# A New Method on Data Clustering Based on Hybrid K-Harmonic Means and Imperialist Competitive Algorithm

MarjanAbdeyazdan

Department of Computer Engineering, Mahshahr branch, Islamic Azad University, Mahshahr, Iran.

E-mail: abdeyazdan87@yahoo.com , m.abdeyazdan@mahshahriau.ac.ir

## Abstract

Data clustering is one of the commonest data mining techniques. The K-means algorithm is one of the most wellknown clustering algorithms thatare increasingly popular due to the simplicity of implementation and speed of operation. However, its performancecouldbe affected by some issues concerningsensitivity to the initialization and getting stuck in local optima. The K-harmonic means clustering method manages the issue of sensitivity to initialization but the local optimaissue still compromises the algorithm. Particle Swarm Optimization algorithm is a stochastic global optimization technique which is a good solution to the above-mentioned problems. In the present article, the PSOKHM, a hybrid algorithm which draws upon the advantages of both of the algorithms, strives not only to overcome the issue of local optima in KHM but also the slow convergence speed of PSO. In this article, the proposed GSOKHM method, which is a combination of PSO and the evolutionary genetic algorithmwithin PSOKHM,has been positedto enhancethe PSO operation. To carry out this experiment, four real datasets have been employed whose results indicate thatGSOKHMoutperforms PSOKHM.

**Key word:** Data Clustering, PSOKHM, Genetic Algorithm

# Council for Innovative Research

## 1. Introduction

Data clustering is one of the most essential methods in data control and management thatcould partition data into classes accordingto their similar features. Data clustering is a process in whichsets of objective data are divided into separate groupings of classes-clusters- in such a way that objects in the same cluster are more similar while theyare dissimilar to objects of other classes. Clustering has multiple applications in various spheres of activities such as pattern recognition, machine learning, data mining, data recovery and bioinformatics. TheK-means is one of the techniques that are being extensively used in clustering.

The principal objective in KM clustering is that the total dissimilarity among objects in one cluster would be less than that of the center of neighboring cluster. The most significantshortcoming of KM is that the results of clustering are sensitive to the initial choice of cluster centers and may converge with local optima (1, 5). The K-harmonic means (KHM), which was proposed in 2002, aims at minimizing the harmonicmeans of all points in a dataset distancing from cluster centers.Although KM solves the initialization problem, it is still wrestling with the issue of getting stuck at local optima. Therefore, to arrive at a better clustering algorithm we need to seek a solution to overcome getting stuck in local optima. Particle swarm optimization (PSO) is an optimization technique based on population that is inspired by collective and cooperative behavior of bird flocks and fish school. This technique could help KHM to evade local optima trap. The PSOKHM attempts to benefit from bothmethods in order to improve clustering process.

Our proposed method is comprised of a combination of PSO and evolutionary genetic algorithm in PSOKHM to improve PSO operation. Moreover, to examine the efficiency of the proposed algorithm four sets of real data have been employed. As the article continues in section 2, PSOKHM algorithm will be discussed in which PSO and KHM will be briefly dealt with. In section 3, the proposed GSOKHM will be introduced.Section 4 will deal with the results of the proposed methods using four real datasets and a comparison will be drawnbetween these results and those of the precursors. Finally, section 5 will present a summary of what has been done in this study.

## 2. The hybrid PSOKHM clustering

In order to explain the above-mentioned hybrid algorithm, PSO and KHM algorithms will be briefly discussed then a discussion of the PSOKHM will follow.

### 2.1. K-harmonic means algorithm

The KM clustering is a simple and rapid method which is widely being used due to the simplicity of implementation and less iteration. In an attempt to find the clusters centers (C1, C2, C3), the KM algorithm behaves in a way that minimizes the sum of squares of the distance for each Xi point from the nearest cluster center (Cj). The KM efficiency depends on the initialization of centers thatis one of the major shortcomings of this algorithm. There has been established a strong connection between data points and the nearest clustering centers which prevents clusters centers from departingthe boundaries of local density of data. The KHM method solves this problem by replacing the minimum distance of a point from centers used in KM with the harmonic means of distance of each point from all centers. The harmonic means give a privilege to every data points according to their proximity to each center which is considered as a feature of harmonic means.

