

Paralelizarea algoritmilor evolutivi pentru instruirea mașinilor cu vectori suport

Gheorghe Radu, Ștefan-Gheorghe Pentiuc, Ionuț Bălan

Abstract—Evolutionary Approach to Support Vector Machines (EASVMs) is a novel technique constructed as an alternative of the standard SVMs architecture (R. Stoean et al, 2009). In this paper we present new parallel algorithms for EASVMs. The experimentation were deployed on a cluster of High Performance Computing Laboratory in “Ștefan cel Mare” University of Suceava. The computational results show the validity of the approach in terms of runtime, accuracy and speedup.

Index Terms—Evolutionary algorithms, parallel computing, supervised learning, support vector machines.

I. INTRODUCERE

SVM-urile (mașinile cu vectori suport) s-au dovedit a fi instrumente deosebit de puternice pentru clasificare și regresie. Totuși, în ciuda originalității și performanței viziunii de învățare, motorul de instruire apare ca fiind restrictiv, rareori transparent, capabil de a converge numai pentru câteva funcții de decizie particulare. Această situație a motivat cercetătorii în a investiga modalități alternative de instruire a SVM. În acest articol vom prezenta o abordare evolutivă a SVM-urilor, EASVM – *evolutionary approach to support vector machines*, urmărind, în principal, lucrarea [1], cu mențiunea că primele idei legate de o astfel de abordare au fost cuprinse în [2].

Algoritmii evolutivi sunt algoritmi robusți, dar care necesită un volum apreciabil de resurse pe parcursul rulării lor. O metodă de paralelizare a algoritmilor evolutivi (EA) pentru instruirea SVM, precum și rezultatele obținute, vor fi deasemenea prezentate în această secțiune, urmărind [3].

Optimizarea evolutivă permite adaptarea hiperplanului de decizie la datele de instruire disponibile, rezolvând direct problema de optimizare primală. Este luată în considerare ideea de bază a SVM, conceptul geometric de instruire, dar abordarea propusă se abate de la tratarea matematică standard. În plus, tehnicile evolute deschid calea generalizărilor ce implică suprafețe de decizie nestructurate,

Gheorghe RADU is with “Henri Coandă” Air Force Academy, #160, M. Viteazul street, 500183 Brașov, gh.radu@gmail.com. He is also PhD student at “Ștefan cel Mare” University of Suceava

Ștefan-Gheorghe PENTIUC, PhD, is with “Ștefan cel Mare” University of Suceava, #13, Universității street, 720229 Suceava, pentiuc@eed.usv.ro

Ionuț BĂLAN is with “Ștefan cel Mare” University of Suceava, #13, Universității street, 720229 Suceava, ionutz19sv@yahoo.com. He is also PhD student at “Ștefan cel Mare” University of Suceava

neliniare. Se propune astfel o nouă abordare: instruirea urmărește SVM-urile standard, în timp ce valorile optime ale coeficienților hiperplanului (w și b) sunt determinate direct de un algoritm evolutiv, cu respectarea echilibrului între acuratețea clasificării și abilitatea de generalizare (vezi și [4]). Deși reprezentarea sugerată pare să fie simplă, determinarea altor detalii ale algoritmului (cum ar fi: interpretarea, operatorii, parametrii) nu este așa ! În abordarea pe care o prezentăm, calculul evolutiv al coeficienților funcției de decizie se face în cadrul conceptului geometric de instruire a SVM.

II. ALGORITMUL EASVM (ALGORITM EVOLUTIV PENTRU INSTRUIREA SVM)

Problema primală generală de găsire a hiperplanului de decizie (vezi și [5], pag. 67) este rezolvată utilizând cu consecvență un algoritm evolutiv. Elementele evolutive ale acestui algoritm sunt prezentate în continuare:

Reprezentarea: Coeficienții hiperplanului, w și b , sunt codificați în structura unui cromozom: $c=(w_1, \dots, w_n, b)$. Inițial, cromozomii sunt generați aleator, astfel încât

$$w_i \in [-1, 1], i=1, 2, \dots, n, b \in [-1, 1]. \quad (1)$$

Asignarea (funcției de) fitness: Funcția de fitness derivă din funcția obiectiv și este supusă restricțiilor problemei de optimizare. Depărtându-se, însă, de SVM standard, pentru funcția de fitness este derivată o formulare neliniară, diferită. Parametrul w este mapat din Φ în H . Rezultă, deci, că norma pătratică implicată în condiția de generalizare devine

