

A Rank-Based K-medoids Clustering Algorithm by a Specific P System

Qian Li¹

School of Management Science and Engineering, Shandong Normal University, 250014, Jinan, Shandong, China
E-mail: shangrilaying@126.com

Xiyu Liu²

School of Management Science and Engineering, Shandong Normal University, 250014, Jinan, Shandong, China
E-mail: sdxylu@163.com

In this paper, a rank-based K-medoids algorithm by a specific P system is proposed, which exhibits novel aspect of applying membrane computing in clustering. The traditional K-medoids clustering result suffers sensitivity to initial medoids selected randomly. To conquer the defect, we introduce the rank based on similarity between pairs of objects for the initialization. As a biological computing model, P system imitates the structure and function of living cells, and the reactions in it process in a maximal parallel and distributed manner. P system is adequate to solve clustering problem for its high parallelism and lower computational time complexity. A specific P system with a sequence of rule is constructed to realize the rank-based K-medoids algorithm. Through test verification, it can improve the clustering quality.

CENet2015
12-13 September 2015
Shanghai, China

¹Speaker

²Corresponding Author

1. Introduction

Clustering analysis is the process of dividing a set of objects into non-overlapping subsets. Each subset is a cluster with intra-cluster data similar and inter-cluster data dissimilar. Clustering is a rapidly developing area which contributes to research field including machine learning, spatial database technology, biology[1-3] and marketing[4,5], and so on [6].

Membrane computing is a new computing model firstly proposed by Romanian scientist Gheorghe Paun, and that is why the model is called P systems. It simulates the function of living cells, abstracting biochemical reactions and material exchanges to perform calculation process on the cellular level.

The traditional K-medoids clustering result suffers sensitivity to initialization of medoids and local optimum [7]. We adopt the principle of ranking the similarity between pairs of objects to avoid the disadvantages. In this paper, we combine the rank-based K-medoids clustering with a specific P system to reduce the computational time complexity employing the great parallelism of P systems.

2. The Rank-Based K-medoids Algorithm

As a well-known algorithm of partitioning approach, K-medoids selects k objects in the dataset as centers for each cluster and form k clusters through iterations. It reduces sensitivity to outliers which the K-means algorithm suffers [8]. Assuming a dataset $D = \{a_1, a_2, \dots, a_n\}$ with n objects, the K-medoids algorithm divides D into $k < n$ clusters. The similarity among objects can be defined by applying the Euclidean distance, Manhattan distance and so on[9,10]. In this paper, the Euclidean distance is applied and it is assumed that all the data is in a two-dimensional space.

First of all, a $n \times n$ dissimilarity matrix $D_{n,n}$ is constructed as follows:

$$D_{nn} = \begin{pmatrix} w_{11}, w_{12}, \dots, w_{1n} \\ w_{21}, w_{22}, \dots, w_{2n} \\ \dots \\ w_{n1}, w_{n2}, \dots, w_{nn} \end{pmatrix} \quad (2.1)$$

Where, w_{ij} represents the value by rounding the dissimilarity between any the i-th and j-th object.

In the rank-based K-medoids algorithm, we introduce the concept of dissimilarity rank and group. We employ novel method that ranks objects according to their similarities other than the direct use of their similarity values [11]. By this method, the more dissimilar object gets higher rank. That is, $\text{rank}(a_i, a_j) = f$ indicates that a_i is the f-th similar object to a_j among n objects in the dataset. The similarity rank matrix (denoted as SimRank in this paper) is constructed by sorting the similarity values between any object a_j and the other objects. The SimRank is expressed as follows:

$$\text{SimRank} = (r_{i,j}), \text{rank}(a_i, a_j) = r_{i,j}, \forall a_i, a_j \in D \quad (2.2)$$

There are two key points to declare about the SimRank. For one thing, SimRank is a matrix that reveals the similarity rank or closeness among objects in the datasets, and also, it shows the extent of dissimilarity through numbers from 1 to n. For another, SimRank is not always a symmetric matrix on account of two points not always being at the same rank of each other. Group is another important concept in our proposed algorithm. The number of members of

a group is determined by a given parameter g . Within a group, it can compute the dissimilarity rank value (denoted as dv) of an object. The dissimilarity rank value of any object a_i in a group G is defined as follows:

$$dv(a_i) = \sum_{a_j \in G} r_{i,j} \quad (2.3)$$

The steps of the proposed algorithm proceeds as follows:

