

UNIVERSITATEA "ALEXANDRU IOAN CUZA" DIN IAȘI
FACULTATEA DE INFORMATICĂ



LUCRARE DE LICENȚĂ

Abordarea clasificării prin metoda programării genetice

propusă de
Constantin-Sebastian Stanciu

Sesiunea: *iulie, 2019*
Coordonator științific
Conf. dr. Mihaela Breabăn

UNIVERSITATEA "ALEXANDRU IOAN CUZA" DIN IAȘI
FACULTATEA DE INFORMATICĂ

Abordarea clasificării prin metoda programării genetice

Constantin-Sebastian Stanciu

Sesiunea: iulie, 2019

Coordonator științific
Conf. dr. Mihaela Breabăn

Avizat,
Îndrumător Lucrare de Licență

Titlul, Numele și prenumele

Data _____ Semnătura _____

DECLARAȚIE privind originalitatea conținutului lucrării de licență

Subsemnatul(a) _____
domiciliul în _____
născut(ă) la data de _____, identificat prin CNP _____,
absolvent(a) al(a) Universității „Alexandru Ioan Cuza” din Iași, Facultatea de
_____ specializarea _____,
promoția _____, declar pe propria răspundere, cunoscând consecințele
falsului în declarații în sensul art. 326 din Noul Cod Penal și dispozițiile Legii Educației
Naționale nr. 1/2011 art.143 al. 4 și 5 referitoare la plagiat, că
lucrarea de licență cu titlul:

_____ elaborată sub îndrumarea dl. _____ / d-na _____,
pe care urmează să o

susțină în fața comisiei este originală, îmi aparține și îmi asum conținutul său în întregime.

De asemenea, declar că sunt de acord ca lucrarea mea de licență să fie verificată prin
orice modalitate legală pentru confirmarea originalității, consimțind inclusiv la introducerea
conținutului său într-o bază de date în acest scop.

Am luat la cunoștință despre faptul că este interzisă comercializarea de lucrări
științifice în vederea facilitării falsificării de către cumpărător a calității de autor al unei
lucrări de licență, de diploma sau de disertație și în acest sens, declar pe proprie răspundere că
lucrarea de față nu a fost copiată ci reprezintă rodul cercetării pe care am întreprins-o.

Data azi, _____ Semnătură student _____

DECLARAȚIE DE CONSIMȚĂMÂNT

Prin prezenta declar că sunt de acord ca Lucrarea de licență cu titlul „Abordarea clasificării prin metoda programării genetice”, codul sursă al programelor și celelalte conținuturi (grafice, multimedia, date de test etc.) care însoțesc această lucrare să fie utilizate în cadrul Facultății de Informatică. De asemenea, sunt de acord ca Facultatea de Informatică de la Universitatea „Alexandru Ioan Cuza” din Iași, să utilizeze, modifice, reproducă și să distribuie în scopuri necomerciale programele-calculator, format executabil și sursă, realizate de mine în cadrul prezentei lucrări de licență.

Iași,

Absolvent Constantin Sebastian Stanciu

CUPRINS

Introducere	6
1.Problema detecției de evenimente în monitorizarea calității apei	8
1.1 Învățarea supervizată	8
1.2 Problematika datelor dezechilibrate	9
1.3 Detectarea anomaliilor pentru calitatea apei potabile	13
1.4 Metode de sampling	14
1.5 Învățare bazată pe cost	15
2. Algoritmi de clasificare	17
2.1 Arbori de decizie	17
2.2 AdaBoost	18
2.3. Regresie logistică	19
2.4. Extra trees	21
3. Programare genetică	22
3.1 Concepte și algoritm	22
3.2. Programare genetică pentru clasificare binară	27
3.3 Funcții fitness pentru clasificare	28
4. Analiză experimentală	29
4.1 Arbori de decizie	30
4.2 Regresie logistică	31
4.3 AdaBoost	31
4.4 Extra trees	32
4.5 Programare genetică	33
4.6 Experimentări pentru clasificare cu programare genetică	37
4.7 Ansamblu de clasificatori bazați pe programare genetică	38
4.8 Interpretarea rezultatelor	40

Introducere

Învățarea automată (Machine Learning) are multiple aplicații, una dintre cele mai importante fiind extragerea de cunoștințe din date (data mining). Aceasta este de două tipuri: învățare supervizată și învățare nesupervizată. Diferența dintre acestea este faptul că în cazul primului tip de învățare setul de date conține pentru fiecare set de atribute și o etichetare a acestora care reprezintă variabila ce trebuie prezisă, care lipsește în schimb în cazul celui de-al doilea tip.

Oamenii sunt pasibili să facă greșeli în analiza sau în încercarea de a stabili relații între atributele dintr-un set de date, fapt ce îngreunează rezolvarea unor probleme, în schimb algoritmi de învățare automată elimină eroarea umană și sunt rapizi, datorită puterii computaționale din ce în ce mai mari.

O mulțime de aplicații de învățare automată pot include sarcini de învățare supervizată, acest domeniu fiind și tema pe care este orientată această lucrare, în special spre clasificare, realizată cu ajutorul programării genetice, care nu este o metodă clasică folosită în această problemă, urmărindu-se dacă această metodă poate fi o soluție viabilă în anumite cazuri de clasificare. Așadar scopul acestei lucrări de licență este studiul programării genetice ca și metodă de clasificare.

Clasificarea este o problemă de actualitate, cu o vastă importanță în domeniul informatic și nu numai. Așa cum este sugerat și de nume, clasificarea reprezintă identificarea unui obiect sau fenomen necunoscut ca membru al unei clase de obiecte sau fenomene, cu ajutorul unor informații anterioare disponibile. O caracteristică principală a clasificării este faptul că algoritmul de rezolvare selectează dintr-un set deja definit de soluții, pe baza atributelor care îl definesc, neîncercând găsirea unor noi soluții. O problemă care poate apărea în acest caz este faptul că uneori regulile de clasificare, care sunt învățate pe baza setului de date de antrenament nu sunt potrivite pentru clasificarea anumitor noi obiecte, ceea ce poate duce la clasificări greșite. Acest lucru poate avea consecințe mari, de exemplu în domeniul medical, un algoritm de detectare a unei boli clasifică un pacient drept sănătos, când de fapt el este bolnav. Această clasificare poate avea un cost mare, fiindcă se întârzie detectarea bolii și începerea tratamentului. În această lucrare accentul va fi pus pe clasificarea binară, elementele putând face parte din două clase de obiecte. Eficiența unui algoritm de clasificare fiind măsurată cu ajutorul a patru concepte: *true positives*, *true negatives*, *false positives* și *false negatives*. Dintre acestea o deosebită importanță le au ultimele două, deoarece suma acestora reprezintă numărul de clasificări greșite ale algoritmului. În numeroase cazuri de clasificare există un “trade off” între acestea, deoarece costurile lor nu sunt aceleași, precum în cazul enunțat anterior, o clasificare false positive este mult mai puțin costisitoare decât una false negative.

Această problemă a clasificării binare are multiple aplicații în probleme din viața reală, precum în domeniul medical, în detectarea fraudelor, a e-mailurilor spam sau în cazul detectării produselor contaminate.

Principala temă a acestei lucrări este de a studia eficiența programării genetice în clasificarea binară, pentru un set de date de antrenament care este dificil și performanțele algoritmilor clasici de clasificare nu oferă rezultate suficient de bune. De asemenea se vor executa un număr ridicat de experimente cu scopul studierii performanței pentru metoda clasificării cu programare genetică. Lucrarea este structurată în patru capitole, fiecare având sub-capitole.

În primul capitol intitulat “Problema detecției de evenimente în monitorizarea calității apei” se va prezenta mai în detaliu în ce constă problema clasificării, cum poate fi afectată performanța clasificatorilor de către structura setului de date de antrenament și explicarea conceptului de set de date dezechilibrat. De asemenea, va fi prezentat setul de date de antrenament care va fi folosit și structura acestuia.

În al doilea capitol “Algoritmi de clasificare” sunt prezentați o parte dintre cei mai cunoscuți și folosiți algoritmi pentru clasificare.

În al treilea capitol este prezentat conceptul de programare genetică, importanța funcției fitness, modul în care va fi abordată problema clasificării cu ajutorul programării genetice și funcțiile de fitness ce au fost folosite în experimentele realizate.

În al patrulea capitol numit în mod sugestiv “Analiză experimentală” sunt analizate și comparate performanțele clasificatorilor clasici și a celor obținuți prin programare genetică. De asemenea este prezentată și construirea unui ansamblu de clasificatori obținuți prin programare genetică și performanțele acestuia.

1.Problema detecției de evenimente în monitorizarea calității apei

1.1 Învățarea supervizată

Învățarea supervizată este reprezentată de algoritmi care reușesc cu ajutorul unui set de date oferit să găsească ipoteze/modele generale prin care să facă predicții despre obiecte noi, să facă mapări de la setul de date input la output. Cu alte cuvinte se urmărește crearea unui clasificator care să asignează etichete unor noi instanțe despre care se știe doar setul de atribute, dar nu și eticheta acestora.

Abordările existente în activitatea denumită generic *mineritul datelor (Data Mining)* se împarte în două tipuri de învățare: supervizată și nesupervizată. Ambele se bazează pe un set de date, însă în cazul învățării supervizate setul de date include și eticheta (outputul corespunzător) (Tabelul 1), iar în cazul învățării nesupervizate acesta lipsește. Învățarea supervizată se împarte în continuare în două categorii: clasificare și regresie, diferență dintre acestea fiind faptul că în cazul clasificării se prezintă o etichetă, așadar se urmărește crearea unei mapări de la variabilele de input la o variabilă categorică, discretă de output, iar în cazul regresiei maparea se face de la variabilele de input la o variabilă de output numerică, continuă. În cazul utilizării algoritmilor de clusterizare, numiți și învățare nesupervizată se urmărește descoperirea unor clase de obiecte neștiute, dar care pot fi importante, care pot surprinde anumite legături, similarități, conexiuni între obiectele respective.

	Atribut 1	Atribut 2	...	Atribut n	Eticheta
1	k	0		xy	1
2	r	1		xx	0
3	s	0		xx	2
...					
m	a	1		xy	6

Tabelul 1.
Instanțe cu etichetă cunoscută.

Există o expansiune continuă a acestui domeniu, care are o importanță deosebită în multe domenii complexe și variate, de la probleme de zi cu zi până la probleme de securitate

națională precum armată, medicină, finanțe. Multe dintre aceste probleme au date nebalansate, adică clasele din setul de date nu sunt reprezentate în mod egal sau cel puțin aproximativ în mod egal. De exemplu în cazul unei probleme de clasificare binară, adică setul de date are două etichete dacă dintr-un număr total de 1000 de instanțe, 800 reprezintă clasa întâi, iar 200 reprezintă clasa a doua, setul de date este considerat nebalansat, din cauza diferenței mari de instanțe reprezentative pentru cele două clase, raportul dintre ele fiind 800:200 sau mai simplu 4:1.

În cazul clasificării binare, în cadrul a numeroase probleme clasele de elemente sunt 0 și 1, clasa 0 fiind denumită clasa negativă, iar clasa 1 clasa pozitivă. În cadrul acestei lucrări se vor folosi de asemenea aceste notații și denumiri, deoarece și în cazul setului de date folosit ca instrument de studiu în această lucrare clasa 0 definește absența unui eveniment, iar clasa 1 prezența acestuia, clasa 0 definește instanțele de apă necontaminată iar clasa 1 pe cele de apă contaminată. Precum am precizat în introducere patru concepte vor fi de folos în evaluarea performanței unui algoritm de clasificare binară: true positive (se prezice în mod corect clasa pozitivă), true negative (se prezice în mod corect clasa negativă), false positive (se prezice în mod greșit clasa pozitivă), false negative (se prezice în mod greșit clasa negativă). Aceste concepte se vor abrevia în această lucrare astfel: true positives ->TP, true negative ->TN, false positive ->FP, false negative ->FN.

1.2 Problematica datelor dezechilibrate

În multe situații una dintre clase este slab reprezentată, adică numărul de exemple aparținând acesteia este mult mai mic decât numărul de exemple din cealaltă clasă (celelalte clase) - aceasta este denumită clasa minoritară. În astfel de situații poate apărea paradoxul acurateții, în sensul că acuratețea poate fi foarte mare deoarece algoritmul învață foarte bine să clasifice clasa majoritară, în exemplul dat anterior o acuratețe de 80% poate părea destul de bună, dar este posibil ca algoritmul să clasifice toate instanțele ca făcând parte din prima clasă, fiind incapabil să clasifice obiectele care ar trebui să facă parte din clasa minoritară. Acuratețea este definită după următoarea formulă:

$$A = \frac{TP + TN}{TP + TN + FP + FN}$$

Rata erorii este definită drept raportul dintre clasificările greșite și numărul total de clasificări, conform formulei:

$$RataErorii = \frac{FP + FN}{TP + TN + FP + FN}$$

Alte doua metrice foarte importante în analiza performanței unui algoritm sunt: *precizia* (*precision* în limba engleza) și *reamintirea* (denumita în engleza *recall*). În această lucrare vor fi folosiți în continuare termenii din limba engleza. Acestea sunt definite astfel:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

Cu ajutorul preciziei și a recallului se calculează scorul F1. Această metrică combină aceste doua concepte având valoarea maxima 1 și cea minima 0. Scorul F1 este în general folosit când numărul *TN* nu are o importanță ridicată și se dorește ca modelul sa nu conțină multe FP sau FN. Din aceste motive această metrica este potrivită în cazul seturilor de date neechilibrate și va fi folosită și în cadrul acestei lucrări pentru a compara algoritmii. O valoare mai ridicată indică o performanță mai buna a algoritmului, formula pentru scorul F1 fiind:

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall}$$

Un alt concept este matricea de confuzie (confusion matrix), care conține în cazul clasificării binare doua linii și doua coloane, iar în cazul unei set de date cu n clase conține n linii și n coloane. Valoarea de pe poziția (i,j) reprezintă numărul de instanțe din clasa i, care au fost clasificate ca făcând parte din clasa j, de exemplu în figura 1, elementul de pe poziția (0,0), reprezintă numărul de instanțe care fac parte din clasa negativă și care au fost clasificate în mod corect.

	Negativ	Pozitiv
Negativ	TN	FP
Pozitiv	FN	TP

Figura 1.

Matrice de confuzie pentru clasificare binară

În cazul clasificării binare, matricea de confuzie ajuta pentru vizualizarea numărului total de clasificări corecte și a celor greșite astfel: suma elementelor de pe diagonala principală reprezintă clasificările corecte, iar cel de pe diagonala secundara clasificările greșite.

Această problemă are un interes major fapt ce se ilustrează prin numeroasele conferințe, articole, al căror număr crește la o rată ridicată, devenind obiectivul principal al acestei lucrări în care se urmărește eficiența unei metode mai puțin clasice pentru o astfel de problemă, programarea genetică. În acest capitol va fi prezentată problema învățării din seturi de date nebalansate în profunzime.

Teoretic orice set de date în care distribuția claselor nu este egală poate fi considerat dezechilibrat, însă de obicei o diferență mică nu influențează eficiența clasificatorilor într-un mod care poate fi considerat semnificativ. De obicei printr-un set de date nebalansat se înțelege o diferență majoră, fiind posibile reprezentări de tipul 1:100, 1:500 sau 1:1000, în fiecare dintre acestea o clasă are mult mai multe instanțe decât cealaltă (Figura 2). Deși din definiția anterioară s-ar putea înțelege că se face referire la problemele de clasificare binară, acest lucru se poate întâmpla și în cazul unor seturi de date care conțin 3 sau mai multe clase, dezechilibrate având loc între anumite clase din setul de date (Figura 3).