The following symbols are used to formulize KHM algorithm:

Data to be clustered:   $X = \{x_1, x_2, \ldots, x_n\}$

The group of cluster centers: $C = \{c_1, c_2, \ldots, c_k\}$

The membership function that defines Xi data share belonging to Cjm$(c_j|x_i)$

The weight function that defines the impact of Xi on the repeated calculation of center parameters at the next iteration.w(xi)

The basic algorithm for KHM clustering is as follows:

1) Initial quantization algorithm with an estimated C centers (centers random selection)
2) The calculation of the value of the objective function is as follows:

1) $$KHM(X, C) = \sum_{i=1}^{n} \frac{1}{\sum_{j=1}^{k} \frac{1}{||x_i - c_j||^p}}$$

where p is an input parameter with the valuep $\geq 2$.

1.     For every Xi data, the membership function m(cj|xi) for each Cj is calculated as the followings:

2)     $$m(c_j|x_i) = \frac{||x_i - c_j||^{p-2}}{\sum_{j=1}^{k} ||x_i - c_j||^{-p-2}}$$

2.     For every Xi data, the weight of respective W(Xi) is calculated as the followings:

3)      $$w(x_i) = \frac{\sum_{j=1}^{k} ||x_i - c_j||^{-p-2}}{(\sum_{j=1}^{k} ||x_i - c_j||^{-p})^2}$$

3.      For every Cj its distance from all Xi points according to their membership functions weights are calculated as the followings:

4)      $$c_j = \frac{\sum_{i=1}^{n} m(c_j|x_i)w(x_i)x_i}{(\sum_{i=1}^{k} m(c_j|x_i)w(x_i))}$$

4.      The steps 2 through 5 are performed either according to the predefined numbers of iterations or until KHM(X,C) stops changing to a considerable extent.

5.      Xi point is allocated to j cluster with the biggest m(cj|xi)

This algorithm indicates that KHM is not necessarily sensitive to initialization of centers but the tendency towards converging with local optima is existent. (1,3,12)

## 2.2.    Particle Swarm Optimization (PSO)

The Particle Swarm Optimization (PSO) was firstly developed by Kennedy and Eberhart in 1995. It has been successfully employed in several scientific and applied fields since then. PSO is an optimization algorithm based on population in which an individual is considered as a particle and every population consists of a number of these particles. In PSO the solution space is regarded as a search space and every position in this search space is a problem-based solution. In this population, particles, working in collaboration, try to find the best position (the best solution) in the search space (solution space).

Moreover, every particle travels according to its velocity. At each iteration, the movement of every particle is calculated using the following formulas:

5)      $$x_i(t+1) \leftarrow x_i(t) + v_i(t)$$

6)      $$v_i(t+1) \leftarrow \omega v_i(t) + c_i \, \text{rand}_1 (\text{pdest}_i(t) -$$

In equations 5 and 6, xi(t) is the position of the lith at the t moment and vi(t) is the velocity of lith at the t moment. Pbesti(t) is the best position that has been found by the lith particle so far. Gbest(t)is the best position that has been found by the whole population so far. ωisthe inertia weight that denotes a proportion of the previous velocity and c1, c2 are the velocity constants that denotes the impact of the particle best position and the global best position.

In addition, rand1 and rand2 are variables ranging from 0 to 1. The procedure of PSO algorithm is shown in figure 1.

Initialize a population of particles with random positions and velocities in the search space.

While (termination conditions are not met)

{

For each particle I do

Update the position of particle I according t equation

(5).

Update the velocity of particle I according to equation

(6).

Map the position of particle I in the solution space and evaluate its fitness value according to the fitness function.

Update pbesti (t) and gbest (t) if necessary.

