

Opera suspicionată (OS)
Suspicious workOpera autentică (OA)
Authentic work

OS	FILIMON, Maria Viorela. Securizarea accesului, interogarea și sincronizarea bazelor de date distribuite utilizând tehnici de inteligență artificială. Teză de doctorat, Universitatea Tehnică din Cluj-Napoca. 2011. Conducător științific: Prof.VĂLEAN, Honoriu (Universitatea Tehnică Cluj-Napoca). Comisia de evaluare a tezei: Prof.NEDEVSCHI, Sergiu (Universitatea Tehnică Cluj-Napoca), Prof.Univ. PROȘTEAN, Octavian (Universitatea Politehnică Timișoara), Prof.ILEANĂ, Ioan (Universitatea "1 Decembrie 1918" Alba Iulia), Prof.ABRUDEAN, Mihail (Universitatea Tehnică Cluj-Napoca).
OA	Rotar, C. <i>Modele naturale si algoritmi evolutivi</i> , Cluj-Napoca: Editura Accent. 2008.

Incidența minimă a suspiciunii / Minimum incidence of suspicion

p.46:5 – p.54:26

p.12:3 – p.26:15

Fișa întocmită pentru includerea suspiciunii în Indexul Operelor Plagiate în România de la
Sheet drawn up for including the suspicion in the Index of Plagiarized Works in Romania at
www.plagiate.ro



UNIVERSITATEA TEHNICĂ
DIN CLUJ-NAPOCA

FACULTATEA DE AUTOMATICĂ ȘI CALCULATOARE

Maria Viorela Filimon

TEZĂ DE DOCTORAT

SECURIZAREA ACCESULUI, INTEROGAREA ȘI SINCRONIZAREA BAZELOR DE DATE DISTRIBUITE UTILIZÂND TEHNICI DE INTELIGENȚĂ ARTIFICIALĂ

Conducător științific,
Prof.dr.ing. **Honoriu VĂLEAN**

Comisia de evaluare a tezei de doctorat:

- PREȘEDINTE: - Prof.dr.ing. *Sergiu Nedevschi* - Decan, Facultatea de Automatică și
Calculatoare, Universitatea Tehnică din Cluj-Napoca;
- MEMBRI: - Prof.dr.ing. *Honoriu Vălean* - Conducător științific, Facultatea de
Automatică și Calculatoare, Universitatea Tehnică din Cluj-Napoca;
- Prof.dr.ing. *Octavian Proștean* - Referent, Universitatea „Politehnica”
din Timișoara;
- Prof.dr. ing. *Ioan Ileană* - Referent, Universitatea „1 Decembrie 1918”
din Alba Iulia;
- Prof.dr.ing. *Mihail Abrudean* – Referent, Facultatea de Automatică și
Calculatoare, Universitatea Tehnică din Cluj-Napoca.

CUPRINS

1. Introducere	5
2. Descrierea procesului de extragere a cunoștințelor din baze de date	7
2.1. Preprocesarea datelor.....	10
2.1.1. Curățarea datelor.....	11
2.1.2. Integrarea datelor.....	12
2.1.3. Transformarea datelor.....	13
2.1.4. Reducția datelor.....	13
2.1.5. Discretizarea datelor.....	16
2.2. Metode pentru extragerea cunoștințelor.....	18
2.2.1. Clasificarea / regresia.....	18
2.2.2. Clusterizarea.....	28
2.2.3. Analiza dependențelor dintre legături.....	32
2.2.4. Identificarea anomaliilor și deviațiilor.....	34
2.3. Metode de evaluare a modelelor descoperite.....	35
2.3.1. Curbe de învățare.....	35
2.3.2. Costul și acuratețea clasificării.....	36
2.3.3. Curbe ROC.....	37
2.3.4. Comparația statistică a performanțelor clasificării.....	38
3. Baze de date distribuite versus baze de date centralizate	39
4. Metode inteligente de securizare a bazelor de date distribuite	43
4.1. Reguli de decizie.....	44
4.2. Algoritmi evolutivi.....	46
4.3. Clasificatori bazați pe reguli de decizie și algoritmi evolutivi pentru detecția intrușilor.....	54
4.4. Proiectarea unor clasificatori bazați pe reguli de decizie și algoritmi evolutivi cu control endocrin pentru detecția intrușilor. Dezvoltarea unor metode de clusterizare bazate pe densitate.....	57
5. Metode inteligente de interogare și sincronizare a bazelor de date distribuite	65
5.1. Mașini cu suport vectorial. Considerații matematice.....	65
5.1.1. Mașini cu suport vectorial pentru clasificarea binară.....	65
5.1.2. Clasificarea în mai multe clase.....	72
5.1.3. Optimizarea minimă secvențială.....	73
5.2. Clasificarea Cost-Sensitive utilizând mașini cu suport vectorial.....	74
5.3. Proiectarea unor metaclasificatori utilizând mașini cu suport vectorial.....	75
6. Rezultate experimentale	79
6.1. Evaluarea metodelor inteligente de securizare pentru baze de date distribuite	79
6.1.1. Descrierea bazelor de date distribuite utilizate.....	79
6.1.2. Lansarea atacurilor în rețelele considerate.....	82
6.1.3. Preprocesarea datelor colectate.....	91
6.1.4. Evaluarea clasificatorului bazat pe reguli de decizie (RIPPER) și algoritmi genetici (AG).....	103
6.1.5. Evaluarea clasificatorului bazat pe reguli de decizie (RIPPER) și algoritmi evolutivi cu control endocrin (ECEA) dezvoltat.....	109

