilitate, este interschimbat cu un element ales aleator dintr-o sub-populație arbitrară. Chiar dacă în cazul static aceste mecanisme împiedică convergența prematură, aceste mecanisme nu sunt suficient de bune în cazul dinamic.

Problema pe care am analizat-o în [ZZ06] este de a găsi cel mai simplu mecanism care să fie suficient de puternic încât să permită urmărirea optimului care se modifică. În cazul modificărilor continue ale poziției optimului (fără salturi mari), ar fi suficient să perturbăm cel mai bun element pentru a ajunge cu populația în zonele unde se află noul optim. Folosirea unui ”nor” de elemente care oscilează în jurul optimului ar asigura posibilitatea algoritmului de a urmări optimul.

În DE cu elemente aleatoare, populația este împărțită în două părți: o mulțime de elemente care urmează regurile DE și o altă mulțime $\{x_1, \dots, x_\mu\}$ de elemente care oscilează în jurul celui mai bun element, x_* , al populației anterioare:

$$y_i^j = x_*^j + K_j \cdot N(0, 1), \quad i = \overline{\mu + 1, m}, \quad j = \overline{1, n} \quad (3.1)$$

unde $N(0, 1)$ este o variabilă normală aleatoare cu o distribuție standard și parametrul $K_j > 0$ controlează amplitudinea perturbării.

Mecanisme de adaptare pentru PSO

Algoritmul PSO are probleme de convergență prematură când viteza, v_i devine prea mică și particulele nu se mai mișcă. În [ZZ06] am analizat două mecanisme de împiedicare a convergenței: creșterea vitezei particulei și perturbarea poziției particulei. În [LA05], este propus un mecanism de perturbare a vitezei particulei: dacă valoarea vitezei este mai mică decât o valoare prag v_c , atunci este înlocuită cu o valoare aleatoare $U(-1, 1)V_{max}/\rho$. Varianta propusă în [LA05] se numește Turbulent Particle Swarm Optimization (TPSO). Noi am propus o variantă mai

simplă, nu de perturbare a vitezei ci de perturbare a poziție particulei staționare:

$$x_i^j(t+1) = \begin{cases} x_i^j(t) + v_i^j(t+1), & \text{if } v_i^j(t+1) \geq v_c \\ x_i^j(t) + K_j \cdot N(0,1), & \text{if } v_i^j(t+1) < v_c \end{cases} \quad (3.2)$$

unde $N(0,1)$ este o variabilă aleatoare normală cu o distribuție standard și parametrul $K_j > 0$ controlează amplitudinea perturbării. Cum orice modificare în poziție influențează viteza particulei, particula poate ieși din optimul local.

3.3.3 Rezultate numerice

Pentru a analiza impactul mecanismelor de păstrare a diversității în medii dinamice am folosit pentru teste o variantă dinamică a funcției Ackley și o variantă a funcției de test Moving Peaks Benchmark [Bra].

Evoluția distanței medii dintre cel mai bun element din populație și optimul curent este ilustrată în Figura 3.1 pentru câteva variante ale algoritmului. Pentru valori mici a lui γ (e.g. $\gamma = 1$), varianta adaptivă de DE nu reușește să urmărească optimul (Figure 3.1 (d)). Pe de altă parte, folosirea unui procent între 25% și 50% de elemente aleatoare îmbunătățește semnificativ comportamentul algoritmului. S-a remarcat că pentru a urmări optimul este necesar să se reevalueze elementele populației curente la fiecare generație.

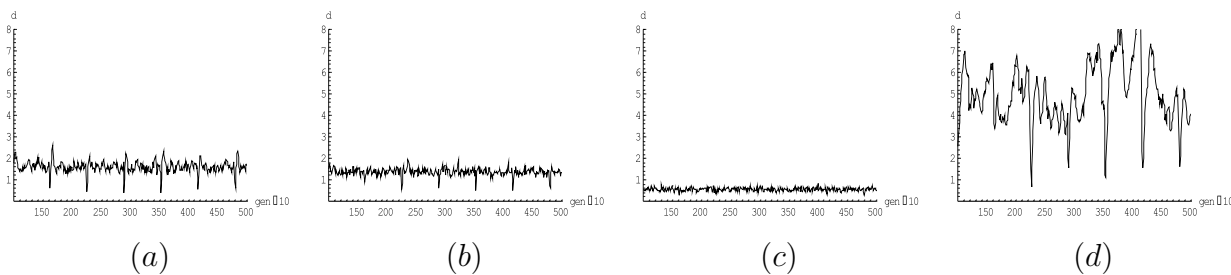


Figura 3.1: Comportamentul DE pentru varianta dinamică a funcției Ackley. (a) DE cu parametri fiși ($p = 0.5$, $F = 0.5$); (b) DE adaptiv cu $\gamma = 1.25$; (c) DE cu parametri fiși D ($p = 0.5$, $F = 0.5$) și 50% elemente aleatoare; (d) DE adaptiv cu $\gamma = 1$

Pentru algoritmul DE, am analizat impactul adaptării parametrilor și elementelor aleatoare. După cum se observă în Figura 3.2, varianta adaptivă DE are un comportament acceptabil dar nu la fel de bun ca în cazul folosirii unui mic procent de elemente aleatoare.

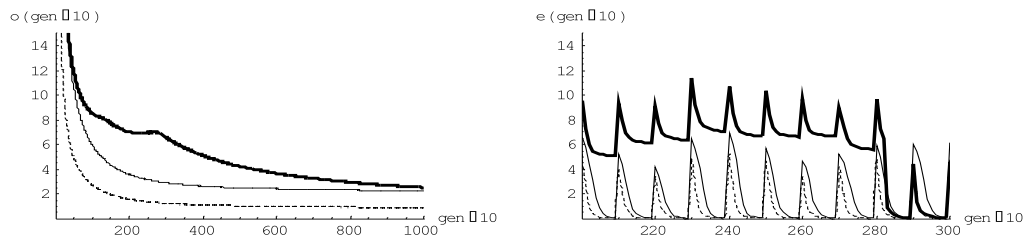


Figura 3.2: Comportamentul DE pentru un optim care își schimbă poziția: eroarea offline (stânga), eroarea curentă (dreapta). Variante DE: DE adaptiv ($\gamma = 1$) + 50% elemente aleatoare (linia groasă), DE adaptiv ($\gamma = 1$) + 25% elemente aleatoare (linia întreruptă), DE adaptiv ($\gamma = 1.25$) fără elemete aleatoare (linia subțire).

Cum convergența algoritmului PSO este mai lentă decât a algoritmului DE, este de așteptat ca algoritmul PSO clasic să poată urmări optimele care se mișcă lin. Figura 3.3 sugerează că pentru probleme dinamice simple nu este nevoie de un algoritm suplimentar de menținere a diversității (precum perturbarea vitezei/poziției sau elemente aleatoare).

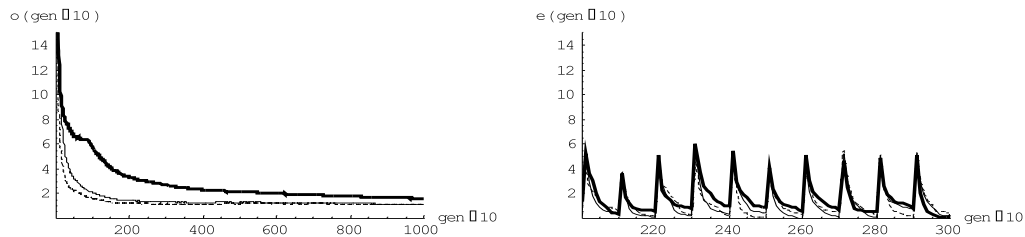


Figura 3.3: Comportamentul PSO pentru un optim care își schimbă poziția: eroarea offline (stânga), eroarea curentă (dreapta). Variante PSO: clasic (linia întreruptă), PSO perturbat (linia subțire), PSO + 25% elemente aleatoare (linie normală).

3.4 Concluzii

Atât Turbulent PSO introdus în [LA05] cât și simpla perturbare pe care am propus-o în [ZZ06] s-au dovedit a fi eficiente în evitarea convergenței premature. În ce privește multipopulațiile, experimentele numerice au arătat că încetinesc convergența prematură dar nu reușesc să o evite.

În cazul optimelor care se schimbă lent s-au obținut rezultate bune prin utilizarea unui mecanism simplu de menținere a diversității (ex. adaptarea parametrilor cu valori mari ale varianței sau adăugarea de elemente aleatoare în cazul DE) sau chiar nefolosirea nici unui

mecanism suplimentar (în cazul PSO). Dar dacă severitatea schimbării este mare, scheme suplimentare ar trebui folosite (ex. excluziunea [MM05], charged PSO [BB02]).

Capitolul 4

Planificarea taskurilor în medii dinamice distribuite

Problema de planificare a taskurilor a fost larg studiată în ultimul timp, drept consecință există mulți algoritmi care rezolvă diferite variante ale problemei de planificare.

4.1 Problema planificării pe grid

Problema de planificare pe grid este o problemă combinatorială NP-complexă. S-au propus diferite variante ale problemei în funcție de proprietățile taskurilor, particularitățile mediului și ale procesului de planificare.

Formalizare. Scopul problemei de planificare este de a găsi o alocare a resurselor la taskuri astfel încât constrângerile taskurilor (hardware, software) și constrângerile de precedență să fie îndeplinite și să se optimizeze criteriile de calitate. Atribuirea taskurilor la resurse se face pe baza unei estimări a timpului de execuție a taskului pe diferite resurse. Fie un set de n taskuri, $\{t_1, \dots, t_n\}$, care urmează să fie planificate pe un set de $m < n$ procesoare, $\{p_1, \dots, p_m\}$. Presupunem că pentru fiecare pereche (t_i, p_j) știm o estimare $ET(i, j)$ a timpului necesar ca taskul t_i să se execute pe procesorul p_j . O planificare este o mapare a taskurilor la resurse, $S = (p_{j_1}, \dots, p_{j_n})$, unde $j_i \in \{1, \dots, m\}$ și p_{j_i} reprezintă procesorul la care taskul t_i este atribuit. Dacă \mathcal{T}_j reprezintă un set de taskuri asignate la procesorul p_j și T_j^0 reprezintă momentul de timp de la care procesorul j este liber atunci timpul de terminare al acestui procesor va fi $CT_j = T_j^0 + \sum_{i \in \mathcal{T}_j} ET(i, j)$. Makespanul este timpul maxim de terminare a tuturor procesoarelor, $makespan = \max_{j=1, \dots, m} CT_j$. Problema care rămâne de rezolvat este de a minimiza valoarea makespanului.