End for

}

**Figure 1: pseudo- code PSO algorithm**

## 2.3. PSOKHM algorithm

The KHM tends to converge faster than the PSO since it needs less function evaluation. However, due to its voracious nature, it would get stuck in a local optima. The PSOKHM hybrid-clustering algorithm attempts to take advantage of both methods through combining PSO and KHM.This hybridalgorithm repeats KHM four times in each generation for which employs 8 generations to improve particles within the population. Furthermore, PSO algorithm repeats 8 times in each generation.

Each particle is a vector of real numbers with K*D dimensions where k is cluster numbers and d is dimensions of the to-be-clustered data. A sample of a particle in population is shown in Figure 2.

The result of its evaluation is the KHM objective function. A summary of the PSOKHM algorithm is illustrated in Figure 3. As the figure shows in each generation, PSO denotes the number of iterations applied on particles. Subsequently, the KHM algorithm applies on the results of PSO iteration again.

| $X_{11}$ | $X_{12}$ | ... | $X_{1d}$ | ... | $X_{k1}$ | $X_{k2}$ | ... | $X_{kd}$ |
|---|---|---|---|---|---|---|---|---|

**Figure 2: a representation of a particle**

Step 1:Set the initial parameters including the maximum iterative count IterCount, the population size Psize,$\omega$, c1 and c2.

Step 2:Initialize a population of aizePsize.

Step 3: Set iterative count Gen 1=0.

Setp 4: Set iterative count Gen2= Gen 3=0.

Setp 5: (pso Method)

Step 5.1: Apply the PSO operator to update the Psize particles.

Step 5.2:Gen2=Gen2+1. If Gen2<8, go to Step 5.1.

Step 6: (KHM Method)For each particle I do

Step 6.1: Take the position of particle I centers of the KHM algorithm.

Step 6.2:Recalculate each cluster center using the KHM algorithm.

Step 6.3: Gen 3=Gen3+1. If Gen3<4, go to Step 6.2.

Step 7:Gen1=Gen+1.If Gen 1<IterCount, go to Step 4.

Step 8:Assign data point i x to cluster j with the biggest m(cj|xi).

**Figure 3: PSOKHM combinatorial clustering algorithm**

## 3. The proposed GSOKHM method

In order to improve the efficiency of the PSO algorithm within PSOKHM, attempts have been madeto combine PSO with another evolutionary algorithm like genetic algorithm so that a more efficient data clustering results. Genetic algorithm is one of the randomized algorithm which draws on the selection, crossover and mutation.