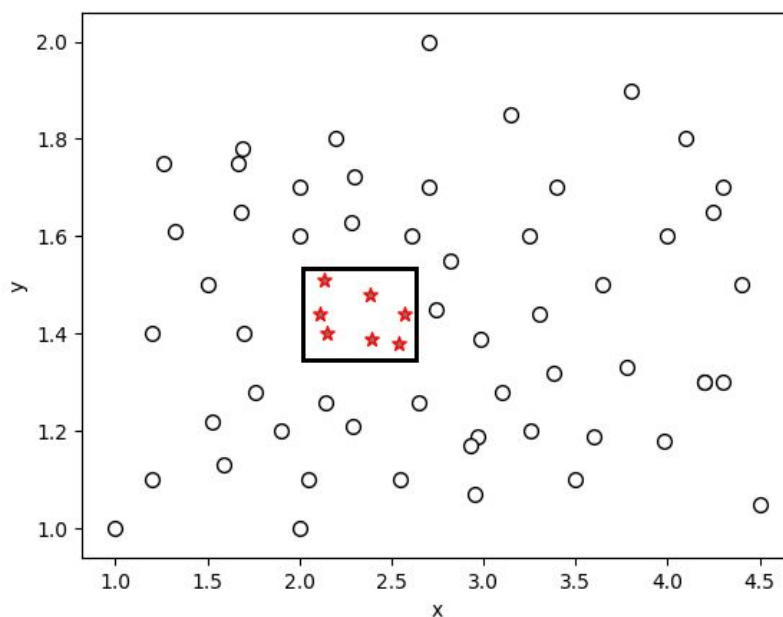


Figura 2.

Set de date în care raportul dintre cele două clase este dezechilibrat, raportul fiind de 1:8.

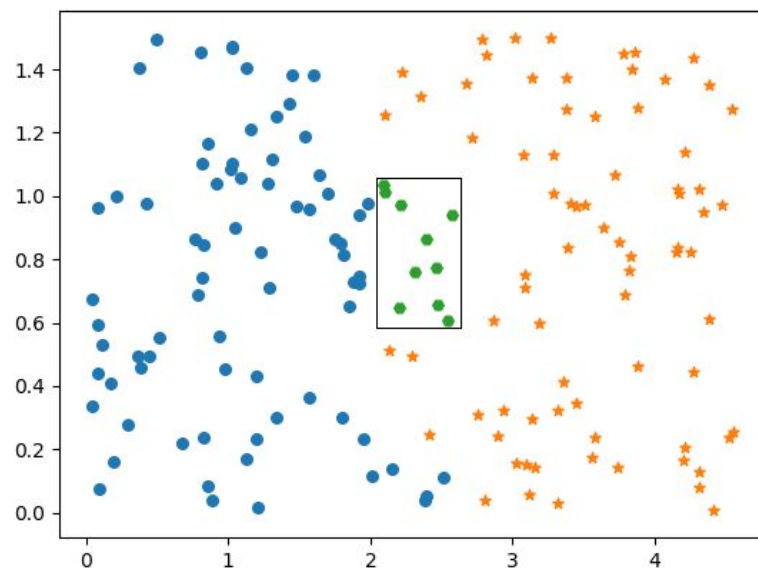


Figura 3.

Set de date dezechilibrat cu mai multe clase. Raportul dintre clasa minoritară și clasele majoritare este 1:5

Dezechilibrele pot fi intrinseci, în cazul în care dezechilibrul este o consecință directă a naturii setului de date, de exemplu un set de date cu ajutorul căruia se încearcă clasificarea pacienților drept bolnavi de cancer sau sănătoși, este normal ca acesta să conțină un număr mai mare de instanțe ale persoanelor sănătoase, decât al celor bolnave. Există și dezechilibre care nu țin de natura setului de date sau a problemei care se încearcă să se rezolve, ci de factori precum timp sau puterea de stocare, acestea se numesc extrinseci. Dezechilibrele extrinseci sunt la fel de interesante ca și cele intrinseci, acestea putând apărea chiar și în cazul unor seturi de date ce nu sunt neapărat inițial dezechilibrate, de exemplu dacă un set de date este obținut treptat iar pe parcursul transmisiei apar întreruperi, în care datele nu sunt trimise, este posibil să obținem un set nebalansat extrinsec, obținut dintr-un set de date balansat.

De asemenea, alte două concepte importante sunt dezechilibrul relativ și dezechilibrul cauzat de instanțe rare, numit și raritate absolută. Un exemplu pentru dezechilibrul relativ este un set de date binar care conține 50.000 de exemple și un raport între clase de 50:1, astfel clasa minoritară având 1000 de instanțe. Dacă se mai obțin date, iar setul de date se dublează, iar raportul dintre clase rămâne același, clasa minoritară va conține 2000 de instanțe, un număr destul de consistent de instanțe, dar scăzut relativ la clasa majoritară. Acest timp de dezechilibru apare în numeroase aplicații din viața reală și va fi studiată în această lucrare. Dezechilibrul cauzat de instanțele rare apare în domenii în care clasa minoritară are număr limitat de exemple reprezentative, de aceea învățarea este dificilă indiferent dacă există dezechilibru între clase sau nu.

1.3 Detectarea anomaliilor pentru calitatea apei potabile

În studiul acestei probleme am folosit un set de date binar dezechilibrat. Am preluat aceste date din cadrul GECCO(“Genetic and Evolutionary Computation Conference“), care a organizat în anul 2018 competiția “Internet of Things: Online Anomaly Detection for Drinking Water Quality”.

Setul de date este oferit de o companie de îmbuteliere a apei(Thüringer Fernwasserversorgung), fiind preluat din viața reală și reprezintă măsurători ale parametrilor realizate de către această companie, făcându-se astfel o legătură directă cu industria. Datele au fost obținute prin trecerea apei prin anumite zone speciale, cu senzori, unde sunt măsurate cei mai importanți indicatori de calitate ai apei. Setul de date reprezintă de fapt, o serie de timp, instanțele fiind măsurători ale acestor atribute o dată pe minut.

Setul de date are următoarele atribute: Time (timpul când a fost făcută măsurătoarea), Tp (temperatura în grade celsius), Cl (cantitatea de dioxid de clor în mg/L), pH (valoarea PH-ului), Redox (potențialul de reducere), Leit (conductibilitatea electrică a apei), Trueb (turbiditatea), Cl_2 (cantitatea de dioxid de clor la a doua măsurătoare), Fm (debitul în conductă 1), Fm_2 (debitul în conductă 2) și EVENT (eticheta, variabilă booleană ce specifică dacă apa este contaminată).

Scopul acestei competiții este realizarea unui clasificator, care să detecteze schimbările de calitate ale apei, contaminarea acesteia. Această sarcină are o importanță deosebită pentru companiile de apă publică, pentru cetățeni, dar nu numai, apa fiind vitală pentru orice formă de viață. Un astfel de clasificator este esențial în menținerea și distribuția apei potabile pe întreaga planetă.

Acest set de date prezintă dezechilibru intrinsec, fiind evident că pe parcursul acestor măsurători instanțele de apă contaminată să fie mult mai slab reprezentate, decât cele ale apei potabile și este de asemenea un dezechilibru relativ, clasa minoritară având un număr de 1726, nefiind slab reprezentată decât relativ la clasa majoritară, care conține 137.840 de instanțe.

Datele sunt foarte dezechilibrate, raportul dintre instanțele clasei minoritare și cele ale clasei majoritare fiind de aproximativ 1:80 (Figura 4).

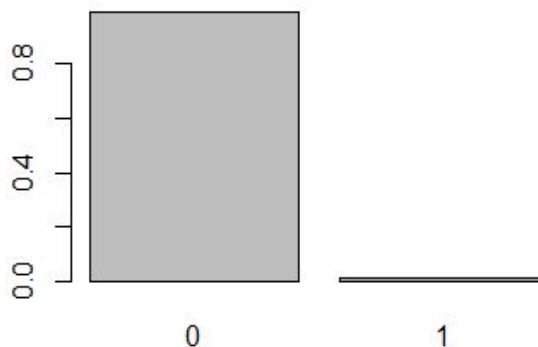


Figura 4.
Vizualizare a raportului dintre instanțele celor doua clase

1.4 Metode de sampling

Aceste metode constau în modificarea setului de date prin diferite mecanisme cu scopul de a micșora dezechilibrul dintre clase sau chiar de a echilibra instanțele ce reprezintă clasele din setul de date, cu scopul de a îmbunătăți performanțele clasificatorilor.

Există numeroase astfel de metode, dar în cadrul acestei prezentări vor fi menționate doar o parte dintre acestea. În primul rând există două concepte fundamentale pentru aceste metode: oversampling, care constă în a adăuga în setul de date instanțe din clasele minoritare, și undersampling, care constă în ștergerea unor instanțe din setul de date, care aparțin claselor majoritare.

Cel mai simplu mod de realizare a acestor două concepte este în mod aleatoriu. În cazul oversampling-ului, se selectează în mod aleatoriu un subset de dimensiune dorită din clasa minoritară și apoi se adaugă în setul de date. În cazul undersampling-ului, se selectează în mod aleatoriu un subset din clasa majoritară, ulterior instanțele din acest subset sunt șterse din setul de date. În mod evident prin aceste metode se micșorează dezechilibrul din setul de date, însă pot cauza o serie de probleme. În cazul undersampling-ului se elimină instanțe din setul de date, fapt ce poate cauza pierderea unor informații, concepte importante pentru clasificarea corectă a instanțelor din clasa majoritară, iar în cazul oversampling-ului se adaugă în setul de date instanțe replicate, ceea ce poate duce la supra specializare.

O metodă mult mai interesantă este metoda SMOTE, ce a fost aplicată cu succes în multiple aplicații, performanțele acesteia fiind studiate și menționate în numeroase articole de specialitate. Această metodă creează în mod artificial instanțe din clasa minoritară cu ajutorul

celor deja prezente în setul de date cu algoritmului KNN, prin găsirea celor mai apropiați k vecini pentru instanțele din clasa minoritară. Apoi se sintetizează o nouă instanță între instanță din clasa minoritară și unul dintre vecinii sai, precum este ilustrat în figura 5

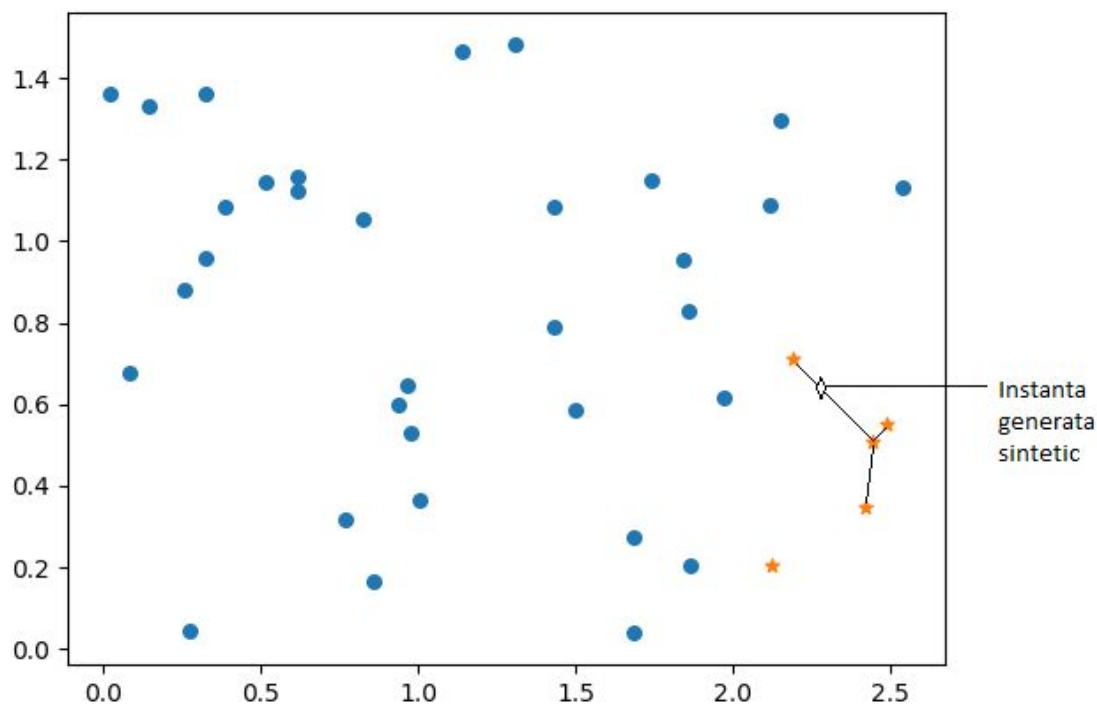


Figura 5.
Generarea unei instanțe în mod sintetic cu SMOTE

1.5 Învățare bazată pe cost

Spre deosebire de metodele de sampling care modifică setul de date pentru a crește performanțele clasificatorilor, învățarea bazată pe cost constă în asocierea unor costuri diferite pentru clasificările greșite, în funcție de clasa din care face parte instanța respectivă. Acest concept este o alternativă viabilă pentru metodele de sampling în încercarea de a îmbunătăți performanțele algoritmului de clasificare.

Modul în care funcționează acest concept este prin atribuirea unor costuri diferite pentru clasificările greșite, de exemplu în cazul unei probleme de clasificare binară costul clasificării greșite a unei instanțe din clasa majoritară va fi mai mic decât costul unei clasificări greșite din clasa minoritară și se încearcă minimizarea costului pentru întreg setul de antrenament. În general nu se atribuie niciun cost pentru clasificările corecte.

Matricea de cost este un concept important în acest tip de învățare, fiind o modalitate de a vizualiza costurile, care sunt furnizate fie de către experți în domeniu sau stabilite prin diferite metode. În cadrul experimentelor din această lucrare costul unei clasificări greșite din clasa majoritară va fi 1, iar costul unei clasificări greșite din clasa minoritară va fi raportul dintre numărul de instanțe din clasa majoritară și numărul de instanțe din clasa minoritară. În matricea de cost pe poziția (i,j) se găsește costul clasificării unei instanțe din clasa i în clasa j , și în acest caz arata astfel:

	Negativ	Pozitiv
Negativ	0	1
Pozitiv	80	0

Figura 6.
Matrice de costuri

2. Algoritmi de clasificare

În acest capitol al acestei lucrări sunt prezentate o parte dintre cele mai cunoscute și folosite metode de clasificare, o parte dintre acestea au fost studiate în cadrul materiei Învățare automata din primul semestru al anului III, care vor fi prezentate și metode care nu au fost studiate în cadrul facultății.

În primul rând acești algoritmi au câteva caracteristici comune, specifice pentru a rezolva problema clasificării, precum: se bazează pe un set de date care conține instanțe de tipul atribut-valoare, cu scopul de a clasifica obiecte noi despre care se cunosc doar atributele, de asemenea atributele valoare din setul de date sunt discrete, adică care au un număr finit de valori spre deosebire de variabilele continue din cazul regresiei, celalalt tip de învățare supervizată.

2.1 Arbori de decizie

Arborii de decizie sunt un tip de clasificator ce se reprezintă drept arbori prin sortarea atributelor în funcție de valoarea acestora. Fiecare nod din arbore reprezintă un atribut, iar fiecare muchie conține o valoare pe care nodul respectiv o poate avea. Clasificarea unei instanțe noi se realizează prin parcurgerea arborelui de la rădăcina la frunza care corespunde acestei parcurgeri și care conține eticheta clasei.

Această metodă se bazează pe împărțirea setului de antrenament în funcție de valoarea atributelor, practic se considera ca nodul rădăcina realizează acest lucru. Apoi aceeași procedură se realizează pe fiecare partiție din setul de date, obținută prin împărțirile anterioare. Există mai multe metode pentru o împărțire cât mai eficientă a setului de date, de exemplu câștigul de informație (Hunt et al., 1966) și indexul gini (Breiman et al., 1984).

Termenul de câștig de informație este strâns legat de cel al entropiei, care este o măsură a incertitudinii, de aceea entropia în cazul aruncării unei monezi este maximă, probabilitatea de a cădea cap sau pajura fiind egale. Entropia este calculată în felul următor:

$$H(X) = - \sum_{i=1}^n p_i \log_2 p_i$$

Valoarea entropiei este maximală atunci când exista o dezordine totală, în sensul ca toate probabilitățile sunt egale. Valoarea entropiei este calculata în acest caz astfel:

$$H = -n \cdot \frac{1}{n} \log_2 \frac{1}{n} = \log_2 n$$

Pentru calculul entropiei se face următoarea convenție, necesara în cazul în care una dintre probabilități este egala cu 0:

$$0 \cdot \log_2 0 = 0$$

În cazul arborilor de decizie în care se folosește câștigul de informație se încearcă maximizarea acestuia, iar atributul cu cel mai mare câștig de informație va fi selectat pentru a partiționa setul de date. Câștigul de informație este calculat cu ajutorul entropie și reprezintă cat de multă informație oferă un atribut despre clasa din care face parte, de aceea un atribut care împarte setul de date în mod perfect oferă un câștig de informație maximal.