- 1 Initialize k medoids
 - 1) Calculate the similarities between any pairs of objects based on their distance,
 - 2) Rank similarity values of any objects and thus construct the SimRank matrix,
 - 3) Select k medoids randomly;
- 2 Optimize medoids
 - 1) Choose g most similar objects to each medoids so as to form a group, employing the SimRank (g is a given parameter that indicates the number of members of a group),
 - 2) Compute the dissimilarity rank values of every object in a group, and then update the medoids by a specific object with the maximum dv ,
 - 3) Go to step 2 2) until the medoids remain unchanged;
- 3 Allocate the common objects
Distribute the rest common objects to the most similar medoids.

The similarities between any pairs of objects in dataset is computed once, and updating the medoids costs $O(k \times g)$ per iteration, where k implies the number of clusters and g is the number of members in a group needed to select the next medoids. While in traditional K-medoids clustering, it distribute every objects in a partition to update the new medoids, and it costs $O(n^2)$ per iteration.

The advantage of the rank based K-medoids algorithm lies in ranking objects according to their similarities other than the direct use of their similarity values. As a consequence, the clustering result gains more accuracy.

3. A Specific P System for the Rank-Based K-medoids Algorithm

3.1 The designed P system

The reader is assumed to be acquainted with the basic prerequisites of P system. It is suggested that the reader refer to papers for further learning[12,13].In this section, a specific P system is designed to realize the rank-based K-medoids algorithm. The P system is a construct incorporating skin membrane and k elementary membranes and especially an output membrane. Object a_i represents the original corresponding point a_i in the data set, and the subscript of a_i means a_i represents the i-th data point [14]. The matrix D_m is applied to compare the distance between the n objects.

The P system for the rank-based K-medoids is designed as follows:

$$\pi = (O, H, \mu, M_0, M_1 \dots M_k, M_{c_0}, R_0, R_1, \dots, R_k, P_i, c_0) \text{ Where:}$$

- 1) $O = (\alpha_{11}, \gamma_{1,1,0}, \theta_1, \xi_1, a_1, a_2, \dots, a_n, s_0)$ denotes the collection of objects in the P system;
- 2) $M_0 = (\alpha_{11}, \gamma_{1,1,0}, \theta_1, \xi_1, a_1, a_2, \dots, a_n)$ denotes the initial objects in membrane 0;
- 3) $M_1 = M_2 = \dots = M_k = (s_0)$ denotes the initial objects in membrane 1, 2...k;
- 4) $M_{c_0} = (\lambda)$ denotes the output membrane.

Rule in R_0 [15,16]:

