# Improved rough *K*-means clustering algorithm based on firefly algorithm

## Tingyu Ye, Jun Ye* and Lei Wang

School of Information Engineering,
Nanchang Institute of Technology,
Nanchang 330000, China
Email: yty01122@163.com
Email: yejun68@sina.com
Email: ezhoulei@163.com
*Corresponding author

**Abstract:** The rough *K*-means clustering algorithm has a strong ability to deal with data with uncertain boundaries. However, this algorithm also has limitations such as sensitivity to initial data selection, as well as it use of fixed weights and thresholds, which results in unstable clustering results and decreased accuracy. In response to this problem, combined with the firefly algorithm, the original algorithm has been improved from three aspects. Firstly, based on the ratio of the number of objects in the dataset to the product of the difference of the objects in the dataset, a more reasonable method of dynamically adjusting the weights of approximation and boundary set is designed. Secondly, a method of adaptively realising the threshold $\varepsilon$ associated with the number of iterations is given. Then, by constructing a new objective function, and take the objective function value as the firefly brightness intensity to perform the search and update iteration of the initial cluster centre point, the optimal solution obtained by each iteration of firefly is taken as the initial centre position of the algorithm. Experiment result shows that the new algorithm has improved the clustering effect.

**Keywords:** rough *K*-means algorithm; firefly algorithm; cluster centre; lower approximation and boundary set; objective function.

**Biographical notes:** Tingyu Ye is currently working towards his MSc at the School of Information Engineering, Nanchang Institute of Technology, China. His research interests include smart computing and machine learning.

Jun Ye is a Professor at the School of Information Engineering, Nanchang Institute of Technology, China. His research interests include rough set, evolutionary computation, knowledge discovery and data mining.

Lei Wang is a Professor at the School of Information Engineering, Nanchang Institute of Technology, China. His research interests include rough set, evolutionary computation, data mining.

## 1   Introduction

Aiming at the limitation of *K*-means clustering algorithm in dealing with uncertain boundary objects, a rough *K*-means clustering algorithm is proposed by using the advantage of rough set in dealing with uncertain boundary data (Lingras and West, 2004). The algorithm solves the problem of clustering objects with uncertain boundaries. However, arbitrarily selecting the initial clustering centre and setting fixed weights and thresholds affect the clustering accuracy. Researchers have proposed many improved algorithms. Zhou (2010) gives a definition of the density of the area where the object is located, and determines the initial cluster centre according to the iteration of the sample density, it improves the stability of cluster centres. Li et al. (2013) combines the theory of granular computing and uses the principle of maximum and minimum distance combination to find the initial cluster centres, avoiding the influence of isolated points. Peters (2014) replaced the absolute distance threshold with the relative distance threshold, and gave a weight-based calculation method for the mean value of cluster centres, and achieved good results. Sun et al. (2016) introduced the concept of fuzzy set and defined a method for measuring boundary set objects. This algorithm adaptively adjusts the influence coefficient of the samples in the boundary area on the cross clusters, and weakens the influence of the boundary area on the central mean. Ofek and Okach (2017) introduced an adaptive measurement of the imbalance of cluster class size, and proposed an improved method based on the imbalance measurement of cluster class size. Liu et al. (2019b) proposed an improved clustering algorithm based on ant colony algorithm, it reduces the sensitivity of the initial centre point and the adverse effects of data differences. Some other methods have also achieved better results (Liu et al., 2019a; Li et al., 2019).

Firefly algorithm is a meta-heuristic algorithm, which is widely used to solve optimisation problems. Some research results show that, in terms of solving clustering optimisation problems, the Firefly algorithm is superior to other swarm intelligence algorithms (Wang et al., 2017). This paper starts with the intelligent optimisation method, regards clustering as a combinatorial optimisation problem, and improves the rough *K*-means from three aspects.

## 2   Rough *K*-means clustering algorithm

Lingras introduced these two operators of rough set into the *K*-means algorithm, and on this basis. He proposed a rough *K*-means clustering algorithm. The main contents are as follows.