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

În cazul arborilor ID3 este folosit câștigul de informație, iar în cazul arborilor CART se folosește indexul gini, care indică cat de mixate sunt clasele în grupurile create de o împărțire, o tăietură a acestuia. Indexul Gini este calculat în cazul unui set de date care conține n clase astfel:

$$gini(X) = 1 - \sum_{i=1}^n p_i^2$$

Indexul gini are valoarea 0 în cazul unei împărțiri perfecte și valoarea maxima în cazul în care rezulta clase cu număr egal de instanțe, calculându-se astfel:

$$gini(X) = 1 - k\left(\frac{1}{k^2}\right) = 1 - \frac{1}{k}$$

2.2 AdaBoost

Boostingul (stimularea) este o metodă bazată pe ideea creării unui predictor puternic format din mai multe reguli mai slabe și cu acuratețe mai slabă. Algoritmul AdaBoost, prescurtare de la “Adaptive Boosting” este o un algoritm bazat pe conceptul boosting care a dovedit o eficiență crescută, fiind unul dintre cei mai populari algoritmi, cu aplicări în multiple domenii. Algoritmii de tip boosting crează predictorul puternic cu ajutorul unei mulțimi de clasificatori mai slabi în mod secvențial.

Modul în care AdaBoost funcționează este prin asignarea inițială a unui cost de $1/n$ tuturor instanțelor din setul de antrenament. De asemenea, după ce clasificatorul a fost antrenat i se asignează o pondere, în funcție de acuratețea sa, iar costul instanțelor clasificate greșit crește cu scopul clasificării corecte a acestora în iterația următoare. Acest proces continuă până când se atinge o acuratețe îndeajuns de bună sau până când se atinge un număr maxim de clasificatori. Clasificarea unui noi instanțe se realizează prin votul tuturor clasificatorilor, ținând cont de ponderea acestora.

Algoritmul AdaBoost se bazează pe câteva formule matematice. Formula următoare reprezintă eticheta asignată de clasificator unui instanțe x :

$$H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

În formula aceasta h_t reprezintă eticheta asignată de clasificatorul slab t instanței x , iar α_t este ponderea asignată clasificatorului t , având valoare negativă în cazul în care acuratețea este mai mică decât 0.5 și pozitivă în caz contrar. Formula de calcul este următoarea, E fiind numărul clasificărilor greșite:

$$\alpha_t = 0.5 * \ln\left(\frac{1 - E}{E}\right)$$

Recalcularea ponderii unei instanțe se realizează conform următoarei formule:

$$D_{t+1}(i) = \frac{D_t(i)e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$$

$D_t(i)$ reprezintă ponderea de la nivelul anterior, iar Z_t reprezintă suma tuturor ponderilor, cu scopul de a normaliza ponderile.

2.3. Regresie logistică

Regresia logistică este un algoritm de clasificare chiar dacă conține cuvântul regresie în denumire, care ar putea dezorienta, indicând spre celalalt tip de învățare supervizată. Așadar regresia logistică caută un model care să estimeze probabilitatea unei variabile dependente cu ajutorul unor variabile independente.

În cazul setului de date folosit în această lucrare variabila dependentă este EVENT, iar variabilele independente sunt restul variabilelor prezentate. Cu ajutorul probabilității obținute se reușește clasificarea unor noi instanțe, prin compararea acestora cu un prag ales, în cazul în care probabilitatea este mai mare decât pragul instanța este clasificată pozitiv, iar în caz contrar este clasificată în mod negativ.

În cazul problemei abordate în această lucrare nu putem folosi regresia liniară deoarece aceasta presupune prezicerea unei variabile dependente cu ajutorul unei alte variabile independente, de asemenea nici regresia multiplă care este asemănătoare cu cea multiplă, dar cu mai multe variabile independente, aceste două metode necesită un set de date cu distribuție normală. Distribuția normală este cel mai important tip de distribuție din cadrul statisticii deoarece este reprezentativă pentru numeroase fenomene naturale, precum înălțimea, greutatea sau scorul IQ.

Abordarea problemei cu astfel de metode nu este posibilă și din cauza faptului că acestea ar putea prezice valori pentru variabila dependentă EVENT care să depășească intervalul $[0,1]$, fapt care ar încălca definiția probabilității. În cazul regresiei logistice se obține următoarea formulă, care nu presupune problema aceasta, prin găsirea unei funcții care asignează mulțimii de variabile independente o valoare care reprezintă probabilitatea de încadrare în valorile variabilei dependente.

$$0 \leq h_{\theta}(x) \leq 1$$

În formula dată, $h_{\theta}(x)$ are următoarea formă:

$$h_{\theta}(x) = g(\theta^T x)$$

$$g(t) = \frac{1}{1 + e^{-t}}$$

Din aceste două relații se obține în mod trivial următoarea relație, care reprezintă probabilitatea estimată ca pentru datele x clasificarea corectă să fie clasa 1 (clasa pozitivă):

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Conform definiției anterioare se calculează probabilitatea apartenenței la clasa pozitivă, iar fiind o problemă de clasificare binară, probabilitatea clasei negative este diferența dintre 1 și această valoare. Pentru clasificarea unor date cu ajutorul modelului creat este necesară stabilirea unui prag care ia în mod evident valori din intervalul (0,1), care dacă este depășit de către valoarea calculată $h_{\theta}(x)$ datele sunt clasificate ca făcând parte din clasa pozitivă sau negativă în caz contrar.

Un rol important în funcționarea buna a acestui algoritm o au parametrii θ , iar aceștia sunt calculați cu ajutorul unei funcții de cost. Scopul acestui procedeu este de a găsi parametrii θ care minimizează această funcție și care vor fi ulterior folosiți în cadrul clasificării. Costul clasificării unei instanțe din setul de date de antrenament este:

$$Cost(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$$

Funcția de cost pentru regresia logistica în cazul unui set de antrenament cu m instanțe este suma costurilor fiecărei instanțe împărțită la numărul de instanțe:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(h(x_i)) + (1 - y_i) \log(1 - h(x_i))]$$

2.4. Extra trees

Extra trees este o altă metodă de clasificare de tip ansamblu, precum algoritmul AdaBoost, dar spre deosebire de acesta elementele ansamblului funcționează în mod paralel. această denumire provine de la prescurtarea numelui inițial “Extremely randomized trees”. Acest algoritm funcționează prin crearea mediei mai multor arbori care sunt creați prin partiționarea setului de date în mod aleator.

Această abordare s-a dovedit foarte eficientă, crescând precizia și având o complexitate computațională scăzută, evitându-se și supra-adaptarea (overfitting), care se datorează unei adaptări prea mari la setul de date de antrenament, în sensul ca ipoteza este influențată și de date irelevante sau de zgomot, ceea ce cauzează o dificultate în clasificarea unor instanțe noi, din afara setului de antrenament.

Ideea algoritmului este asemănătoare ca în cazul clasificării cu ajutorul algoritmului Random Forest, dar există și anumite diferențe, una dintre cele mai importante fiind faptul ca în cazul Extra trees punctul de tăiere în cazul proprietăților selectate este ales în mod random, spre deosebire de Random Forest în cazul căruia se caută cel optim. Din acest motiv

algoritmul Extra trees este mult mai rapid din punct de vedere computațional, iar din punct de vedere al performanței algoritmului sunt asemănători.

Spre deosebire de metoda arborilor de decizie prezentați anterior în cadrul acestei lucrări, Extra trees propune un ansamblu de arbori, nu doar un arbore, de aceea se reduce supra-adaptarea și de aceea în cele mai multe cazuri această metodă funcționează mai bine pentru clasificarea unor noi instanțe.

La fel ca în cazul arborilor de decizie calitatea împărțirii setului de date de către un atribut este calculată cu ajutorul câștigului de informație sau a indexului gini. Aceste metode au fost prezentate în cadrul acestui capitol în cadrul prezentării arborilor de decizie și nu vor mai fi reluate.

3. Programare genetică

3.1 Concepte și algoritmi

Programarea genetică reprezintă obiectul de studiu principal al acestei lucrări. Acest concept este împrumutat din natura în sensul că structurile biologice care reușesc o adaptare buna supraviețuiesc și se reproduc la o scară mai ridicată decât cele ce sunt inadaptate, conform teoriei evoluționiste a lui Darwin, prin selecție naturală. Adaptabilitatea unui individ poate fi măsurată prin conceptul de fitness.

Programarea genetică este o tehnică ce face parte din grupul algoritmilor evolutivi, din care fac parte și algoritmi genetici, ce are scopul de a rezolva probleme în mod automat indiferent de domeniul acestora, fără a fi necesar să se specifice forma sau dimensiunea soluției, ci doar în ce constă rezolvarea problemei, de exemplu în cazul unei probleme de optimizare, găsirea unei valori de maxim sau de minim.

Pentru o prezentare cât mai completă a programării genetice mai sunt necesare prezentarea anumitor concepte. Rezolvarea problemelor pe calculator se efectuează cu ajutorul algoritmilor, pentru care nu există o definiție standard, definiția din dicționarul Cambridge, considerat o sursă cu grad mare de încredere este: “Un set de instrucțiuni matematice care trebuie urmat într-o ordine fixată și care cu ajutorul unui calculator va ajuta la calcularea răspunsului pentru o problemă matematică”. Există două strategii cu ajutorul cărora se rezolvă probleme cu ajutorul calculatorului, în mod determinist și în mod probabilist.

Algoritmi determiniști sunt de două tipuri: algoritmi exacți și euristici. Caracteristica comună a acestora este că pentru un anumit input produc mereu același output, urmând mereu aceeași secvență de stări. În cazul algoritmilor determiniști exacți se obțin rezultate exacte, dar nu pot fi mereu folosiți în practică în timp util, în special în cazul problemelor ce fac parte din clasa NP-complete, pentru care nu se cunoaște încă un algoritm determinist exact care să le rezolve în timp polinomial. Euristicile sunt algoritmi care produc soluții aproximative, care nu sunt de cele mai multe ori optime, dar sunt suficient de bune în cazul în care se dorește o rezolvare într-un timp limitat pentru o problemă dificilă precum cele NP-complete. Exemple de probleme ce fac parte din categoria NP-complete sunt: problema comisului voiajor, problema colorării unui graf, problema rucsacului.

Algoritmi probabilști reprezintă o metodă de rezolvare a problemelor prin utilizarea noțiunii “random”, adică la cel puțin un pas din cadrul algoritmului decizia se realizează în

mod aleatoriu. Din acest motiv nu se poate prezice valoarea soluției ce se va obține, la rulări succesive este posibil ca rezultatele pentru un același input să difere, reieșind astfel ca soluțiile nu pot fi exacte, ci sunt aproximative.

Modelul tradițional de rezolvare al problemelor se numește “hard computing” și are trăsături precum faptul că se bazează pe principii ca precizie și rigoare și nu are posibilitate de adaptare (exemple : calculul integral, calculul diferențial). Aceste metode nu pot fi folosite în multe probleme din viața reală în care apar incertitudini, inexactități sau ambiguități, ci doar în probleme relativ simple.

O tehnică mai modernă sunt metodele “soft computing” care au posibilitatea de autoadaptare. Considerăm ca cel mai intuitiv mod de a prezenta acest principiu este prin antiteză cu conceptul prezentat anterior “hard computing”, care este descris prin rigiditate, inflexibilitate. Tehnica soft computing este tolerantă la incertitudini și aproximări, imitând creaturile naturale flexibile și adaptive. Această paradigmă se bazează pe un feedback intern cu ajutorul căruia proprietățile individuale ale problemei sunt exploatate la maxim.

Programarea genetică face parte din categoria calculului evolutiv, ce este o sub-ramură a Inteligenței Artificiale, ce cuprinde paradigme soft computing. Așadar programarea genetică aparține paradigmei soft-computing, fiind algoritmi probabiliști inspirați din biologie și natură. În cadrul paradigmei programării genetice structurile supuse adaptării sunt programe pe calculator de dimensiune și forme variabile. Încercarea de a rezolva o problemă cu ajutorul programării genetice reprezintă practic o căutare a unui program pe calculator cu un fitness foarte bun în spațiul programelor posibile. Această paradigmă oferă posibilitatea căutării celui mai adaptat program pe calculator pentru problema ce se încearcă să se rezolve, prin anumite metode inspirate din natura și teoria evoluționistă a lui Charles Darwin.

Execuția programării genetice începe cu o populație inițială de programe generate random, care poate fi considerată o căutare oarbă în spațiul programelor pe calculator. Programele generate sunt compuse din funcții și terminali aleși în funcție de domeniul problemei.

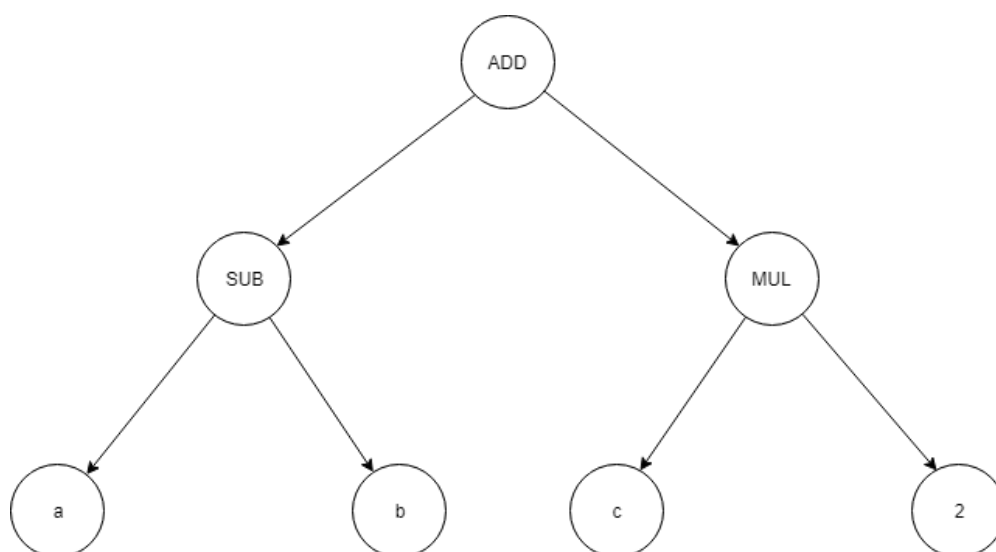


Figura 7. Reprezentarea unei soluții candidat în Programarea genetică

În cadrul programării genetice programele sunt reprezentate sub forma de arbori. În arborele din Figura 7 este reprezentat programul $\text{ADD}(\text{SUB}(a,b), \text{MUL}(c,2))$. Frunzele arborelui sunt terminalii, iar nodurile interne sunt denumite primitive. Terminalii sunt împărțiți în două subtipuri, constante și argumente, în exemplul anterior există o singură constantă (numărul 2) și trei argumente (a, b, c). Constantele își păstrează valoarea pe întreaga execuție, iar argumentele sunt datele de intrare ale programului.

Funcția fitness reprezintă un termen important în terminologia programării genetice, fiind modul în care se testează performanța programelor din setul de programe. Cu ajutorul funcției fitness se pot determina cele mai bune programe, informație care este folosită ulterior pentru modificarea și îmbunătățirea acestora. Astfel se poate ușor intui faptul că structura programelor pe calculator care reies din cadrul paradigmei programării genetice sunt o consecință a funcției de fitness, care reprezintă și în cadrul natural factorul dominant al evoluției.