6.2. Evaluarea metodelor inteligente de interogare și sincronizare pentru baze de date distribuite.....	115
6.2.1. Descrierea bazelor de date utilizate.....	115
6.2.2. Replicarea și sincronizarea bazelor de date distribuite considerate.....	115
6.2.3. Preprocesarea datelor.....	124
6.2.4. Găsirea subsetului optim de atribute.....	127
6.2.5. Influența parametrului de complexitate al clasificatorului SMO.....	128
6.2.6. Influența gradului nucleului polinomial.....	130
6.2.7. Evaluarea metaclasificatorului Cost-Sensitive utilizând mașini cu suport vectorial.....	131
6.2.8. Evaluarea metaclasificatorului dezvoltat.....	135
7. Aplicații dezvoltate care utilizează metodele propuse.....	141
7.1. Aplicația MEDIS.....	141
7.2. Aplicația Monitorizarea parametrilor vântului.....	145
7.3. Devierea traficului de rețea în cazul detecției unor atacuri.....	146
8. Concluzii și contribuții proprii.....	149
8.1. Concluzii.....	149
8.2. Contribuții proprii ale lucrării.....	151
8.3. Direcții viitoare de cercetare.....	154
Bibliografie.....	155

1. Introducere

Lumea modernă trece printr-o transformare fundamentală, de la societatea industrială care a marcat secolul al XX-lea, la societatea informațională a secolului al XXI-lea. Acest proces dinamic a produs schimbări fundamentale în toate aspectele vieții noastre (diseminarea cunoștințelor, interacțiunea socială, economie și practici de afaceri, angajament politic, media, sănătate, odihnă și agrement), dând naștere unei societăți care se hrănește cu informație. Din nefericire, forma recreativă și de divertisment a informației, vine mai degrabă spre noi într-o formă crudă, de date înregistrate în depozite de mari dimensiuni.

Pentru că dimensiunea datelor stocate în depozite crește (cercetătorii de la Universitatea din Berkeley au calculat că în fiecare an sunt generate și salvate aproximativ un exabyte de date), ne găsim copleșiți de date, și realizăm că ne lipsește puterea procesării acestora pentru a extrage informații folositoare din ele. Chiar și cu ajutorul calculatorului, echipat cu cele mai recente sisteme de gestiune a bazelor de date și cu tehnologii noi de comunicare, suntem neajutorați când ne confruntăm cu diversitatea și complexitatea pe care această lume o prezintă. Ritmul rapid de creștere, cantitatea enormă de date, adunate și stocate în mari și numeroase baze de date, a depășit cu mult capacitatea noastră umană de înțelegere. Abundența de date, cuplată cu nevoia de putere a instrumentelor de analiză a datelor a fost descrisă ca o bogăție de date, dar ca o sărăcie a informațiilor.

Ca și rezultat, datele colectate în mari baze de date devin „morminte de date”, date arhivate care arareori sunt vizitate. În consecință, deciziile importante sunt bazate nu pe informațiile bogate stocate în arhive de date, ci mai degrabă pe intuiție, pe factorul de decizie a acestuia. În plus, luându-se în considerare tehnologiile sistemelor expert, care de obicei se bazează pe utilizatori sau experți, se introduc manual cunoștințele în baze. Din nefericire, această procedură este predispusă deviațiilor și erorilor fiind o mare consumatoare de timp și costuri. Instrumentele de extragere a cunoștințelor efectuează analiza datelor și pot descoperi modele de date importante, care contribuie foarte mult la strategiile de afaceri, baze de cunoștințe, cercetări științifice și medicale. Adâncirea diferențelor dintre date și informații cheamă pentru o dezvoltare sistematică a instrumentelor de extragere a cunoștințelor, care vor transforma mormintele de date în „pepite” ale cunoașterii.

O notă recentă a Gartner Group Advanced Technology Research, a poziționat extragerea cunoștințelor din date și inteligența artificială în topul celor cinci domenii cheie ale tehnologiei, și că acest lucru „va avea un impact major într-o arie largă de industrii în următorii 3-5 ani”. Gartner, de asemenea a prezentat arhitecturile paralele și data mining ca două dintre cele 10 noi tehnologii în care companiile vor investi în următorii 5 ani. În conformitate cu o notă recentă a Gartner HPC Research, „cu înaintarea rapidă în captarea datelor, transmisia și stocarea acestora, utilizatorii marilor sisteme vor avea din ce în ce mai multă nevoie să pună în aplicare metode noi și inovatoare pentru extragerea cunoștințelor, care angajează procesarea masivă paralelă (*massively parallel processing - MPP*) a sistemelor pentru a crea noi surse de avantaj în afaceri”.