**Definition 1** Assuming $U = \{Y_1, Y_2, \cdots Y_N\}$, for category $C_k$, the upper approximation set is $\overline{C_k}$, the lower approximation set is $\underline{C_k}$, and the boundary set $C_k^B = \overline{C_k} - \underline{C_k}$, then the update formula of cluster centre $m_k$ of category $C_k$ is:

$$
m_k = \begin{cases}
W_l \displaystyle\sum_{x_i \in \underline{C_k}} \frac{x_i}{|\underline{C_k}|} + W_b \displaystyle\sum_{x_j \in C_k^B} \frac{x_j}{|C_k^B|} & C_k^B \neq \phi \\
\displaystyle\sum_{x_i \in \underline{C_k}} \frac{x_i}{|\underline{C_k}|} & C_k^B = \phi
\end{cases}
\tag{1}
$$

Among them, the $W_l$ is the weight of the lower approximation, the $W_b$ represent the weights of the boundary, and $W_l + W_b = 1$. $\left|\underline{C_k}\right|$ is the number of data in the lower approximate, the $\left|C_k^B\right|$ represents the number of data in the boundary.

The algorithm is based on whether the difference between the distance between other cluster centres and the object and the minimum distance is less than the threshold $\varepsilon$, and assigns the object to be classified to the upper approximation or the lower approximation set. Update the position of the cluster centre by formula (1), and repeat this process until each cluster centre is unchanged.

## 3   Firefly algorithm

The size of the firefly population is $N$, the $i$th firefly is represented as $X_i = (x_{i1}, x_{i2} \cdots x_{iD})$. The attraction of fireflies to each other depends on two factors, brightness and attractiveness.

**Definition 2:** The light intensity of fireflies:

$$
I = I_0 e^{-\gamma r_{ij}}, \; I \infty f(x_i),
\tag{2}
$$

Among them, the $x_i$ is the $i$th firefly, $f(x_i)$ represents the value of the objective function $I \infty f(x_i)$ of the specific problem. The $I_0$ is the light intensity of the firefly at $\gamma = 0$.

**Definition 3:** Firefly attraction:

$$
\beta = \beta_0 e^{-\gamma r_{ij}^2}
\tag{3}
$$

Among them, $\beta_0$ is the attraction at $\gamma = 0$. The Euclidean distance between two fireflies is determined by the following formula:

$$
r_{ij} = \left\| x_i - x_j \right\| = \sqrt{\sum_{d=1}^{D} (x_{id} - x_{jd})^2}
\tag{4}
$$

Among them, $D$ represents the dimensionality of the problem to be solved, $d = 1, 2, \ldots D$.

**Definition 4:** Location update formula:

$$x_i(t+1) = x_i(t) + \beta(x_j(t) - x_i(t)) + \alpha(\text{rand} - 0.5) \tag{5}$$

Equation (5) represents the position update formula for the movement of firefly $i$ to firefly $j$ ($i \neq j$), $\alpha$ is the step factor in [0, 1], rand is a random number, rand $\in$ [0, 1], $t$ is the number of iterations, the sub-item $\alpha(\text{rand} - 0.5)$ is a perturbation item.

## 4 Improved new algorithm

### 4.1 Improvement of lower approximation and boundary weight

In the original algorithm, the $W_l$ and $W_b$ in cluster centre update formula (1) adopt fixed weight values, that is, the values of and remain unchanged throughout the clustering process. Since the search for cluster centres is a dynamic process of continuous iteration and update, the values of $W_l$ and $W_b$ will change as the iteration changes. Obviously, taking the same value in different periods will cause the clustering accuracy to drop significantly. A reasonable measure of the importance of the approximate and boundary regions to the update of the cluster centre position. And dynamic allocation of the weights of $W_l$ and $W_b$ is one of the effective ways to improve the clustering accuracy. For this reason, Zhou (2010) proposed improved methods, after analysing the changes in the initial and later positions of the cluster centres. They constructed a Logistic growth curve:

$$W_l = \frac{1}{(k + ae^{-bt})} \qquad W_b = 1 - W_l \tag{6}$$

Among them, the $t$ is the number of times, the $k$, $a$, and $b$ are the function adjustment parameters. From equation (6), we can see that as the number of iterations increases, the weight of $W_l$ gradually increases, while $W_b$ gradually decreases. To a certain extent, the curve dynamically reflects its importance to the clustering centre. However, the curve has shortcomings. Because the adjustment parameters $k$, $a$ and $b$ on the curve are given manually in advance and are subjective. When the number of iterations increases, the weight of the lower approximation set hardly changes. Ofek and Okach (2017) defines a comparison between the number of data in the lower approximation set and the number of data in the upper approximation set, and uses this as an adaptive weighting formula:

$$\frac{W_l}{W_b} = \frac{|C_k|}{|\overline{C_k}|}, \; W_l + W_b = 1 \tag{7}$$

Liu et al. (2019b) gives an adaptive weighting formula that compares the number of data in the lower approximate set with the number of data in the boundary set.