4.2 Abordări într-un mediu static

Deoarece soluțiile suboptimale ale acestei probleme sunt acceptate în practică, noi am testat două euristici inspirate din natură: algoritmi evolutivi (Simple Population-based Scheduler - SPS) și unul inspirat de euristica Ant Colony Optimization. Experimentele realizate au fost făcute în medii statice și dinamice.

4.2.1 Simple Population-based Scheduler

În [ZFZ11] am propus un planificator numit *SPS - SimplePopulationScheduler*. Elementul cheie al acestui algoritm este strategia de perturbare (mutație), care este modulul care definește strategia de îmbunătățire a planificării. Un alt element care influențează performanțele planificatorului este pasul de inițializare.

Inițializarea populației. Construirea de planificări (sub)-optimale se bazează pe crearea unei planificări inițiale care este îmbunătățită iterativ. Când se construiește o planificare inițială se iau două decizii: (i) ordinea în care taskurile sunt atribuite la procesoare; (ii) criteriul folosit pentru a selecta procesorul corespunzător fiecărui task. În funcție de aceste criterii, există câteva strategii care pot fi aplicate: aleator, OLB, MET, MCT, Max-Min și Min-Min. Fiecare din aceste strategii generează planificări, care ar putea suferi îmbunătățiri. Din această cauză, ar fi util să nu folosim doar o singură strategie ci o populație de strategii construite folosind diferite strategii. Adăugarea de planificări generate prin diferite euristici în populația inițială a fost des utilizată [BSB⁺01, PKN10, Xha07].

Perturbarea. Planificările inițiale construite de euristicele de planificare sunt de obicei ne-optimale și pot fi îmbunătățite prin mutarea sau interschimbarea de taskuri între resurse. În funcție de criteriul utilizat, pentru a selecta resursele sursă, destinație și taskul care va fi realocat, se pot construi o mulțime de strategii de perturbare a planificării [XA10]. Majoritatea operatorilor de perturbare care se regăsesc în euristicele de planificare se bazează pe două operații: mutarea ("move") unui task de pe o resursă pe alta și interschimbarea ("swap") a două taskuri între resursele lor. Strategiile de perturbare prezentate în Tabelul 4.1 au fost selectate pe baza simplității, eficienței și balansului între caracterul greedy și aleator. Perturbarea "random move" corespunde operatorului "local move" [Xha07] și este asemănător cu operatorul de mutație folosit de algoritmi evolutivi. Operatorul "greedy move" este asemănător cu operatorul "steepest local move" [Xha07], dar implică un grad mai mare de lăcomie deoarece implică întotdeauna cel mai încărcat procesor. Operatorul "greedy swap" este similar cu "steepest local

Sursă		Destinație		Strategie
Procesor	Task	Procesor	Task	
Aleator	Aleator	Aleator	-	<i>Random Move</i>
Cel mai încărcat (max CT)	Aleator	Cea mai bună îmbunătățire	-	<i>Greedy Move</i>
Cel mai încărcat (max CT)	Aleator	Cel mai puțin încărcat (min CT)	Random	<i>Greedy Swap</i>

Tabelul 4.1: Caracteristicile strategiilor utilizate pentru a perturba planificarea

swap” [Xha07], dar este mai puțin lacom și costisitor deoarece nu implică o căutare peste întregul set de taskuri. Dacă aplicăm la fiecare pas doar un tip de perturbare (random move, greedy move or greedy swap), vom numi strategia de perturbare *SimplePerturbation*, dacă aplicăm o variantă care combină celor trei strategii perturbarea va fi numită *HybridPerturbation*.

Deoarece doar un pas de aplicare a perturbației nu duce neapărat la o îmbunătățire a calității perturbației, am optat pentru o aplicare iterativă a pasului de perturbare, până când fie n iterații sunt executate (fiecare task are șansa să fie mutat) sau până când un număr, g_p , maxim de perturbați s-au executat cu eșec.

4.2.2 Ant Colony Optimization

Această secțiune prezintă particularitățile planificatorului inspirat de comportamentul coloniilor de furnici Ant Scheduler, inspirat de algoritmul introdus în [RL04].

Inițializarea Feromonului. Pentru a ajuta construirea de planificări bune am utilizat ideea de încorporare a informației corespunzătoare unei planificări generate cu euristica Min-Min. Ideea de a adăuga informații obținute de euristici greedy este abordată comună în cazul planificatoarelor evolutive [CX06] și îmbunătățește și comportamentul planificatorului bazat pe ACO.

Căutarea locală. Pe baza concluziilor oferite de studiile anterioare [RL04], căutarea locală poate îmbunătăți semnificativ comportamentul planificatorului bazat pe ACO. Astfel, am aplicat un pas de îmbunătățire a soluției optime la fiecare epocă a algoritmului. Pasul de îmbunătățire este bazat pe un operator de rebalansare, care este asemănător cu operatorul de mutație (*HybridPerturbation*), utilizat de algoritmul evolutiv. Operația este repetată până când nu se mai observă o îmbunătățire a soluției.

4.2.3 Rezultate numerice

În cazul algoritmului Simple Population-based Scheduler (SPS) am analizat influența strategiilor de perturbare [ZFZ11]. De asemenea am analizat abilitatea SPS și a planificatorului bazat pe ACO (AS) de a construi planificări într-un mediu static pentru trei tipuri de medii: consistent (C), semi-consistent (S) și inconsistent (I) [ZZC10a, ZYC10b].

Mediu de test. Datele de test utilizate au fost introduse în [BSB⁺01], și oferă matrici care conțin valori pentru o estimare a timpului de execuție (ET) generate pe baza diferitelor presupuneri asupra heterogenității taskurilor, heterogenității și consistenței resurselor.

Analiza strategiilor de perturbare. Scopul studiului numeric în cazul mediilor statice în [ZFZ11] a fost de a analiza influența diferitelor strategii de perturbare asupra performanțelor algoritmului Simple Population-based Scheduler (SPS). De asemenea, am analizat câteva strategii de inițializare: (i) aleatoare; (ii) introducerea în populație a unor euristici de planificare; (iii) utilizarea unor euristici de planificare perturbate aleator; (iv) utilizarea euristicii Min-Min și a elementelor perturbate ale ei. După cum ne așteptam, au fost obținute cel mai bun comportament atunci când populația inițială conține elemente generate folosind alte euristici de planificare, iar cel mai rău comportament s-a observat în cazul elementelor aleatoare.

Am comparat patru variante de perturbare (aleatoare, bazată pe mutarea unui task, bazată pe interschimbarea a două taskuri și una hibridă) cu algoritm memetic hibridizat cu tabu search (MA+TS) [Xha07]. Chiar dacă se bazează pe operatori mai simpli, algoritmul SPS găsește soluții calitativ similare cu algoritmul MA+TS, în cazul datelor inconsistente ("u_i_**" problems), algoritmul propus găsește rezultate mai bune.

Analiza planificatoarelor SPS și AS în mediu static. În [ZZC10a, ZYC10b] am analizat comportamentul planificatoarelor Simple Population-based Scheduler (SPS) și a planificatorului Ant Scheduler (AS) într-un mediu static. În cazul algoritmului SPS, am folosit pentru teste o strategie simplă de perturbare - GreedyMove. Ca date de test am folosit datele caracterizate prin taskuri și procesoare cu un grad mare de heterogenitate în diferite tipuri de medii: consistent ("u_c_hihi"), inconsistent ("u_i_hihi") și semi-consistent ("u_s_hihi"). Rezultatele obținute în [CX06] sunt puțin mai bune decât cele obținute de planificatoarele noastre SPS și AS în cazurile consistent și semiconsistent, în timp ce în cazul inconsistent planificatoarele noastre se comportă mai bine.

4.3 Analiza robusteții euristiciilor

În [ZZC10b], am studiat robustețea algoritmilor de planificare prezentați anterior. O planificare bună ar trebui să satisfacă trei proprietăți: (i) ar trebui să se obțină relativ repede; (ii) ar trebui să se adapteze ușor într-un mediu cu modificări line; (iii) ar trebui să fie robustă, adică să fie puțin sensibilă la modificările din timpul rulării [CJSZ08]. Rezultatele obținute sugerează, în cazul mediilor consistente și semiconsistente, că euristica Min-Min duce la rezultate mai robuste decât planificatoarele SPS și AS, în timp ce în cazul inconsistent, planificatoarele SPS și AS dau rezultate mai bune.

4.4 Planificare Online

În cazul planificării online, taskurile vin secvențial și sunt planificate atunci când ajung; această abordare diferă față de abordările anterioare unde taskurile vin sub formă de grupuri (batch). În [ZFZ11] câteva euristici dinamice care se bazează pe îmbătrânire au fost testate în comparație cu variantele lor bazate de populații (euristica specifică a înlocuit operatorul de perturbare din SPS). Comportamentul acestor algoritmi a fost de asemenea comparat cu o variantă a algoritmului SPS care utilizează o variantă de aplicare neiterată a perturbației (pentru fiecare generație, perturbarea este aplicată o singură dată). Printre algoritmi testați se află o variantă a algoritmului DMECT descris în [Fr9], algoritmul MinQL descris în [FMC09] și euristicele Min-Min și Max-Min la care s-a adăugat un mecanism de îmbătrânire.

Mediu de test. Pentru planificare online, am considerat un model de simulare în care timpii de execuție a taskurilor (ET) urmăresc o distribuție Pareto. Rata de venire este simulată folosind un polinom de gradul 8 determinat pe baza activității unei rețele reale [Fei10]. A fost generat un număr de 500 de taskuri pentru fiecare test. Replanificarea a fost executată la fiecare 250 unități de timp utilizând un min de 1000 de unități pentru ET. Toate testele au fost repetate de 20 de ori.