Utilizând aceste noțiuni, acuratețea se definește ca fiind raportul dintre numărul instanțelor clasificate corect și numărul total al instanțelor, atât clasificate corect, cât și clasificate eronat [30]:

$$\text{acuratețea} = \frac{AP + AN}{AP + AN + FP + FN} \quad (41)$$

4.2. Algoritmi evolutivi

Istoric

Primele semnale privind posibilitatea simulării proceselor evolute cu ajutorul calculatorului au fost înregistrate în anii '50. Unul dintre pionierii acestui domeniu, biologul Alex S. Fraser (1923-2002), prin lucrările sale (ex. "*Simulation of genetic systems by automatic digital computers*" – 1957), a reușit să aducă în atenția cercetătorilor ideea simulării selecției artificiale a organismelor. Interesul asupra fenomenelor evoluției, privitye din perspectivă computațională, a crescut în decadele următoare, culminând cu lucrarea de referință a lui John Holland („*Adaptation in Natural and Artificial Systems*”, 1975) în care este descris algoritmul genetic standard. Holland elaborează *Teorema schemelor* în încercarea de a explica forța algoritmilor genetici în rezolvarea problemelor de căutare și optimizare. În paralel cu dezvoltarea algoritmilor genetici se conturează și celelalte direcții majore ale Calculului Evolutiv: strategiile evolute, programarea genetică, programarea evolutivă. Cu toate acestea, cea mai prolifică direcție rămâne cea a algoritmilor genetici, popularitatea lor fiind justificată de simplitatea și succesul înregistrat în rezolvarea multor probleme dificile.

Un alt punct de referință în istoria algoritmilor inspirați de fenomene naturale îl constituie lucrările lui David E. Goldberg (vezi de ex. *Genetic Algorithms in Search, Optimization and Machine Learning*, 1989) în care sunt revăzute și formulate principiile algoritmilor genetici, dintr-o perspectivă matură asupra noului domeniu conturat. Goldberg furnizează o ipoteză interesantă dar controversată asupra abilității algoritmului genetic de a găsi soluțiile bune a multor probleme. Cu toate că Goldberg afirmă că *teorema schemelor* (Holland, 1973) susține *ipoteza blocurilor constructive*, teoria formulată de Goldberg nu se verifică în totalitate și nu reușește să explice potențialul algoritmilor genetici în rezolvarea problemelor.

Studiul teoretic și aplicabilitatea algoritmilor genetici în rezolvarea multor probleme grele se datorează cercetătorilor amintiți, astfel încât aceștia sunt recunoscuți pe drept ca fiind “părinții” Algoritmilor Genetici.

Modelul natural al algoritmilor genetici

Principala sursă de inspirație a algoritmilor genetici este în mod cert teoria evoluției naturale enunțată de Charles Darwin (1809-1882). Cu toate acestea, descoperirile din genetică au influențat, de asemenea, structura algoritmilor genetici. Teoria Darwinistă reprezintă un ansamblu de observații asupra evoluției privită la nivel macroscopic și asupra identificării factorilor responsabili de acest fenomen. Adaptarea unei populații este interpretată prin două principii majore: *selecția naturală* – supraviețuirea celor mai adaptați indivizi și *mutația* – variație a caracteristicilor indivizilor apărută ca răspuns la mediul înconjurător. Cercetările și lucrările lui Darwin (ex. *Originea speciilor*, 1959) nasc controverse și astăzi, însă afirmațiile conform cărora evoluția este călăuzită de principiul supraviețuirii celor mai performanți indivizi și a diversificării prin mutațiile apărute, sunt fundamentele algoritmilor evolutivi.

Teoria neo-darwinistă, sau *teoria eredității* este un alt punct de vedere asupra evoluției, de data acesta interpretat prin transferul de material genetic de la părinți la descendenți. La sfârșitul secolului al XIX-lea, un preot și om de știință austriac, Gregor Mendel (1822-1884) descoperă mai mult sau mai puțin accidental teoria eredității prin experimentele efectuate în vederea încrucișării speciilor diferite de mazăre. Acesta lasă posterității câteva însemnări în caietele sale de observații și o lucrare publicată care a fost ignorată la vremea

aceea. Mai târziu, aceste observații enunțate de Mendel sunt recunoscute ca fiind legile eredității și autorul acestora este considerat părintele geneticii moderne. Multe aspecte ale algoritmilor genetici sunt împrumutate din genetică. Dintre acestea remarcăm moștenirea trăsăturilor (genelor) părinților de către descendenți, codificarea indivizilor prin secvențe de *gene*, dar și caracterul aleator al aportului trăsăturilor provenite de la părinți.

Algoritmii genetici nu respectă strict ingredientele celor două teorii amintite, iar terminologia împrumutată nu acoperă semnificația originală; aceștia se formează ca o rețetă practică inspirată de evoluție și ereditate în care factorul uman a strecurat secvențe de calcule și parametri artificiali în scopul obținerii unor rezultate mai bune. Privit la nivel superficial, un algoritm genetic simulează evoluția unei populații, însă, în esență, fenomenul evoluției artificiale este ghidat de funcții matematice și controlat prin parametri suplimentari pentru a genera o soluție cât mai apropiată de optimul problemei. Avantajul imediat al algoritmilor genetici constă în gradul mare de generalitate și paleta largă a problemelor abordate. Originalitatea, simplitatea și sursa inedită de inspirație poate justifica parțial atenția acordată dezvoltării acestor metaeuristici. Succesul real înregistrat în probleme de dificultate ridicată este argumentul forte al algoritmilor genetici [36].

