

# Algoritmi paraleli de procesare a imaginilor

Carmen PUGHINEANU

**Abstract**—In this article we aim to use the power calculation offered by a cluster to process the images. We achieved and proved the use of the parallel algorithms to process the images. We studied and parallelised two clustering algorithm, namely the k-means algorithm and ISODATA. The experimental results compare the two algorithms from the point of view of the execution time, obtained by the algorithm implementation using the MPI standard.

**Index Terms**—grayscale, cluster, segmentation, parallelization, threshold value

## I. INTRODUCERE

În prima parte a articolului imaginile color vor fi convertite în imagini grayscale, apoi se vor paraleliza [1] cei doi algoritmi, k-means și ISODATA. În partea a doua are loc compararea timpilor de execuție și concluziile obținute în urma segmentării imaginilor.

Segmentarea este procesul de clasificare [2] a pixelilor dintr-o imagine. Este unul dintre procedeele cele mai complexe în procesele de vizualizare.

Prin segmentarea imaginilor digitale obținem partiționarea acestora în mai multe segmente (seturi de pixeli) necesare pentru a localiza obiecte, a identifica zone ce sunt uniforme în raport cu un anumit criteriu.

Procesul de segmentare a unei imagini constă în partiționarea completă a imaginii într-un ansamblu de mulțimi disjuncte, nevide și conexe, ce satisfac un anumit criteriu [3].

Algoritmul k-means [4], model propus de MacQueen (1967) este considerat cel mai simplu algoritm de segmentare a imaginilor.

Algoritmul ISODATA propus este în principiu asemănător cu algoritmul k-means dar acesta permite ca numărul de clustere care urmează a se determina, să se modifice automat în timpul iterației prin fuziunea grupurilor similare și divizarea celor cu deviații standard mari [5].

Acești algoritmi lucrează cu un set mare de date și necesită o putere de calcul destul de mare. Din acest motiv implementarea secvențială a acestor algoritmi nu este foarte eficientă.

În acest articol se urmărește implementarea algoritmului k-means pe un cluster format din 9 calculatoare (master + 8 noduri). Sistemul de operare folosit pentru acest cluster este Fedora 12(GNU/Linux). Pentru realizarea CLUSTER-ului s-a utilizat două switch-uri, unul de 100MB(MPI) și unul de 1GB(folosit de NFS).

## II. TRANSFORMAREA IMAGINILOR DIN RGB ÎN IMAGINI GRAYSCALE

Pentru a putea aplica algoritmul este nevoie de o fază de preprocesare și anume transformarea imaginii color în imagini cu niveluri de gri.

Fișierele alese pentru crearea setului de date au format .bmp și chiar dacă au o mărime mult mai mare decât altele, ele nu au pierderi de date, iar reprezentarea datelor este deosebit de simplă.

La o imagine pe 24 biți, sunt necesari 3 octeți pentru definirea unui pixel. Fiecare din cei trei octeți asociați unui pixel codifică cantitatea de roșu, verde și albastru (imaginea RGB - Red-Green-Blue) cu valori cuprinse între 0 și 255.

## III. DESCRIEREA SETULUI DE DATE DE INTRARE PENTRU ALGORITMUL K-MEANS SI ALGORITMUL ISODATA

Setul de date de intrare are o formă proprie în funcție de structura fiecărui punct al imaginii. Fiecărui punct i se asociază un vector de trăsături caracteristice. Aceste trăsături sunt nivelul de gri (grayscale) al pixelului și coordonatele sale în imagine.

Mulțimea vectorilor de caracteristici astfel obținuți (numărul lor este evident egal cu numărul de puncte din imagine), formează setul de date de intrare și este partiționată cu ajutorul algoritmului descris în acest articol.

Imaginea este astfel reprezentată printr-o matrice în care fiecare punct corespunde unui pixel. Pentru a realiza procesul de grupare a pixelilor este nevoie de o funcție de similitudine.

Fiecărui cluster (fiecarei clase) i se atribuie o anumită etichetă și toate punctele imaginii al căror vector de trăsături aparține aceluiași cluster vor primi eticheta clusterului respectiv.

## IV. SEGMENTAREA IMAGINILOR

Segmentarea imaginilor [6] este una din cele mai importante componente ale procesului de analiză și interpretare a imaginilor. Erorile care apar în procesul de segmentare vor afecta toate componentele următoare.