$$\frac{W_l}{W_b} = \frac{|C_k|}{|C_k^B|}, \; W_l + W_b = 1 \tag{8}$$

Equations (7) and (8) dynamically determine the weight value according to the change in the number of sample objects, it objectively reflects the dynamic change of the number of

sample objects in the upper and lower approximate sets and boundary sets during the clustering process. However, only the number of objects in the above and below approximate sets and boundary sets are used to determine the weights, which can neither reflect the differences in the distribution of objects in the same category, nor the distribution of samples in different categories. In fact, within the same category, the distance distribution of the objects in the dataset relative to the cluster centre is not the same, and its effect on the cluster centre is different. Among different categories, the importance of the objects in the dataset to the cluster centre is also different. Therefore, determine the weight values of the two datasets, not only the influence of changes in the number of objects in the lower approximation and boundary concentration must be considered, but also the influence of the objects on the cluster centre due to the difference in distance distribution. Combining these factors, this paper proposes an adaptive and dynamic adjustment of the approximate and boundary set weight method.

**Definition 5:** Let $X_i$ be the object in the lower approximation set of category $C_k$ and $m_k$ is the centre. Definition the distance distribution from the data object of the lower approximate set to the cluster centre is:

$$d(\underline{C_k}, m_k) = \sum_{x_i \in \underline{C_k}} d(x_i, m_k) \tag{9}$$

**Definition 6:** Let $X_j$, be the object in the boundary set of category $C_k$ and $m_k$ is the centre. Definition the distance distribution from the boundary set object to the cluster centre is:

$$d(C_k^B, m_k) = \sum_{x_i \in C_k^B} d(x_j, m_k) \tag{10}$$

Among them, $d(x_i, m_k)$ and $d(x_j, m_k)$ represent the Euclidean distance from the object and to the cluster centre respectively. From the above analysis of equations (9) and (10), we can see that the more data in the dataset, the greater the distribution distance, and the more important the influence of the dataset on the location of the cluster centre. Conversely, the less data in the dataset, the smaller the distribution distance, and the smaller the impact of the dataset on the location of the cluster centre. Therefore, we can use the ratio of the product of the number of objects in the dataset and the distribution distance to adaptively adjust the weight value, assuming that $W_l^{'}$ is the weight of the lower approximation, $W_b^{'}$ is the weight of the boundary set, the formula is defined as:

$$\frac{W_l^{'}}{W_b^{'}} = \frac{|\underline{C_k}| . d(\underline{C_k}, m_k)}{|C_k^B| . d(C_k^B, m_k)}, |C_k^B| \neq \phi \tag{11}$$

Among them, $W_l^{'} + W_b^{'} = 1$, $\left|\underline{C_k}\right|$ and $\left|C_k^B\right|$ respectively represent the number of objects. In the initial stage of clustering, most of the objects are not divided and are in the boundary concentration. That is, there are many objects in $\left|C_k^B\right|$, the distribution distance of $d(C_k^B, m_{k)})$ is large, at this time, the influence of the boundary set on the cluster centre is greater, and the weight that we can get from the ratio equation (11) is greater. At this time, the influence of the boundary set on centre is greater, and the weight value of $W_b^{'}$

obtained from the ratio equation (11) is greater. In the later stage of clustering, most objects are divided into the lower approximation set, that is, there are more objects in $\left|\underline{C_k}\right|$, and the distribution distance of $d(\underline{C_k}, m_k)$ is larger. At this time, the importance of the lower approximate set to the cluster centre is greater. Similarly, the weight of $W_l^{'}$ obtained from the ratio equation (11) is greater. Therefore, we redefine the new cluster centre update calculation formula as follows:

$$m_k = \begin{cases} W_l^{'} \sum_{x_i \in \underline{C_k}} \dfrac{x_i}{|\underline{C_k}|} + W_b^{'} \sum_{x_j \in C_k^B} \dfrac{x_j}{|C_k^B|} & C_k^B \neq \phi \\ \sum_{x_i \in \underline{C_k}} \dfrac{x_i}{|\underline{C_k}|} & C_k^B = \phi \end{cases} \qquad (12)$$