$$\begin{aligned}
 r1 &= \left\{ \alpha_{ij} a_i a_j \rightarrow \alpha_{i(j+1)} a_i a_j U_{ij}^{\omega_{ij}} \mid 1 \leq i, j \leq n \right\} \cup \left\{ \alpha_{i(n+1)} \rightarrow \alpha_{(i+1),1} \mid 1 \leq i \leq n \right\} \cup \left\{ \alpha_{n(n+1)} \rightarrow \lambda \right\} \\
 r2 &= \left\{ U_{ij}^{\omega_{ij}} \gamma_{i,j,q} \rightarrow \gamma_{i,j,q+1} (\zeta_{ij}^{q+1})_{in_1} (\zeta_{ij}^{q+1})_{in_2} \dots (\zeta_{ij}^{q+1})_{in_k} \zeta_{ij}^{q+1} \mid \omega_{ij} = 0, 1 \leq i, j \leq n, 0 \leq q \leq n \right\} \\
 &\cup \left\{ \gamma_{i,j,n+1} \rightarrow \gamma_{(i+1),1,0} \mid 1 \leq i, j \leq n \right\} \cup \left\{ \gamma_{n,n+1} \rightarrow e \right\} \\
 r3 &= \left\{ U_{i1}^{\omega_{i1}} U_{i2}^{\omega_{i2}} \dots U_{in}^{\omega_{in}} \rightarrow U_{i1}^{\omega_{i1}-1} U_{i2}^{\omega_{i2}-1} \dots U_{in}^{\omega_{in}-1} \mid 1 \leq i \leq n \right\} \\
 r4 &= \left\{ e a_i \theta_t \rightarrow e A_{it} \theta_{t+1} \mid 1 \leq i \leq n, 1 \leq t \leq k \right\} \cup \left\{ e \theta_{k+1} \rightarrow \chi_1 \right\} \\
 r5 &= \left\{ \chi_t A_{it} a_j \zeta_{ij}^q \xi_p \rightarrow \chi_t A_{it} G_{pj,in_t} \xi_{p+1} \mid q = 0, 1 \leq i, j \leq n, 1 \leq t \leq k, 1 \leq p \leq m \right\} \\
 r6 &= \left\{ A_{it} \zeta_{ij_1}^{q_1} \zeta_{ij_2}^{q_2} \dots \zeta_{ij_n}^{q_n} \rightarrow A_{it} \zeta_{ij_1}^{q_1-1} \zeta_{ij_2}^{q_2-1} \dots \zeta_{ij_n}^{q_n-1} \mid 1 \leq i, j_t, q_t \leq n \right\} \\
 r7 &= \left\{ \chi_{k+1} \xi_m \rightarrow \xi_1 \eta_{1,in_1} \eta_{1,in_2} \dots \eta_{1,in_k} \right\} \\
 r8 &= \left\{ d^k a_i A_{1j_1} A_{2j_2} \dots A_{kj_k} U_{ij_1}^{\omega_{ij_1}} U_{ij_2}^{\omega_{ij_2}} \dots U_{ij_k}^{\omega_{ij_k}} \rightarrow d^k a_{i,in_p} A_{1j_1,in_1} A_{2j_2,in_2} \dots A_{kj_k,in_k} \right\} \\
 &\left\{ \omega_{ij_p} = 0, 1 \leq j_1, j_2, \dots, j_k \leq n \right\} \\
 r9 &= \left\{ U_{ij_1}^{\omega_{ij_1}} U_{ij_2}^{\omega_{ij_2}} \dots U_{ij_k}^{\omega_{ij_k}} \rightarrow U_{ij_1}^{\omega_{ij_1}-1} U_{ij_2}^{\omega_{ij_2}-1} \dots U_{ij_k}^{\omega_{ij_k}-1} \mid 1 \leq j_1, j_2, \dots, j_k \leq n \right\} \\
 r10 &= \left\{ b^i d^j \rightarrow \phi_{in_1} \phi_{in_2} \dots \phi_{in_k} \mid 1 \leq i \leq n, 0 \leq j \leq k \right\} \cup \left\{ d^k \rightarrow \beta_{in_1} \beta_{in_2} \dots \beta_{in_k} \right\} \\
 r10 &= \left\{ \tau \beta \omega \rightarrow (\beta \omega)_{in_n} \mid \omega \subseteq O_{ij} \cup \{a_p \mid 1 \leq p \leq n\} \right\}
 \end{aligned}$$

Rules in membrane t ($1 \leq t \leq k$) :