This algorithm is one of the most well known evolutionary algorithms which is widely used in problem-solving optimization. The genetic algorithm could be very efficient in solving local optima within KHM and improving the efficiencyof PSO algorithm. To use the combination of PSO and GA for this specific application, GSO algorithm is used in a way that is shown in figure 4.
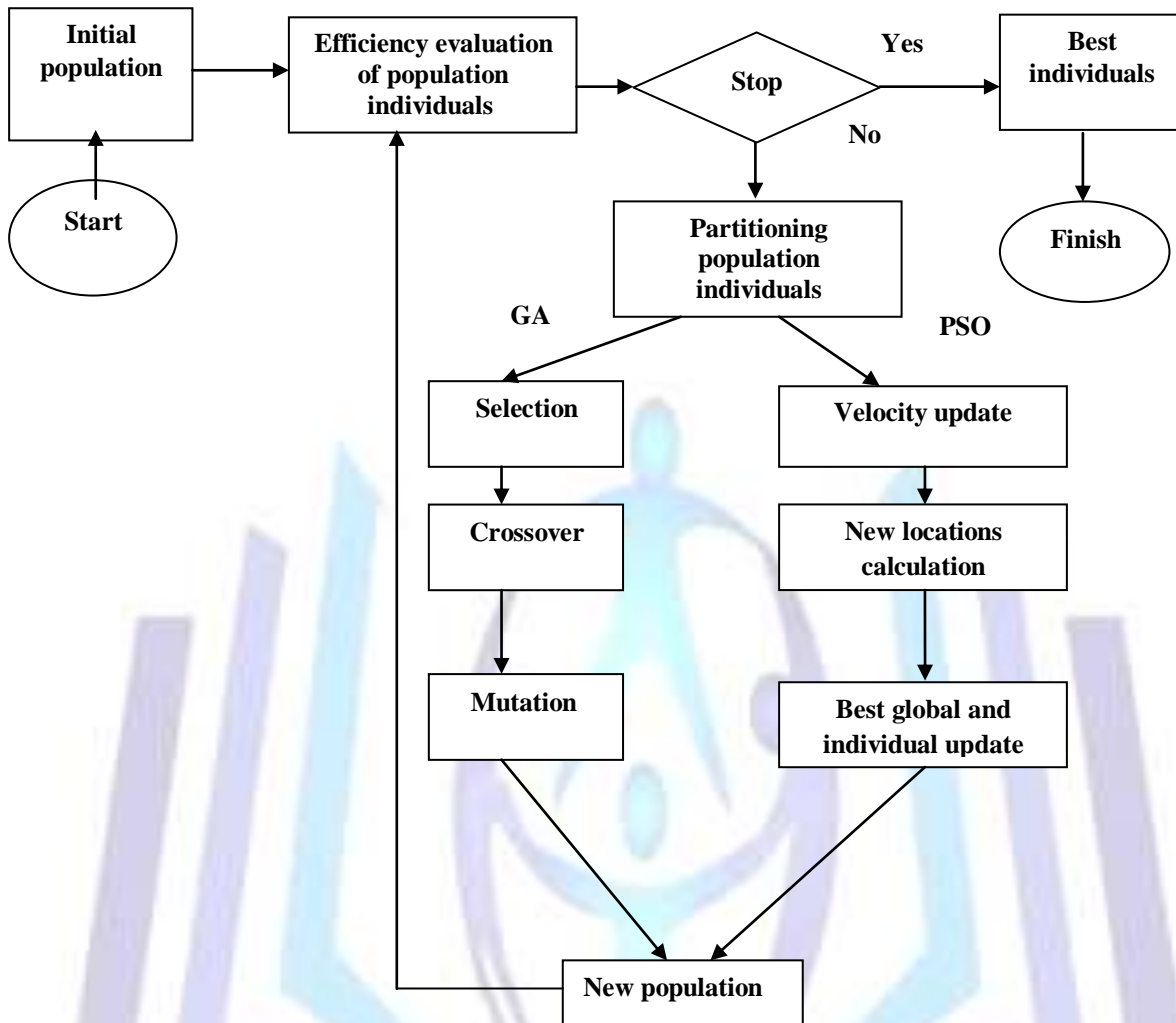
**Figure 4: GSO hybrid algorithm**

As it isevidentin the figure above, members of the population are partitioned into two equal classes at each iteration and PSO and GA operators are directly applied on each class which will eventually be combined to evaluate changes. This procedures proceeds until arrived at a favorable conditions. Furthermore, Roulette wheel is employed for selection in AG algorithm and crossover is carried out as depicted in figure 5. To perform mutation points of random particles of each generation are randomly selected and will be replaced by another random value.
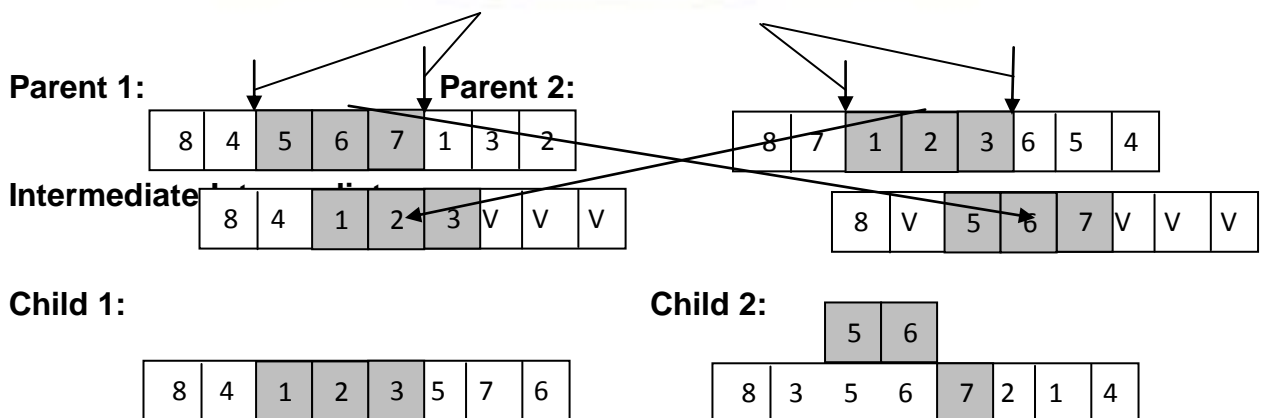