Among them, $W_l^{'} + W_b^{'} = 1$, the $W_l^{'}$ is the weight of the lower approximation, the $W_b^{'}$ represent the weights of the boundary, and satisfies the above formula (11). Adaptive adjustment threshold

## 4.2   Threshold adaptive improvement

The threshold $\varepsilon$ determines whether the sample object is divided into the upper approximation or the lower approximation set in the category. Therefore, the reasonable selection of the threshold $\varepsilon$ is very important. In the classic rough $K$-means algorithm, the $\varepsilon$ is artificially given a fixed value, and this value does not change with iteration. In fact, looking at the changes of clustering objects from the clustering process, the beginning of clustering, and the attribution relationship of data objects is not clear, and $\varepsilon$ should be larger, so that most of the data objects are classified into the upper approximation set. In the later stage of clustering, the number of iterations continues to increase, and the belonging relationship of the objects becomes clear, and more and more data objects are classified into the lower approximation set of the class, and $\varepsilon$ should be smaller. This paper designs an adaptive threshold $\varepsilon$ realisation method:

$$\varepsilon = \varepsilon - \frac{1}{t^2} \qquad (13)$$

Among them, $t$ is the number of iterations, and the initial value of $t$ is 2. Obviously, initial value of the $\varepsilon$ cannot be too large or too small, too large will increase the number of iterations, and computationally expensive; If $\varepsilon$ is too small, the initial clustering will result in an empty approximate area, which will affect the cluster centre update.

## 4.3   Design objective function

Designing the objective function is the key step of the firefly improved algorithm used in this paper. It directly determines the clustering direction of the firefly, the number of iterations and the pros and cons of the solution, which is related to the clustering accuracy and anti-noise ability of the algorithm. This paper defines two functions of the degree of aggregation within a category and the degree of dispersion between categories to

construct the objective function, and use this to find the optimal initial cluster centre and perform clustering.

**Definition 7:** Suppose the object set $U = \{x_i, i = 1, 2, \ldots, N\}$ has $N$ samples and $K$ cluster centres $C_K = \{m_1, m_2 \cdots m_k\}$, then the internal aggregation function is:

$$J(C_K) = \sum_{k=1}^{K} (W_l' \sum_{x_i \in \underline{C_k}} \sqrt{|x_i - m_k|^2} + W_b' \sum_{x_i \in C_k^B} \sqrt{|x_i - m_k|^2}) \tag{14}$$

Among them, the $x_i$ is the object to be classified, the $m_k$ is the cluster centre of category $C_k$, the $W_l'$ is the weight of the lower approximate set of the $k$th category, the $W_b'$ is the weight of the boundary set of the $k$th category, it satisfying the formula (11). The goal of the cohesion function is to minimise the sum of the distances from the object $x_i$ to the cluster centre $m_k$, that is, the maximum degree of aggregation of objects in the same category, which reflects the degree of aggregation between objects in the same category.

**Definition 8:** Suppose the object set $U$ has $N$ samples and $K$ cluster centres $K$, then the inter-class dispersion function of each cluster centre is:

$$D(C_K) = \sum_{i=1}^{K} \omega \sum_{j=i+1}^{k} \sqrt{|m_i - m_j|^2} , \quad \omega = \frac{1}{\sqrt{k}} \tag{15}$$

Among them, $m_i$ and $m_j$ represent the cluster centres of class $C_i$ and class $C_j$. As the number of iterations increases, the degree of internal aggregation of objects in the same category continues to decrease, and the dispersion between different categories continues to expand. In order to avoid clustering centres falling in sparse areas or isolated points, this paper uses a weight coefficient $\omega = \frac{1}{\sqrt{k}}$ to balance the distance between clustering within clusters and dispersion between clusters, the $k$ represents the number of categories, and its initial value is 1. It makes the results more consistent with the actual distribution of the data.

**Definition 9:** Objective function:

$$f(t) = \frac{D(C_K)}{J(C_K)} \tag{16}$$

Among them, $J(C_K)$ is the value of the aggregation function within the class, and $D(C_K)$ is the value of the inter-class dispersion function. In formula (16), the smaller the class cohesion distance and the larger the separation distance between classes, the larger the objective function value. And the better the clustering effect obtained.