$$\begin{aligned}
 r1' &= \left\{ (\beta \omega \rightarrow \beta \omega a_i)_{a_i} \mid 1 \leq i \leq n, \omega \subseteq O_{ij} \cup \{a_p \mid 1 \leq p \leq n\} \right\} \\
 r2' &= \left\{ \beta \omega \rightarrow (\tau \beta \omega)_{out} \mid \omega \subseteq O_{ij} \cup \{a_p \mid 1 \leq p \leq n\} \right\} \\
 r3' &= \left\{ A_{tp} G_{1j_1} G_{2j_2} \dots G_{mj_m} \zeta_{pj_1}^{q_1} \zeta_{pj_2}^{q_2} \dots \zeta_{pj_m}^{q_m} \rightarrow A_{tp} G_{1j_1} G_{2j_2} \dots G_{mj_m} \zeta_{pj_1}^{q_1} \zeta_{pj_2}^{q_2} \dots \zeta_{pj_m}^{q_m} \zeta_p^{Q_p} \right\} \\
 &\left\{ 1 \leq j_1, j_2, \dots, j_m, p \leq n, 1 \leq i \leq m \right\} \\
 r4' &= \left\{ \eta_i G_{ij_i} G_{1j_1} G_{2j_2} \dots G_{(i-1)j_{(i-1)}} G_{(i+1)j_{(i+1)}} \dots G_{mj_m} A_{tp} \zeta_{j_i j_1}^{q_1} \zeta_{j_i j_2}^{q_2} \dots \zeta_{j_i j_{(i-1)}}^{q_{(i-1)}} \zeta_{j_i j_{(i+1)}}^{q_{(i+1)}} \dots \zeta_{j_i j_m}^{q_m} \zeta_{j_i p}^{q_p} \rightarrow \right. \\
 &\left. \eta_i G_{ij_i} G_{1j_1} G_{2j_2} \dots G_{(i-1)j_{(i-1)}} G_{(i+1)j_{(i+1)}} \dots G_{mj_m} A_{tp} \zeta_{j_i j_1}^{q_1} \zeta_{j_i j_2}^{q_2} \dots \zeta_{j_i j_{(i-1)}}^{q_{(i-1)}} \zeta_{j_i j_{(i+1)}}^{q_{(i+1)}} \dots \zeta_{j_i j_m}^{q_m} \zeta_{j_i}^{Q_i} \right\} \\
 &\left\{ 1 \leq j_1, j_2, \dots, j_m, p \leq n, 1 \leq i \leq m \right\} \\
 r5' &= \left\{ s_h A_{tp} G_{ij_i} \zeta_p^{Q_p} \zeta_{j_i}^{Q_{j_i}} \rightarrow s_{h+Q_{j_i}-Q_p} A_{tp} G_{ij_i} \zeta_p^{Q_p} \zeta_{j_i}^{Q_{j_i}} \mid 1 \leq t \leq k, 1 \leq j_i, p \leq n, 1 \leq i \leq m \right\} \\
 r6' &= \left\{ s_{h'} A_{tp} G_{ij_i} \rightarrow s_0 A_{ij_i} G_{ip} \mu \mid h' > 0, 1 \leq t \leq k, 1 \leq j_i, p \leq n, 1 \leq i \leq m \right\} \\
 &\cup \left\{ s_{h'} \rightarrow s_0 \nu \mid h' \leq 0 \right\} \cup \left\{ \eta_i \rightarrow \eta_{i+1} \right\} \cup \left\{ \eta_{m+1} \rightarrow \lambda \right\} \\
 r7' &= \left\{ \zeta_p^{Q_p} \zeta_{j_i}^{Q_{j_i}} \mid 1 \leq j_i, p \leq n \right\} \\
 r8' &= \left\{ \mu^i \nu^j G_{1j_1} G_{2j_2} \dots G_{mj_m} A_{tp} \rightarrow (b a_{j_1} a_{j_2} \dots a_{j_m} A_{tp} \chi_1)_{out} \mid 1 \leq j_1, j_2, \dots, j_m, p \leq n, 1 \leq t \leq k \right\} \\
 &\cup \left\{ \nu^i A_{tp} \rightarrow (A_{tp} d)_{out} \right\}_{-\mu^j} \mid 1 \leq i, j \leq n \\
 \rho &= \left\{ r_i < r_j \mid 1 < i < j \leq 11 \right\} \cup \left\{ r'_i < r'_j \mid 1 \leq i, j \leq 8 \right\}
 \end{aligned}$$

POS (CENet2015) 009

3.2 The computations in P system

In this section, a comprehensive introduction of the computations and responses in the specific P system designed to realize the ranked K-medoids is presented. We make the clarification in the way the proposed algorithm processing.