Programarea genetică dispune de doi operatori genetici de baza, încrucișarea și mutația, cu ajutorul cărora se realizează căutarea în spațiul posibilelor soluții. Mutația este o operație care este controlată de un parametru probabilistic și se realizează pe un singur individ din populație. În urma mutației se obține un nou individ prin modificarea structurii părintelui. Un exemplu simplu de mutație este următorul:

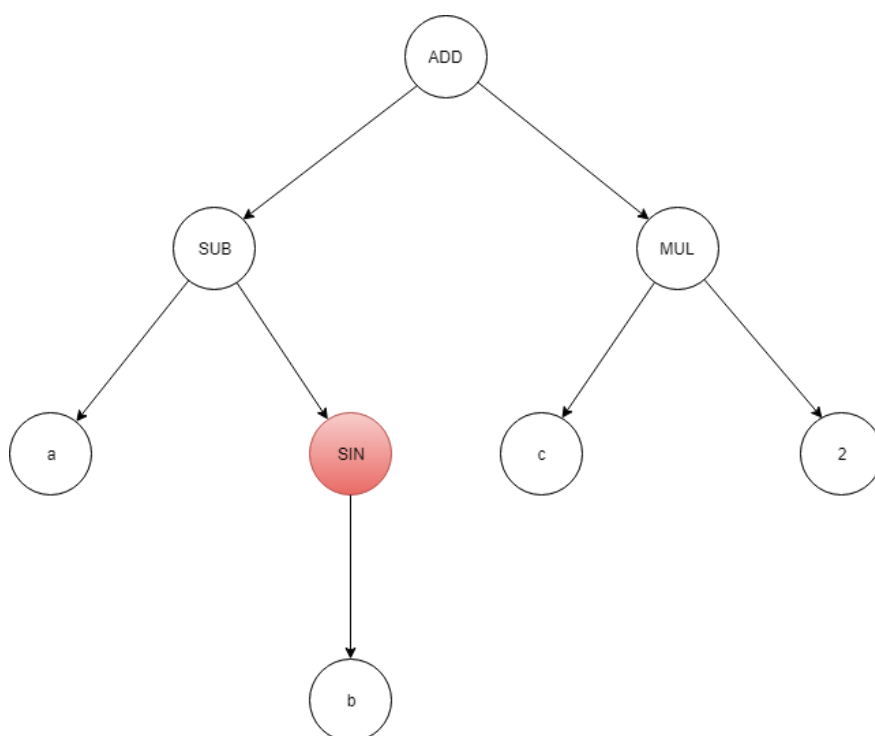


Figura 8.Individ înainte de realizarea mutației

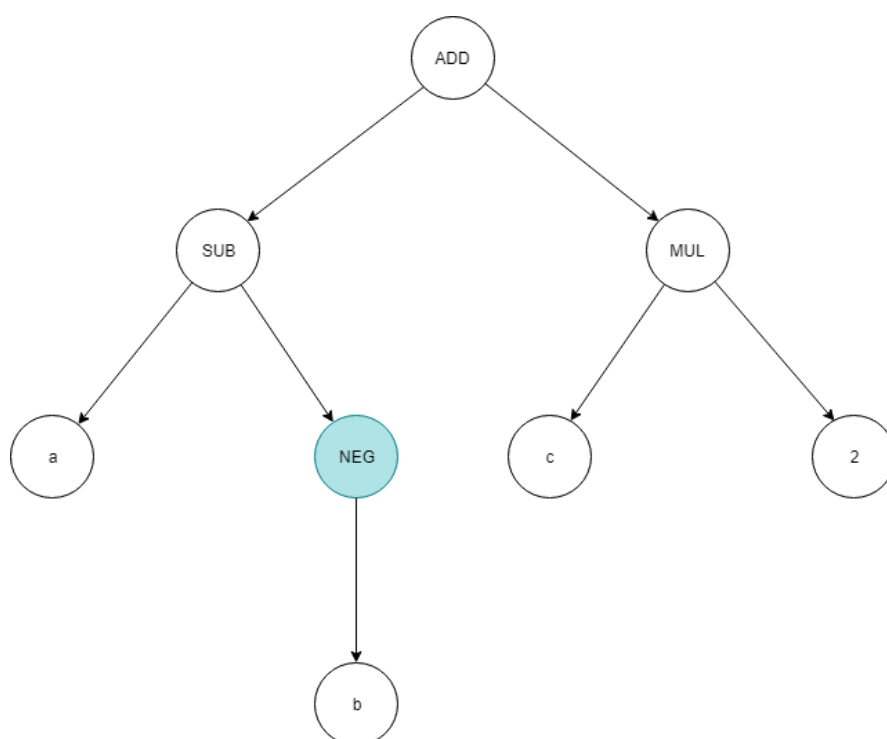


Figura 9.Individ după realizarea mutației

În exemplul dat se observă ca mutația a constant în înlocuirea primitivei SIN cu cea NEG, astfel pornind de la structura $\text{ADD}(\text{SUB}(a, \text{SIN}(b)), \text{MUL}(c, 2))$ obținem prin mutație structura $\text{ADD}(\text{SUB}(a, \text{NEG}(b)), \text{MUL}(c, 2))$.

Încrucișarea se realizează cu ajutorul a doua programe și se obțin doua noi programe, care sunt în general diferite între ele și diferite față de părinții lor. Operația începe prin selectarea în mod random în cadrul fiecărui părinte a unui punct de încrucișare, astfel se obțin doua fragmente de încrucișare reprezentate de sub-arborii ce au drept rădăcina punctul de încrucișare ales anterior. Primul copil se produce prin ștergerea fragmentului de încrucișare din cadrul primului părinte și adăugarea în locul acestuia a fragmentului de încrucișare din al doilea părinte, al doilea copil se produce în mod simetric.

În exemplul următor se ilustrează modul în care încrucișarea funcționează. Nodurile colorate în galben vor fi considerate punctele de încrucișare. Părinții au următoarele expresii: $\text{AND}(\text{NOT}(x), \text{OR}(y, z))$ și $\text{OR}(\text{AND}(x, \text{OR}(x, y)), \text{NOT}(z))$, iar copii rezultați au expresiile $\text{AND}(\text{NOT}(x), \text{AND}(x, \text{OR}(x, y)))$ și $\text{OR}(\text{OR}(y, z), \text{NOT}(z))$.

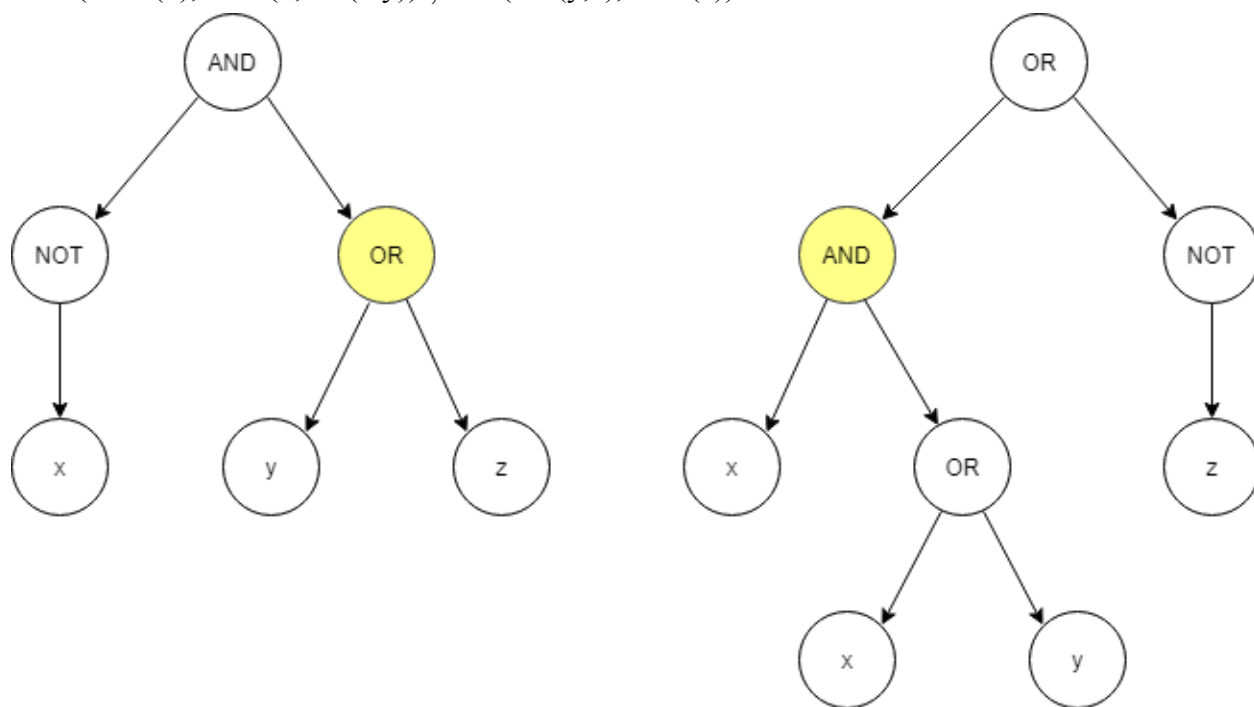


Figura 10.
Alegerea punctelor de încrucișare pentru doi indivizi

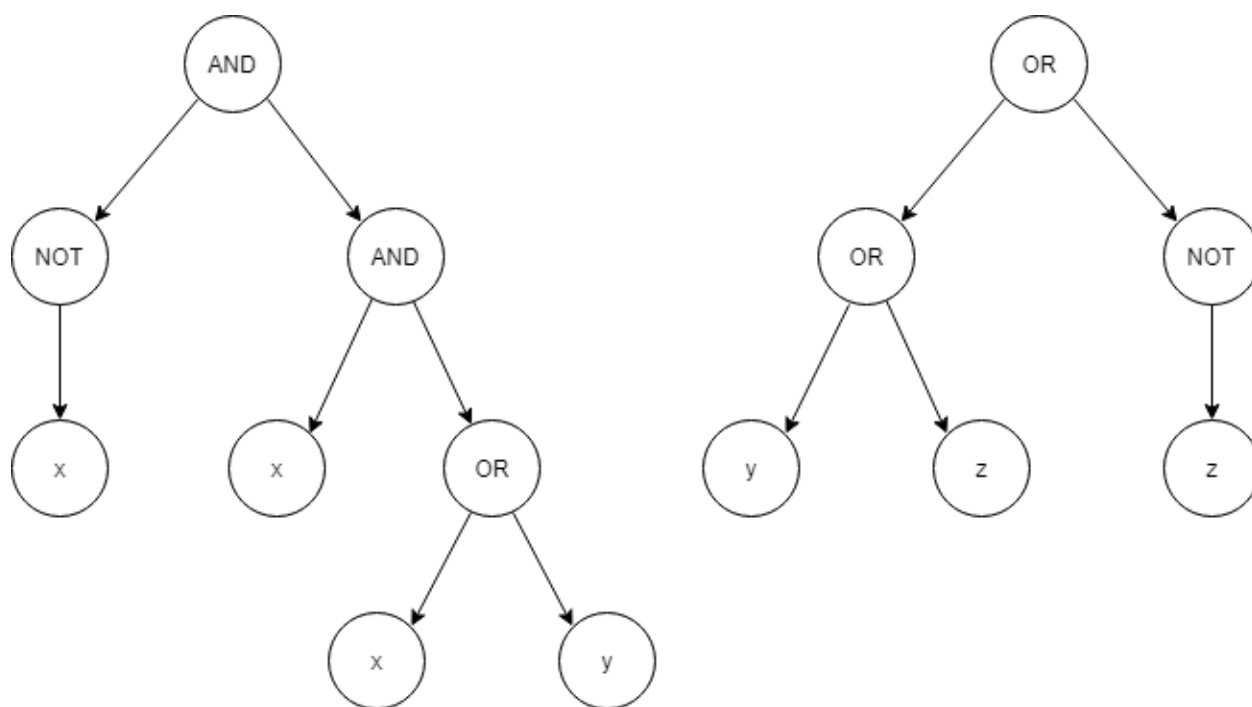


Figura 11.
Indivizii rezultați în urma încrucișării

3.2. Programare genetică pentru clasificare binară

În această lucrare vom studia dacă programarea genetică poate fi folosită cu succes în cazul clasificării binare. Modul în care acest studiu este realizat este prin căutarea unor expresii matematice. Argumentele problemei vor fi datele din setul de antrenament. Pentru fiecare expresie din setul de expresii se va parcurge întreg setul de date și se va verifica dacă pentru acele date clasificarea realizată este corectă sau nu.

Modul în care este abordată această problemă este prin găsirea unei probabilități în funcție de rezultatul dat de expresia matematică găsită. Acest lucru este realizat prin următoarea formulă matematică:

$$predictie(x) = \frac{1}{1 + e^x}$$

Argumentul funcției este chiar rezultatul expresiei matematice. Cu ajutorul acestei funcții se obține o valoare între 0 și 1. Valoarea returnată va fi comparată cu un prag și în funcție de rezultatul comparației se va face clasificarea. Dacă predicția este mai mare sau egală decât valoarea pragului datelor vor fi clasificate ca făcând din clasa pozitivă, iar în caz contrar din clasa negativă.

Graficul funcției este ilustrat în figura 12. Se poate observa ca pentru un prag cu valoarea de 0.5 rezultatul expresiei trebuie sa aibă valoare mai mica sau egala cu 0 pentru ca datele sa fie clasificate în clasa pozitivă și valoare mai mare de 0.5 pentru ca asignarea sa fie în clasa negativă.

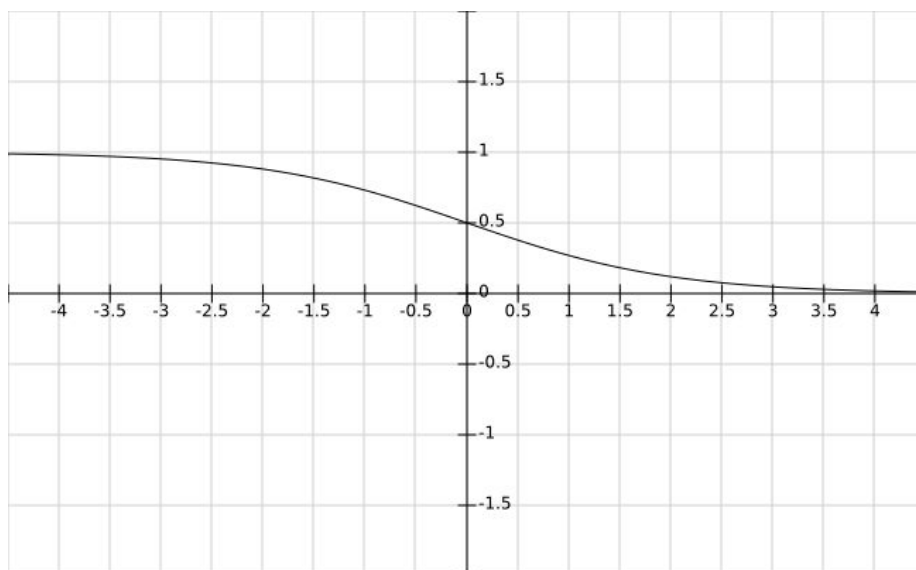


Figura 12.
Graficul funcției utilizate în clasificarea binară

3.3 Funcții fitness pentru clasificare

Pentru găsirea unor expresii care sa fie folosite pentru clasificare prin metoda programării genetice este necesară definirea unei funcții de fitness, cu ajutorul căreia sa se calculeze cat de bună este expresia în rezolvarea problemei de clasificare propusă, structura acestora fiind o consecință a funcției de fitness.

În studierea performantei programării genetice din această lucrare vor fi încercate mai multe funcții de fitness. Prima metodă încercata este cross-entropia, definita astfel:

$$H(p, q) = - \sum_x p(x) \log q(x)$$

În formula data $p(x)$ reprezintă probabilitatea vrută și $q(x)$ probabilitatea adevărată. În cazul acestei funcții de fitness scopul va fi minimizarea fitness-ului.

O alta funcție de fitness este numărarea numărului de erori. De asemenea și în acest caz se va urmări scăderea fitness-ului.

Metrica de măsurare a performantei unui clasificator din această lucrare va fi scorul F1, așadar o alta funcție de fitness folosită va fi chiar calcularea scorului F1 obținut pe setul de antrenament, așadar scopul va fi maximizarea fitness-ului.

Una dintre metodele care au fost folosite cu succes pentru îmbunătățirea performanței unui clasificator în cazul învățării dintr-un set de date dezechilibrat este învățarea bazată pe cost. O altă funcție de fitness este construită cu ajutorul acestui concept, calculându-se numărul de erori, însă costul clasificărilor greșite nu va mai același. Costurile au fost stabilite cu ajutorul raportului dintre instanțele clasei minoritare și a celei majoritare(1:80). Costul unei clasificări greșite a unei instanțe din clasa majoritare va avea cost 1, iar costul unei clasificări greșite a unei instanțe din clasa minoritară va fi 80.