**Figure 5: crossover of genetic algorithm within GSOKHM**

## 4. Experiments and Results

Four real datasets are employed to measure the proposed method which include Iris, Wine, Glass and Contraceptive Method Choice (CMC) with small, medium and large dimensions. These datasets are presented in 15. Table 1 shows a summary of thefutures of these datasets. Additionally, table 1 illustrates the parameters values employed in the algorithm.

**Table 1: datasets features**

| Name of data set | No. of classes | No. of features | Size of dataset (size of classes in parentheses) |
|---|---|---|---|
| Iris | 3 | 4 | 150(50,50,50) |
| Glass | 6 | 9 | 214(70,17,76,13,9,29) |
| CMC | 3 | 9 | 1473(629,334,510) |
| Wine | 3 | 13 | 178(9,71,48) |

**Table 2: GSOKHM algorithm parameters**

| Parameter | Value |
|---|---|
| Psize | 18 |
| | 0.73 |
| $C_1$ | 1.5 |
| $C_2$ | 1.5 |
| Pmutation | 0.02 |
| IterCount | 0.5 |

## 5. Results

In this section, considering the objective algorithm KHM, the efficiency of KHM, PSOKHM and GSOKHM methods are evaluated and compared. Besides, the intended clustering quality is being investigated using the two criteria below:

The sum over all data points based on the harmonic average of the distance from a data point to all the centers as is shown in equation (2-10). It is evident the smaller the values of this set, the better the clustering quality would be.

F-measure criterion which employs precision and recall to recover data.

Every iclass, as shown using the class labels in the evaluated datasets,is considered as a set of nithat is favorable for a search. Everyj cluster, generated by the algorithm, is regarded as the sum of ni of the recovered section by a search. nij denotes the number of objects of i class within j cluster.Precision and recall criteriafor everyi class and j cluster are defined as follows:

7) $\quad r(i,j) = \frac{n_{ij}}{n_i}$

8) $\quad p(i,j) = \frac{n_{ij}}{n_j}$

Neighboring F-measure value is calculated as follows:

9) $\quad F(i,j) = \frac{(b^2+1).p(i,j).r(i,j)}{b^2.p(i,j)+r(i,j)}$

We consider b=1 to have a trade-off for p(o,j) and r(I,j). The global F-measure value for datasets about the size of n is shown below:

10) $\quad F = \sum_i \frac{n_i}{n} max_i\{F(i,j)\}$

It is clear that the more the F-measure value, the better the clustering quality would be.The reported results are averages of runs of the program. The proposed algorithms are implemented using MATLAB 7.6.0 (R2008a) installed on a Vista Home Premium OS with 2.4 GHz CPU and 6 GB RAM. So far, the experiments carried out on KHM algorithm indicate that p is a key parameter to arrive at the values of the objective function.

To this end, our experiments were carried out on a variety of p values and the results are presented in the form of tables for comparison. These tables are the results of the objective function KHM (X,C) which are in accordance with different p values P=2, P=2.5 and P=3. Moreover, not only the objective function KHM (X,C) and F-measure were calculated but the runtime of the proposed algorithms were also calculated and added to the tables.Finally, as the major results of the evaluation,the average independent runs of the algorithms are presented and compared in the tables.