$\|\Phi(w)\|^2$, iar ecuația hiperplanului de separare este

$$\langle \Phi(w), \Phi(x_i) \rangle - b = 0.$$

Este utilizată forma $\langle u, w \rangle = u^T w$ și nucleul este folosit pentru a transforma norma într-un produs scalar. Formularea condiției (funcției) de fitness (de a fi minimizată) încorporează funcția obiectiv și restricțiile sunt formulate prin penalizarea cromozomilor eronat clasificați printr-o funcție t care returnează valoarea argumentului, dacă acesta e negativ și valoarea 0, în caz contrar. Pentru clasificare, expresia funcției de fitness este :

$$f(w, b) = \frac{1}{2} K(w, w) + C \sum_{i=1}^m \xi_i +$$

$$\sum_{i=1}^m [t(y_i(K(w, x_i) - b) - 1 + \xi_i)]^2 \quad (2)$$

Selecția și operatorii de variație : Sunt aplicate cele mai uzuale scheme folosite în codificarea reală. Acestea sunt: selecția turnir, încrucișarea intermediată și mutația cu perturbare (distribuție) normală.

Condiția de oprire : Algoritmul se oprește după un număr predefinit de generații.

Îndată ce sunt găsite cele mai apropiate valori optime ale coeficienților hiperplanului de decizie, o nouă formă poate fi clasificată direct, prin utilizarea ecuației:

$$\text{class}(x_i) = \text{sgn}(K(w, x_i) - b). \quad (3)$$

Acuratețea clasificării este definită ca fiind raportul dintre numărul cazurilor/formelor corect etichetate și numărul total de exemple de test.

A. O construcție naivă a EASVM

Vrem să evaluăm dacă algoritmul EASVM poate da rezultate bune atunci când este comparat cu abordarea standard a SVM. Mai întâi, trebuie să determinăm parametrii potriviți ai algoritmului. Astfel, dacă indicatorii de erori ξ_i , $i = 1, 2, \dots, m$ apar în restricțiile formulate pentru optimalitatea planului, aceștia trebuie tratați (corespunzător) și în algoritmul EASVM. În acest caz, algoritmul EASVM se depărtează de înțelesul strict geometric (al unei SVM canonice) al unei deviații și, pur și simplu, evoluează indicatorii de erori. Prin urmare, algoritmul EASVM include acești indicatori în structura unui cromozom, codificarea acestuia devenind :

$$c = (w_1, \dots, w_n, b, \xi_1, \dots, \xi_m), \quad (4)$$

unde $w_i \in [-1, 1]$, $i = 1, 2, \dots, n$, $b \in [-1, 1]$ și

$$\xi_j \in [0, 1], \quad j = 1, 2, \dots, m.$$

Construcția naivă a EASVM produce, de asemenea, bune rezultate și când se compară cu SVM canonic. În orice caz, abaterile standard mai mici dovedesc o stabilitate mai înaltă a abordării EASVM. Trebuie de asemenea să remarcăm că, pentru nuclee standard, nu ne putem aștepta ca algoritmul EASVM să fie considerabil mai bun decât SVM standard, din moment ce transformarea nucleu ce induce instruirea este aceeași. În orice caz, flexibilitatea algoritmilor evolutivi ca instrumente de optimizare face din algoritmul EASVM o alegere atractivă din perspectiva performanței, datorită abilității prospective a acestuia de a evolua (în plus) nuclee adaptate la problemă, indiferent dacă aceste nuclee sunt pozitiv (semi-)definite sau nu, ceea ce este imposibil pentru SVM canonice.

Totuși, pentru mulțimi mari de date (spam filtering, de exemplu) timpul de execuție necesar instruirii este foarte mare. Aceasta rezultă din faptul că se folosesc indivizi/cromozomi foarte mari, cu foarte multe caracteristici

(gene), cum ar fi indicatorii de eroare ai tuturor formelor din mulțimea de instruire, care sunt incluși în reprezentare. Această problemă poate fi rezolvată în mai multe moduri :

- prin aplicarea procedurii chunking ([Perez04]), timpul de execuție se reduce de cca opt ori, cu prețul unei mici pierderi în ceea ce privește acuratețea; în afară de rezolvarea problemei lungimii unui cromozom, mecanismul propus reduce numărul mare de calcule necesare referirii formelor de instruire în expresia funcției de fitness;

- prin paralelizarea algoritmului EASVM; această metodă va fi prezentată în capitolul următor.