4. Analiză experimentală

În acest capitol sunt analizate rezultatele obținute prin metodele de învățare automată prezentate și prin metoda programării genetice, astfel se determină eficiența celei din urmă. Acest lucru este realizat cu ajutorul unui set de date de test, din cadrul aceleiași competiții (“Internet of Things: Online Anomaly Detection for Drinking Water Quality”), care a fost folosit pentru a verifica performanțele metodelor trimise de către concurenți.

Astfel încât comparația dintre programarea genetică și metodele de învățare automată să fie cât mai obiectivă și relevantă în cazul unor metode de învățare automată a fost realizată optimizarea hiper-parametrilor, adică se caută parametri pentru aceste metode astfel încât rezultatele să fie cât mai bune. Un exemplu de hiper-parametru este adâncimea maximă a unui arbore de decizie, așadar hiper-parametrii sunt parametrii ai metodelor, ce sunt setați înainte de începerea propriu zisă a procesului de învățare.

O metodă foarte folosită pentru analiza performanței unui model de învățare automată este K-fold Cross Validation, care constă în împărțirea setului de date de antrenament în k subseturi, dintre care $k-1$ sunt folosite pentru antrenament și un subset este folosit pentru testare, procesul se repetă până când fiecare dintre subseturi a fost folosit pentru testare. Analiza performanței se va realiza prin media rezultatelor obținute. În figura 13 este ilustrat 3-fold Cross Validation, sub-seturile colorate în culoarea portocaliu fiind setul de test pentru iterația respectivă.



Figura 13.
Vizualizare a 3-fold Cross Validation

Căutarea grid se realizează prin definirea unui set de valori posibile pentru parametrii algoritmilor. Se realizează toate asignările posibile ale parametrilor cu aceste valori și se calculează performanța cu ajutorul Cross Validation. Se găsește astfel cea mai buna combinație de parametri și se îmbunătățește performanța algoritmului. Acest tip de căutare este folosit pentru găsirea hiper-parametrilor împreună cu 10-fold Cross Validation, setul de date fiind împărțit în 10 sub-seturi, conform principiilor prezentate anterior.

Pentru a compara performanțele algoritmilor se va folosi scorul F1, care a fost și criteriul oficial de departajare în cadrul competiției, de unde au fost preluate setul de date de antrenament și cel de test. această metrică este potrivită deoarece nu se ia în considerare numărul clasificărilor corecte pentru elementele din clasa negativă, dar este sensibilă la clasificările greșite. Un scor F1 ridicat indica faptul ca instanțele din clasa pozitivă sunt identificate într-o proporție ridicată în mod corect, fără a produce un număr ridicat de alarme false.

4.1 Arbori de decizie

Pentru arborii de decizie s-a realizat optimizarea hiper-parametrilor. Parametrul pentru evaluarea unei taieturi a fost ales criteriul Gini, parametrul înălțime maxima nu este limitata la o anumită valoare, iar argumentul `min_samples_split`, care reprezintă numărul minim de instanțe conținut de către o frunza a fost setat la numărul 10.

Rezultatele clasificării se regăsesc în următoarea matrice de confuzie:

	Negativ	Pozitiv
Negativ	104937	7820
Pozitiv	784	1545

Cu ajutorul acestei matrice putem observa și calcula cu ușurința următoarele aspecte ce ilustrează performanța clasificatorului:

$$TN = 104937 \quad TP = 1545$$

$$FN = 784 \quad FP = 7820$$

$$Precision = \frac{1545}{1545 + 7820} = 0.16497$$

$$Precision = \frac{1545}{1545 + 784} = 0.66337$$

$$F1 = \frac{2 \cdot precision \cdot recall}{precision + recall} = 0.26423$$

4.2 Regresie logistică

În cazul acestei metode de clasificare nu a fost cazul găsirii și setării unor hiper-parametri. Rezultatele obținute de acest clasificator sunt prezentate în următoarea matrice de confuzie:

	Negativ	Pozitiv
Negativ	112757	0
Pozitiv	1918	411

O prezentare mai în detaliu a performanței acestui clasificator sunt următoarele rezultate.

$$TN = 112757 \quad TP = 411$$

$$FN = 1918 \quad FP = 0$$

$$Precision = \frac{411}{411 + 0} = 1$$

$$Recall = \frac{411}{411 + 1918} = 0.17647$$

$$F1 = \frac{2 \cdot precision \cdot recall}{precision + recall} = 0.3$$

4.3 AdaBoost

Pentru această metodă de clasificare a fost identificat hiper-parametrul ce reprezintă numărul maxim de estimatori. Valoarea acestuia a fost setată la 50. Matricea de confuzie corespunzătoare acestui clasificator este:

	Negativ	Pozitiv
Negativ	109245	3512
Pozitiv	1351	978

Cu ajutorul matricei obținem cu ușurință următoarele rezultate:

$$TN = 109245 \quad TP = 978$$

$$FN = 1351 \quad FP = 3512$$

$$Precision = \frac{978}{978 + 3512} = 0.21781$$

$$Recall = \frac{978}{978 + 1351} = 0.41992$$

$$F1 = \frac{2 \cdot precision \cdot recall}{precision + recall} = 0.28683$$

4.4 Extra trees

În cazul acestei metode de clasificare au fost identificați următorii hyper-parametri: numărul de arbori a fost setat la valoarea de 40, criteriul de evaluare a unei tăieturi este indexul Gini, înălțimea nu este limitată la o valoare maximă, iar argumentul `min_samples_split`, care reprezintă numărul minim de instanțe conținute de către o frunză a fost setat la numărul 10.

Matricea de confuzie corespunzătoare acestui clasificator este:

	Negativ	Pozitiv
Negativ	112616	141
Pozitiv	1636	693

Rezultatele pentru acest clasificator sunt:

$$TN = 112616 \quad TP = 693$$

$$FN = 1636 \quad FP = 141$$

$$Precision = \frac{693}{693 + 141} = 0.83093$$

$$Recall = \frac{693}{693 + 1636} = 0.29755$$

$$F1 = \frac{2 \cdot precision \cdot recall}{precision + recall} = 0.43818$$

4.5 Programare genetică

Pentru programarea genetică a fost definit următorul set de primitive: minim, maxim, adunare, scădere, înmulțire, împărțire, negație, sinus și cosinus. Cu ajutorul acestor primitive se vor crea expresii matematice care iau ca date de intrare valorile din setul de date. În urma experimentelor a fost aleasa o funcție de fitness ce reprezenta numărul de clasificări greșite din setul de antrenament, iar în acest caz în mod evident se urmărește minimizarea fitness-ului.

O serie de alți parametri cu importanță deosebită au fost stabiliți astfel: dimensiunea populației are valoarea de 200, numărul de iterații realizate pentru căutarea expresiilor este 55 și modalitatea de selecție ale expresiilor este de tip turneu. Acest mod de selecție este realizat prin alegerea în mod aleatoriu a unui număr k de indivizi din populație și returnarea celui cu cel mai bun fitness. Numărul de indivizi pentru selecție a fost stabilit la 25. În cazul mutației punctul de mutație poate fi înlocuit cu un arbore produs în mod aleator, ce are înălțime cuprinsă între 0 și 3, iar probabilitatea de mutație a unui individ este 0.2. Încrucișarea a doi

indivizi este prin alegerea în mod aleator a cate unui punct de încrucișare pentru fiecare dintre ei, în modul prezentat vizual în capitolul anterior. Probabilitatea de încrucișare a doi indivizi a fost stabilita la valoarea de 0.5 .

Pentru acești parametri au fost obținute rezultatele din tabelul următor, care conține 30 de expresii. Este important de specificat ca pentru obținerea acestuia au fost necesare mai mult de 30 de iterații, deoarece nu au fost incluse rezultatele care aveau probabilitate mare de overfitting, adică au fost eliminate soluțiile în cadrul cărora execuția programului a înregistrat pentru un număr relativ mare de iterații scăderi mici ale valorii funcției fitness, fapt ce poate indica o adaptare a expresiei pentru zgomotul și detaliile din setul de date, ceea ce afectează calitatea clasificatorului pentru date noi.

Expresia:	Fitness	Scorul F1
mul(sub(pH, add(neg(mul(Fm, Cl_2)), add(cos(Redox), sub(Fm, mul(sub(Fm, add(neg(mul(Fm, Cl_2)), add(sub(sub(Fm, mul(neg(neg(mul(add(neg(mul(Fm, sub(Cl, Cl_2))), mul(sub(Fm, add(neg(mul(Fm, Cl_2)), add(sub(sub(Fm, mul(neg(neg(mul(add(neg(mul(Fm, mul(Cl_2, Cl_2))), add(sub(sub(Fm, mul(neg(neg(mul(Fm, Cl_2))), pH)), -1, Redox))), Cl_2))), pH)), -1, Redox))), pH)), Cl_2))), pH)), -1, Redox))), add(Redox, sub(Fm, mul(neg(pH), neg(neg(mul(add(neg(mul(Fm, mul(Cl_2, Cl_2))), add(sub(sub(Fm, mul(neg(neg(mul(Fm, Cl_2))), add(neg(mul(Fm, sub(Cl, Cl_2))), mul(sub(Fm, add(neg(mul(Fm, Cl_2)), add(sub(sub(Fm, mul(neg(neg(mul(add(neg(mul(Fm, mul(Cl_2, Cl_2))), add(sub(sub(Fm, mul(neg(neg(mul(Fm, Cl_2))), pH)), -1, Cl)), Cl_2))), pH)), -1, Redox))), pH))), -1, Redox))), mul(Fm, mul(Cl_2, Cl_2))))))))), Redox)	1101	0.26544
mul(add(mul(add(mul(min(safeDiv(add(Redox, Fm), sub(Fm_2, add(mul(add(mul(min(Fm, add(neg(pH), Cl)), Redox), sub(Fm, Redox)), cos(max(Cl, pH))), sin(pH))), mul(max(mul(Cl, pH), safeDiv(Fm, Redox)), sub(Fm, Redox))), Redox), min(sub(neg(cos(Redox)), Redox), mul(add(neg(pH), sub(Fm, Redox)), cos(max(Cl, pH))))), cos(max(sub(pH, min(add(Cl, Cl_2), add(min(add(Cl, Cl_2), add(neg(pH), neg(Cl_2))), neg(Cl_2))), sub(safeDiv(sin(Cl_2), add(Leit, Tp)), sin(cos(Cl))))), Cl), safeDiv(sub(sub(Fm, Redox), min(sub(Fm, Redox), Cl)), Cl_2))	1132	0.3811
mul(add(mul(add(mul(min(safeDiv(add(Redox, Fm), sub(Fm_2, add(mul(add(mul(min(Fm, add(neg(pH), Cl)), Redox), sub(Fm, Redox)), cos(max(Cl, pH))), sin(pH))), mul(max(mul(Cl, pH), safeDiv(Fm, Redox)), sub(Fm, Redox))), Redox), min(sub(sub(sub(Fm, Redox), min(sub(Fm, Redox), Cl)), Redox), mul(add(mul(min(Fm, add(neg(pH), Cl)), Redox), max(Leit, Fm)), cos(max(Cl, pH))))), cos(max(sub(pH, min(add(Cl, Cl_2), add(min(add(Cl, Cl_2), add(neg(pH), neg(Cl_2))), neg(Cl_2))), sub(safeDiv(sin(Cl_2), add(Leit, Tp)), sub(Fm, Redox))))), Cl), safeDiv(sub(sub(Fm, Redox), min(sub(Fm, Redox), Cl)), Cl_2))	1138	0.37574
safeDiv(max(safeDiv(min(Cl_2, Cl_2), Fm), cos(max(Cl_2, pH))), mul(neg(cos(sub(pH, min(Leit, Cl_2)))), safeDiv(sub(Fm, Redox), min(add(cos(max(Cl_2, pH)), sin(sin(pH))), Fm)))	1189	0.40356
mul(sub(safeDiv(add(cos(pH), Leit), add(Fm_2, Leit)), Cl_2, sub(mul(sin(Cl_2), sub(sub(pH, Cl_2), sin(add(Cl_2, sub(sub(pH, Cl_2), sin(add(sin(sub(pH, sin(sub(sub(pH, Cl_2), cos(pH))))), Tp))))))), cos(sin(sin(sub(pH, sin(sub(sub(pH, Cl_2), cos(pH))))))))	1040	0.34673
mul(neg(sin(sin(neg(pH)))), sub(sub(mul(Fm, Trueb), safeDiv(Fm, pH)), mul(sin(add(sub(safeDiv(Trueb, cos(sub(Redox, neg(sin(Cl))))), mul(sin(pH), neg(pH))), Trueb)), add(add(Tp, neg(sub(mul(add(add(Leit, neg(sub(mul(Fm, Trueb), neg(sub(mul(Fm, Trueb), safeDiv(Fm_2, pH))))), sub(Redox, Fm)), Trueb), safeDiv(Fm, pH))), sub(Redox, Fm))))	1114	0.21975
mul(cos(safeDiv(cos(sin(neg(safeDiv(mul(neg(sin(neg(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(Redox, Fm))), Fm), sub(Leit, Fm))))), Fm), sub(sub(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(pH, Fm))), Fm), sub(sub(sub(Tp, pH), pH), Cl_2), sub(Cl_2, neg(pH))), Fm))))), cos(sub(sub(Cl_2, neg(pH)), Cl))), add(sub(Redox, Fm), mul(neg(sin(neg(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(pH, Fm))), Fm), sub(Leit, Fm))))), Fm)))	1137	0.44187
mul(add(neg(neg(neg(sin(safeDiv(Fm, Fm_2))))), mul(Cl_2, add(safeDiv(mul(Redox, pH), mul(Leit, Tp)), min(sub(sub(sub(Fm, max(safeDiv(Tp, cos(Redox)), add(Tp, mul(pH, Cl_2))), max(safeDiv(Fm, Fm_2), Cl)), Trueb), max(Tp, Cl))))), min(mul(max(add(Trueb, Cl), max(Leit, Redox)), Redox), safeDiv(Redox, add(sub(sub(Fm, Redox), safeDiv(mul(Leit, Tp), add(sub(sub(Fm, Redox), cos(Redox)), neg(sin(safeDiv(mul(safeDiv(Trueb, Fm_2), max(max(Tp, Cl), Redox)), Fm_2))))), Tp)))	1190	0.14204