**Table 3: the results for p=3**

|  | Iris | Glass | Wine | CMC |
|---|---|---|---|---|
| **KHM** | | | | |
| **F-Measure** | 0.8923 | 0.48311 | 0.6900 | 0.4491 |
| **KHM(X,C)** | 74.95 | 376.33 | 7,479,216 | 150.950 |
| **Runtime (sec)** | 0.1811 | 0.3244 | 0.2496 | 0.7720 |
| **PSOKHM** | | | | |
| **F-Measure** | 0.8990 | 0.4245 | 0.7023 | 0.4436 |
| **Runtime(sec)** | 2.19 | 5.43 | 2.55 | 18.68 |
| **GSOKHM** | | | | |
| **F-Measure** | 0.9129 | 0.4354 | 0.7090 | 0.4510 |
| **KHM(X,C)** | 11,72 | 105,25 | 64,490 | 9,424 |
| **Runtime(sec)** | 2.73 | 6.89 | 3.48 | 21,45 |

**Table 4: the results for p=2.5**

|  | Iris | Glass | Wine | CMC |
|---|---|---|---|---|
| **KHM** | | | | |
| **F-Measure** | 0.8853 | 0.4130 | 0.6694 | 0.4496 |
| **KHM(X,C)** | 44.07 | 633.40 | 194,607,300 | 687,737.3 |
| **Runtime (sec)** | 0.1331 | 0.2898 | 0.1554 | 1.5656 |
| **PSOKHM** | | | | |
| **F-Measure** | 0.8990 | 0.4245 | 0.7023 | 0.4447 |
| **KHM(X,C)** | 23.159 | 89.98 | 8,442,950 | 82,307.2 |
| **Runtime(sec)** | 2.19 | 5.88 | 2.78 | 19.45 |
| **GSOKHM** | | | | |
| **F-Measure** | 0.9017 | 0.4100 | 0.6902 | 0.4446 |
| **KHM(X,C)** | 3,687 | 70,89 | 3,572,228 | 12,921 |
| **Runtime(sec)** | 2.93 | 6.27 | 3.71 | 22,45 |

**Table 5: the results for p=3**

|  | Iris | Glass | Wine | CMC |
|---|---|---|---|---|
| **KHM** |  |  |  |  |
| **F-Measure** | 0.8853 | 0.4130 | 0.6694 | 0.4469 |
| **KHM(X,C)** | 44.07 | 633.40 | 194,607,300 | 687,737.3 |
| **Runtime (sec)** | 0.1331 | 0.2898 | 0.1554 | 1.5656 |
| **PSOKHM** |  |  |  |  |
| **F-Measure** | 0.8951 | 0.4180 | 0.6835 | 0.4447 |
| **KHM(X,C)** | 23,159 | 89,98 | 8,442,950 | 82,307,2 |
| **Runtime(sec)** | 2.19 | 5.88 | 2.78 | 19.45 |
| **GSOKHM** |  |  |  |  |
| **F-Measure** | 0.9017 | 0.4100 | 0.6902 | 0.4446 |
| **KHM(X,C)** | 3,687 | 70,89 | 3,572,228 | 12,921 |
| **Runtime(sec)** | 2.93 | 7.27 | 3.71 | 22,45 |

The results demonstrate that for all p values the mean of KHM(X,C) function within the proposed GSOKHM was smaller than that of KHM and PSOKHM resulting in more optimized data.On the other hand, we concluded that, except in the case of CMC data in other cases, the value is more than the other two previous samples in GSOKHM. Therefore, this results in more efficiency. From runtime perspective, this algorithm demands much more time compared to KHM but it is comparable with the PSOKHM combinatorial algorithm.

Finally, due to the considerable reduction of the value of KHM(X,C) function and the increase of F-measure, this could be concluded that GSOKHM algorithm generates better clustering quality than that of its precursors.