Rezultate numerice. Au fost testate câteva euristici dinamice de planificare la care s-a adăugat un mecanism de îmbătrânire (DMECT, MinQL) versus variantele lor care folosesc scheletorul algoritmului SPS dar înlocuiesc pasul de perturbare cu propriul comportament (pDMECT, pMinQL). Comportamentul lor a fost de asemenea comparat cu o versiune a algoritmului SPS în care pasul de perturbare se aplică o singură dată.

Tabelul 4.2 prezintă beneficiile abordării bazate pe populații (pDMECT, pMinQL) în cazul planificării online. Variantele bazate pe populații ale algoritmilor DMECT și MinQL s-au dovedit mai bune decât variantele simple, iar varianta pDMECT are un comportament similar cu SPS (cele mai bune valori din tabelul 4.2 au fost validate folosind un test t-test cu un nivel de semnificație egal cu 0.05). Singura diferență notabilă este legată de viteza de execuție. pDMECT a avut nevoie de aproape 30 de secunde pentru a construi o planificare în timp ce SPS doar 3 secunde în medie.

	DMECT	pDMECT	MinQL	pMinQL	SPS	Max-Min	Min-Min
MS	66556.20	49409.11	76564.40	54332.89	46996.76	61165.15	68774.87
	± 15097.85	± 9522.13	± 18114.51	± 9891.15	± 8812.87	± 11936.19	± 15101.05
Timp	66.56	28343.04	3.06	2254.64	2777.70	684.49	669.21
(ms)	± 15.50	± 10702.15	± 2.52	± 314.45	± 578.22	± 242.15	± 209.99

Tabelul 4.2: Valoarea makespan (MS) obținută de euristicele de planificare în cazul planificatorului online

4.5 Adaptarea în mediul dinamic

În această secțiune vom prezenta analiza făcută în [ZZC10a, ZCC10b] referitoare la utilizarea unor mecanisme de păstrare a informațiilor din pași anteriori de planificare pentru algoritmi inspirați din natură prezentați anterior.

4.5.1 Simple Population-based Scheduler

În cazul mediului dinamic, procesul de planificare constă în evenimente de planificare consecutive. La fiecare eveniment de planificare, lista de procesoare diferă de lista de procesoare de la pasul anterior. Dacă diferența dintre listele de procesoare nu este prea mare, putem reutiliza informațiile găsite în evenimentul anterior de planificare.

Pentru a putea reutiliza informațiile găsite în evenimentele anterioare de planificare în [ZZC10a, ZCC10b] am analizat comportamentul a două variante: utilizarea planificărilor bune din evenimentul anterior de planificare și folosirea unei arhive care conține cele mai bune planificatoare din evenimentele anterioare.

4.5.2 Ant Colony Optimization

Matricea de feromon asigură comunicare între furnici, deoarece este accesată de toate furnicile în fiecare epocă a fiecărui eveniment de planificare. Din această cauză, pare a fi utilizată pentru a fi folosită ca metodă de comunicare între evenimente de planificare diferite. Acest lucru presupune că matricea de feromon nu este reinițializată la fiecare eveniment de planificare când lista de procesoare se schimbă. Valorile obținute în evenimentul anterior de planificare sunt menținute. Particularitatea acestei abordări este că valorile matricii de feromon care aparțin procesoarelor care nu sunt prezente în lista de procesoare de la evenimentul curent de planificare rămân neschimbate în pasul curent.

4.5.3 Experimente în mediul dinamic

Analiza experimentală este realizată pe datele de test prezentate în [BSB⁺01] pentru estimările timpilor de execuție și caracterul dinamic al mediului, prin faptul că unele procesoare devin neutilizabile în unele momente ale procesului de planificare.

Mediu de test. Pentru a simula comportamentul dinamic al mediului, am generat pentru fiecare eveniment de planificare o listă nouă de procesoare disponibile prin înlăturarea aleatoare a unui procent de procesoare din lista inițială de 16 procesoare. Procentele utilizate în experimente au fost de 10%, 20% și 40%.

Rezultate numerice. În cazul mediului dinamic am analizat comportamentul algoritmilor dinamici SPS și AS. Pentru algoritmul evolutiv am folosit o strategie de perturbare simplă, *Greedy Move*.

Scopul experimentelor a fost de a analiza impactul mecanismelor de memorie (refolosirea planificărilor anterioare sau matricii de feromon) asupra abilității planificatorului evolutiv și a celui bazat pe comportamentul furnicilor de a se adapta în mediul dinamic. Tablelul 4.3 conține raportul dintre numărul de evenimente de planificare în care varianta dinamică s-a comportat mai bine ca varianta statică și de asemenea raportul dintre numărul de evenimente în care varianta statică s-a comportat mai bine decât cea dinamică. Putem observa că, comportamentul celor doi algoritmi este diferit. În cazul SPS, atât timp cât diferența dintre mașinile indisponibile este mai mică decât 12%, folosirea arhivei de memorie aduce beneficii pentru toate tipurile de mediu (consistent, semi-consistent, inconsistent). Totuși, beneficiul este mai mare în cazul mediilor consistente și mai mic în cazul mediilor inconsistente. Pentru algoritmul AS, mecanismul de memorie bazat pe conservarea matricii de feromon nu aduce îmbunătățiri în cazul

mediului consistent. Pe de altă parte, dacă mediul este inconsistent, varianta dinamică de AS este mai bună decât varianta statică și de asemenea depășește performanțele variantei dinamice obținute cu SPS.

Problema	Procesoare neutilizabile (%)	SPS		AS	
		Dinamic	Static	Dinamic	Static
C	10	42/50	0/50	0/50	49/50
C	20	5/50	3/50	2/50	39/50
C	40	1/50	3/50	0/50	30/50
S	10	35/50	1/50	0/50	49/50
S	20	5/50	2/50	5/50	43/50
S	40	0/50	3/50	12/50	35/50
I	10	31/50	0/50	49/50	0/50
I	20	4/50	1/50	39/50	1/50
I	40	0/50	1/50	20/50	2/50

Tabelul 4.3: Raportul dintre numărul de evenimente de planificare când există o diferență statistică între variantele din mediul static și dinamic ale algoritmilor AS și SPS

4.6 Concluzii

În acest capitol am prezentat doi algoritmi de planificare: un planificator evolutiv (Simple Population-based Scheduler) și unul inspirat din comportamentul furnicilor. Am analizat comportamentul planificatoarelor în medii statice [ZZC10a, ZCC10b, ZFZ11], dinamice [ZZC10a, ZCC10b] și în cazul planificării online [ZFZ11].

Capitolul 5

Algoritmi de grupare a datelor inpirați din natură

Pe lângă tehnicile clasice de grupare ierarhică și partițională, tehnicile inspirate din natură, ca gruparea de date evolutivă sau gruparea de date bazată pe swarms, sunt aplicate cu succes și în gruparea datelor.

5.1 Coloniile de furnici și gruparea datelor

Comportamentul coloniilor de furnici a inspirat dezvoltarea diferitelor tehnici de grupare. Diferite abordări sunt bazate pe diferite aspecte ale comportamentului furnicilor: (i) organizarea cimitirelor și sortarea larvelor; (ii) recunoașterea chimică a vecinilor de mușuroi [LMV02]; (iii) abilitatea de a construi punți din furnici; (iv) strategiile de vânare a prădătorilor la specia *Pachycondyla Apicalis*.

5.2 Algoritmul AntClust

Algoritmul `AntClust` propus în [LMV02], simulează așa numita ”apropiere colonială” bazată pe feromon în coloniile de furnici. Feromonul este o substanță chimică pe care furnicile o posedă și care le permite să facă diferența dintre colegii de colonie și intruși. Fiecare furnică deține o informație despre ”mirosul” coloniei, care este mereu modificată. Pornind de la această idee, Labroche et al. a propus un model de furnică artificială care reușește să clasifice date.

Fiecare furnică, i , are următoarele caracteristici: (1) o dată asociată, x_i , care este un element unic ce nu se modifică în timpul procesului de grupare; (2) o etichetă, L_i , care este un număr ce identifică clusterul; (3) un prag de similaritate, T_i , care este utilizat pentru a decide dacă două furnici sunt suficient de similare ca să aparțină aceleiași colonii (4) vârsta, A_i , care de fapt este egală cu numărul de întâlniri la care a participat furnica (5) un parametru adaptiv, M_i , care reprezintă dimensiunea coloniei așa cum o percepe furnica, (6) un parametru adaptiv, M_i^+ , care reprezintă gradul de acceptanță al furnicii în colonie. Cu cât gradul de acceptanță este mai mare cu atât furnica este mai bine integrată în colonie.

Procesul de grupare se realizează în trei pași: pasul de învățare a pragului de similaritate, pasul de întâlniri aleatoare și pasul de rafinare a clusterelor mici. În pasul în care se realizează întâlnirile se aplică cinci reguli în funcție de proprietățile furnicilor: (i) crearea de noi colonii, (ii) acceptarea unei furnici într-o colonie, (iii) întâlnirea pozitivă între două furnici care aparțin aceleiași colonii, (iv) întâlnirea dintre furnici care nu aparțin aceleiași colonii.

5.3 Tratarea zgomotelor în tehnicile de grupare inspirate de comportamentul furnicilor

În [ZZ05b] am analizat comportamentul algoritmului `AntClust` pentru două seturi de date generate artificial. Primul set este format din 6 cluster elipsoidale, generate utilizând o distribuție normală bidimensională, peste care s-a suprapus un zgomot cu o distribuție uniformă (Figura 5.1(a)). Al doilea set este format din 2050 de puncte grupate în 4 cluster care au forme geometrice diferite (punctele au fost generate uniform în interiorul formelor) și 750 de puncte uniform distribuite în afara formelor geometrice, reprezentând zgomotul (Figura 5.1(d)).