Descrierea algoritmului genetic standard

Pe scurt, un algoritm genetic operează cu o mulțime de indivizi (denumiți în mod uzual *cromozomi*) asupra cărora se aplică *operatorii genetici*. Fiecare individ reprezintă o soluție posibilă din spațiul de căutare. Maniera de codificare a unei soluții este binară, reală sau specifică, în funcție de opțiunea programatorului și de contextul problemei. Mulțimea de cromozomi formează o *populație*. Evoluția populației este condusă de două elemente: *operatorii genetici* și *funcția de evaluare* a calității cromozomilor. *Selecția*, *încrucișarea* și *mutația* sunt operatorii genetici uzuali, fiind inspirați de procesele naturale care stau la baza evoluției.

Indiferent de natura problemei considerate, se poate construi cel puțin o funcție de evaluare a soluțiilor posibile. Prin aceasta, fiecărui individ al populației i se atribuie o valoare numerică reprezentând performanța sa în raport cu cerințele problemei. Ulterior, măsura calității indivizilor populației curente este folosită în procesul de *selecție* a acelor indivizi, denumiți *părinți*, asupra cărora se aplică operatorii de încrucișare și mutație pentru obținerea noii generații. Principiul este simplu: cu cât părinții selectați sunt mai performanți, cu atât șansa ca descendenții obținuți să fie calitativ superiori este mai mare. Procesul se reia având ca populație curentă noua generație de cromozomi. Se observă de-a lungul evoluției populației o creștere a calității indivizilor săi. După un număr considerabil de generații soluția globală a problemei poate fi aproximată suficient de bine printr-un individ al ultimelor generații.

Conceperea unui algoritm genetic de rezolvare a unei probleme concrete presupune evidențierea următoarelor componente:

1. individul
2. populația
3. funcția de evaluare
4. selecția
5. operatorii de variație (încrucișare și mutație)
6. condiția de oprire a algoritmului

1) Individul

Analizând specificațiile problemei, se poate identifica spațiul de căutare. Orice punct al acestuia constituie o soluție posibilă. În funcție de specificul soluției posibile, se poate determina o manieră inspirată de codificare numerică sau nenumerică a acesteia. Codificarea unei soluții formează un *cromozom* și identifică un individ al

populației curente. Fiecare cromozom este format dintr-un șir de valori ale unui alfabet dat. Valoarea de pe o poziție oarecare a șirului codificării se numește *genă*. Alfabetul A al valorilor genelor este stabilit în prealabil și poate fi format din:

- valori binare: $A = \{0,1\}$
- valori reale: $A \subseteq R$
- alt alfabet (*codificare specifică*)

Lungimea șirului de gene din codificarea cromozomului poate fi constantă sau variabilă și depinde de numărul de trăsături definitorii ale unei soluții posibile a problemei:

Fie c un cromozom oarecare și A alfabetul genelor din reprezentarea sa. Cromozomul c se descrie prin vectorul de n elemente:

$$c = (\alpha_1, \alpha_2, \dots, \alpha_n), \quad (42)$$

$$\text{unde } \alpha_i \in A, \forall i \in \{1, 2, \dots, n\}. \quad (43)$$

În funcție de alfabetul genelor, codificarea cromozomială se clasifică în:

- codificare binară, $A = \{0,1\}$
- codificare reală, $A \subseteq R$
- codificare specifică – alta decât cea binară sau reală.

2) Populația

O mulțime finită de cromozomi alcătuiește populația. Dimensiunea populației (numărul de indivizi care o formează) este, în general, o valoare constantă stabilită în prealabil și reprezintă un parametru important al algoritmului dezvoltat. Într-o abordare mai flexibilă, dimensiunea populației poate fi variabilă, un exemplu în acest sens ar fi descreșterea numărului de indivizi de-a lungul evoluției, o dată cu creșterea performanței acestora, fapt care ar ușura determinarea soluțiilor finale din ultima generație produsă. Nu este exclusă utilizarea mai multor populații cooperante sau a unei populații suplimentare cu rol de memorie, în care sunt reținute soluții bune găsite în cadrul generațiilor intermediare. Modelul standard al algoritmului genetic presupune exploatarea unei populații unice de dimensiune constantă. Formal, o astfel de populație se descrie astfel:

$$P = \{c_1, c_2, \dots, c_m\}, \quad (44)$$

unde: $\forall j \in \{1, 2, \dots, m\}$, c_j reprezintă un *cromozom*.

Prima generație se numește *populație inițială* și construirea acesteia se face în principal prin generarea aleatoare a unor soluții posibile în domeniul de căutare. Dacă sunt cunoscute zone ale spațiului care conțin soluții bune ale problemei, populația inițială poate fi îmbogățită cu indivizi care codifică puncte din aceste zone. În această manieră se accelerează convergența populației, oferindu-se șansa ca încă de la primele generații să fie obținute soluții optime ale problemei. Există puține situații în care procedeul descris mai sus poate fi aplicat deoarece din specificațiile problemei rareori putem deduce zonele promițătoare ale spațiului de căutare. Modelul standard al algoritmului genetic presupune generarea aleatoare a populației inițiale.

Cu cât dimensiunea populației este mai mare, cu atât șansa obținerii rapide a unei bune aproximări a soluțiilor problemei crește. De asemenea, distribuția indivizilor populației joacă un rol important în economia construirii algoritmului genetic: o distribuție echilibrată a populației inițiale mărește substanțial viteza de identificare a zonelor promițătoare ale spațiului de căutare. Acest fapt aduce cu sine o convergență mai bună înspre soluțiile problemei.