Segmentarea unei imagini digitale pe nivele de gri în funcție de nivelul de gri și de coordonatele pixelilor din imagine presupune o partiție a imaginii în mulțimi conexe, astfel încât fiecare mulțime are toți pixelii din imagine cu nivel de gri apropiat.

Pixelii din imagine sunt asociați uneia sau alteia dintre mulțimi în funcție de proprietatea lor de similitudine. În urma alocării tuturor pixelilor din imagine la o mulțime sau alta,

rezultă segmentarea imaginii în numărul de mulțimi, obținându-se o partiție a imaginii.

#### V. MASURI ALE SIMILARITĂȚII DINTRE DOUA PUNCTE

Spre deosebire de clasificările clasice, care au un scop bine stabilit și criterii clare de grupare, clasificările automate au drept unic criteriu similaritatea dintre pixelii ce urmează a fi grupați. Există o varietate de metode ce pot fi utilizate pentru măsurarea similitudinii.

În data clustering cea mai folosită metodă de măsurare a similarității o reprezintă distanța Euclidiană dată de următoarea formulă:

$$d(x, y) = \|x - y\| = \sqrt{\sum_{i=1}^3 (x_i - y_i)^2} \quad (1)$$

unde:  $x=(x_1, x_2, x_3)$  și  $y=(y_1, y_2, y_3)$  sunt două puncte date în coordonatele carteziene.

Practic, o clasificare automată realizează o grupare optimă de puncte sub raportul similarității dintre ei.

#### VI. ALGORITMUL K-MEANS

Algoritmul k-means (MacQueen, 1967) este unul din cei mai populari algoritmi de partiționare și urmărește identificarea a  $k$  grupe (clustere) distincte, astfel încât datele din fiecare clasă să fie suficient de similare [7].

Fiecare clasă va avea un reprezentant și va fi considerat centrul clasei. Aceste centre se definesc aleatoriu și pe la o distanță cât mai mare una de alta, deoarece alegerea acestor puncte poate influența foarte mult rezultatul dat de algoritm. Următorul pas o reprezintă determinarea apartenenței fiecărui punct la un cluster. Zonele din imagine cu aceeași similaritate pot să nu fie conexe, deoarece punctele imaginii depind de trăsăturile utilizate, și deci orice segmentare prin clustering va fi urmată de etichetarea imaginii segmentate, pentru identificarea obiectelor individuale. Când punctele nu își schimbă clusterul din care fac parte atunci gruparea este terminată. În caz contrar trebuie să se recalculeze centrele pentru fiecare cluster în funcție de punctele care fac parte din el. După calcularea acestor centre se realizează din nou determinarea apartenenței fiecărui punct la un cluster. Astfel este generată o buclă. La fiecare iterație din buclă, clusterelor își schimbă centrul. Ieșirea din buclă se poate realiza în momentul în care clusterelor nu își mai schimbă centrele. Setul de date de intrare ce urmează a fi prelucrat este mare. Dacă setul de date ce conține  $n$  puncte  $(x_1, x_2, \dots, x_n)$ , în care fiecare punct este considerat ca un vector de observații de dimensiune  $d$  ( $d=3$ ), atunci clusterizarea k-means presupune împărțirea acestui set în  $k$  partiții (clase) disjuncte și nenule ( $k < n$ )  $S = \{S_1, S_2, \dots, S_k\}$  astfel încât să se minimizeze (suma pătratelor distanțelor din interiorul clusterului) funcția de similaritate:

$$E = \arg \min_s \sum_{i=1}^k \sum_{x_j \in S_i} \|x_j - c_i\|^2 \quad (2)$$

unde:  $c_i$  este media lui  $S_i$

Algoritmul secvențial de grupare a  $n$  puncte în  $k$  clase este: Aleator se stabilesc  $k$  puncte (pixeli) ca centrele inițiale a claselor  
repetă

Pentru fiecare punct se caută cea mai apropiată clasă

Dacă punctele și-au schimbat clasa din care face parte

Recalculează centrele claselor în funcție de punctele care fac parte din fiecare clasă

Până când punctele care și-au schimbat clasa este mai mic decât o valoare prestabilită

Algoritmul serial k-means are timpul de execuție egal cu  $O(R_s kn)$  unde  $k$  este numărul dorit de clustere (grupe sau clase) și  $R_s$  este numărul de iterații.