Algoritmul `AntClust` (aplicat pentru $k_M = m^2/10$ pentru primul set de date și $k_M = m^2/2$ pentru al doilea set) identifică 6 și respectiv 4 cluster cum se observă în Figura 5.1(d). Cum M^+ este o măsură a acceptanței furnicii în colonie, este natural să o interpretăm ca pe un nivel de separare a datelor de zgomote. În acest caz, datele care au valori mici pentru M^+ pot fi considerate ca zgomote. Problema care apare în acest caz este de a găsi un prag care face separarea dintre date și zgomote. Deoarece toate valorile M^+ se găsesc în intervalul $[0, 1)$ am putea folosi o valoare absolută (ex. 0.1). De asemenea se poate aplica un criteriu bazat pe o ordonare statistică. Rezultatele care s-au obținut prin considerarea unui sfert din datele M^+ asociate ca fiind zgomote sunt prezentate în Figurile 5.1 (b), (c) și 5.1 (e), (f). În cazul

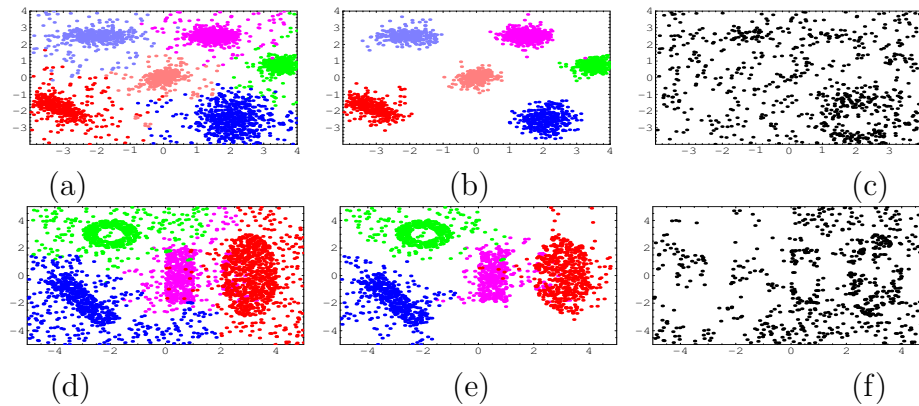


Figura 5.1: (a),(d) Rezultate obținute prin aplicarea algoritmului *AntClust* inițial; (b),(e) Clusterelor identificate prin ignorarea punctelor pentru care parametrul de acceptanță (M^+) este mai mic decât un sfert din toate valorile lui M^+ ; (c),(f) Punctele ignorate (estimarea zgomotului)

clusterelor elipsoidale, rezultatele sunt acceptabile, dar în cazul datelor geometrice, zgomotul nu este bine definit. Acest lucru ne sugerează folosirea unei măsuri legate de densitate.

5.3.1 Adăugarea informației de densitate la algoritmul *AntClust*

Pentru a introduce informația legată de densitate în algoritmul *AntClust* am propus adăugarea unui nou parametru D_i pentru fiecare furnică. Acest parametru este o estimare a densității regiunii în care se află furnica. Acest parametru este inițializat cu 0 și este ajustat la fiecare întâlnire a furnici i cu o altă furnică j cu $D_i := D_i + \Delta_i$ unde

$$\Delta_i = \begin{cases} \exp\left(\frac{(1-S(i,j))^2}{2\sigma_i^2}\right) & \text{if } 1 - S(i,j) \leq \sigma_i \\ 0 & \text{if } 1 - S(i,j) > \sigma_i \end{cases} \quad (5.1)$$

unde σ_i sunt parametrii care controlează aria de influență a fiecărei date. De obicei $\sigma_i = \sigma$ pentru toți i . Această măsură de densitate este similară cu cea folosită de algoritmul *DENCLUE*, dar în loc de a o calcula folosind o căutare sistematică a vecinilor, ne folosim de întâlnirile aleatoare la care participă furnicile.

După terminarea pasului de întâlniri, parametru D_i se împarte la vârsta furnici, A_i , informația de densitate fiind păstrată în intervalul $[0, 1)$. Prima întrebare care apare este dacă informația stocată în D_i este diferită de cea din M^+ .

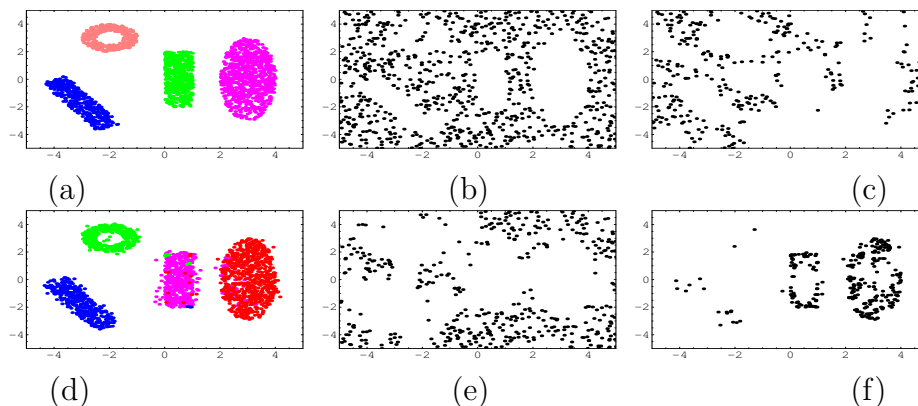


Figura 5.2: Clustere geometrice (a) clusterelor originale; (b) zgomotul suprapus pe cluster; (c) datele în a treia categorie; (d) clusterelor identificate (a patra categorie); (e) zgomotul identificat (prima categorie); (f) datele în a doua categorie.

Un prim mod de utilizare a valorilor D_i este în pasul de rafinare a clusterelor: elementele fără cluster sunt asignate unui cluster doar dacă valoarea de densitate este mai mare decât un sfert din valorile de densitate. Când căutăm elementul cel mai similar utilizat pentru a fi asignat la un cluster, se iau în considerare doar elementele care aparțin clusterului și au densitatea mai mare decât un sfert din valorile densităților. În acest mod apare o nouă clasă, cea a valorilor neclasificate care sunt considerate a fi zgomote.

Parametru σ joacă un rol important în definirea clusterelor. Am analizat influența acestui parametru asupra comportamentului algoritmului `AntClus` pentru setul de date care conține cluster sub formă de figuri geometrice (Figura 5.1 (d)) și zgomot distribuit uniform în exteriorul figurilor geometrice. Algoritmul ar trebui să identifice 5 clase, patru corespunzând datelor utile și una zgomotului. Pentru a evalua calitatea grupării am calculat eroarea de măsură introdusă în [LMV02].

Pe de altă parte, rolurile diferite pe care D_i și M_i^+ le joacă ne sugerează folosirea ambelor măsuri pentru a separa zgomotul de datele utile. Acest lucru presupune împărțirea datelor în patru categorii pe baza valorii parametrilor D_i și M_i^+ . Separarea se realizează cu ajutorul a două praguri T_M (prag pentru gradul de acceptanță) și T_D (prag pentru densitate). Cele patru categorii sunt: (1) *Prima categorie*. Conține toate datele i pentru care $M_i^+ \leq T_M$ și $D_i \leq T_D$ și corespunde datelor care pot fi considerate zgomote (Figura 5.2(e)); (2) *A doua categorie*. Conține toate datele i pentru care $M_i^+ \leq T_M$ și $D_i > T_D$ și corespunde datelor care au o estimare mare a densității dar un grad mic de acceptanță. Aceste date nu au putut fi clasificate chiar dacă aparțin unor regiuni dense. De obicei aceste date care se află pe marginea clusterelor (Figura 5.2(f)) și sunt similare cu punctele granița ale algoritmului DBSCAN; (3) *A*

treia categorie. Conține toate datele i pentru care $M_i^+ > T_M$ și $D_i \leq T_D$ și care corespunde la date cu un grad mare de acceptanță la cluster, estimate a se afla în regiuni rare. De obicei aceste date corespund unor regiuni rare, dar au participat la multe întâlniri (Figura 5.2(c)). Existența acestei categorii se explicată prin diferitele praguri utilizate în estimarea lui M_i^+ (T_i) și prin modul în care se calculează D_i (σ); (4) *A patra categorie.* Conține toate datele i pentru care $M_i^+ > T_M$ and $D_i > T_D$ și corespunde datelor care aparțin unor regiuni dense și au fost acceptate în clusterelor lor (Figura 5.2(d)).

5.4 AntClust aplicat în cazul datelor medicale

În [ZZ05a] am analizat comportamentul algoritmului AntClust propus de Labroche et al. [LMV02] în cazul datelor medicale caracterizate atât de attribute numerice cât și categoriale, dar și de date lipsă. Algoritmii de grupare bazați pe comportamentul furnicilor a fost aplicați cu succes în gruparea documentelor, web mining, dar comportamentul în cazul datelor medicale care conțin attribute numerice și categoriale precum și elemente lipsă nu a fost analizat. În [LMV02], AntClust a fost aplicat în cazul datelor medicale care conțin date numerice. Toate datele conțin doar attribute numerice și nu conțin valori lipsă.

5.4.1 Măsuri de similaritate/disimilaritate potrivite în cazul datelor medicale

Un element cheie în procesul de clasificare este alegerea unor măsuri adecvate de similaritate/disimilaritate. Când alegem astfel de măsuri, trebuie să luăm în considerare natura datelor care vor fi procesate, deoarece o măsură potrivită datelor numerice nu ar fi potrivită și datelor categoriale. Pe de altă parte, datele medicale sunt de multe ori caracterizate prin valori lipsă. Există diferite metode de corectare a acestor măsuri: înlăturarea datelor lipsă, strategii bazate pe calcularea unor distanțe parțiale sau cel mai apropiat prototip.

Datele medicale conțin frecvent attribute categoriale, exprimând nivele (grade) de prezență a unor proprietăți (ex. simptom). Codificarea clasică în acest caz este: 0-dacă simptomele nu sunt prezente și 1,2,3 ... pentru nivelele de prezență ale simptomului. Uneori, o secvență de valori consecutive nu este cea mai potrivită codare deoarece diferența dintre 0 (lipsa simptomului) și primul nivel (cel mai mic nivel al prezenței simptomului) ar trebui să fie mai mare decât diferența dintre restul de nivele. Acest lucru se întâmplă deoarece aceste nivele codează niște informații subiective care nu sunt reflectate corect de valorile care le codează.