neg(mul(Cl, cos(add(pH, mul(Cl, pH))))), sub(mul(Cl_2, safeDiv(mul(sin(neg(cos(pH))), pH), sin(max(neg(cos(pH)), 0))), neg(cos(add(pH, mul(Cl, pH))))), pH), sin(max(neg(cos(pH)), 0))), neg(cos(add(pH, Cl_2))), 0)))))))))		
mul(mul(Cl_2, safeDiv(add(safeDiv(sin(Trueb), safeDiv(sin(add(add(sin(sin(sin(safeDiv(sin(sin(pH)), Cl_2))), Cl_2, Cl_2)), add(neg(mul(sin(add(add(sin(safeDiv(sin(sin(sin(add(add(neg(safeDiv(sin(sin(safeDiv(sin(sin(pH)), Cl_2))), Cl_2))), Cl_2, Cl_2))), mul(sin(sin(sin(safeDiv(sin(sin(sin(sin(safeDiv(sin(sin(pH)), Cl_2))), Cl_2, Cl_2))), Redox))), Redox), safeDiv(sin(add(add(sin(safeDiv(sin(sin(pH)), Cl_2), Cl_2, add(add(sin(sin(sin(sin(sin(sin(sin(safeDiv(sin(sin(pH)), Cl_2))))))))) Cl_2, Cl_2))), add(neg(mul(safeDiv(sin(sin(pH)), Cl_2, mul(safeDiv(sin(sin(pH)), Cl_2, safeDiv(add(sin(add(add(sin(sin(sin(safeDiv(sin(sin(pH)), Cl_2))), Cl_2, Cl_2))), Redox))), sub(Fm, add(safeDiv(sin(Trueb), safeDiv(sin(sin(add(add(sin(sin(sin(safeDiv(sin(sin(pH)), Cl_2))), Cl_2, Cl_2))), add(neg(mul(Fm, safeDiv(sin(sin(sin(sin(sin(safeDiv(sin(sin(pH)), Cl_2))), Cl_2))), Redox))), Redox)))	859	0.32385
min(add(safeDiv(Cl, neg(add(Cl_2, Cl))), mul(Cl_2, pH)), mul(sub(Fm, Redox), neg(safeDiv(safeDiv(Cl, min(pH, Fm)), safeDiv(neg(min(cos(pH), add(pH, Fm))), cos(add(pH, pH))))))	1101	0.3624
mul(add(mul(add(mul(min(safeDiv(add(Redox, Fm), sub(Fm_2, add(mul(add(mul(min(Fm, add(neg(pH), Cl)), Redox), sub(Fm, Redox)), cos(max(Cl, pH))), Cl))), mul(max(mul(Cl, pH), safeDiv(Fm, Redox)), min(Cl, Fm))), Redox), min(sub(Fm, Redox), Cl)), cos(max(sub(pH, min(add(Cl, Cl_2), add(min(add(Cl, Cl_2), add(neg(pH), neg(Cl_2))), neg(Cl_2))), sub(safeDiv(sin(Cl_2), add(Leit, Tp)), sin(cos(Tp))))), Cl), safeDiv(sub(sub(Fm, Redox), min(sub(Fm, Redox), max(sub(pH, min(add(Cl, Cl_2), add(min(add(Cl, Cl_2), add(neg(pH), neg(Cl_2))), neg(Cl_2))), sub(safeDiv(sin(Cl_2), add(Leit, Tp)), sin(cos(Cl))))), Cl_2))	1154	0.44996
mul(add(neg(neg(neg(sin(safeDiv(Fm, Fm_2))))), mul(Cl_2, add(safeDiv(mul(Redox, pH), mul(Leit, Tp)), min(sub(sub(sub(Fm, max(safeDiv(Tp, cos(Redox)), add(Tp, mul(pH, Cl_2))), max(safeDiv(Fm, Fm_2), Cl)), Trueb), max(Tp, Cl))), min(mul(max(add(Trueb, Cl), max(Leit, Redox)), Redox), safeDiv(Redox, add(sub(sub(Fm, Redox), safeDiv(mul(Leit, Tp), add(sub(sub(Fm, Redox), cos(Redox)), neg(sin(safeDiv(mul(safeDiv(Trueb, Fm_2), max(max(Tp, Cl), Redox)), Fm_2))))), Tp)))	1190	0.14204
neg(min(max(cos(min(Fm_2, Fm))), max(add(add(pH, min(max(add(Cl_2, mul(pH, add(safeDiv(mul(Cl_2, Fm_2), Redox), mul(pH, Cl_2))), min(safeDiv(mul(Redox, pH), neg(pH))), Cl)), add(add(mul(pH, Cl_2, Cl_2, Cl_2))), pH), max(cos(Cl_2), sin(Cl_2))), cos(add(pH, min(max(add(min(max(Fm_2, max(add(safeDiv(Cl_2, Redox), pH), max(cos(Cl_2), sin(Cl_2))), cos(add(pH, min(max(add(safeDiv(Cl_2, sin(Redox)), pH), add(pH, Cl_2)), add(add(pH, Cl_2, Cl_2))), mul(pH, add(safeDiv(mul(Cl_2, Fm_2), Redox), mul(pH, Cl_2))), min(Cl, Cl)), add(add(pH, Cl_2, Cl_2))))))	1005	0.40762
mul(add(mul(add(mul(min(safeDiv(add(Redox, Fm), sub(Fm_2, add(mul(add(mul(min(Fm, add(neg(pH), Cl)), Redox), sub(Fm, Redox)), cos(max(Cl, pH))), sin(pH))), mul(max(mul(Cl, pH), safeDiv(Fm, Redox)), sub(Fm, Redox))), Redox), min(sub(sub(sub(Fm, Redox), min(sub(Fm, Redox), Cl)), Redox), mul(add(mul(min(Fm, add(neg(pH), Cl)), Redox), max(Leit, Fm)), cos(max(Cl, pH))), cos(max(sub(pH, min(add(Cl, Cl_2), add(min(add(Cl, Cl_2), add(neg(pH), neg(Cl_2))), neg(Cl_2))), sub(safeDiv(sin(Cl_2), add(Leit, Tp)), sub(Fm, Redox))), Cl), safeDiv(sub(sub(Fm, Redox), min(sub(Fm, Redox), Cl)), Cl_2))	1138	0.37574
neg(safeDiv(cos(mul(add(-1, mul(Cl, Trueb)), pH)), mul(Fm, mul(safeDiv(safeDiv(Fm, Cl_2), Cl_2), mul(safeDiv(safeDiv(Fm_2, Cl_2, Cl_2), add(mul(safeDiv(pH, mul(add(-1, mul(Cl, pH))), pH)), Tp), mul(safeDiv(Tp, mul(safeDiv(-1, mul(Cl, pH)), Tp)), Tp))))))	1052	0.25608
mul(add(mul(add(mul(min(safeDiv(add(Redox, Fm), sub(Fm_2, add(mul(add(mul(min(Fm, add(neg(pH), Cl)), Redox), sub(Fm, Redox)), cos(max(Cl, pH))), sin(pH))), mul(max(mul(Cl, pH), safeDiv(Fm, Redox)), sub(Fm, Redox))), Redox), min(sub(neg(cos(Redox)), Redox), mul(add(neg(pH), sub(Fm, Redox)), cos(max(Cl, pH))), cos(max(sub(pH, min(add(Cl, Cl_2), add(min(add(Cl, Cl_2), add(neg(pH), neg(Cl_2))), neg(Cl_2))), sub(safeDiv(sin(Cl_2), add(Leit, Tp)), sin(cos(Cl))))), Cl), safeDiv(sub(sub(Fm, Redox), min(sub(Fm, Redox), Cl)), Cl_2))	1132	0.3811
min(sin(safeDiv(max(pH, Trueb), safeDiv(sin(safeDiv(Cl, Cl_2), cos(Cl))), min(safeDiv(safeDiv(Cl, sin(safeDiv(sin(add(pH, Cl_2), Cl_2))), sin(min(Fm, pH))), sub(safeDiv(Trueb, cos(add(pH, min(safeDiv(Cl, Cl_2), sub(safeDiv(sin(safeDiv(Cl, Cl_2), cos(add(pH, Cl_2))), safeDiv(max(pH, Trueb), safeDiv(sub(safeDiv(sin(add(min(mul(Cl_2, Tp), safeDiv(Fm, Fm_2), cos(add(Redox, Cl))), cos(add(pH, Cl_2))), safeDiv(max(pH, Trueb), sin(max(mul(Cl, Fm_2), cos(Leit))))), cos(add(pH, Cl_2))))))))) safeDiv(cos(sin(add(Cl, Cl_2))), safeDiv(add(pH, min(Cl_2, sub(safeDiv(sin(safeDiv(Cl, Cl_2), cos(add(pH, Cl_2))), safeDiv(max(pH, Trueb), Cl_2))), cos(add(pH, Cl_2))))))	726	0.30679
mul(cos(safeDiv(cos(sin(neg(safeDiv(mul(neg(sin(neg(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(Redox, Fm))), Fm), sub(Leit, Fm))), Fm), sub(sub(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(pH, Fm))), Fm), sub(sub(sub(Tp, pH), pH), Cl_2), sub(Cl_2, neg(pH))), Fm))), cos(sub(sub(Cl_2, neg(pH)), Cl))), add(sub(Redox, Fm), mul(neg(sin(neg(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(pH, Fm))), Fm), sub(Leit, Fm))), Fm)))	1137	0.44187

add(neg(cos(mul(sub(Cl, pH), Cl))), mul(min(sin(cos(mul(pH, max(sub(max(sin(sin(Cl)), pH), max(Fm, Cl)), sin(sin(sin(add(pH, Cl))))))), min(sin(cos(mul(pH, Cl))), Cl)), mul(max(max(Cl, Leit), min(sin(min(Cl_2, cos(mul(sin(add(safeDiv(add(Leit, Cl_2), mul(pH, Tp)), neg(Cl))), sin(cos(mul(pH, Cl))))))), Cl)), mul(Cl, min(Cl_2, Cl))))	1120	0.13705
mul(mul(Cl_2, safeDiv(mul(max(mul(mul(Cl_2, safeDiv(mul(Redox, pH), Cl)), Cl), mul(safeDiv(mul(Redox, pH), Cl), safeDiv(mul(Cl, pH), Redox))), pH), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), add(Cl, sub(Redox, max(safeDiv(safeDiv(Redox, pH), Cl), max(mul(mul(Cl_2, safeDiv(mul(safeDiv(mul(mul(Cl_2, safeDiv(mul(mul(mul(Cl_2, safeDiv(mul(max(mul(mul(Cl_2, safeDiv(mul(max(mul(mul(Cl_2, safeDiv(mul(Redox, pH), Cl)), Cl), mul(safeDiv(mul(Fm, pH), Cl), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), Tp), Cl_2))), Cl), mul(safeDiv(mul(Cl, pH), Cl), safeDiv(max(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), pH), Cl)), mul(mul(Cl_2, Cl), pH)), pH), Cl)), Cl), mul(Fm, safeDiv(Cl_2, mul(mul(mul(Cl_2, Cl), pH), Redox))), pH), Cl)), Cl_2, mul(safeDiv(mul(Redox, pH), Cl), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))))))	995	0.29119

Se pot face câteva observații cu ajutorul acestui tabel. În primul rând un fitness scăzut, adică un număr mic de erori pentru setul de date de antrenament nu înseamnă neapărat ca expresia data va clasifica bine și instanțele din setul de testare, acest lucru fiind cauzat de overfitting. Pentru expresiile din tabel s-au obținut următoarele rezultate:

$$F1Minim = 0.13705$$

$$F1Maxim = 0.4774$$

$$MediaF1 = 0.31582$$

Deviația standard măsoară dispersia unui set de date relativ la medie și are următoarea formula:

$$DevStandard = \sqrt{\frac{\sum_{i=1}^n (x_i - medie)^2}{n - 1}} \Rightarrow DevStandard = 0.09661$$

4.6 Experimentări pentru clasificare cu programare genetică

În cadrul acestui sub-capitol vor fi prezentate câteva modalități prin care s-a încercat creșterea performanțelor clasificării prin metoda programării genetice. În primul rând au fost încercate mai multe funcții fitness. Acestea au fost prezentate anterior în sub-capitolul “Funcții fitness pentru clasificare”, din capitolul anterior.

În cazul cross-entropiei rezultatele au fost destul de slabe în sensul ca evoluția programului consta în apropierea probabilității instanțelor din clasa majoritara de 0, fără a realiza într-o măsură semnificativă apropierea instanțelor din clasa minoritara de 1, deci clasificările realizate cu expresiile care au fost evaluate în acest mod aveau un număr mare de false negative.

În cazul funcției fitness reprezentată chiar de valoarea metricii F1, rezultatele au fost asemănătoare cu cazul în s-a minimizat numărul de erori, din acest motiv nu s-ar justifica o prezentare mai detaliată a rezultatelor.

O alta funcție fitness utilizată a fost una inspirată de învățare bazată pe cost, deci costul unei clasificări incorecte nu a mai fost considerat egal pentru ambele clase din setul de date. Costurile au fost următoarele:

	Negativ	Pozitiv
Negativ	0	1
Pozitiv	80	0

În cazul acestei funcții fitness s-a constatat că clasificatorii obținuți au înregistrat un recall bun, mai mare decât în cazul celorlalte metode, însă precizia a fost foarte slabă, din acest motiv scorul F1 a fost scăzut.

Un alt mod prin care s-a încercat creșterea performanței a fost prin metoda de sampling SMOTE. Cu ajutorul acestei metode a fost dublat numărul de instanțe din clasa minoritară, fiind folosită valoarea 5 pentru determinarea celor mai apropiați vecini. Cu folosirea noului set de date nu s-au obținut îmbunătățiri semnificative pentru scorurile F1 obținute, diferențele fiind irelevante având în vedere nedeterminismul presupus de programarea genetică. Consider însă că metodele de sampling reprezintă o direcție importantă care trebuie luată în considerare în încercarea de a aprofunda studiul acestei probleme și de a îmbunătăți performanțele clasificării.

4.7 Ansamblu de clasificatori bazați pe programare genetică

Analizând rezultatele obținute ce au fost introduse în tabelul prezentat se poate constata că apare supra-specializarea expresiilor, în sensul că o valoare mică a fitnessului pentru un anumit individ nu presupune că acesta va clasifica bine și instanțele din setul de testare, acest fapt fiind cauzat de faptul că modelul este influențat și de zgomot și detaliile din setul de antrenament, afectând astfel performanța clasificării de instanțe noi.

Pentru a micșora supra-specializarea se va încerca o metodă inspirată de modul în care metoda Extra trees funcționează, adică se va crea un ansamblu de clasificatori mai slabi, decizia unei clasificări pentru o instanță fiind dată de către majoritatea clasificărilor oferite de către clasificatorii din ansamblu. Se poate observa că expresiile create au formă, structură și dimensiune diferită, de aceea un astfel de ansamblu reprezintă un model promițător pentru scăderea supra-specializării. De asemenea, acest tip de ansamblu este diferit față de cel propus de AdaBoost, în care construcția clasificatorilor din ansamblu este influențată de performanțele celor deja aflați în ansamblu și li se atribuie ponderi, acest tip de ansamblu este de tip boosting. Ansamblul creat pentru acest experiment este de tip bagging, fiecare model

rulând în mod independent, iar apoi decizia clasificării se realizează în funcție de majoritatea clasificărilor ale elementelor din ansamblu.

O parte dintre expresiile folosite pentru clasificare au fost preluate din tabelul prezentat anterior, dar au fost introduse și altele obținute prin modificarea parametrilor pentru programarea genetică, pentru a obține expresii cat mai variate, scăzându-se astfel probabilitatea ca acestea să fie afectate de supra-specializare în mod asemănător.