#### VII. ALGORITMUL ISODATA

Un rol important în procesul de segmentare îl constituie alegerea inițială a numărului de clase. O alegere greșită determină erorile de supra-segmentare sau sub-segmentare. Supra-segmentarea se poate corecta prin reunirea ulterioară a obiectelor deja determinate și este deci preferabilă alegerea unui  $k$  mai mare decât necesar.

Inițial se dorește segmentarea imaginii în  $k$  clustere, dar algoritmul permite ca numărul de clustere să se modifice automat în timpul iterației prin fuziunea grupurilor similare și divizarea celor cu deviații standard mari.

Parametrii utilizați în acest algoritm sunt [8]:

- 1 -  $K$  = numărul de clustere dorit;
- 2 -  $I$  = numărul maxim de iterații permise;
- 3 -  $P$  = numărul maxim de perechi de cluster care se pot uni (fuziona);
- 4 -  $\theta_N$  = o valoare de prag pe care o poate avea un număr minim de probe din fiecare cluster (utilizată pentru divizare);
- 5 -  $\theta_S$  = o valoare de prag pentru deviația standard (utilizată pentru operația de separare);
- 6 -  $\theta_C$  = o valoare de prag pentru distanțe perechi (folosite pentru operația de îmbinare).

Pașii algoritmului

1. alegerea arbitrară a celor  $k$  centre,  $M_1, M_2, \dots, M_k$
2. fiecare probă  $N$  se atribuie celui mai apropiat cluster:

$$X \sim w_j \text{ dacă}$$

$$D_L(X, M_j) = \max\{D_L(X, M_i), i=1, 2, \dots, k\}$$

3. se îndepărtează clusterelor cu mai puțin de  $\theta_N$  membri
4. se recalculează fiecare centru:

$$M_j = \frac{1}{N_j} \sum_{x \sim w_j} X, j=1, \dots, k$$

5. se calculează distanța medie  $D_j$  între probele din clusterul  $w_j$  și centrul corespunzător:

$$D_j = \frac{1}{N_j} \sum_{X \sim w_j} D_L(X, M_j), j=1, \dots, k$$

6. Se calculează distanța medie generală a probelor din centrele clusterelor lor respective

$$M_j = \frac{1}{N} \sum_{j=1}^k N_j D_j$$

7. Dacă sunt prea puține clusterse se merge la pasul 8;

Altfel

dacă sunt prea multe clusterse se merge la pasul 11;

altfel se merge la pasul 14.

Pașii 8 - 10 pentru operația de separare, Pașii 11-13 sunt pentru operația de fuziune

8. Primul pas pentru separare.

Se găsește vectorul deviațiilor standard  $\sum |j = [\sigma_1^{(j)}, \dots, \sigma_n^{(j)}]$  pentru fiecare cluster

$$\sigma_i^{(j)} = \sqrt{\frac{1}{N_j} \sum_{X \sim w_j} (x_i - m_i^{(j)})^2}, i=1, \dots, n, j=1, \dots, k$$

9. se află componenta maximă notată cu  $\sigma_{\max}^{(j)}$  pentru orice  $j=1, \dots, k$

10. Dacă pentru fiecare  $\sigma_{\max}^{(j)}$  ( $j=1, \dots, k$ ) tot ce urmează este adevărat:

- $\sigma_{\max}^{(j)} > \theta_s$
- $D_j > D$ ,
- $N_j > 2\theta_N$

atunci  $M_j$  este divizat în două noi centre de cluster  $M_j^+$  și  $M_j^-$  prin adăugarea lui  $\pm \delta$  la componenta lui  $M_j$  corespunzând lui  $\sigma_{\max}^{(j)}$  unde  $\delta$  poate fi  $\alpha \sigma_{\max}^{(j)}$  pentru vreun  $\alpha > 0$ . Apoi se șterge  $M_j$  și  $k \leftarrow k+1$  și apoi se merge la pasul 2

altfel se merge la pasul 14.