În concluzie, o populație de dimensiune suficient de mare și având un grad de diversitate considerabil reprezintă unul dintre factorii majori ce conduc la obținerea soluțiilor dorite.

3) Funcția de evaluare

Una dintre condițiile fundamentale ale funcționării eficiente ale unui algoritm genetic o reprezintă alegerea inspirată a funcției de evaluare. În cazul unei probleme de optimizare numerică (ex. determinarea maximului/minimului global al unei funcții matematice), funcția de evaluare se poate alege ca fiind chiar funcția criteriu sau o funcție construită pe baza funcției criteriu.

Pentru o mai bună înțelegere prezentăm în continuare un exemplu:

Fie $f : D \rightarrow R$, $D = [-\pi, \pi]$, dată de formula:

$$f(x) = \sin(x) \quad (45)$$

Se dorește determinarea maximului funcției f .

Criteriul care ghidează căutarea soluției globale este maximizarea funcției date, respectiv, determinarea punctului x pentru care $\sin(x)$ este maxim.

Orice valoare reală din intervalul $[-\pi, \pi]$ poate fi reprezentată în limbaj binar, cu o anumită precizie impusă. Un cromozom al populației corespunzător unei valori reale x va fi construit ca un vector de cifre binare ale codificării punctului x . Se impune observația că o abordare mai flexibilă permite ca reprezentarea cromozomială a unui punct x al spațiului de căutare să fie dată de un vector particular cu o singură componentă. Unica genă a codificării reprezintă chiar valoarea reală a punctului corespunzător.

În acest caz particular, funcția de evaluare este dată de însăși funcția criteriu. Fie cromozomul c reprezentarea (binară sau reală) a punctului x . Funcția de evaluare este construită astfel:

$$\text{performanta}(c) = \sin(x) \quad (46)$$

Valorile calculate ale funcției de performanță califică și diferențiază indivizii populației. Astfel, este posibilă identificarea indivizilor performanți ai populației, cei cărora li se oferă ulterior o șansă mai mare de a produce descendenți.

Există situația în care este dificil de exprimat matematic criteriul de optimizare al problemei date. În aceste situații, utilizatorul va construi funcții de evaluare prin care se va încerca calificarea cât mai coerentă a soluțiilor candidat, fapt pentru care rezultatele oferite de algoritmul dezvoltat sunt dependente de factorul uman. Mai mult, probleme de optimizare de dificultate ridicată sunt cele în care se dorește determinarea optinelor multiple ale unei funcții (optimizare multimodală) sau cele în care se impun anumite restricții (optimizare cu restricții). Pentru rezolvarea acestor probleme, funcția de evaluare va fi atent construită pentru a respecta restricțiile impuse, respectiv, pentru a permite populației să convergă înspre toate optimele funcției criteriu. Un alt caz în care construirea funcției de evaluare necesită un efort mai mare este cel al problemelor de optimizare multiobiectiv. În discuția referitoare la funcția de evaluare ne-am referit strict la probleme de optimizare. Justificarea acestui fapt rezidă din afirmația că orice problemă de căutare a soluțiilor în spațiul soluțiilor posibile, căutare ghidată de unul sau mai multe criterii, poate fi reinterpretată ca o problemă de optimizare: din mulțimea soluțiilor posibile se caută acelea care optimizează cel mai bine criteriile problemei. Deși principala aplicabilitate a algoritmilor genetici constă în rezolvarea problemelor de optimizare, afirmația precedentă ne permite extinderea ariei de aplicabilitate a algoritmilor genetici și înspre probleme de alt gen.

4) Selecția

Evoluția populației, obținerea generațiilor de indivizi mai performanți decât predecesorii lor, este cauzată de operatorul de selecție și de operatorii de variație (recombinarea și mutația) aplicați.

Unul dintre factorii majori ai evoluției este selecția naturală. O interpretare simplificată a acestui fenomen ne permite înțelegerea procesului prin care o populație poate converge înspre anumite puncte ale spațiului de căutare, respectiv, înspre soluțiile unei probleme date. Pornind de la observația acceptată că indivizii performanți ai unei specii dau naștere unor descendenți asemănători și de performanțe apropiate și, în schimb, indivizii slab calificați produc indivizi noi ale căror performanțe sunt, în general, slabe, putem considera două scenarii:

- în primul scenariu, indivizii unei populații suferă încrucișări și mutații pur aleatoare, fără a se aplica principiul selecției;
- al doilea scenariu implică o selecție prealabilă a părinților noii generații, selecție făcută pe baza valorilor de performanță a indivizilor populației curente.

Cele două populații ale scenariilor descrise se monitorizează de-a lungul evoluției. Se observă că, dacă în primul caz absența selecției induce o stagnare în procesul evoluției populației înspre indivizi de calitate superioară, în a doua situație, selecția favorizează o creștere accelerată a performanței medii a generațiilor produse. Comparând ultimele generații ale ambelor scenarii putem deduce importanța selecției în evoluția unei populații.