### 4.4   Algorithm implementation steps

The core of the new algorithm is to use the objective function value designed by formula (16) to represent the brightness of the firefly, namely:

$$I = f(t) \tag{17}$$

The value of $I$ is determined by $f(t)$. The larger the value of $f(t)$, the larger the value of $I$, that is, the brighter the firefly. According to the objective function we designed, every search for cluster centres performed by Firefly is also a sub-process of clustering sample objects, This not only ensures the maximum aggregation of sample objects in each category, and the distance between different categories is as discrete as possible, And can avoid the influence of isolated points, effectively reducing the number of iterations. The best clustering result is obtained while finding the best cluster centre. The main steps of the algorithm:

**Step 1:** Given the number of categories $K$, the number of fireflies $N$, and the initial values of $T_{max}$, $\gamma$, $\beta_0$, $\varepsilon$.

**Step 2:** Select $K$ sample points as the location of the fireflies. Then calculate the distance between the object to be classified $x_i$ and each cluster centre, at the same time, divide $x_i$ into the upper approximate set $C_k$ of the category corresponding to the nearest centre $\overline{C_k}$.

**Step 3:** For any clustered object $x_i$, find the cluster centre with the smallest distance from it. If there are other cluster centres $C_i$ such that $\left| C_i - C_j \right| \le \varepsilon$, then the object $x_i$ is classified into the corresponding category $\overline{C_k}$, otherwise it is classified into the corresponding $\overline{C_k}$.

**Step 4:** According to the results obtained in step 3, the values of $W_l'$ and $W_b'$ are calculated by equation (11), and the position of the cluster centre is updated by equation (12), and the value of v is adaptively determined by equation (13).

**Step 5:** Calculate the value of the objective function from equations (14)–(16), and get the brightness of the firefly from equation (17).

**Step 6:** If $I_j > I_i$ is present, firefly $i$ will move to the position of firefly $j$. The size of the movement is determined by equation (3), and the position of firefly $i$ will be updated through equation (5).

**Step 7:** If the maximum number of times set by the algorithm is reached, go to step 8, otherwise go to step 3.

**Step 8:** Output $K$ categories.

## 5   Analysis of experimental results

In the environment of Win 7 operating system and application software Matlab9.0, this paper uses three datasets of Iris, Wine and Balance-scale in the UCI library to verify the effectiveness and effect of the algorithm. We analyse from the three aspects of clustering

accuracy, number of iterations and clustering effect. In order not to lose generality, we choose the other three algorithms for comparison (Lingras and West, 2004; Li et al., 2013; Liu et al., 2019b).

- *Comparison of clustering accuracy.* Set $K = 3$, $N = 120$, $T_{max} = 150$, $\gamma = 1$, $\beta = 1$, $\alpha = 0.06$. Since Lingras and West (2004) and Li et al. (2013) use a fixed weight method, different values will change the position of the cluster centre. Therefore, in the experiment, take the best values for literature 1 and 3, set $W_l' = 0.8$, $W_b' = 0.2$, $\varepsilon = 0.05$. Liu et al. (2019b) and this paper use adaptive weights and $\varepsilon \in (0, 0.5]$ and perform 20 experiments on the 3 selected UCI datasets to get the average value. The experimental results are shown in Table 1.

**Table 1**     Comparison of clustering accuracy

| Accuracy | Iris | Wine | Balance-Scale |
|---|---|---|---|
| Lingras and West (2004) | 85.32 | 69.11 | 52.26 |
| Li et al. (2013) | 89.33 | 73.60 | 56.96 |
| Liu et al. (2019b) | 90.12 | 73.70 | 57.34 |
| This paper | 90.79 | 74.15 | 57.35 |

From Table 1, we know that on the Iris and Wine datasets, the clustering accuracy of this algorithm has been improved, which is better than the other three algorithms, on the Balance-Scale dataset, it is better than the Lingras and West (2004) and Li et al. (2013) algorithms, and slightly better than Liu et al. (2019b).

- *Comparison of iteration times.* Set the same parameters as the above, and perform 20 experiments on the selected 3 UCI datasets to get the average value. The results are shown in Figure 1.