11. Primul pas spre fuziune. Se calculează distanțele, în perechi,  $D_{i,j}$  între fiecare două clusterse:  $D_{i,j} = D_L(M_i, M_j)$ , (pentru toți  $i \neq j$ ) și se aranjează aceste  $k(k-1)/2$  distanțe în ordine ascendentă.

12. Se găsesc cei mai mici  $D_{i,j}$  nu mai mulți ca  $P$ , care sunt de asemenea mai mici decât  $\theta_c$  și se mențin în ordine ascendentă:  $D_{i_1 j_1} \leq D_{i_2 j_2} \leq \dots \leq D_{i_p j_p}$ .

13. Se execută fuzionare în perechi: pentru  $l=1, \dots, P$ , se execută următoarele:

Dacă nici unul din  $M_{i_1}, M_{j_1}$  nu au fost folosite în această iterație, atunci se fuzionează pentru a forma un nou centru:

$$M = \frac{1}{N_{i_1} + N_{j_1}} [N_{i_1} M_{i_1} + N_{j_1} M_{j_1}]$$

se șterge  $M_{i_1}$  și  $M_{j_1}$ , și  $k \leftarrow k-1$ . Megeți la pasul 2.

14. Algoritmul se termină dacă numărul maxim de iterații  $I$  este atins. Altfel se merge la Pasul 2.

## VIII. IMPLEMENTAREA PARALELĂ A ALGORITMILOR K-MEANS ȘI ISOADATA

Procedura de clasificare în cazul paralel se derulează astfel: se pleacă de la un aranjament arbitrar al punctelor într-un număr prestabilit de grupe, după care punctele sunt mutate dintr-o grupă în alta în vederea minimizării varianței din interiorul grupelor și deci maximizării varianței dintre grupe. Numărul de mutări poate fi, de asemenea, specificat [9].

Ideea de bază este a defini  $k$  centri de greutate, unul pentru fiecare cluster, atât cât este posibil mai îndepărtate unele de altele. Următorul pas este să luăm fiecare punct de-a lungul setului datelor de intrare și să-l asociem celui mai apropiat centru de greutate. Când nu există puncte nedecise, primul pas este complet, iar o grupare inițială s-a făcut. În acest moment trebuie să recalculăm  $k$  noi centre de greutate ale clusterelor rezultați din pasul anterior.

În funcție de anumite valori de prag, în algoritmul ISODATA există posibilitatea ca centrul unui cluster să fie divizat în alte două noi centre, sau ca aceste centre să fuzioneze pentru a forma un nou centru, astfel numărul de clusterse crește sau scade cu 1.

În continuare se vor partiționa calculele care vor fi efectuate prin asocierea fiecărei operații cu datele cu care operează. Această partiționare produce un număr de task-uri, fiecare task conținând date și un set de operații pentru ele. Pot exista cazuri în care o operație trebuie efectuată asupra datelor din task-uri diferite. În aceste cazuri este necesară realizarea comunicației între aceste task-uri.

Pentru cazul de față se va folosi strategia realizării aceluiași set de calcule pe seturi diferite de date. Astfel datele sunt împărțite în părți egale, numărul acestora fiind egal cu numărul de procese care se vor executa în paralel.

Implementarea acestor algoritmi se realizează cu ajutorul standardului MPI.

## IX. REZULTATE EXPERIMENTALE

Imaginea în formatul RGB din figura 1, este transformată în imagini grayscale, figura 2:



Fig. 1. Imaginea în format RGB.



Fig.2. Imaginea grayscale.

Urmează prelucrarea imaginii din figura 2 cu algoritmul k-means paralelizat. Imaginea a fost segmentată în 5 clase și a fost extrasă imaginea ce face parte din clasa numărul 1, figura 3:

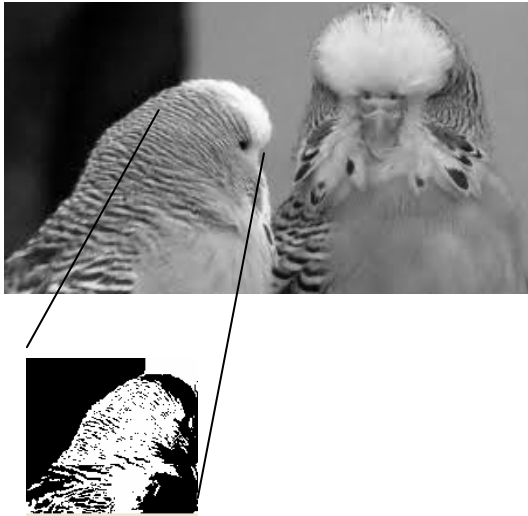


Fig. 3. Extragerea zonei (clasa 1) de interes din imaginea din fig. 2.