Am studiat aceste probleme pentru o bază de date de test - Dermatology - de la UCI Machine Learning Repository. Am studiat influența pasului de învățare a măsurii de similaritate pentru diferite formule de calcul și adecvarea unei codări sau măsurii de similaritate.

5.5 Concluzii

În acest capitol am analizat abilitatea unei tehnici de grupare inspirate din natură, **AntClust**, de a se descurca în cazul datelor afectate de zgomot și a datelor medicale incomplete. Algoritmul **AntClust** a fost inspirat de abilitatea furnicilor de a realiza noi mușuroaie, de acceptare a furnicilor în mușuroaie și de recunoaștere a veciniilor din colonie. Secțiunea 5.2 conține descrierea algoritmului și unele recomandări de setare a parametrilor.

A fost analizată abilitatea algoritmului **AntClust** de a separa datele de zgomot. A fost realizată o analiză a utilității parametrilor M^+ dar și a parametrului de densitate D introdus. Calculul acestui nou parametru și faza de post procesare introdusă nu au modificat semnificativ timpul de execuție al algoritmului. Cum parametrul σ poate fi calculat pe baza pragului de similaritate, nu este nevoie de setarea unui nou parametru. Rezultatele preliminare sunt încurajatoare, dar multe lucruri legate de informațiile stocate de furnică nu sunt complet clarificate. Valorile folosite pentru setarea pragurilor T_M și T_D au fost determinate pe baza unor experimente, nu pe baza unui rezultat analitic. De asemenea, se doresc rezultate referitoare la estimarea parametrilor M , M^+ și D .

De asemenea, am analizat abilitatea algoritmului **AntClust** de a identifica clustere în datele medicale. Importanța codificării datelor categoriale a fost de asemenea analizată și comparată cu algoritmul K-Means. În acest caz, am observat că dacă parametrii și codificarea este aleasă cu grijă algoritmul **AntClust** se aplică cu succes în gruparea de date medicale. Pe de altă parte, când sunt implicate attribute care conțin valori clasificate pe nivele, o codificare adecvată a valorilor îmbunătățește rezultatele grupării.

Capitolul 6

Abordări bazate pe algoritmi evolutivi pentru extragerea de reguli

Procesul de extragere de reguli (data mining) este unul din pașii cheie al KKD, care presupune aplicarea câtorva metode de procesare pentru a facilita aplicarea algoritmilor de data mining pentru a îmbunătăți și rafina procesul de descoperire a cunoștințelor.

6.1 Extragerea de reguli din date

O regulă poate fi prezentată în următoarea formă: *IF "câteva condiții aplicate atributelor prezise sunt adevărate" THEN "unele condiții aplicate atributelor țintă sunt adevărate"*. Dacă există doar un atribut țintă și acesta specifică o clasă, atunci problema descrisă este una de *clasificare*, adică data care satisface condiția antecedent (IF) aparține unei clase specificate de consecvent (THEN). Când atributele țintă nu exprimă o clasă, vom avea de a face cu reguli de *predicție*, exprimând ipoteze despre dependențele dintre antecedentul și consecventul regulii. Dacă setul de atribute consecvent și antecedent nu sunt predefinite, vorbim despre reguli de *asociere* care exprimă asocierea diferitelor atribute. Descoperirea și selectarea de reguli este un proces de căutare ghidat de măsuri de calitate: *acuratețea*, *sensul* și *înteresul* regulii. Aceste măsuri de multe ori intră în conflict, de exemplu regulile care au o acuratețe mare nu sunt neapărat interesante și ușor de citit, din acest motiv procesul de căutare trebuie să fie multicriterial. Pe de altă parte, datele care se regăsesc în practică (de ex. datele medicale studiate în [LZZ08]) conțin un număr mare de valori lipsă (MVs), care implică diferite tehnici de corecție pentru a putea fi evaluate.

6.1.1 Măsuri de evaluare

Calitatea regulilor de clasificare, predicție sau asociere poate fi contorizată utilizând diferite măsuri statistice, fiecare dintre ele specificând anumite caracteristici ale regulilor. Regurile extrase din date ar trebui să caracterizeze datele, să fie inteligibile și să conțină informații noi și interesante, astfel măsurile de calitate se pot împărți în următoarele clase:

- măsuri de acuratețe - cuantifică gradul în care regulile exprimă datele: suport (*Supp*), încredere (*Conf*), acuratețe (*Acc*), specificitate (*Spec*) și sensitivitate (*Sens*);
- măsuri de înțelegere - cuantifică gradul în care regula este ușor înțeleasă;
- măsuri de interes - cuantifică gradul în care regula conține noi informații sau informații deja știute (coeficientul Phi (Φ), procentul de neînțelegere (*OR*) și măsura cosinus (*cos*)).

Regurile au următoarea structură: *IF* (AT_1, AT_2, \dots, AT_k) *THEN* (CT_1, CT_2, \dots, CT_l), unde AT_i este antecedentul regulii iar CT_j este consecventul regulii. Fiecare termen conține cel puțin un atribut și este un triplet $\langle a, op, valoare/interval \rangle$ unde a este un atribut, op este un operator (egal, diferit, în, not în, mai mic, mai mare) și $valoare$ este o posibilă valoare sau set de valori ale atributului.

6.1.2 Valori lipsă în date

Când evaluăm o regulă pentru un set incomplet de date, trebuie să decidem cum vom trata valorile lipsă. Cea mai radicală abordare de a trata valorile lipsă (MV) este de a elimina observațiile care conțin valori lipsă (de ex. datele incomplete) din analiză. Variantele pe care le-am analizat în [LZZ08] se bazează pe ideea de a trata datele incomplete în timpul procesului de extragere de cunoștințe. În acest caz, valorile lipsă intră în discuție în momentul când verificăm dacă o înregistrare se potrivește cu o anumită regulă.

În cazul datelor incomplete, când o regulă R este evaluată, valorile lipsă pot interfera cu procesul de evaluare numai dacă attribute sunt implicate în regulă. Va trebui să decidem cum trebuie modificate calculele în acest caz. Am analizat două metode.

Metoda 1. În această variantă, am încercat să limităm penalizarea regulii care implică attribute incomplete.

Metoda 2. În această variantă, nu am ignorat datele incomplete ci le-am penalizat utilizând valori non-crisp. În loc de valoare 0 de potrivire pentru attributele lipsă, fiecărei valori atribut lipsă a atributului i s-a atribuit o probabilitate de a satisface termenul regulii.

Aceasta este o variantă de estimare a probabilității care arată că o valoare lipsă va satisface termenul regulii, bazată pe presupunerea că valorile atributelor sunt distribuite uniform în domeniul lor de definiție. Prin utilizarea altor modele de distribuție pentru valorile atributelor, s-ar putea obține diverse valori de potrivire. Diferența dintre prima variantă și cea curentă constă în faptul că în primul caz $cover(R, S) \in \{1, \dots, cardS\}$, în timp ce în al doilea caz $cover(R, S) \in [1, cardS]$.

6.2 Un algoritm evolutiv pentru extragerea de reguli din date

În [ZLZ08] ne-am propus să implicăm utilizatorul în procesul de căutare, deoarece găsirea unei combinații de măsuri de calitate predefinite este dificilă; de asemenea utilizatorii se pot răzgândi asupra calității regulii în timpul procesului de optimizare [OAT⁺99]. Abordarea propusă este bazată pe un algoritm evolutiv multi-obiectiv (MOEA).

Am aplicat tehnica de extragere de reguli pentru datele de test oferite de repositoryul UCI și pentru un set de date obstetrice colectate timp de un an de la un spital de Obstetrică-Ginecologie. Deoarece scopul studiului nostru a fost de a ajuta medicii specialiști în a lua decizii, trebuie să luăm în considerare feedbackul oferit de specialiști și să le permitem să intervină în procesul de extragere de reguli.

Codarea. Fiecare element (cromozom) al populației corespunde la o regulă și este format dintr-o listă de componente (gene) care corespund la toate attributele din setul de date. Fiecare componentă este constituită din trei câmpuri: $\langle \text{flag de prezență}, \text{operator}, \text{valoare} \rangle$. Un element este o listă de valori mixte (binare, întregi, reale și interval). Diferența dintre attributele antecedent și consecvent este realizată doar în pasul de evaluare a elementului.

Operatori evolutivi. În fiecare generație, o nouă populație este creată din cea curentă folosind câțiva operatori evolutivi.

Prin *încrucișare*, o nouă regulă este construită pornind de la două reguli selectate aleator din populația curentă.

Mutația deține rolul de a modifica regulile obținute prin încrucișare. Pentru fiecare atribut, operatorul de mutație este aplicat cu o probabilitate (de ex. $p_m = 1/n$) și poate afecta doar un câmp (de ex. flag-ul de prezență, operatorul sau valoarea) și doar unul la fiecare pas de mutație.

Selecția și arhiva. După ce o nouă populație este creată prin intermediul mutației și încrucișării, este aplicat un pas de selecție (tipic pentru MOEA). Strategia de selecție este similară cu cea utilizată de algoritmul NSGA-II [DK07], ceea ce presupune că elementele din populația reunită (părinții și copiii) sunt aranjate în funcție de o relație de nedominanță. Pentru a stimula diversitatea în frontul Pareto rezultat, o distanță de crowding este utilizată ca un criteriu secund de selecție: pentru două elemente care au același rang, este selectat cel cu o distanță de crowding mai mare (sugerând că aparține la zone mai puțin dense). Distanța de crowding poate fi definită în spațiul variabilelor de decizie sau obiectiv. O caracteristică particulară a abordării noastre este legată de distanța de crowding dintre reguli.

Am analizat două tipuri de distanțe, una care exprimă diferența structurală dintre reguli și cealaltă exprimă diferența dintre dată și subsetul de date acoperit de regulă.