Implementarea operatorului de selecție se face în diverse maniere. Unul dintre cei mai populari algoritmi de selecție este cel al *selecției proporționale* inspirat de algoritmul ruletei. Selecția proporțională presupune alegerea părinților noii generații în mod proporțional cu valoarea de performanță a acestora. Altfel spus, probabilitatea de selectare a unui individ performant este mai mare decât a unuia slab calificat, și această probabilitate este direct proporțională cu valoarea performanței individului. Acest procedeu favorizează o convergență rapidă a populației, însă în unele cazuri prezintă dezavantaje majore, fapt pentru care se optează pentru implementarea altor variante de selecție.

5) Operatorii de variație (încrucișare și mutație)

Selecția propriu-zisă nu este suficientă pentru a induce dinamica unei populații. Indivizii selectați vor suferi modificări prin operatorii genetici de variație pentru a permite obținerea unor descendenți diferiți, respectiv pentru a permite o explorare bună a spațiului de căutare. Principalii operatori de variație sunt încrucișarea și mutația.

Varianta standard a operatorului de încrucișare se aplică asupra a doi părinți selectați anterior și produce unul sau doi descendenți care moștenesc trăsăturile părinților în proporții diferite. Prezentăm în continuare un operator de încrucișare aplicabil pentru codificarea binară. Tipul operatorului este de 2:2, respectiv, din 2 părinți se produc 2 descendenți. Datorită specificului său, operatorul descris mai jos este cunoscut sub denumirea de *încrucișare cu un punct de tăietură*.

Fie c_1 și c_2 doi indivizi oarecare ai populației curente:

$$c_1 = (\alpha_1, \alpha_2, \dots, \alpha_n) \quad (47)$$

$$c_2 = (\beta_1, \beta_2, \dots, \beta_n) \quad (48)$$

Se generează aleator o valoare naturală $t \in \{1, 2, \dots, n\}$, având semnificația punctului de tăietură. Ambii părinți se vor diviza, obținându-se secvențele:

$$c_1^1 = (\alpha_1, \alpha_2, \dots, \alpha_t), \quad c_1^2 = (\alpha_{t+1}, \alpha_{t+2}, \dots, \alpha_n) \quad (49)$$

$$c_2^1 = (\beta_1, \beta_2, \dots, \beta_t), \quad c_2^2 = (\beta_{t+1}, \beta_{t+2}, \dots, \beta_n) \quad (50)$$

Descendenții d_1 și d_2 ai părinților c_1 și c_2 se obțin prin concatenarea secvențelor obținute:

$$d_1 = c_1^1 + c_2^2 \quad (51)$$

$$d_2 = c_2^1 + c_1^2 \quad (52)$$

Observație: Operatorul + are în această descriere formală semnificația concatenării secvențelor.

Se obțin astfel noii indivizi:

$$d_1 = (\alpha_1, \dots, \alpha_t, \beta_{t+1}, \dots, \beta_n) \quad \text{și} \quad (53)$$

$$d_2 = (\beta_1, \dots, \beta_t, \alpha_{t+1}, \dots, \alpha_n) \quad (54)$$

Prin aplicarea încrucișării, noii indivizi moștenesc genele părinților, fapt pentru care aceștia sunt asemănători indivizilor care îi produc. Se întâlnește deseori situația în care o anumită soluție a spațiului de căutare nu poate fi obținută prin aplicarea operatorului de încrucișare descris. Există, de asemenea, suficiente argumente în favoarea aplicării altor tipuri de operatori de încrucișare (exemplu: *încrucișarea cu puncte multiple de tăietură*). Pentru asigurarea unei bune explorări a spațiului de căutare, un alt operator genetic este introdus: *mutația*. Teoria evoluției naturale menționează faptul că selecția nu este unicul factor responsabil al fenomenului evoluției. Mutăția naturală, definită prin micile modificări la nivelul codificării cromozomiale, este cea care, în combinație cu mecanismul selecției, asigură adaptarea la mediu. Implementarea procesului natural al mutației se face prin *operatorul de mutație*, care, la rândul său, cunoaște variate forme descrise în lucrări de specialitate.

Prezentăm în continuare o variantă simplă a operatorului de mutație. Acest operator este de tipul 1:1, respectiv, dintr-un părinte selectat se obține un unic descendent prin alterarea valorii unei singure gene. De asemenea, aplicabilitatea operatorului descris este strict limitată la *codificarea cromozomială binară*.

Fie c un cromozom al populației curente:

$$c = (\alpha_1, \alpha_2, \dots, \alpha_k, \dots, \alpha_n) \quad (55)$$

Se generează aleator o poziție k din secvența binară a codificării.

Individul mutant $c' = (\alpha_1, \alpha_2, \dots, \beta_k, \dots, \alpha_n)$ se obține prin copierea nealterată a genelor părintelui c , cu excepția genei de pe poziția k . Valoarea genei mutate se obține prin inversare:

Dacă $\alpha_k = 1$ atunci

$$\beta_k \leftarrow 0$$

Altfel

$$\beta_k \leftarrow 1$$

SfDacă

Alegerea operatorilor folosiți joacă un rol decisiv în economia dezvoltării algoritmilor genetici. Utilizatorul optează pentru anumite forme ale operatorilor genetici pe baza unor criterii care sunt, de cele mai multe ori, subiective. Specificațiile problemei sunt cele care ne pot oferi informațiile necesare pentru a decide maniera de codificare a indivizilor. Odată stabilită forma de codificare (*binară, reală, specifică*), paleta operatorilor se restrânge semnificativ. Cu toate acestea, diversitatea procedurilor de încrucișare și mutație, clasificate pe tipul codificării pre-fixate, fac dificilă alegerea inspirată a celor mai eficiente variante. Preferințele se îndreaptă înspre acei operatori care se dovedesc a fi eficienți în rezolvarea unor probleme similare. Lipsa constrângerilor

legate de implementarea operatorilor genetici, permite utilizatorului să își construiască operatori specifici problemei sau clasei de probleme pe care le abordează.