B. Algoritmul EASVM simplificat

Chiar dacă algoritmul EASVM ne oferă o alternativă viabilă la SVM, el poate fi încă îmbunătățit în ceea ce privește simplitatea. Problema curentă de optimizare necesită tratarea indicatorilor de erori care, în prezenta variantă de algoritm evolutiv, sunt incluși în reprezentare. Aceasta poate complica sever problema prin creșterea lungimii cromozomului (numărul de variabile) cu numărul de forme de instruire. Mai mult decât atât, o astfel de metodologie se depărtează puternic de conceptul de SVM canonic. În această secțiune este prezentată o procedură de reprezentare numai a coeficienților hiperplanului și de calcul al indicatorilor de erori (în loc de a-i evolua), procedură care își menține performanța.

Din moment ce algoritmul EASVM furnizează direct și interactiv coeficienții hiperplanului în orice moment, se propune scoaterea indicatorilor de eroare din reprezentarea algoritmului evolutiv și, în schimb, calcularea acestor valori. În consecință, reprezentarea unui individ conține numai w și b . În plus, toți indicatorii ξ_i , $i = 1, 2, \dots, m$, vor putea fi calculați în vederea referirii lor în funcția de fitness. În ceea ce privește clasificarea, se ia individul curent (hiperplanul de separare), iar hiperplanele suport sunt determinate prin mecanismul propus de Bosch și Smith, în lucrarea [Bosch98]. Mai întâi calculăm :

$$\begin{aligned} m_1 &= \min \{K(w, x_i) | y_i = +1\} si \\ m_2 &= \max \{K(w, x_i) | y_i = -1\}, \quad i = 1, 2, \dots, m; \\ p &= |m_1 - m_2|; \\ w' &= \left(\frac{2}{p}\right)w; \quad b' = \left(\frac{1}{p}\right)(m_1 + m_2). \end{aligned} \quad (5)$$

Pentru fiecare formă de instruire x_i , $i = 1, 2, \dots, m$, deviația față de hiperplanul suport corespunzător este obținută astfel :

$$\begin{aligned} \delta(x_i) &= k(w', x_i) - b' - 1, \quad y_i = +1; \\ \delta(x_i) &= k(w', x_i) - b' + 1, \quad y_i = -1. \end{aligned} \quad (6)$$

Dacă semnul deviației unei forme este egal cu clasa, atunci indicatorul său este 0; altfel, deviația absolută (normalizată) este returnată ca indicator de eroare. Forma funcției de fitness rămâne ca în formula (2), evident fără a lua ξ_i ca argumente ale funcției.

III. PARALELIZAREA ALGORITMULUI EASVM

Prin utilizarea algoritmilor evolutivi în rezolvarea unei probleme, rezultatele sunt obținute într-un timp destul de mare, cu un consum important de resurse. Datorită faptului că acest tip de algoritmi pot fi relativ ușor paralelizați, am recurs la această variantă, în primul rând cu scopul de a obține rezultatele din ce în ce mai bune, într-un timp rezonabil de mic.

A. Structura paralelă folosită

Testele au fost efectuate pe un cluster (RedPower) cu procesoare XEON Quad Core, din dotarea laboratorului de Calcul de Înaltă Performanță al Universității „Ștefan cel Mare” din Suceava. Această arhitectură este formată din 28 de noduri, fiecare nod având 8 nuclee, viteza ceasului fiind de 2,3 GHz, iar memoria pe nod de 1,96 GB.

Comunicațiile dintre nodurile folosite în testele pe care le-am efectuat se realizează folosind biblioteca OpenMPI [6]. Paralelizarea acestui algoritm a fost realizată la nivel de date, în sensul că populația este împărțită în mod egal între procesele MPI ce rulează pe fiecare nod al clusterului, alegându-se un model insular ce permite acest lucru.

B. Datele de test

Pentru testarea algoritmului propus s-au folosit patru seturi de date [7]:

- Pima Indians Diabetes Data Set;
- Iris Data Set;
- Spambase Data Set;
- Soybean (Small) Data Set.

În fiecare caz considerat s-au făcut câte 5 rulări, luând-se în considerare valorile medii ale rezultatelor obținute. În Tabelul 1 sunt prezentate principalele caracteristici ale seturilor de date, așa cum apar în cadrul algoritmului genetic în fiecare caz în parte.