At the very beginning, rule r1 is performed accordingly to the priority relationship. It calculates the distance as dissimilarity between any two points and produce object U_{ij} . The multiplicity of object U_{ij} , that is object ω_{ij} , represents the distance distinctly. Moreover, the subscript of object α_{ij} is utilized to control the cyclic process until the distance between any two points is acquired. Then it executes rule r2, r3 to construct the SimRankR(\mathcal{V}_{ij}) where $\mathcal{V}_{i,j,q}$ means the similarity rank of α_i towards α_j is q. Meanwhile, object ζ_{ij}^{q+1} is generated and transferred to membranes labeled from 1 to k. This course ends with yielding an object e to stimulate the rule r4 of selecting the initial k medoids randomly. When it accomplishes the last loop, the object e and θ_{k+1} are converted to \mathcal{X}_1 to impel the rule r5. It performs rule r5 and r6 to select the group of the most preferred objects. Once it completes the loop of selecting the group of the most preferred objects for all medoids, it starts to deliver an object η_i to each membrane labeled from 1 to k to initiate the computations in the corresponding membrane.

Now, the responses in membranes labeled from 1 to k are activated. Rule r3' is performed to calculate the total preference factors of the medoid A_p by adding the preference factors it holds toward other members in the same group. When responses in membranes labeled from 1 to k are accomplished, the membrane 0 with triggering rules repeats the process of clustering until there is no b in it, which indicates that the medoids in all k membranes remain unchanged and the clustering adjustment process terminates [16]. And an object β is generated to stimulate rules in membranes labeled from 1 to k to output the clustering result.

Finally, the result in the form of string enters the output membrane c_0 . One computation process is accomplished.

3.3 Experiment and analysis

In order to give a better interpretation of our P system model for the ranked K-medoids clustering, we take an example to simulate the procedure of the P system. There are 7 points:

$$a_1 (1, 2), a_2 (2, 2), a_3 (3, 0), a_4 (3, 3), a_5 (6, 1), a_6 (7, 1), a_7 (8, 3)$$

The P system is supposed to distribute the points into 2 clusters with the given parameter m of the value 3. Diagram 3 depicts the original state of the 7 points.

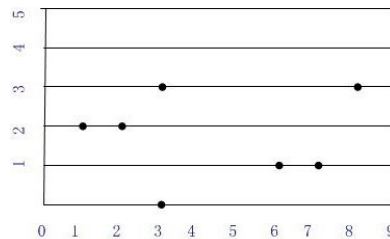


Figure 3: The initial state of the points

First of all, the dissimilarity matrix $D_{7,7}$ is presented. We calculate the square of the distance between any two points as dissimilarity, thus $D_{7,7}$ is the same to $D_{7,7}'$. The SimRank

matrix R is constructed by sorting the dissimilarity values between any object a_i and the other objects.

$$D_{7,7} = \begin{pmatrix} 0 & 1 & 8 & 5 & 26 & 37 & 50 \\ 1 & 0 & 5 & 2 & 17 & 26 & 37 \\ 8 & 5 & 0 & 9 & 10 & 17 & 34 \\ 5 & 2 & 9 & 0 & 13 & 20 & 25 \\ 26 & 17 & 10 & 13 & 0 & 1 & 8 \\ 37 & 26 & 17 & 20 & 1 & 0 & 5 \\ 50 & 37 & 34 & 25 & 8 & 5 & 0 \end{pmatrix} \quad R_{7,7} = \begin{pmatrix} 1 & 2 & 4 & 3 & 5 & 6 & 7 \\ 2 & 1 & 4 & 3 & 5 & 6 & 7 \\ 3 & 2 & 1 & 4 & 5 & 6 & 7 \\ 3 & 2 & 4 & 1 & 5 & 6 & 7 \\ 7 & 6 & 4 & 5 & 1 & 2 & 3 \\ 7 & 6 & 4 & 5 & 2 & 1 & 3 \\ 7 & 6 & 5 & 4 & 3 & 2 & 1 \end{pmatrix}$$

(3.1)

Then, the process of the clustering with the designed P system begins.

Eventually, the clustering result was attained that the 7 points were classified into 2 clusters. Consequently, the ranked K-medoids algorithm certified effective. And the clustering effect sketch is shown in Fig.4.