După un anumit număr de generații, se creează o arhivă care conține elemente nedominate. Nu toate elementele nedominate din populația curentă sunt transferate în arhivă, ci ele sunt filtrate astfel încât distanțele structurale și bazate pe cover dintre două elemente ale arhivei să fie mai mari decât un prag (în analiza noastră am folosit 0.01).

6.3 Interacțiunea utilizatorului cu procesul de extragere de informații

O căutare iterativă permite utilizatorului să interfereze cu procesul evolutiv pentru a-l ghida spre zone interesante din spațiul de căutare. Ideea de a permite utilizatorului să interfereze cu procesul evolutiv a mai fost explorată în contextul optimizării multi-obiectiv evolutive [DK07] prin faptul că cere utilizatorului să furnizeze un așa numit punct de referință în spațiul obiectiv.

În varianta interactivă, procesul de căutare este alcătuit din câțiva pași; la fiecare pas, populația evoluează un număr de generații și o arhivă cu regulile nedominate este oferită utilizatorului împreună cu toate măsurile obiectiv utilizate pentru setul testat (măsuri nu neaparat limitate la cele utilizate ca criterii în procesul de optimizare). În implementarea noastră am folosit următoarele măsuri: suportul, încrederea, acuratețea, specificitate, sensibilitate,

înțelegere, gradul de înțelegere, lift, uncovered negative și riscul relativ [OAT⁺99]. Pe baza acestor criterii și a evaluării subiective, utilizatorul poate decide dacă o regulă este neinteresantă sau neînțeleasă. Apoi utilizatorul poate marca aceste reguli și se poate trece la noul pas al căutării. Efectul marcării regulilor ca nenesesare este următorul: în primul rând, elementele populației marcate sunt înlocuite cu reguli generate aleator. Acest lucru permite algoritmului să exploreze noi regiuni ale spațiului de căutare, asemănătoare cu tehnica imigranților aleatori; în al doilea rând, regulile marcate sunt adăugate unei liste de reguli interzise (L_p), creându-se o arhivă care conține părți ale spațiului de căutare care sunt considerate a fi neinteresante de către utilizator.

Un posibil efect secundar al utilizării unui criteriu suplimentar de optimizare este că de obicei duce la un număr mare de elemente nedominate. Pentru a evita acest efect, lista de reguli neinteresante poate fi utilizată pentru a introduce unele constrângeri, de ex. toate regulile similare cu cele din listă sunt considerate nefezabile. Fezabilitatea regulii intervine când se verifică relația de dominanță dintre două reguli: dacă o regulă este fezabilă și cealaltă nu este, prima regulă o domină pe cea de a doua, ignorând valorile criteriului de calitate; dacă ambele sunt fezabile sau nefezabile, dominanța este decisă pe baza criteriului de calitate.

6.4 O abordare evolutivă de reducere a regulilor extrase din date

Algoritmul NNGE (Non-Nested Generalized Exemplars) este un algoritm hibrid de învățare bazat pe instanțe care deduce din date reguli de clasificare reprezentate ca hiper-dreptunghiuri neimbricate și nesuprapuse. După ce un set de hiperdreptunghiuri \mathcal{H} se generează folosind algoritmul NNGE, acesta este postprocesat utilizând un algoritm evolutiv pentru a încerca reducerea dimensiunii și a acurateții clasificării [ZPNZ11].

6.5 Experimente numerice

Experimentele numerice sunt realizate pe două tipuri de date: date medicale și date despre caracteristicile taskurilor care urmează a fi planificate. În cazul datelor medicale, este analizat un algoritm evolutiv multi-obiectiv. Pentru datele de test de planificare, am analizat comportamentul unui algoritm evolutiv de postprocesare utilizat în selecția regulilor generate cu algoritmul NNGE.

6.5.1 Setul de date medicale

Rezultatele numerice conțin rezultate obținute pentru date medicale de test de la repository-ul UCI și un set de date colectat de la un spital de Obstetrică Ginecologie pe parcursul unui an. În [ZLZ08] scopul experimentelor a fost de a valida abilitatea abordării evolutive de a descoperi reguli precise, și de a analiza impactul intervenției utilizatorului în procesul de căutare. Pe de altă parte, în [LZZ08] am analizat influența câtorva metode de tratare a datelor incomplete în procesul de găsire a regulilor de clasificare din date medicale.

Date de test

Setul de date *ObGyn* este format din date colectate timp de un an (2006) într-un spital de Obstetrică-Ginecologie. Setul de date conține 2686 de înregistrări care sunt împărțite în două clase: nașteri înainte de termen (370 înregistrări, reprezentând 13.77%) și nașteri la termen (2316 înregistrări, reprezentând 86.23%). Fiecare înregistrare conține 63 de atribute corespunzătoare diferitelor caracteristici ale mamei și copilului. Procentul de valori lipsă din date este de 23% și sunt distribuite neuniform printre atribute.

Extragerea de reguli de asociere

Pentru a valida abilitatea algoritmului multi-obiectiv propus de a extrage reguli valoroase din date, l-am testat prima dată în cazul problemelor de clasificare. Rezultatele au fost obținute pentru seturile de date *Pima Indians Diabetes* și *Wisconsin Breast Cancer* (1991), bazate pe două criterii de optimizare; acuratețea (*Acc*) și *uncovered negative (UN)*. Rezultatele MOEA sunt comparabile cu cele obținute prin aplicarea altor clasificatori implementați în WEKA: clasificatori simpli (ZeroR, OneR), clasificatori de reguli conjunctive (CR), clasificatori bazați pe tabele de decizie (DT), clasificatori bazați pe învățarea incrementală regulilor (JRIP), clasificatori bazați pe cel mai apropiat vecin îmbinat cu o tehnică de neimbricare a prototipurilor (NNge), arbori de decizie parțială (PART).

Impactul interacțiunii cu utilizatorul

Două variante de includere a evaluării utilizatorului în procesul de căutare au fost implementate și testate: criterii utilizator și constrângeri utilizator. *Criteriile utilizator*: pentru fiecare element, un criteriu suplimentar de optimizare este adăugat între distanțele minimal structural și

cover-based pentru toate elementele din lista marcată. *Constrângerile utilizatorului*: un element care este structural sau semantic similar cu cel puțin o regulă marcată este considerat nefezabil. Fezabilitatea regulii intervine când se verifică relația de dominanță dintre două reguli: o regulă nefezabilă este întodeauna domină o regulă fezabilă; dacă ambele sunt fezabile sau nefezabile, dominanța este decisă pe baza criteriului de calitate.

Impactul acestor două variante asupra numărului de reguli nedominate este ilustrat în tabelul Table 6.1, pentru setul de date *Wisconsin Breast Cancer* (fiecare pas execută 100 de generații). Intervenția utilizatorului constă în marcarea calității tuturor regulilor cu valori de calitate mai mici decât un prag (de ex. încrederea, sensibilitatea, specificitatea sau acuratețea mai mici decât 0.75). Rezultatele din Tabelul 6.1 demonstrează că intervenția utilizatorului a condus la găsirea unor reguli mai sensibile și precise. După cum ne-am așteptat, varianta II *criteriile utilizatorului* au permis unui număr mai mare de reguli să evolueze în comparație cu varianta III (*constrângerile utilizatorului*).

Var/ Pas	Nr. reguli	Coef.Marcare reguli	Conf. interval	Sens. interval	Spec. interval	Acc. interval	Lift interval
I/1	16	–	[0.42, 1]	[0.25, 0.98]	[0.31, 1]	[0.54, 0.93]	[1.24, 2.9]
I/2	6	–	[0.23, 1]	[0.25, 0.55]	[0.88, 1]	[0.81, 0.93]	[2.88, 12.4]
I/3	4	–	[0.97, 1]	[0.28, 0.53]	[0.99, 1]	[0.94, 0.96]	[12.1, 12.4]
II/1	16	10	[0.42, 1]	[0.25, 0.98]	[0.31, 1]	[0.54, 0.93]	[1.24, 2.9]
II/2	36	16	[0.75, 1]	[0.07, 0.98]	[0.97, 1]	[0.92, 0.98]	[10.06, 12.4]
II/3	32	19	[0.21, 1]	[0.21, 0.99]	[0.88, 1]	[0.84, 0.98]	[2.6, 12.4]
III/1	16	10	[0.42, 1]	[0.25, 0.98]	[0.31, 1]	[0.54, 0.93]	[1.24, 2.9]
III/2	15	6	[0.23, 1]	[0.28, 0.98]	[0.83, 1]	[0.81, 0.99]	[2.8, 12.4]
III/3	11	2	[0.75, 1]	[0.27, 0.98]	[0.97, 1]	[0.96, 0.99]	[9.3, 12.4]

Tabelul 6.1: Breast data set: intervalele măsurilor de calitate ale regulilor găsite pentru trei scenarii ale algoritmului MOEA. Varianta I: fără interațiunea utilizatorului. Varianta II (*criterii utilizator*) și III (*constrângeri utilizator*): la fiecare pas, utilizatorul marchează regulile care au măsurile de confidență, sensibilitate, specificitate, sau acuratețe mai mici de 0.75.

De asemenea am folosit variantele interactive de MOEA pentru a explora spațiul de căutare în cazul datelor obstretice (setul de date *ObGyn*), scopul experimentului fiind găsirea potențialilor factori de risc pentru nașteri premature (de ex. nașterea înainte de 37 de săptămâni gestaționale). Modelele de risc pot fi exprimate ca reguli de clasificate (IF ”condițiile antecedent sunt sat-

isfăcute” THEN clasa=prematuro) sau ca reguli de predicție (IF ”condițiile antecedente sunt satisfăcute” THEN ”vârsta gestației este mai mică decât o valoare”).

Strategia propusă de implicare a utilizatorului în procesul de căutare este doar un prim pas în dezvoltarea unui sistem interactiv a cărui scop este de a ajuta medicul să exploreze datele și să extragă noi, posibil neașteptate, cunoștințe. Rezultatele obținute în cazul datelor obstetice nu au fost suficient de relevante din diferite motive: pe de o parte, particularitățile datelor (multe valori lipsă și eronate); pe de altă parte, limitările strategiei evolutive. Folosirea de valori numerice pentru valorile continue a fost utilizată pentru a evita o discretizare preliminară, dar acest lucru a condus la un spațiu de căutare foarte mare și descoperirea de reguli care nu sunt ușor de interpretat. Folosirea de variabile fuzzy în loc de variabile crisp ar putea îmbunătăți calitatea regulilor, mai ales în cazul datelor medicale.