Ansamblul creat este format din expresiile din următorul tabel:

Expresii:
$\min(\text{mul}(\text{mul}(\text{Cl}, \text{neg}(\cos(\text{sub}(\text{Cl}, \text{pH}))))), \text{mul}(\text{mul}(\text{Cl}, \text{mul}(\sin(\cos(\text{sub}(\cos(\sin(\sin(\text{mul}(\text{neg}(\text{Redox}), \cos(\text{sub}(\cos(\sin(\text{pH})), \text{pH})))))), \text{pH}))), \text{pH})), \text{add}(\text{neg}(\text{Trueb}), \text{Cl}))), \text{mul}(\text{mul}(\text{Cl}, \text{mul}(\sin(\cos(\text{sub}(\cos(\sin(\sin(\text{mul}(\text{neg}(\text{mul}(\text{neg}(\text{Redox}), \text{neg}(\text{Cl}_2))), \text{neg}(\text{Cl}_2))))), \text{pH}))), \text{Cl}_2)), \text{add}(\text{neg}(\text{Cl}_2), \text{Cl})))$
$\text{mul}(\text{sub}(\text{safeDiv}(\text{add}(\cos(\text{pH}), \text{Leit}), \text{add}(\text{Fm}_2, \text{Leit})), \text{Cl}_2), \text{sub}(\text{mul}(\sin(\text{Cl}_2), \text{sub}(\text{sub}(\text{pH}, \text{Cl}_2), \sin(\text{add}(\text{Cl}_2, \text{sub}(\text{sub}(\text{pH}, \text{Cl}_2), \sin(\text{add}(\sin(\text{sub}(\text{pH}, \sin(\text{sub}(\text{sub}(\text{pH}, \text{Cl}_2), \cos(\text{pH}))))), \text{Tp}))))), \cos(\sin(\sin(\text{sub}(\text{pH}, \sin(\text{sub}(\text{sub}(\text{pH}, \text{Cl}_2), \cos(\text{pH}))))))))))$
$\text{mul}(\text{neg}(\sin(\sin(\text{neg}(\text{pH})))), \text{sub}(\text{sub}(\text{mul}(\text{Fm}, \text{Trueb}), \text{safeDiv}(\text{Fm}, \text{pH})), \text{mul}(\sin(\text{add}(\text{sub}(\text{safeDiv}(\text{Trueb}, \cos(\text{sub}(\text{Redox}, \text{neg}(\sin(\text{Cl}))))), \text{mul}(\sin(\text{pH}), \text{neg}(\text{pH}))), \text{Trueb})), \text{add}(\text{add}(\text{Tp}, \text{neg}(\text{sub}(\text{mul}(\text{add}(\text{Leit}, \text{neg}(\text{sub}(\text{mul}(\text{Fm}, \text{Trueb}), \text{neg}(\text{sub}(\text{mul}(\text{Fm}, \text{Trueb}), \text{safeDiv}(\text{Fm}_2, \text{pH}))))), \text{sub}(\text{Redox}, \text{Fm})), \text{Trueb}), \text{safeDiv}(\text{Fm}, \text{pH}))), \text{sub}(\text{Redox}, \text{Fm}))))$
$\text{mul}(\text{add}(\text{mul}(\text{add}(\text{mul}(\min(\text{safeDiv}(\text{add}(\text{Redox}, \text{Fm}), \text{sub}(\text{Fm}_2, \text{add}(\text{mul}(\text{add}(\text{mul}(\min(\text{Fm}, \text{add}(\text{neg}(\text{pH}), \text{Cl})), \text{Redox}), \text{sub}(\text{Fm}, \text{Redox})), \cos(\max(\text{Cl}, \text{pH}))), \text{Cl}))), \text{mul}(\max(\text{mul}(\text{Cl}, \text{pH}), \text{safeDiv}(\text{Fm}, \text{Redox})), \min(\text{Cl}, \text{Fm}))), \text{Redox}), \min(\text{sub}(\text{Fm}, \text{Redox}), \text{Cl})), \cos(\max(\text{sub}(\text{pH}, \min(\text{add}(\text{Cl}, \text{Cl}_2), \text{add}(\min(\text{add}(\text{Cl}, \text{Cl}_2), \text{add}(\text{neg}(\text{pH}), \text{neg}(\text{Cl}_2))), \text{neg}(\text{Cl}_2))), \text{sub}(\text{safeDiv}(\sin(\text{Cl}_2), \text{add}(\text{Leit}, \text{Tp}))), \sin(\cos(\text{Tp}))))), \text{Cl}), \text{safeDiv}(\text{sub}(\text{sub}(\text{Fm}, \text{Redox}), \min(\text{sub}(\text{Fm}, \text{Redox}), \max(\text{sub}(\text{pH}, \min(\text{add}(\text{Cl}, \text{Cl}_2), \text{add}(\min(\text{add}(\text{Cl}, \text{Cl}_2), \text{add}(\text{neg}(\text{pH}), \text{neg}(\text{Cl}_2))), \text{neg}(\text{Cl}_2))), \text{sub}(\text{safeDiv}(\sin(\text{Cl}_2), \text{add}(\text{Leit}, \text{Tp}))), \sin(\cos(\text{Cl}))))), \text{Cl}_2))$
$\text{safeDiv}(\text{mul}(\text{add}(\text{neg}(\text{neg}(\sin(\text{safeDiv}(\text{Cl}, \cos(\text{mul}(\text{add}(\text{neg}(\text{neg}(\text{pH})), \text{Cl}), \text{neg}(\sin(\text{Leit}))))))), \text{add}(\text{neg}(\text{neg}(\sin(\text{pH}))), \text{Cl})), \text{neg}(\text{Cl})), \text{safeDiv}(\cos(\text{safeDiv}(\max(\text{safeDiv}(\cos(\text{safeDiv}(\max(\text{pH}, \text{pH}), \text{sub}(\text{neg}(\sin(\text{pH})), \text{Fm}))), \cos(\text{mul}(\text{add}(\text{neg}(\text{neg}(\sin(\text{pH}))), \text{Cl}), \text{neg}(\text{add}(\text{add}(\text{sub}(\text{pH}, \text{Cl}), \text{safeDiv}(\text{Trueb}, \text{Redox})), \min(\sin(\text{Cl}_2), \text{add}(\text{Tp}, \text{Cl}_2))))))), \text{safeDiv}(\text{add}(\max(\text{Trueb}, \text{pH}), \text{add}(\text{pH}, \text{Leit})), \cos(\text{neg}(\text{Cl}_2))), \text{sub}(\text{Tp}, \text{Fm}))), \cos(\text{mul}(\text{add}(\text{neg}(\text{neg}(\sin(\text{pH}))), \text{Cl}), \text{neg}(\text{pH}))))$
$\text{mul}(\cos(\text{safeDiv}(\cos(\sin(\text{neg}(\text{safeDiv}(\text{mul}(\text{neg}(\sin(\text{neg}(\text{safeDiv}(\text{mul}(\text{neg}(\text{safeDiv}(\text{mul}(\text{Cl}, \text{Fm}), \text{sub}(\text{Redox}, \text{Fm}))), \text{Fm}), \text{sub}(\text{Leit}, \text{Fm}))))), \text{Fm}), \text{sub}(\text{sub}(\text{safeDiv}(\text{mul}(\text{neg}(\text{safeDiv}(\text{mul}(\text{Cl}, \text{Fm}), \text{sub}(\text{pH}, \text{Fm}))), \text{Fm}), \text{sub}(\text{sub}(\text{sub}(\text{Tp}, \text{pH}), \text{pH}), \text{Cl}_2)), \text{sub}(\text{Cl}_2, \text{neg}(\text{pH}))), \text{Fm}))), \cos(\text{sub}(\text{sub}(\text{Cl}_2, \text{neg}(\text{pH}))), \text{Cl}))), \text{add}(\text{sub}(\text{Redox}, \text{Fm}), \text{mul}(\text{neg}(\sin(\text{neg}(\text{safeDiv}(\text{mul}(\text{neg}(\text{safeDiv}(\text{mul}(\text{Cl}, \text{Fm}), \text{sub}(\text{pH}, \text{Fm}))), \text{Fm}), \text{sub}(\text{Leit}, \text{Fm}))))), \text{Fm})))$
$\text{mul}(\cos(\text{pH}), \text{sub}(\text{mul}(\text{sub}(\text{Fm}, \text{Redox}), \text{neg}(\text{add}(\sin(\sin(\text{pH})), \cos(\text{pH})))), \text{mul}(\text{pH}, \text{mul}(\text{pH}, \sin(\text{add}(\text{Tp}, \text{Tp}))))))$
$\text{mul}(\text{mul}(\sin(\sin(\sin(\text{add}(\sin(\sin(\text{add}(\text{add}(\text{Cl}_2, \text{pH}), \text{add}(\text{Cl}_2, \text{Cl}_2))), \text{add}(\text{Cl}_2, \text{Cl}_2))), \cos(\text{neg}(\text{add}(\text{add}(\text{Cl}_2, \text{pH}), \text{add}(\text{Cl}_2, \text{pH})))), \text{safeDiv}(\text{add}(\text{add}(\text{mul}(\text{safeDiv}(\text{neg}(\text{add}(\text{add}(\text{Cl}_2, \text{pH}), \text{add}(\text{Cl}_2, \text{pH}))), \cos(\text{add}(\text{add}(\text{Cl}_2, \text{pH}), \text{add}(\text{Cl}_2, \text{Cl}_2))), \text{pH}), \text{add}(\text{mul}(\text{Cl}_2, \text{Tp}), \text{neg}(\text{sub}(\text{Redox}, \text{add}(\text{Tp}, \text{neg}(\text{sub}(\sin(\text{add}(\text{mul}(\text{Cl}, \text{pH}), \text{Cl}_2))), \text{add}(\text{add}(\text{add}(\sin(\cos(\text{sub}(\text{Fm}_2, \text{Cl}_2))), \text{mul}(\text{Cl}_2, \text{Redox})), \text{pH}), \text{Fm}))))))), \text{Fm}), \text{sub}(\text{Redox}, \text{add}(\text{add}(\text{add}(\text{Cl}_2, \text{mul}(\text{Cl}_2, \text{Redox})), \text{add}(\text{Cl}_2, \text{pH})), \text{Fm})))$
$\text{neg}(\text{safeDiv}(\sin(\cos(\text{mul}(\text{safeDiv}(\text{mul}(\text{Tp}, \text{Cl}), \text{neg}(\cos(\text{pH}))), \text{Cl}))), \text{sub}(\text{safeDiv}(\text{pH}, \text{sub}(\text{sub}(\cos(\cos(\cos(\text{neg}(\cos(\text{add}(\text{mul}(\text{safeDiv}(\text{mul}(\text{Tp}, \text{Cl}), \text{pH}), \text{Cl}), \text{add}(\text{Tp}, \text{Tp}))))), \text{neg}(\text{mul}(\text{safeDiv}(\text{mul}(\text{Tp}, \text{Cl}), \text{pH}), \text{mul}(\text{safeDiv}(\text{mul}(\text{Tp}, \text{mul}(\text{safeDiv}(\text{mul}(\text{Tp}, \text{Cl}), \text{neg}(\cos(\text{pH}))), \text{Cl})), \text{neg}(\cos(\text{pH}))), \text{Cl}))), \text{safeDiv}(\text{safeDiv}(\text{mul}(\text{safeDiv}(\text{mul}(\text{Tp}, \text{Cl}), \text{neg}(\cos(\text{pH}))), \text{Cl}), \text{neg}(\cos(\text{pH}))), \text{neg}(\cos(\text{pH}))), \text{mul}(\text{safeDiv}(\cos(\cos(\text{Trueb})), \text{neg}(\cos(\text{pH}))), \text{Cl})))$
$\text{mul}(\text{mul}(\text{mul}(\text{mul}(\text{Leit}, \text{Cl}_2), \text{sub}(\text{mul}(\text{mul}(\text{Leit}, \text{Cl}_2), \text{sub}(\text{mul}(\text{mul}(\text{sub}(\text{mul}(\text{Leit}, \text{add}(\text{mul}(\text{Leit}, \text{Cl}_2), \text{neg}(\text{Tp}))), \text{Redox}), \text{Cl}_2), \text{add}(\text{mul}(\text{sub}(\text{mul}(\text{Leit}, \text{add}(\text{mul}(\text{Leit}, \text{Cl}_2), \text{neg}(\text{Tp}))), \text{Redox}), \text{Cl}_2), \text{neg}(\text{mul}(\text{Leit}, \text{Cl}_2))), \text{Fm})), \text{Fm})), \text{safeDiv}(\text{sub}(\text{Fm}, \text{Redox}), \sin(\text{pH}))), \text{add}(\text{mul}(\text{Leit}, \text{add}(\text{mul}(\text{Leit}, \text{Cl}_2), \text{neg}(\text{Tp}))), \text{neg}(\text{safeDiv}(\text{sub}(\text{mul}(\text{Leit}, \text{add}(\text{mul}(\text{Leit}, \text{Cl}_2), \text{neg}(\text{Tp}))), \text{Redox}), \sin(\text{pH}))))$
$\text{safeDiv}(\max(\text{Fm}, \cos(\max(\cos(\text{Fm}), \text{pH}))), \text{mul}(\text{neg}(\cos(\text{sub}(\text{sub}(\text{pH}, \min(\text{Leit}, \text{Cl}_2)), \min(\text{Leit}, \text{Cl}_2))), \text{safeDiv}(\text{sub}(\max(\text{Redox}, \text{Fm}), \text{Redox}), \min(\text{add}(\cos(\max(\text{Cl}_2, \text{pH})), \sin(\sin(\text{pH}))), \text{Fm})))$
$\text{mul}(\text{mul}(\text{Cl}_2, \text{safeDiv}(\text{add}(\text{safeDiv}(\sin(\text{Trueb}), \text{safeDiv}(\sin(\text{add}(\text{add}(\sin(\sin(\sin(\sin(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2))), \text{Cl}_2), \text{Cl}_2))), \text{Cl}_2), \text{Cl}_2), \text{Cl}_2)), \text{add}(\text{neg}(\text{mul}(\sin(\text{add}(\text{add}(\sin(\text{safeDiv}(\sin(\sin(\sin(\text{add}(\text{add}(\text{neg}(\text{safeDiv}(\sin(\sin(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2))), \text{Cl}_2), \text{Cl}_2, \text{Cl}_2))), \text{Cl}_2))), \text{Cl}_2), \text{Cl}_2, \text{Cl}_2))), \text{Cl}_2))), \text{Redox})), \text{Redox}), \text{safeDiv}(\sin(\text{add}(\text{add}(\sin(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2))), \text{Cl}_2), \text{add}(\text{add}(\sin(\sin(\sin(\sin(\sin(\sin(\sin(\sin(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2))))))))), \text{Cl}_2), \text{Cl}_2))), \text{add}(\text{neg}(\text{mul}(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2), \text{mul}(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2), \text{safeDiv}(\text{add}(\sin(\text{add}(\text{add}(\sin(\sin(\sin(\sin(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2))))), \text{Cl}_2), \text{Cl}_2, \text{Cl}_2))), \text{Redox}))), \text{sub}(\text{Fm}, \text{add}(\text{safeDiv}(\sin(\text{Trueb}), \text{safeDiv}(\sin(\sin(\text{add}(\text{add}(\sin(\sin(\sin(\sin(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2))))), \text{Cl}_2), \text{Cl}_2, \text{Cl}_2))), \text{add}(\text{neg}(\text{mul}(\text{Fm}, \text{safeDiv}(\sin(\sin(\sin(\sin(\sin(\sin(\text{safeDiv}(\sin(\sin(\text{pH}))), \text{Cl}_2))))), \text{Cl}_2))), \text{Redox}))), \text{Redox})))$