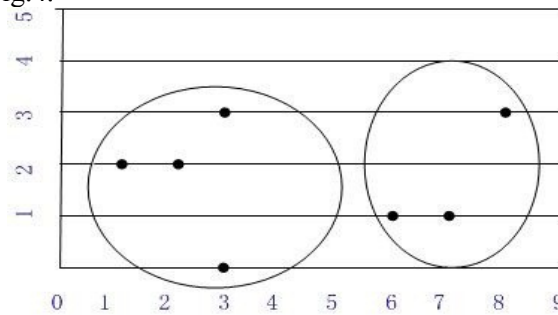


Figure 4: The final clustering result

4. Conclusion

In this paper we propose a specific P system to solve the clustering problem in the framework of the cellular computing with membranes. Tests prove that the P system is adequate to implement the ranked K-medoids algorithm for its high parallelism. However, some questions remain to be discussed furthermore. On the one hand, the feasibility and the efficiency of the model in large database need to be studied. On the other hand, the number of members in a group is determined by a given parameter g which is determined by data experiment. Lastly, it is of great significance to realize other clustering methods by a P system.

References

- [1] J. Ponomarenko, T. Merkulova, G. Orlova, O. Fokin, E. Gorshkov, M. Ponomarenko, *Mining DNA sequences to predict sites which mutations cause genetic diseases*, Knowledge-Based Systems. 15 (4) 225–233.
- [2] J. Shi, Z. Luo, 2010, *Nonlinear dimensionality reduction of gene expression data for visualization and clustering analysis of cancer tissue samples*, Computers in Biology and Medicine 40 (8) 723–732.
- [3] D. Sebiskveradze, V. Vrabie, C. Gobinet, A. Durlach, P. Bernard, E. Ly, M. Manfait, P. Jeannesson, O. Piot, 2011, *Automation of an algorithm based on fuzzy clustering for analyzing tumoral heterogeneity in human skin carcinoma tissue sections*, Laboratory Investigation. 91 (5) 799–811.
- [4] P. C. Chang, C. H. Liu, C.Y. Fan, 2009, *Data clustering and fuzzy neural network for sales forecasting: a case study in printed circuit board industry*, Knowledge-Based Systems. 22 (5) 344–355.

- [5] E. Hadavandi, H. Shavandi, A. Ghanbari, 2010, *Integration of genetic fuzzy systems and artificial neural networks for stock price forecasting*, *Knowledge-Based Systems*. 23 (8) 800–808.
- [6] J. Han and M. Kambr, 2012, *Data Mining Concepts and Techniques*. USA: Elsevier Inc., ch.8
- [7] R. Xu, D. Wunsch, 2005, *Survey of clustering algorithms*, *IEEE Transactions on Neural Networks*. 16 (3) 645–678.
- [8] S. Khan, A Ahmadb, 2004, *Cluster center initialization algorithm for K-means clustering*, *Pattern Recognition Letters*, 25(1):1293–1302.
- [9] Park, H.S., Jun, C.H. 2009, *A simple and fast algorithm for K-medoids clustering*. *Expert Systems with Applications*. 36(2), 3336–3341.
- [10] Q. Zhang, I. Couloigner, 2005, *A new and efficient K-medoid algorithm for spatial clustering*, *Computational Science and Its Applications_ICCSA*. Springer 207–224.
- [11] S. Mohammad ,R. Zadegan, M. Mirzaie, F. Sadoughi, 2013, *Ranked k-medoids: A fast and accurate rank-based partitioning algorithm for clustering large datasets*. *Knowledge-Based Systems*.133–143.
- [12] G. Paun, 2010, *A quick introduction to membrane computing*, *The Journal of Logic and Algebraic Programing*, 79(3):291-294.
- [13] G. Paun, G. Rozenberg and A. Salomaa. 2010, *Membrane Computing*. New York: Oxford University press, 282-301.
- [14] J.D. Martin-Guerrero, A. Palomares, E. Balaguer-Ballester, E. Soria-Olivas, J. Gomez-Sanchis, A. Soriano-Asensi, 2006, *Studying the feasibility of a recommender in a citizen web portal based on user modeling and clustering algorithms*, *Expert Systems with Applications*. 30 (2) 299–312.
- [15] Y. Zhao, X. Liu, J. Qu, 2012, *The K-medoids clustering algorithm by a class of P system*, *Information and Computational Science*, J., 9(18):215-22.
- [16] L. Han, L. Xiang , X. Liu, J. Luan, 2014, *The K-medoids Algorithm with Initial Centers Optimized Based on a P System*, *Information and Computational Science J.*, 11(6):132-144.