Aplicarea algoritmului ISODATA paralelizat pentru setul de intrare adică imaginea din figura 2, segmentează pixelii tot în 5 clase, fiind extrasă imaginea care face parte din clasa 0, figura 4:

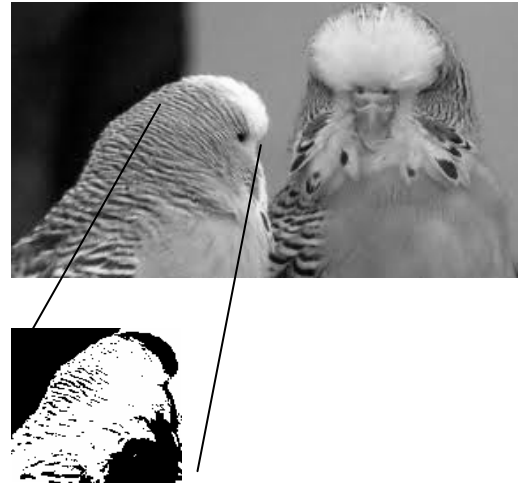


Fig.4. Extragerea zonei (clasa 1) de interes din imaginea din fig. 2.

Prin compararea zonelor de interes extrase, cu cei doi algoritmi paralelizați k-means și ISODATA folosind imaginea din figura 2 observăm că, gruparea pixelilor cu algoritmul ISODATA paralelizat oferă rezultate mult mai bune.

Acum se vor compara timpii de execuție obținuți cu cei doi algoritmi k-means și ISODATA paraleli.

Pentru setul de intrare adică imaginea din figura 2, s-a realizat gruparea pixelilor în 10, 25, 50, 100 clase. Pentru fiecare grupare aplicația a fost executată folosind 2, 4, 8, 16, 32 procese MPI. Rezultatele testelor sunt date în graficele de mai jos:

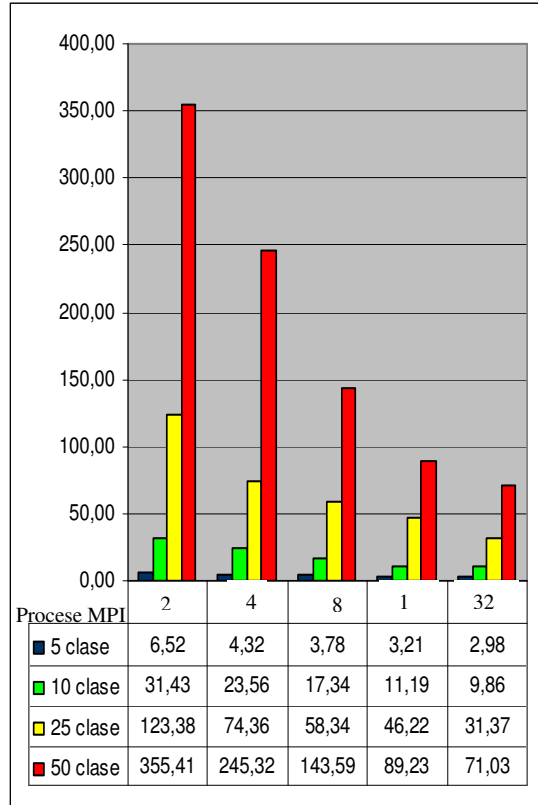


Fig. 5. Rezultatele experimentale ale algoritmului k-means având setul de date de intrare imaginea din fig.2.