mul(sub(pH, add(neg(mul(Fm, Cl_2))), add(cos(Redox), sub(Fm, mul(sub(Fm, add(neg(mul(Fm, Cl_2))), add(sub(sub(Fm, mul(neg(neg(mul(add(neg(mul(Fm, mul(Cl_2, Cl_2))), add(sub(sub(Fm, mul(neg(neg(mul(Fm, Cl_2))), pH)), -1, Redox)), Cl_2))), pH)), -1, Redox))), pH)), Cl_2))), pH)), -1, Redox))), add(Redox, sub(Fm, mul(neg(pH), neg(neg(mul(add(neg(mul(Fm, mul(Cl_2, Cl_2))), add(sub(sub(Fm, mul(neg(neg(mul(Fm, Cl_2))), add(neg(mul(Fm, sub(Cl, Cl_2))), mul(sub(Fm, add(neg(mul(Fm, Cl_2))), add(sub(sub(Fm, mul(neg(neg(mul(add(neg(mul(Fm, mul(Cl_2, Cl_2))), add(sub(sub(Fm, mul(neg(neg(mul(Fm, Cl_2))), pH)), -1, Cl)), Cl_2))), pH)), -1, Redox))), pH))), -1, Redox))), mul(Fm, mul(Cl_2, Cl_2))))))))) Redox)
neg(min(max(cos(min(Fm_2, Fm))), max(add(add(pH, min(max(add(Cl_2, mul(pH, add(safeDiv(mul(Cl_2, Fm_2), Redox), mul(pH, Cl_2)))), min(safeDiv(mul(Redox, pH), neg(pH))), Cl)), add(add(mul(pH, Cl_2), Cl_2), pH))), pH), max(cos(Cl_2), sin(Cl_2))), cos(add(pH, min(max(add(min(max(Fm_2, max(add(safeDiv(Cl_2, Redox), pH), max(cos(Cl_2), sin(Cl_2))), cos(add(pH, min(max(add(safeDiv(Cl_2, sin(Redox)), pH), add(pH, Cl_2))), add(add(pH, Cl_2), add(safeDiv(mul(Cl_2, Fm_2), Redox), mul(pH, Cl_2))))))), mul(pH, add(safeDiv(mul(Cl_2, Fm_2), Redox), mul(pH, Cl_2))), min(Cl, Cl)), add(add(pH, Cl_2, Cl_2))))))
mul(mul(Cl_2, safeDiv(mul(max(mul(mul(Cl_2, safeDiv(mul(Redox, pH), Cl)), Cl), mul(safeDiv(mul(Redox, pH), Cl), safeDiv(mul(Cl, pH), Redox))), pH), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), add(Cl, sub(Redox, max(safeDiv(safeDiv(Redox, pH), Cl), max(mul(mul(Cl_2, safeDiv(mul(safeDiv(mul(mul(Cl_2, safeDiv(mul(mul(mul(Cl_2, safeDiv(mul(max(mul(mul(Cl_2, safeDiv(mul(Redox, pH), Cl)), Cl), mul(safeDiv(mul(Fm, pH), Cl), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), Tp, Cl_2))), Cl), mul(safeDiv(mul(Cl, pH), Cl), safeDiv(max(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), pH, Cl)), mul(mul(Cl_2, Cl), pH)), pH, Cl)), Cl), mul(Fm, safeDiv(Cl_2, mul(mul(mul(Cl_2, Cl), pH), Redox))), pH, Cl)), Cl_2, mul(safeDiv(mul(Redox, pH), Cl), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))))))
safeDiv(neg(safeDiv(sin(min(Fm, pH))), safeDiv(mul(min(mul(sub(Tp, Leit), cos(pH))), Trueb), mul(Cl_2, pH))), sub(add(Cl_2, mul(min(add(add(Cl_2, pH), add(cos(add(pH, Redox))), safeDiv(Cl_2, sin(max(Trueb, Leit))))), pH), mul(Cl_2, Fm))), Redox))), cos(add(add(min(mul(Cl_2, min(sin(add(Tp, add(pH, Redox))), pH)), add(Cl_2, pH)), pH), add(Cl_2, pH)))
mul(mul(Cl_2, safeDiv(mul(max(mul(mul(Cl_2, safeDiv(mul(Redox, pH), Cl)), Cl), mul(safeDiv(mul(Redox, pH), Cl), safeDiv(mul(Cl, pH), Redox))), pH), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), add(Cl, sub(Redox, max(safeDiv(safeDiv(Redox, pH), Cl), max(mul(mul(Cl_2, safeDiv(mul(safeDiv(mul(mul(Cl_2, safeDiv(mul(mul(mul(Cl_2, safeDiv(mul(max(mul(mul(Cl_2, safeDiv(mul(Redox, pH), Cl)), Cl), mul(safeDiv(mul(Fm, pH), Cl), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), Tp, Cl_2))), Cl), mul(safeDiv(mul(Cl, pH), Cl), safeDiv(max(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))), pH, Cl)), mul(mul(Cl_2, Cl), pH)), pH, Cl)), Cl), mul(Fm, safeDiv(Cl_2, mul(mul(mul(Cl_2, Cl), pH), Redox))), pH, Cl)), Cl_2, mul(safeDiv(mul(Redox, pH), Cl), safeDiv(mul(Cl, pH), mul(mul(mul(Cl_2, Cl), pH), Redox))))))
safeDiv(neg(add(min(min(mul(Cl, Redox), neg(add(pH, add(safeDiv(cos(safeDiv(min(Trueb, Trueb), mul(pH, Redox))), cos(max(Redox, Tp))), neg(sub(Cl_2, Tp))))), min(neg(max(safeDiv(max(Cl, max(Redox, Fm)), sub(mul(Fm, Cl_2), min(pH, Trueb))), min(cos(add(max(Cl, Leit), sin(Cl_2))), Cl))), Cl)), add(mul(Cl, Redox), sub(add(mul(Cl, Redox), sub(neg(add(add(Tp, Cl_2), safeDiv(Tp, sub(add(mul(Cl, Redox), neg(safeDiv(mul(Fm_2, Leit), safeDiv(Redox, Trueb))), safeDiv(pH, Cl_2))), safeDiv(pH, Cl_2))), safeDiv(pH, Cl_2))), sub(add(mul(Cl, Redox), sub(neg(add(pH, add(safeDiv(sin(Leit), mul(pH, Leit))), neg(sub(Cl_2, pH))), safeDiv(pH, Cl_2))), safeDiv(pH, Cl_2)))
neg(safeDiv(neg(safeDiv(Cl_2, add(neg(neg(mul(add(Cl_2, mul(Cl, Cl))), neg(mul(Cl_2, pH))))), Cl))), mul(safeDiv(Fm_2, Cl_2), mul(safeDiv(Fm_2, Cl_2), add(neg(safeDiv(add(mul(safeDiv(Fm_2, Cl_2), add(neg(safeDiv(add(neg(mul(add(neg(mul(add(Cl_2, mul(Cl, Cl))), neg(mul(Cl_2, pH))), mul(Cl, Cl))), neg(mul(Cl_2, pH))), sub(Cl_2, Cl))), neg(Cl_2))), neg(mul(add(Cl_2, mul(Cl, mul(safeDiv(Fm_2, Cl_2), mul(safeDiv(Fm_2, Cl_2), add(neg(safeDiv(add(neg(mul(add(neg(mul(add(Cl_2, mul(Cl, Cl))), neg(mul(Cl_2, pH))), mul(Cl, Cl))), neg(mul(Cl_2, pH))), sub(Cl_2, Cl))), neg(Cl_2))), neg(mul(add(Cl_2, mul(Cl, Cl))), neg(mul(Cl_2, pH))))))))) neg(Cl))))))
safeDiv(neg(safeDiv(Cl, mul(safeDiv(sub(safeDiv(sub(Cl, Cl_2), Cl_2), Cl), sub(add(add(Cl_2, add(safeDiv(neg(safeDiv(add(pH, Tp), sin(Redox))), Leit), pH)), pH), pH)), mul(cos(add(add(Cl_2, add(cos(sub(safeDiv(sub(Cl, Cl_2), Cl_2), min(safeDiv(safeDiv(sub(Cl, Cl_2), Cl_2), neg(safeDiv(safeDiv(Fm, Leit), sin(Redox))), safeDiv(sub(Cl, Cl_2), Cl_2))), pH)), pH)), Leit))), cos(add(add(Cl_2, add(Cl_2, pH)), pH)))
neg(safeDiv(cos(mul(add(-1, mul(Cl, Trueb))), pH)), mul(Fm, mul(safeDiv(safeDiv(Fm, Cl_2), Cl_2), mul(safeDiv(safeDiv(Fm_2, Cl_2), Cl_2), add(mul(safeDiv(pH, mul(add(-1, mul(Cl, pH))), pH)), Tp), mul(safeDiv(Tp, mul(safeDiv(-1, mul(Cl, pH)), Tp)), Tp))))))
mul(cos(safeDiv(cos(sin(neg(safeDiv(mul(neg(sin(neg(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(Redox, Fm))), Fm), sub(Leit, Fm))))), Fm), sub(sub(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(pH, Fm))), Fm), sub(sub(sub(Tp, pH), pH), Cl_2)), sub(Cl_2, neg(pH))), Fm))), cos(sub(sub(Cl_2, neg(pH))), Cl))), add(sub(Redox, Fm), mul(neg(sin(neg(safeDiv(mul(neg(safeDiv(mul(Cl, Fm), sub(pH, Fm))), Fm), sub(Leit, Fm))))), Fm)))
min(mul(mul(Cl, neg(cos(sub(Cl, pH)))), mul(mul(Cl, mul(sin(cos(sub(cos(sin(sin(mul(neg(Redox), cos(sub(cos(sin(pH), pH))))), pH))), pH))), add(neg(Trueb), Cl))), mul(mul(Cl, mul(sin(cos(sub(cos(sin(sin(mul(neg(mul(neg(Redox), neg(Cl_2))), neg(Cl_2))), pH))), Cl_2)), add(neg(Cl_2, Cl)))

Rezultatul scorului F1 obținut de către acest ansamblu demonstrează faptul ca acesta metoda poate îmbunătăți performanțele clasificării prin metoda programării genetice într-o măsura considerabila.

$$F1Ansamblu = 0.52447$$

4.8 Interpretarea rezultatelor

În acest capitol au fost evaluați o serie de clasificatori clasici cu scopul de a compara performanțele programării genetice folosită în problema clasificării cu aceștia. Pentru algoritmi de clasificare au fost setați hiper-parametrii cu algoritmul de optimizare al acestora grid search. Acest lucru este esențial pentru o comparație cât mai corectă, deoarece ei influențează într-o proporție ridicată performanțele algoritmilor.

Modul prin care se efectuează compararea performanțelor este scorul F1 obținut de clasificatori pentru setul de testare. Dintre algoritmi de clasificare metoda Extra trees a obținut cel mai bun scor: 0.43818, urmat de regresia logistică cu 0.3, AdaBoost cu 0.28683 și arborele de decizie cu 0.26423. Media scorului F1 pentru clasificatorii obținuți prin programare genetică dintr-un grup format din 30 de expresii a fost de 0.31582, mai slabă decât scorul obținut prin metoda Extra trees, dar mai bună decât al celorlalți algoritmi de clasificare.

Pentru îmbunătățirea performanței programării genetice în cazul acestei probleme sau încercat o serie de metode. Un rezultat a fost notabil, crescând în mod semnificativ performanțele clasificării. A fost creat un ansamblu de clasificatori obținuți prin programare genetică. Crearea unui astfel de model de tip ansamblu presupune în mod evident mai mult efort computațional, dar performanțele sunt îmbunătățite foarte mult. Scorul F1 obținut cu această metodă este 0.52447, mult mai mare decât media scorurilor F1 obținute în mod individual.

Așadar cea mai bună performanță a fost obținută de către ansamblul de clasificatori cu programare genetică, ceea ce arată faptul că această metodă poate fi utilizată cu succes în practică. Vizualizarea comparării rezultatelor arată astfel:

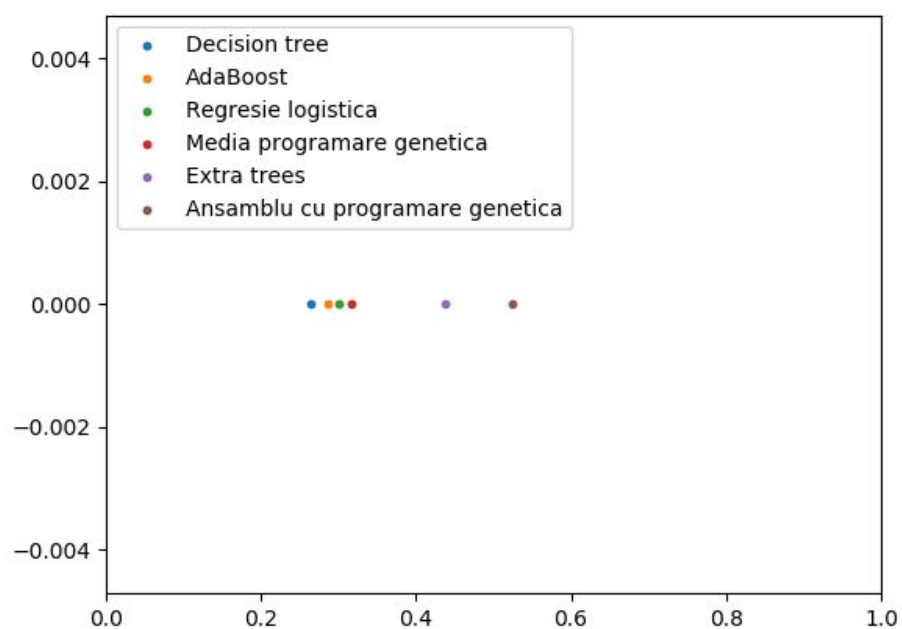


Figura 14.
Vizualizarea rezultatelor obținute

Rezultatele competiției sunt publicate, așadar se poate face o comparație între rezultatele ansamblului și rezultatele obținute de către competitori. Cele mai bune 5 rezultate obținute de către concurenți sunt:

Autori	Scor F1
MINH NGUYEN, AFFAN HAMDANI, KALAICHELVI RAJENDRAN	0.56123041
FITORE MUHARREMI	0.47693977
VICTOR H. A. RIBEIRO (SVM)	0.45098039
VICTOR H. A. RIBEIRO (DT)	0.39030920
ABDUL_RASHEEQ, JAGADESH_DEVAN, JITHI_VARGHESE	0.35265018

Ansamblul nu obține o valoare mai bună decât concurentul de pe primul loc, însă are un scor mai bun decât tot restul concurenților. Consider că rezultatul obținut evidențiază faptul că programarea genetică poate fi folosită cu succes în problema clasificării. Studiul acestei metode poate fi continuat, prin încercarea mai multor moduri de îmbunătățire a performanței, în special prin metode de sampling.

Concluzii

În această lucrare s-a realizat un studiu al performanțelor ce pot fi obținute pentru problema clasificării cu ajutorul programării genetice și compararea cu o serie de algoritmi populari pentru această problema.

Compararea performanțelor s-a realizat în funcție de scorul F1, care este o metrică potrivită pentru învățarea din seturi de date dezechilibrate deoarece nu ia în calcul numărul de clasificări corecte din clasa majoritară, care nu este atât de relevant în acest caz fapt ilustrat prin paradoxul acurateții. Un scor F1 mare indică un număr mare de clasificări true positive și un număr scăzut de clasificări false positives și false negatives.

Seturile de date de antrenament și testare au fost preluate din cadrul competiției “Internet of Things: Online Anomaly Detection for Drinking Water Quality”, acestea fiind măsurători reale ale parametrilor importanți în stabilirea calității apei. Unele variații ale acestor parametrilor sunt normale, însă altele pot indica faptul că aceasta este contaminată, nu este potabilă. Scopul clasificării este de a detecta dacă o instanță definită prin diverse valori pentru atributele măsurate este potabilă sau nu. Așadar, performanțele obținute cu clasificarea prin programare genetică în aceste experimente sunt reprezentative deoarece indică faptul că această metodă poate fi folosită în probleme de clasificare din viața reală cu succes.

Modul în care a fost abordată această problemă a fost prin căutarea unor expresii matematice cu ajutorul programării genetice. Parametrii expresiei sunt chiar atributele unei instanțe din setul de date. Cu ajutorul rezultatului obținut pentru expresie se calculează ulterior o probabilitate pentru clasificare, cu ajutorul căreia se stabilește dacă instanța este clasificată în clasa pozitivă sau în cea negativă.

Media scorurilor F1 pentru un set de 30 de expresii generate prin programare genetică a fost apropiată și puțin mai mare decât rezultatele obținute pentru metodele de clasificare AdaBoost, regresie logistică și arbori de decizie, dar mai mică decât scorul obținut prin metoda Extra trees, care este de tip ansamblu.

Au fost încercate mai multe modalități de îmbunătățire a performanței clasificării cu programarea genetică. Cel mai important rezultat a fost obținut prin crearea unui ansamblu de clasificatori de tip bagging, obținuți prin programare genetică.

Acest ansamblu a oferit cel mai bun scor dintre cele obținute în această lucrare, fiind la o diferență relativ mică față de cea mai bună submisie din cadrul competiției. Aceste rezultate indică faptul că programarea genetică poate oferi rezultate bune în problema clasificării datelor.

Concluzia principală a acestei lucrări este că programarea genetică poate fi folosită cu succes pentru problema clasificării, însă este o metodă costisitoare din punct de vedere al timpului. Crearea unui ansamblu de tipul celui prezentat în această lucrare presupune un timp mai mare decât restul metodelor, algoritmilor cu care s-au făcut comparații.

Bibliografie

1. William J. Clancey - Heuristic Programming Project Computer Science Department Stanford University Stanford. CA 94305 Classification problem solving
2. Robert E. Schapire - Explaining AdaBoost
3. Tom M. Mitchell - Machine Learning, 1997
4. S. B. Kotsiantis - Supervised Machine Learning: A Review of Classification Techniques, 2007
5. Louis Wehenkel, Damien Ernst, Pierre Geurts - Ensembles of extremely randomized trees and some generic applications
6. Pierre Geurts, Damien Ernst, Louis Wehenkel - Extremely randomized trees
7. Dmitry Bystrov - Introduction to soft computing
8. Charles X. Ling, Victor S. Sheng - Cost-Sensitive Learning and the Class Imbalance Problem
9. John R. Koza - Genetic Programming: On the Programming of Computers by Means of Natural Selection, MIT Press, 1992, paginile 162-169
10. Henri Luchian, Mihaela Breabăn - Algoritmi genetici
11. Riccardo Poli, William B. Langdon, Nicholas F. McPhee - A field guide to genetic programming
12. John R. Koza - Genetic Programming: On the Programming of Computers by Means of Natural Selection
13. Haibo He, Edwardo A. Garcia - Learning from Imbalanced Data
- Site-uri sau alte resurse:
14. Sursa originală a setului de date:
<http://www.spotseven.de/gecco/gecco-challenge/gecco-challenge-2018/>
15. <https://deap.readthedocs.io/en>
16. <https://towardsdatascience.com/>