TABEL 1 CARACTERISTICILE SETURILOR DE DATE.

Baza de date	Kernel	Numărul formelor de antrenare	Număr formelor de testare
Pima	polinomial	576	192
Iris	radial	105	45
Spambase	polinomial	3451	1150
Soybean	polinomial	30	17

Trebuie să precizăm că dimensiunea populației în fiecare caz în parte este de 200 indivizi, numărul de generații în care vor evolua aceste populații este de 250, în timp ce atât probabilitatea de recombinare, cât și cea de mutație este de 0,4.

IV. REZULTATELE EXPERIMENTELOR

Rezultatele obținute în urma rulărilor pe diferite arhitecturi sunt prezentate în Tabelul 2.

TABEL 2 ACURATEȚEA MEDIE OBTINUTĂ.

Numărul de noduri	Iris(%)	Spam(%)	Pima(%)	Soybean(%)
1	91,99	74,83	72,93	90,72
2	94,66	77,15	77,83	92,71
4	97,32	78,13	78,32	97,13
8	98,11	79,13	79,15	100
20	99,97	79,45	79,23	100

Datele din tabelul 2 sunt reprezentate graphic în Figura 1, pentru a scoate în evidență câștigurile obținute prin paralelizarea algoritmului evolutiv din cadrul problemei rezolvate.

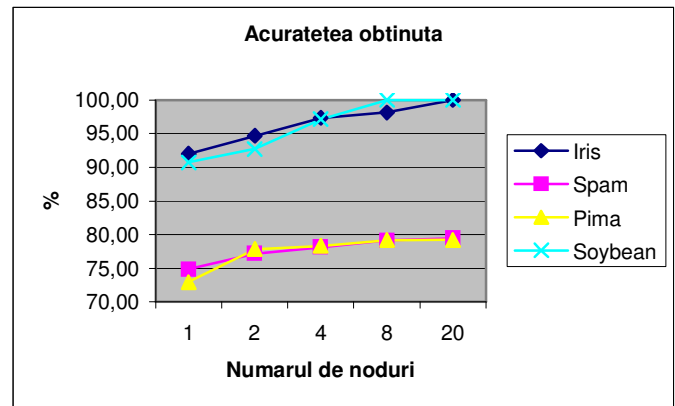


Fig.1 Acuratețea obținută în cele 4 cazuri.

Pe seturile de date studiate se observă că introducerea paralelizării algoritmilor genetici a condus la îmbunătățirea rezultatelor în marea majoritate a cazurilor avute în vedere. În cazul bazei de date Soybean se observă că la 20, respectiv 8 noduri folosite pentru rulare se atinge acuratețea de 100%. Acest rezultat apare ca urmare a numărului mic de forme ce trebuie incluse într-o anumită categorie, aceasta realizându-se cu probabilitate mult mai mare decât în cazul seturilor de date (mult) mai mari.

Un anumit cluster dispune de resurse suficiente pentru rezolvarea anumitor probleme. Prin urmare, nu are rost să supraîncărcăm cu sarcini anumite procesoare, doar pentru obținerea unor rezultate mai bune, în schimbul creșterii duratei de rulare. Pentru a confirma acest lucru am recurs la utilizarea aceluiași număr de indivizi într-o anumită generație de evoluție, dar pe configurații diferite. Rezultatele obținute prin rularea algoritmului asupra setului de date Iris, în condițiile precizate mai sus, sunt prezentate în Tabelul 3.

TABEL 3 REZULTATE OBTINUTE PRIN RULAREA, ÎN DIFERITE CONFIGURAȚII, A SETULUI DE DATE IRIS.

Dimensiunea populației	4000	2000	800	400	200
Numărul de procesoare	1	2	5	10	20
Timpul (s)	129,3	63,93	26,33	12,86	6,69
Acuratețea(%)	92,33	100	98,33	97,77	100

Din Tabelul 3 putem trage concluzia că nu este necesar să folosim populații imense, ce pot duce la sufocarea unei anumite unități de calcul în momentul în care avem la dispoziție mai multe unități de procesare ce sunt legate în paralel. Prin recurgerea la astfel de metode, folosind același număr de indivizi/generație, obținem rezultate mult mai bune și timpi de rulare semnificativ mai mici.

Pentru a putea la răspunde la întrebarea: „este benefică paralelizarea algoritmilor genetici pentru rezolvarea problemelor propuse?”, vom calcula accelerarea obținută pe baza datelor din Iris Data Set. Valorile obținute sunt reprezentate în Figura 2.