Influența tratării valorilor lipsă din date pentru regulile de clasificare

Deoarece datele de test de la UCI nu conțin valori lipsă, am pregătit datele pentru teste prin eliminarea aleatoare a $mv\%$ din valorile atributelor ($mv \in \{10, 20, 30\}$ și reprezintă un procent din numărul total al atributelor din setul de date). Pentru a realiza validarea încrucișată, datele au fost împărțite în patru exemplare. Valorile lipsă au fost introduse doar în datele utilizate pentru antrenare, setul de validare conținând valorile originale ale atributelor.

Pentru a analiza influența valorilor lipsă din date am utilizat trei metode: *Metoda 1* și *Metoda 2* sunt descrise în secțiunea 6.1.2 și *Metoda 3* se bazează pe o simplă strategie de imputare (datele sunt pre-procesate astfel încât valorile lipsă sunt înlocuite cu mediana valorilor existente pentru atributele corespunzătoare). Prima metodă este varianta care încearcă să limiteze penalizarea regulilor și a doua variantă folosește valori non-crisp pentru potrivire.

Experimentele au ilustrat că prin utilizarea de metode diferite de tratare a valorilor lipsă din date s-au obținut diferite seturi de reguli. Pe de altă parte, diferența statistică dintre calitatea regulilor obținute cu diferite metode nu este semnificativă. Rezultate puțin mai bune s-au obținut pentru metodele care ajustează calculul valorilor de calitate pentru a trata valorile lipsă din date (Metoda 1 și Metoda 2).

De asemenea am realizat o analiză și asupra datelor obstetice (obGyn). Scopul a fost găsirea de reguli care corespund nașterilor premature. Numărul mare de atribute duce la utilizarea unui spațiu mare de căutare, creând dificultăți algoritmului evolutiv în căutarea de reguli de calitate ridicată. Din această cauză, am ales diferite seturi de atribute care sunt implicate în reguli. Rezultatele prezentate în Tablul 6.2 au fost obținute când am folosit în antecedentul regulii de

clasificare informația despre sarcinile anterioare, avorturi provocate, avorturi naturale și despre înălțimea fătului. Am utilizat produsul dintre specificitate și sensibilitate ca criteriu unic de optimizare. Din cauză că am folosit un singur criteriu de optimizat la fiecare rulare, a rezultat o singură regulă. Toate regulile obținute au conținut termeni care corespund la înălțimea ”fundus uter” (de ex. înălțimea ”fundus uter” ≤ 29) chiar dacă acest atribut conține valori lipsă în cazul unor date. Acest lucru poate fi explicat prin faptul că atributul care conține înălțimea ”fundus uter” apare de obicei în cazul nașterilor premature. Pe de altă parte, după cum sugerează Tabelul 6.2, rezultatele obținute de cele trei metode nu diferă semnificativ. Cu toate acestea, prima metodă oferă rezultate puțin mai bune decât celelalte două.

	acc	spec	sens	UN	lift
Method 1	0.70±0.09	0.73±0.13	0.52±0.17	0.58±0.12	1.92±0.63
Method 2	0.64±0.08	0.64±0.12	0.59±0.18	0.55±0.1	1.57±0.34
Method 3	0.64±0.05	0.66±0.07	0.54±0.11	0.57±0.06	1.54±0.29

Tabelul 6.2: Rezultate pentru setul de date obstretice

6.5.2 Setul de date utilizat pentru planificarea taskurilor

O alternativă la algoritmi de interschimbare este *selecția utilizând forța brută* (BS). În acest caz, fiecare strategie de planificare (SA) este testată pe setul de date, ceea ce este o abordare care consumă timp. O alternativă ar fi aplicarea strategiei BS doar pe un set de date de antrenare. Setul de date de antrenare conține câteva caracteristici ale platformei pentru scenariul de planificare (taskurile și resursele relaționate) împreună cu cea mai bună SA (eticheta clasei) pentru o configurare specifică, găsită de BS. Acest set de date poate fi folosit pentru a antrena un sistem de clasificare. Apoi, când apare o nouă configurare, clasificatorul generat în pasul anterior este utilizat pentru a găsi SA corespunzător. În [ZF11a, ZF11b] am testat câteva strategii de clasificare pentru a găsi strategia care asigură cea mai bună acuratețe de clasificare și identificarea caracteristicilor evenimentelor de planificare care influențează cel mai mult alegerea algoritmului de planificare.

Date de test

Setul de date a fost generat artificial utilizând modelele descrise în [Fei10]. Fiecare instanță a setului de antrenare este formată din următoarele atribute: *timpul când planificarea se termină*,

media timpului estimat de execuție (ETT); media deviațiilor standard a EET; media timpului estimat de terminare (ECT); media deviației standard a ECT; media dimensiunii taskurilor; media deviațiilor standard a mărimii taskurilor; numărul total de taskuri; și numărul de taskuri lungi utilizate în experiment. Pe lângă informația pentru fiecare configurație, cea mai bună strategie de planificare găsită de BS a fost adăugată pentru a identifica clasa. Șapte SA au fost utilizate pentru a determina cea mai bună strategie: Max-Min [MAS+99], Min-Min [MAS+99], Suffrage [CLZB00], MinQL [FMC09], MinQL-Plain [FMC09], DMECT [Fr9] și DMECT2 [Fr9].

Clasificatorii utilizați în teste. În experimentele noastre am utilizat câțiva clasificatori implementați în WEKA data mining toolkit (o rețea neuronală multi-perceptron (MPL), Radial Basis Function (RBF) network, Non-Nested Generalized Exemplars (NNGE)), un clasificator ne-supervizat Fuzzy C-Mean și un clasificator hibrid care combină algoritmul NNGE cu un algoritm evolutiv de selecție a atributelor și exemplarelor (EP-NNGE descris în Secțiunea 6.4).

Rezultatele testării. O primă analiză a fost realizată pentru a determina timpul necesar strategiilor de planificare pentru selecția celui mai potrivit planificator. Timpul de execuție pentru euristicele de clasificare a fost sub 2.5s (antrenare+clasificare), în timp ce în cazul strategiei BS pentru 7 SA timpul mediu a fost de aproximativ 6 secunde pentru fiecare eveniment de planificare. Acuratețea mare a clasificării și timpul scurt de învățare fac din tehnicile de clasificare un candidat potrivit pentru a înlocui politica de interschimabare a planificatoarelor și strategia BS.

O a doua analiză a fost realizată pentru a determina cea mai bună tehnică de clasificare. Comportamentul clasificatorilor propuși EP-NNGE și EPA-NNGE este similar pentru cele trei seturi de date și mai bun decât restul de euristicilor de clasificare.

6.6 Concluzii

În acest capitol am prezentat un algoritm evolutiv multi-obiectiv pentru extragerea de reguli din datele medicale și o strategie evolutivă de selecție pentru determinarea de reguli de clasificare care asigură o bună alegerea a unei strategii de planificare pentru un eveniment de planificare.

Capitolul 7

Concluzii și direcții viitoare

Această teză se bazează pe aplicarea euristiciilor inspirate din natură în medii incerte; sunt discutate trei probleme de optimizare concretă (planificare pe grid, gruparea datelor, extragerea de reguli din date) care pot fi formulate ca probleme de optimizare în medii dinamice și incerte și câteva funcții de test.

În cazul problemelor de optimizare dinamică am analizat doi algoritmi: Particle Swarm Optimization și Differential Evolution. Pentru a îmbunătăți comportamentul euristiciilor în medii dinamice am introdus un mecanism de perturbare a poziției particulei în cazul algoritmului PSO (Secțiunea 3.3.2) și elemente aleatoare care oscilează în jurul optimului în cazul DE (Secțiunea 3.3.2).

În cazul problemei de planificare pe grid am propus un planificator simplu bazat pe populații și am analizat robustețea lui (Secțiunea 4.3), în cazul planificatoarelor batch comportamentul într-un mediu static (Secțiunea 4.2.3) și dinamic (Secțiunea 4.5.3), în cazul planificării online (Secțiunea 4.4). Deasemenea am studiat influența a diferite strategii de inițializare (Secțiunea 4.2.1) și perturbare (Secțiunea 4.2.3) a planificatorului bazat pe populații.

În cazul problemei de grupare de date am adăugat un parametru de densitate (Secțiunea 5.3.1) algoritmului `AntClust` pentru a-i îmbunătăți performanțele în cazul grupării datelor afectate de zgomot și am analizat performanțele euristicii în cazul datelor medicale care conțin atât atribute numerice cât și categoriale afectate de valori lipsă (Secțiunea 5.4).

În cazul problemei de extragere de reguli din date, am propus introducerea expertului uman (Secțiunea 6.3) pe lângă evaluarea automată a regulilor algoritmului evolutiv. O particularitate a acestei abordări constă în adăugarea unei distanțe de crowding între reguli, care se comportă ca un criteriu de diversitate, iar mulțimea Pareto este stimulată să aleagă elemente din zone mai puțin dense în arhivă (Secțiunea 6.2). De asemenea, am propus metode de tratare a valorilor

lipsă de date și am analizat influența acestora în procesul de extragere de reguli (Secțiunea 6.5.1).

O altă problemă abordată în cazul extragerii de reguli din date se referă la clasificarea aplicațiilor care sunt trimise planificatoarelor astfel încât, atunci când apare un nou eveniment de planificare selectarea algoritmului potrivit să se facă pe baza caracteristicilor taskurilor. Contribuția în acest caz este legată de introducerea planificatorului hibrid bazat pe exemplare neimbricate și aplicarea unei selecții evolutive asupra atributelor și exemplarelor (Secțiunea 6.4).