**Figure 1**    Comparison of running times (see online version for colours)
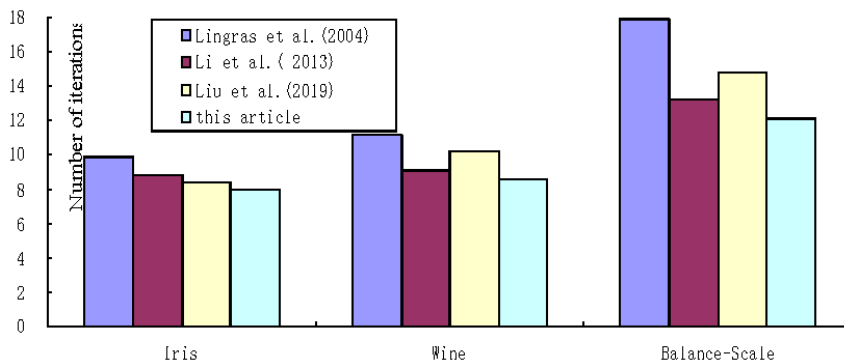


Figure 1 shows that the number of iterations of the algorithm in the three datasets is better than the other three algorithms. Because every iteration of the algorithm in this paper. It needs to calculate the distribution distance of the objects in the two datasets relative to the cluster centre. Therefore, the computational workload is larger than that of Lingras and West (2004) and Liu et al. (2019b), and roughly the same as Li et al. (2013).

- *Comparison of clustering effects.* This paper uses the ratio of the distribution of objects in the category relative to the entire data centre to the distance of the objects in the category in the Sun et al. (2016) to measure the clustering effect. It is an effective indicator to measure the pros and cons of the clustering effect. The formula is:

$$h = \frac{\sum\limits_{i=1}^{K}\sum\limits_{j=1}^{|C_i|}|x_{ij} - \overline{m}|^2}{\sum\limits_{i=1}^{K}\sum\limits_{j=1}^{|C_i|}|x_{ij} - m_i|^2} \qquad (18)$$

Among them, $\overline{m} = \sum\limits_{i=1}^{K}\sum\limits_{J}^{N} x_{ij} \Big/ N$ is the centre of all data objects, and $N$ is the total number of objects.

The four algorithms obtain the distance ratio $h$ on the three datasets as shown in Figures 2–4.

**Figure 2**   The object distance ratio on the Iris data (see online version for colours)
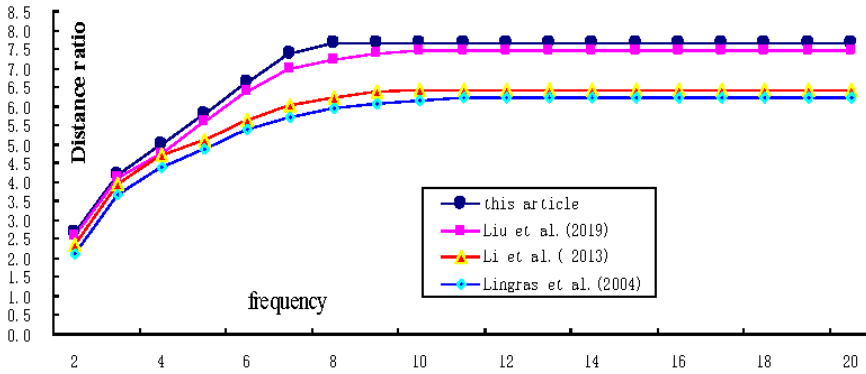


**Figure 3**   The object distance ratio on the Wine data (see online version for colours)
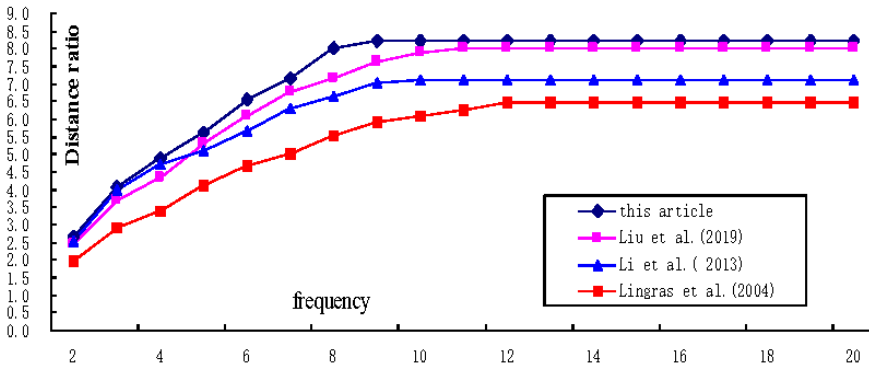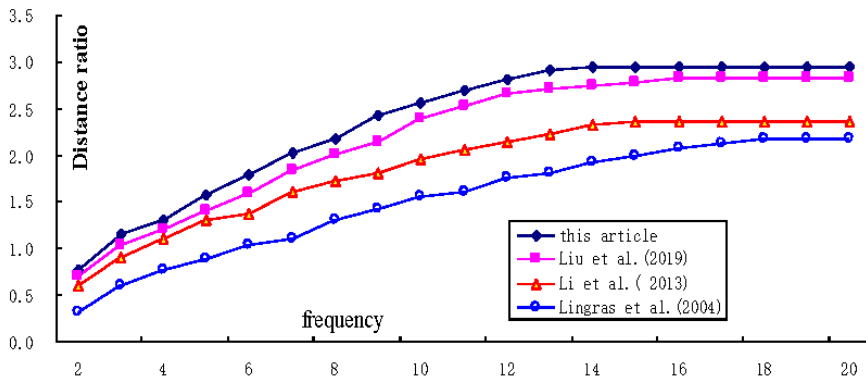
**Figure 4** The object distance ratio on the Balance-Scale data (see online version for colours)