## 6. Summary

This article examines the hybrid algorithm PSOKHM based on advantages of both PSO and KHM algorithms. In fact, this combination not only improves the converging velocity of PSO algorithm but also prevents KHM from falling into local optima traps. In the present article, theproposed method is SOKHM which combines evolutionary genetic algorithm with PSO on the PSOKHM hybrid algorithm. Four sets of real data have been employed to carry out this experiment. These algorithmscalculate data cluster centers through a sum of all data points based on the harmonic means of a point distance from all centers. Thus, this method has brought about better results compared to KHM and PSOKHM. Furthermore, from F-measure criterion perspective, it has also had much more favorable results. Although that this algorithm is very efficient in clustering, it demands more runtime than KHM. Therefore, this method is not applicable when time is a vital factor in systems.

**References**

[1] Yang, F., Sun, T., and Zhang, C., "An efficient hybrid data clustering method based on K-harmonic means and Particle Swarm Optimization", Expert Systems with Applications, 36(9) 847-852, 2009.

[2] Cui, X., and PotokT. E., "Document clustering using Particle Swarm Optimization", IEEE swarm intelligence symposium, Pasadena, california, 2005.

[3] Güngör, Z., and Ünler, A., "K-harmonic means data clustering with tabu-search method", Applied Mathematical Modelling, 32, 1115–1125, 2008.

[4] Hu, G., Zhou, S., Guan, J., and Hu, X., "Towards effective document clustering: A constrained K-means based approach", Information Processing and Management, 44(4), 1397–1409, 2008.

[5] Hammerly, G., and Elkan, C., "Alternatives to the k-means algorithm that find better clusterings", Proceedings of the 11th international conference on information and knowledge management, pp. 600–607, 2002.

[6] Liu, B., Wang, L., and Jin, Y. H., "An effective hybrid PSO-based algorithm for flow shop scheduling with limited buffers". Computers and Operations Research, 35(9), 2791–2806, 2008.

[7] Maitra, M., and Chatterjee, A., "A hybrid cooperative–comprehensive learning based PSO algorithm for image segmentation using multilevel thresholding", Expert Systems with Applications, 34, 1341–1350, 2008.

[8] Pan, H., Wang, L., and Liu, B., "Particle swarm optimization for function optimization in noisy environment", Applied Mathematics and Computation, 181,908–919, 2006.

[9] Tan, P. N., Steinbach, M., and Kumar, V., "Introduction to data mining", pp. 487– 559, Boston: Addison-Wesley, 2005.

[10] Tjhi, W. C., and Chen, L. H., "A heuristic-based fuzzy co-clustering algorithm for categorization of high-dimensional data", Fuzzy Sets and Systems, 159(4),371–389, 2008.

[11] Ünler, A., and Güngör, Z., "Applying K-harmonic means clustering to the partmachine classification problem". Expert Systems with Applications, pp. 361–406, 2008.

[12] Zhang, B., Hsu, M., and Dayal, U., "K-harmonic means – a data clustering algorithm", Technical Report HPL-1999-124. Hewlett-Packard Laboratories, 1999.

[13] Zhang, B., Hsu, M., and Dayal, U., "K-harmonic means". International workshop on temporal, spatial and spatio-temporal data mining, TSDM2000. Lyon, France, September 12, 2000.

[14] Zhou, H., and Liu, Y. H., "Accurate integration of multi-view range images using k-means clustering". Pattern Recognition, 41(1), 152–175, 2008.

[15] ftp://ftp.ics.uci.edu./pub/machine-learning-databases.

[16] Anderson E 1935 The irises of the Gaspe Peninsula. Bulletin of the American Iris Society 59: 2-5.

[17] Atashpaz-Gargari E, Lucas C 2007b Designing an optimal PID controller using Colonial Competitive.

[18] Fathian M, B Amiri, Maroosi Ali 2008 A honey-bee mating approach on clustering. The International Journal of Advanced Manufacturing Technology 43(9-10): 809-821.

[19] Figueiredo MAT, Jain A K 2002 Unsupervised learning of finite mixture models. IEEE Transactions on Pattern Analysis and Machine Intelligence 24(3): 381-396 Fisher R A 1936 The use of multiple measurements in taxonomic problems. Annals of Eugenics 7: 179-188.