6) Condiția de oprire a algoritmului genetic

În esență, algoritmul genetic este o procedură repetitivă prin care o populație de soluții posibile, prin aplicarea operatorilor genetici (selecție, încrucișare, mutație) produce noua generație. Procesul se reia pentru noua populație obținută până când o condiție de terminare este îndeplinită. Condiția de terminare a algoritmului se referă în general la atingerea unui număr maxim prestabilit de generații. Există și situații în care utilizatorul poate să impună o cu totul altă condiție de încheiere a evoluției, condiție independentă de numărul populațiilor generate. Un exemplu în acest sens ar fi înregistrarea unui anumit grad de uniformitate în cadrul populației. În situația în care indivizii populației devin asemănători putem considera că diversitatea slabă a populației nu mai permite o explorare eficientă a spațiului de căutare. Mai mult, putem concludiona că indivizii populației s-au stabilit în zona cea mai promițătoare, corespunzătoare soluției globale. Experimentele arată că nu întotdeauna această concluzie este adevărată. Populația poate fi atrasă de un punct de optim local, producându-se fenomenul de convergență prematură. Convergența prematură este unul dintre fenomenele nedorite ale unui algoritm genetic și evitarea acestui neajuns comportă câteva modificări în structura algoritmului: alegerea unei alte scheme de selecție, reinițializarea aleatoare unei secțiuni din populație sau accentuarea mutațiilor produse. O altă soluție oferită este cea prin care sunt monitorizate ultimele k generații obținute. Pentru a nu cădea în posibilă capcană a scăderii temporare a gradului de diversitate a populației putem considera un număr mai mare de generații analizate din punct de vedere al asemănării indivizilor. În această manieră obiectivitatea deciziei de terminare a evoluției crește substanțial. Parametrul k reprezintă în acest caz o valoare numerică naturală prestabilită.

Revenind la condiția de terminare, ca și în cazul celorlalte componente ale algoritmului genetic, alegerea acesteia se face de către proiectantul algoritmului după o atentă observare a comportamentului populației de-a lungul generațiilor. O bună experiență a proiectantului, cât și concluziile experimentelor efectuate, sunt factorii principali ai construirii unui algoritm genetic valoros.

Înainte de a furniza schema algoritmului genetic standard, este important să afirmăm că acest algoritm nu funcționează ca un șablon de rezolvare a oricărei probleme. Detalii legate de implementare, alegerea operatorilor, construirea funcției de evaluare, precum și stabilirea parametrilor implicați reprezintă aspecte care diferențiază algoritmi obținuți din punct de vedere al aplicabilității acestora pe anumite clase de probleme.

Presupunem că cele 6 elemente enumerate anterior au fost stabilite. În acest caz putem formula structura algoritmului genetic:

Schema algoritmului genetic standard:

Algoritmul AG_standard este:

Fie $P(0)$ - populația inițială; $t = 0$;

Cât timp (NOT condiție de terminare)

***Evaluare*($P(t)$)**

$P_S(t)$ = *Selectie* ($P(t)$)

$P_R(t)$ = *Încrucișare*($P_S(t)$)

$P(t+1)$ = *Mutație* ($P_R(t)$)

$$t:=t+1$$

SfCâtimp

SfAlgoritm

Notăm:

t – numărul generației curente.

P_S - mulțimea indivizilor selectați

P_R - mulțimea indivizilor obținuți prin încrucișare

Soluția finală a algoritmului genetic este dată de cel mai performant individ al ultimei generații produse. Se impune observația că această variantă de extragere a rezultatului nu este cea mai eficientă. Considerăm următorul scenariu: un individ performant, extrem de apropiat de soluția globală a problemei) este obținut într-una dintre generațiile intermediare. Principiul selecției ne conduce la supoziția că probabilitatea acestuia de a fi selectat în vederea aplicării operatorilor de variație este mare. Operatorul de încrucișare nu asigură păstrarea părinților performanți în noua generație. Astfel, individul considerat în scenariul nostru va fi distrus, neavând certitudinea că descendenții săi sunt cel puțin la fel de performanți ca și el. Observăm astfel că indivizii calificați ai generațiilor intermediare pot dispărea pe parcursul evoluției populației.

Remediarea acestui fenomen poate fi făcută extrem de simplu prin suplimentarea algoritmului inițial cu o formă de *elitism*. Prin elitism, cel mai bun, sau cei mai buni k indivizi ai generației curente vor fi copiați nemodificați în noua generație. Se evită în acest mod pierderea soluției globale dacă aceasta este obținută într-o etapă intermediară a evoluției.