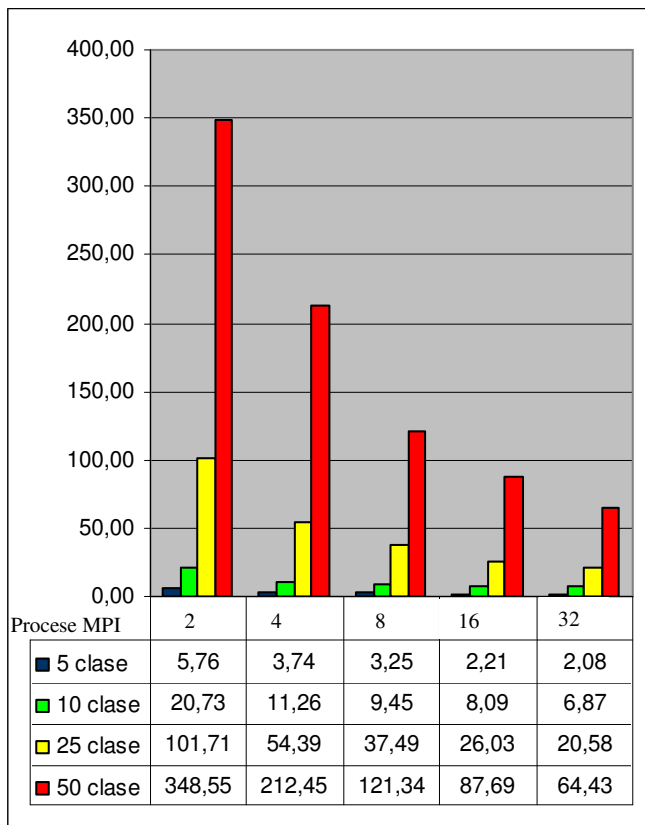


Fig. 6. rezultatele experimentale ale algoritmului ISODATA având setul de date de intrare imaginea din fig.2.

Din aceste rezultate se poate observa că prin creșterea numărului de procese MPI se realizează o scădere semnificativă a timpul de execuție. De asemenea se observă că timpul de execuție crește dacă și numărul de clase în care sunt grupați pixelii crește. Timpul de execuție obținut cu algoritmul ISODATA este mai bun față de algoritmul k-means.

## X. CONCLUZII

În cadrul acestui articol au fost prezentați doi algoritmi paraleli pentru segmentarea imaginilor.

S-a demonstrat că algoritmul paralel ISODATA utilizat în segmentarea imaginii este mult mai bun față de algoritmul paralel k-means, regiunea extrasă cu ISODATA este mult mai fidelă cu regiunea din imagine originală. Mai mult timpul de execuție în algoritmul ISODATA este mai mic față de timpul de execuție în algoritmul k-means.

S-a mai demonstrat că în ambii algoritmi timpul de execuție scade atunci când numărul de procese crește și crește atunci când numărul de clase crește.

## REFERINȚE

- [1] Zhonghui Feng, Bing Zhou, Junyi Shen, A parallel hierarchical clustering algorithm for PCs cluster system, *Neurocomputing*, Volume 70, Issues 4-6, January 2007, Pages 809-818
- [2] Michael W. Berry, Malu Castellanos, *Survey of Text Mining II: Clustering, Classification, and Retrieval*, Springer, 2008, ISBN 1848000456, 9781848000452
- [3] V. Neagoie, O. Stanasila "Teoria recunoașterii formelor", Editura Academiei, 1992
- [4] Kanungo, T.; Mount, D.M.; Netanyahu, N.S.; Piatko, C.D.; Silverman, R.; Wu, A.Y.; An efficient k-means clustering algorithm: analysis and implementation, *Pattern Analysis and Machine Intelligence*, IEEE Transactions on Volume 24, Issue 7, July 2002 Page(s):881 – 892
- [5] R. Vancea, Ștefan Holban, Dan Ciubotariu, "Recunoașterea formelor.
- [6] Pentiuc Ștefan-Gheorghe: *Aplicații ale recunoașterii formelor în diagnosticul automat*, Editura Tehnică, București, 1997.
- [7] T.M.Cover, P.E.Hart, *Nearest Neighbor Pattern Classification*, IEEE Trans. on Information Theory, IT-13, 1, pag. 21-26, (1967)
- [8] <http://fourier.eng.hmc.edu/e161/lectures/classification/node13.html>
- [9] Kanungo, T.; Mount, D.M.; Netanyahu, N.S.; Piatko, C.D.; Silverman, R.; Wu, A.Y.; An efficient k-means clustering algorithm: analysis and implementation, *Pattern Analysis and Machine Intelligence*, IEEE Transactions on Volume 24, Issue 7, July 2002 Page(s):881 – 892 .