[20] Huan Min Xu, Dong Bo Li 2008 A clustering-based modeling scheme of the manufacturing resources for process planning. The International Journal of Advanced Manufacturing Technology 38(1-2): 154—162

[20] Jasour A M, AtashpazGargari E, Lucas C 2008 Vehicle fuzzy controller design using imperialist competitive algorithm. Second Iranian Joint Congress on Fuzzy and Intelligent Systems.Tehran,Iran.

[21] Kao Y T, Zahar E, I. Kao W 2008 A hybridized approach to data clustering. Expert Systems with Applications 34(3): 1754-1762 Krishna K, Murty M 1999 Genetic k -means algorithm. IEEE Transactions on Systems. Man and Cybernet, Part B: Cybernet 29: 433-439.

[22] Laszlo M, Mukherjee S 2007 A genetic algorithm that exchanges neighboring centers for k-means clustering. Pattern Recognition Letters 28(16): 2359-2366 Lloyd 1982 Least square quantization in PCM. IEEE Transactions on Information Theory 28(2): 129-137.

[23] MacQueen J B 1967 Some Methods for classification and Analysis of Multivariate Observations. Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability 281-297.

[24] Morales A K, Erazo F R 2009 A search space reduction methodology for data mining in large data bases. Engineering Application of Artificial Intelligence 22(1): 92-100.

[25] Mualik U, Bandyopadhyay S 2000 Genetic algorithm-based clustering technique. Journal of Pattern Recognition Letters 33: 1455-1465.

[26] Ng M K, Wong J C 2002 Clustering categorical data sets using tabu search techniques. Journal of Pattern Recognition Letters 35(12): 2783-2790 .

[27] Niknam T, Olamaie J, Amiri B 2008a A hybrid evolutionary algorithm based on ACO and SA for cluster analysis. Journal ofApplied Science 8(15): 2695-2702 .

[28] Niknam T, BahmaniFirouzi B, Nayeripour M 2008b An efficient hybrid evolutionary algorithm for cluster analysis. World Applied Sciences Journal 4(2): 300-307.

[29] Niknam T, Amiri B, Olamaie J, Arefi A 2009 An efficient hybrid evolutionary optimization algorithm basedon PSO and SA for clustering. Journal ofZhejiang University of Science A 10(4):512-519.

[30] Niknam T, Amiri B 2010 An efficient hybrid approach based on PSO, ACO and k-means for clusteranalysis. Journal ofApplied Soft Computing 10(1): 183-197.

[31] Rajabioun R, Hashemzadeh F, Atashpaz-Gargari E 2008a Colonial competitive algorithm: a novel approachfor PID controller design in MIMO distillation column process. Int. J. Intelligent Computing and Cybernetics 1(3): 337-355.

[32] Rajabioun R, Hashemzadeh F, Atashpaz-Gargari E, Mesgari B, RajaieeSalmasi F 2008b Identification of a MIMO evaporator and Its decentralized PID controller tuning using colonial competitive algorithm. The International Federation ofAutomatic Control Congress. Seoul Korea:9952-9957.

[33] Roshanaei M, Atashpaz-Gargari E, Lucas C 2008 Adaptive beamforming using colonial competitive algorithm. 2nd International Joint Conference on Computational Engineering. Vancouver. Canada.

[34] Atashpaz-Gargari E, Lucas C 2007a Imperialist competitive algorithm: an algorithm for optimization inspired by imperialistic competition. IEEE Congress on Evolutionary Computation 4661-4667.

[35] Shelokar P S, Jayaraman V K, Kulkarni B D 2004 An ant colony approach for clustering. AnalyticaChimicaActa 509(2): 187-195.

[36] Sung C S, Jin H W 2000 A tabu-search-based heuristic for clustering. Pattern Recognition Letters 33(5): 849-858.

[37] Tibshirani R, Walther G, Hastie T 2001 Estimating the number of clusters in a data set via the gap statistic.

[38] J. Statistical Soc., Series B 63(2):411-423 Zalik K R 2008 An efficient k-means clustering algorithm. Pattern Recognition Letters 29: 1385-1391.