Funcționarea algoritmului genetic

Intuitiv, algoritmul genetic va conduce populația de soluții posibile înspre soluțiile optime ale problemei de rezolvat. De cele mai multe ori acest fenomen decurge în mod firesc, fără anomalii, și putem spune că algoritmul converge înspre optimele problemei. Există și situația în care indivizii populației se blochează în zone sub-optimale ale spațiului de căutare, evoluția lor ulterioară fiind îngreunată sau imposibilă. În aceste situații rezultatul oferit de algoritm este eronat. Cauzele acestui anomalii sunt multiple. Una dintre cele mai frecvente surse de eșec este alegerea neinspirată a funcției de evaluare și a operatorului de selecție, rezultând pierderea diversității populației și incapacitatea de a determina soluțiile globale. Fenomenul nedorit, descris ca obstrucționarea explorării spațiului de căutare și blocarea populației în zone corespunzătoare optinelor locale, este denumit *convergență prematură*.

În esență, convergența prematură este cauzată de pierderea diversității genetice a indivizilor într-o etapă timpurie a evoluției, fapt care generează o stagnare a evoluției ulterioare. Probabilitatea ca indivizii asemănători ai populației curente să genereze descendenți diferiți este redusă, ceea ce induce o rezistență sporită la dislocarea indivizilor din zonele de blocaj.

În scopul evitării convergenței premature au fost dezvoltate diferite strategii de menținere a diversității populației. Dintre acestea enumerăm:

- stabilirea unei dimensiuni mai mari a populației;
- aplicarea unor operatori genetici care s-au dovedit benefici evoluției și menținerii diversității (ex. încrucișarea uniformă);
- strategii de împerechere care evită obținerea de indivizi similari prin aplicarea încrucișării;

- aplicarea mutației cu o mai mare probabilitate;
- evitarea folosirii operatorului de selecție proporțională care, în anumite situații, va favoriza evoluția doar a unei mici fracțiuni de indivizi din populație, conducând la omogenizarea populației.

Un alt aspect al funcționării algoritmului genetic este cel al *elitismului*. În multe situații, o soluție bună a problemei (codificată ca individ al populației) poate să apară într-una dintre generațiile intermediare; având o performanță mare, șansa de a fi selectată în vederea supunerii acesteia la acțiunea operatorilor de variație (mutație și încrucișare) este mare. Aplicarea operatorilor genetici poate să distrugă soluția performantă, fapt pentru care evoluția populației este încetinită. Remedierea acestui fenomen constă în aplicarea unor strategii de salvare a indivizilor performanți determinați în populațiile intermediare. În mod sugestiv, acești indivizi cu performanțe ridicate se numesc *elită* și procedeul prin care elita este menținută de-a lungul generațiilor se numește *elitism*. Cea mai simplă strategie de elitism constă în copierea elitei în noua generație. Acest transfer garantează protejarea soluțiilor calificate de-a lungul evoluției și o mai bună convergență a algoritmului genetic.

Funcționarea algoritmului genetic este strâns legată de două concepte duale: *explorarea* și *exploatarea* spațiului de căutare. Primul concept se referă la capacitatea populației de a acoperi spațiul soluțiilor posibile prin căutare globală și este strâns legat de proprietatea de diversitate a populației. Cel de-al doilea concept se referă la capacitatea populației de a asigura căutarea locală în zonele promițătoare ale spațiului și corespunde unei îmbunătățiri a indivizilor în vederea obținerii unei aproximări mai bune a soluțiilor finale. Cele două fenomene definite de conceptele sus-menționate sunt duale în sensul în care, accentuarea unuia conduce în mod gradual la inhibarea apariției celuilalt. Astfel, o *explorare* asiduă a spațiului de căutare, marcată prin menținerea unei populații diverse, îngreunează procesul de rafinare a soluțiilor bune și cauzează o stagnare în evoluția populației. Contrar, accentuarea *exploatării* zonelor promițătoare conduce la o explorare defectuoasă a spațiului de căutare și poate conduce la fenomenul de convergență prematură. Ambele situații descrise trebuie evitate. Pentru buna funcționare a algoritmului genetic din perspectiva celor două concepte, este necesară stabilirea unui echilibru *explorare-exploatare* obținut, în principal, prin aplicarea unui mecanism de evaluare și selecție bine ales.

Algoritmii genetici au fost folosiți atât în clasificare [64], [99], cât și în alte probleme de optimizare [100], [101], [102]. În *data mining* aceștia sunt utilizați pentru a evalua fitness-ul altor algoritmi.

4.3. Clasificatori bazați pe reguli de decizie și algoritmi evolutivi pentru detecția intrușilor

Întrucât regulile de decizie pot fi reprezentate cu ușurință ca și *bit-string*-uri de lungime fixă, se pot utiliza împreună cu algoritmi genetici cu scopul creșterii performanței clasificatorului. În contextul *data mining*-ului, populația reprezintă spațiul de căutare a soluției, unde o soluție este reprezentată de o regulă DACĂ - ATUNCI. Un cromozom este o regulă DACĂ – ATUNCI, iar o genă reprezintă o sub-condiție a regulii. Bazându-se pe noțiunea de supraviețuire a celui mai performant individ, o nouă populație este formată din cele mai bune reguli ale populației curente, dar și din urmașii acestor reguli. În mod tipic, fitness-ul unei reguli este reprezentat de acuratețea clasificării instanțelor de antrenare.

Algoritmii propus funcționează în felul următor [28], [62], [103]:

1. Încarcă setul de date care urmează a fi procesat;
2. Parsează regula DACĂ – ATUNCI în forma antecedentă și forma consecventă;
3. Calculează distribuția clasei pentru fiecare regulă;