Câteva direcții viitoare de cercetare sunt prezentate la sfârșitul fiecărui capitol. În cazul problemei de planificare pe grid, studiile viitoare vor trata cazul taskurilor relaționate sau folosirea altor metrici ca Total Processing Consumption Cycle care este o alternativă la makespan și este independentă de hardware. În cazul problemei de grupare a datelor, se dorește o analiză a comportamentului algoritmului pentru alte date dar și pentru datele reale. În cazul problemei de descoperire a regulilor din date, se poate realiza un studiu mai amănunțit a diferitelor metode de tratare a valorilor lipsă din date și influența lor asupra designului algoritmului evolutiv. Deoarece datele utilizate pentru testare în cazul algoritmilor de planificare sunt nebalansate, se dorește realizarea unui studiu al diferitelor metode de balansare asupra tehnici de căutare propuse.

Deoarece problemele de optimizare în medii incerte și dinamice se regăsesc printre problemele reale iar euristicele inspirate din natură sunt candidate bune pentru a rezolva aceste probleme, studiile viitoare se vor axa pe analiza unor algoritmi din natură având ca scop rezolvarea unor probleme de optimizare în medii incerte și dinamice.

Bibliografie

- [BB02] T. M. Blackwell and Peter J. Bentley. Dynamic search with charged swarms. In *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO '02*, pages 19–26, San Francisco, CA, USA, 2002. Morgan Kaufmann Publishers Inc.
- [Bra] J. Branke. Moving peaks benchmark. www.aifb.uni-karlsruhe.de/jbr/MovPeaks/movpeaks/.
- [BSB⁺01] Tracy D. Braun, Howard Jay Siegel, Noah Beck, Lasislau L. Bölöni, Muthucumara Maheswaran, Albert I. Reuther, James P. Robertson, Mitchell D. Theys, Bin Yao, Debra Hensgen, and Richard F. Freund. A comparison of eleven static heuristics for mapping a class of independent tasks onto heterogeneous distributed computing systems. *J. Parallel Distrib. Comput.*, 61:810–837, June 2001.
- [CJSZ08] Louis-Claude Canon, Emmanuel Jeannot, Rizos Sakellariou, and Wei Zheng. Comparative evaluation of the robustness of dag scheduling heuristics. In Sergei Gorlatch, Paraskevi Fragopoulou, and Thierry Priol, editors, *Grid Computing*, pages 73–84. Springer US, 2008. 10.1007/978-0-387-09457-1-7.
- [CLZB00] H. Casanova, A. Legrand, D. Zagorodnov, and F. Berman. Heuristics for scheduling parameter sweep applications in grid environments. In *Procs. 9th Heterogeneous Computing Workshop (HCW)*, pages 349–363, Cancun, Mexico, May 2000.
- [CX06] J. Carretero and F. Xhafa. Using genetic algorithms for scheduling jobs in large scale grid applications. *Journal of Technological and Economic Development - A Research Journal of Vilnius Gediminas Technical University*, 12:11–17, 2006.
- [DK07] Kalyanmoy Deb and Abhishek Kumar. Interactive evolutionary multi-objective optimization and decision-making using reference direction method. In *Proceedings of the 9th annual conference on Genetic and evolutionary computation, GECCO '07*, pages 781–788, New York, NY, USA, 2007. ACM.
- [Fei10] D.G. Feitelson. Workload modeling for computer systems performance evaluation.

URL <http://www.cs.huji.ac.il/~feit/wlmod/>, 2010.

- [FMC09] M. Frîncu, G Macariu, and A. Cârstea. Dynamic and adaptive workflow execution platform for symbolic computations. *Pollack Periodica*, 4(1):145–156, 2009.
- [Fr9] Marc E. Frîncu. Dynamic scheduling algorithm for heterogeneous environments with regular task input from multiple requests. In *Proceedings of the 4th International Conference on Advances in Grid and Pervasive Computing*, GPC '09, pages 199–210, Berlin, Heidelberg, 2009. Springer-Verlag.
- [JB05] Y. Jin and J. Branke. Evolutionary optimization in uncertain environments - a survey. *Evolutionary Computation, IEEE Transactions*, 9(3):303–317, 2005.
- [KE95] J. Kennedy and R.C. Eberhart. Particle swarm optimization. *Proceedings of IEEE International Conference on Neural Networks*, pages 1942–1948, 1995.
- [LA05] Hongbo Liu and Ajith Abraham. Fuzzy adaptive turbulent particle swarm optimization. In *Proceedings of the Fifth International Conference on Hybrid Intelligent Systems*, pages 445–450, Washington, DC, USA, 2005. IEEE Computer Society.
- [LMV02] N. Labroche, N. Monmarche, and G. Venturini. A new clustering algorithm based on the chemical recognition system of ants. *Harmelen, F. van (Ed.) Proc. of the 15th European Conference on Artificial Intelligence, Lyon, France*, pages 345–349, 2002.
- [LZZ08] Diana Lungeanu, Daniela Zaharie, and **Flavia Zamfirache**. Influence of missing values handling on classification rules evolved from medical data. In *Industrial Conference on Data Mining - Posters and Workshops'08*, pages 86–95, 2008.
- [MAS⁺99] M. Maheswaran, S. Ali, H. J. Siegel, D. Hensgen, and R. F. Freund. Dynamic mapping of a class of independent tasks onto heterogeneous computing systems. *Journal of Parallel and Distributed Computing*, 59:107–131, 1999.
- [MM05] R. Mendes and A. Mohais. Dynde: a differential evolution for dynamic optimization problems. *Proc. of CEC2005. Edinburgh, Scotland, September 2-5*, pages 2808–2815, 2005.
- [OAT⁺99] M. Ohsaki, H. Abe, S. Tsumoto, H. Yokoi, and T. Yamaguchi. Discovering interesting prediction rules with a genetic algorithm. *Proc. of of Congress on Evolutionary Computing (CEC 1999)*, pages 1322–1329, 1999.
- [PKN10] Andrew J. Page, Thomas M. Keane, and Thomas J. Naughton. Multi-heuristic

- dynamic task allocation using genetic algorithms in a heterogeneous distributed system. *J. Parallel Distrib. Comput.*, 70:758–766, July 2010.
- [RL04] G. Ritchie and J. Levine. A hybrid ant algorithm for scheduling independent jobs in heterogeneous computing environments. *Proceedings the 23rd Workshop of the UK Planning and Scheduling Special Interest Group*, pages 1–7, 2004.
- [SP95] R. Storn and K. Price. Differential evolution - a simple and efficient adaptive scheme for global optimization over continuous spaces. *Technical Report TR-95-012, ICSI, March*, 1995.
- [XA10] Fatos Xhafa and Ajith Abraham. Computational models and heuristic methods for grid scheduling problems. *Future Gener. Comput. Syst.*, 26:608–621, April 2010.
- [Xha07] F. Xhafa. A hybrid evolutionary heuristic for job scheduling on computational grids. In Ajith Abraham, Crina Grosan, and Hisao Ishibuchi, editors, *Hybrid Evolutionary Algorithms*, volume 75 of *Studies in Computational Intelligence*, pages 269–311. Springer Berlin / Heidelberg, 2007. 10.1007/978-3-540-73297-6-11.
- [ZF11a] **Flavia Zamfirache** and Marc Frîncu. Automatic selection of scheduling algorithms based on classification models. *STUDIA*, (2), 2011.
- [ZF11b] **Flavia Zamfirache** and Marc Frîncu. Automatic selection of scheduling algorithms based on classification models. *Proc. of International Conference on Knowledge Engineering: Principles and Techniques, Cluj-Napoca*, 2011.
- [ZFZ11] **Flavia Zamfirache**, Marc Frîncu, and Daniela Zaharie. Population-based meta-heuristics for tasks scheduling in heterogeneous distributed systems. In *Proceedings of the 7th international conference on Numerical methods and applications, NMA'10*, pages 321–328, Berlin, Heidelberg, 2011. Springer-Verlag.
- [ZLZ08] Daniela Zaharie, Diana Lungeanu, and **Flavia Zamfirache**. Interactive search of rules in medical data using multiobjective evolutionary algorithms. In *Proceedings of the 2008 GECCO conference companion on Genetic and evolutionary computation, GECCO '08*, pages 2065–2072, New York, NY, USA, 2008. ACM.
- [ZPNZ11] D. Zaharie, L. Perian, V. Negru, and F. Zamfirache. Evolutionary pruning of non-nested generalized exemplars. *Proc. of 6th IEEE International Symposium on Applied Computational Intelligence and Informatics (SACI 2011) to be held on May 19-21, 2011 in Timisoara, Romania*, pages 57–62, 2011.

- [ZZ05a] Daniela Zaharie and **Flavia Zamfirache**. Ant-based clustering of medical data. *HCMC 2005, First East European Conference on Health Care Modelling and Computation* (eds: F. Gorunescu, E. El-Darzi, M. Gorunescu), pages 332–343, 2005.
- [ZZ05b] Daniela Zaharie and **Flavia Zamfirache**. Dealing with noise in ant-based clustering. In *Congress on Evolutionary Computation'05*, pages 2395–2401, 2005.
- [ZZ06] Daniela Zaharie and **Flavia Zamfirache**. Diversity enhancing mechanisms for evolutionary optimization in static and dynamic environments. *Proc. of 3rd Romanian-Hungarian Joint Symposium on Applied Computational Intelligence*, pages 460–471, 2006.
- [ZZC10a] **Flavia Zamfirache**, Daniela Zaharie, and Ciprian Craciun. Evolutionary task scheduling in static and dynamic environments. *Proc. ICCO-CONTI, Timisoara, Romania*, pages 619–624, 2010.
- [ZZC10b] **Flavia Zamfirache**, Daniela Zaharie, and Ciprian Craciun. Nature inspired meta-heuristics for task scheduling in static and dynamic computing environments. *Scientific Buletin of Politehnica University of Timisoara, Romsnis, BS-UPT TACCS*, 55(69)(3):133–142, 2010.
- [ZZN⁺07] Daniela Zaharie, **Flavia Zamfirache**, Viorel Negru, Daniel Pop, and Horia Popa. A comparison of quality criteria for unsupervised clustering of documents based on differential evolution. *Proc. of International Conference on Knowledge Engineering: Principles and Techniques*, pages 114–121, 2007.