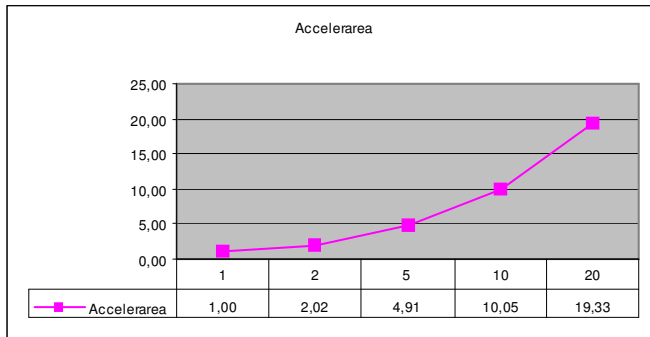


Fig. 2 Accelerearea obținută prin paralelizarea algoritmilor genetici.

Teoretic, accelerarea nu poate depăși numărul de procesoare pe care s-a realizat paralelizarea, dar este de dorit și de așteptat să se apropie cât mai mult de acesta. În cazul nostru se observă anumite diferențe între accelerarea obținută și numărul de procesoare pe care s-a făcut rulare. Principala cauză care a dus la această situație fiind existența unor porțiuni mari de program (diferite de cele rezervate algoritmilor genetici) ce nu au fost (sau nu au putut fi) paralelizate. O altă cauză, minoră, pentru această diferență ar putea fi timpii necesari comunicațiilor dintre procesoare.

V. CONCLUZII

Abordarea EASVM materializează viziunea de învățare a SVM, dar are un mare avantaj în ceea ce privește transparența procesului de instruire, inerent oricărei probleme de optimizare prin intermediul unui algoritm evolutiv.

Sunt utilizate două reprezentări ale algoritmului evolutiv: una este simplă, dar rapidă, cealaltă este mai complicată, dar dovedește o acuratețe mai bună.

Paralelizarea algoritmului EASVM se dovedește a fi o cale extrem de promițătoare, atât în ceea ce privește performanța demonstrată în rezolvarea problemelor (foarte) mari (chiar și fără să fie utilizat mecanismul de *chunking*), cât și acuratețea soluțiilor, superioară în toate cazurile.

REFERINȚE

- [1] R. Stoean, M. Preuss, C. Stoean, E. El-Darzi, D. Dumitrescu, *Support vector machines learning with an evolutionary engine*, Journal of the Operational Research Society **60**, doi:10.1057/jors.2008.124, pp. 1116-1122, 2009.
- [2] Gh. Radu, *Tehnici evolutive de instruire a SVM*, Universitaria Ropet 2003: lucrările Simpozionului Internațional, vol. Matematică-Informatică și Fizică, pp. 97-102, Editura UNIVERSITAS, Petroșani, 2003 (ISBN 973-8260-37-X).
- [3] Gh. Radu, I. Bălan, Șt.-Gh. Pentiuc, *Learning Support Vector Machine Using Parallelized Evolutionary Algorithms*, in Proceedings of the International Conference on Theory and Applications of Mathematics and Informatics, ICTAMI 2011, Alba Iulia, special issue of ACTA UNIVERSITATIS APULENSIS, Mathematics-Informatics, 2011, pp. 331-340, ISSN 1582-5329.
- [4] R. Stoean, M. Preuss, C. Stoean, D. Dumitrescu, *Evolutionary Support Vector Machines and their Application for Classification*, Technical Report Nr. CI-212/06, Collaborative Research Center on Computational Intelligence, University of Dortmund, June, 2006.
- [5] Gh. Radu, Șt.-Gh. Pentiuc, *Mașini cu vectori suport nuanțate – o nouă abordare în clasificarea datelor*, Sisteme distribuite ISSN 1842-6808, pp. 65-70, 10 Dec 2010, Suceava, <http://eed.usv.ro/SistemeDistribuite/2010/Masini%20cu%20vectori.pdf>.
- [6] G. Burns, R. Daoud, J. Vaigl, *LAM: An Open Cluster Environment for MPI*, in Proceedings of Supercomputing Symposium, pp. 379-386, 1994.
- [7] A. Frank, A. Asuncion, *UCI Machine Learning Repository*, University of California, Irvine, School of Information and Computer Science, <http://archive.ics.uci.edu/ml>, 2010.