It can be seen from equation (18) that the tighter the distribution of objects in the same category, the smaller the distance; the more discrete the objects in different categories, the greater the distance to the entire data centre, and the greater the value of $h$. It can be seen from Figures 2–4 that the value of $h$ obtained by the algorithm in this paper is larger than the other three algorithms, and the convergence speed is significantly faster than the other three algorithms. It has achieved a good clustering effect. Based on the comparison of these three aspects, the clustering results obtained by combining the Firefly algorithm in this paper are basically consistent with the actual distribution of the data in the three datasets.

## 6 Conclusion

As the application field of clustering continues to expand, its research value is also prominent. Improving the quality of clustering and improving the performance of clustering algorithms are the goals that the majority of researchers have been working hard on. In this paper, combined with the Firefly algorithm, the original algorithm has been improved from the three aspects of optimising the initial clustering centre, dynamic adjustment of the approximation, boundary weight and adaptive adjustment threshold $\varepsilon$. From another angle, the method of improving the original algorithm is discussed, which provides ideas for improving the adverse effect of the algorithm due to the sensitivity of the initial data, at the same time, the new algorithm provides a more reasonable way to dynamically adjust the lower approximation and boundary weights. Taking into account the impact of changes in the number of data objects in the two datasets on the centre, it also considers the impact of data objects on the cluster centres due to the difference in distance distribution. The overall performance of the rough $K$-means clustering algorithm is improved.

## Acknowledgement

# References

Li, F.J., Qian, Y.H. and Wang, J.T. (2019) 'Clustering ensemble based on sample's stability', *Artificial Intelligence*, No. 273, pp.37–55.

Li, L., Luo, K. and Zhou, B.X. (2013) 'Rough clustering algorithm based on granular computing', *Application Research of Computers,* Vol. 30, No. 10, pp.2916–2919.

Lingras, P. and West, C. (2004) 'Interval set clustering of web user with rough *K*-means', *Journal of Intelligent Information Systems*, Vol. 23, No. 1, pp.5–16.

Liu, J.B., Li, H.X. and Zhou, X.Z. (2019a) 'An optimization-based formulation for three-way decisions', *Information Sciences*, Vol. 495, No.3, pp.185–214.

Liu, Y., Wang, H.Q. and Zhang, X.H. (2019b) 'An improved rough *K*-means clustering algorithm combining ant colony algorithm', *Journal of Data Acquisition and Processing*, Vol. 34, No. 2, pp.341–348.

Ofek, N.R. and Okach, L. (2017) 'Fast-CBUS: a fast clustering-based under sampling method for addressing the class imbalance problem', *Neuro Computing*, Vol. 243, No. 21, pp.88–102.

Peters, G. (2014) 'Rough clustering utilizing the principle of indifference', *Information Sciences*, Vol. 277, No. 1, pp.358–374.

Sun, Z.P., Qian, X.Z. and Wu, Q. (2016) 'Self-adaptive rough *K*-means algorithm based on weighted distance', *Application Research of Computers*, Vol. 33, No. 7, pp.1987–1991.

Wang, H., Wang, W.J. and Sun, H. (2017) 'Firefly algorithm with adaptive control parameters', *Soft Computing*, Vol. 21, No. 17, pp.5091–5102.

Zhou, T. (2010) 'Adaptive rough *k*-means clustering algorithm', *Computer Engineering and Applications*, Vol. 46, No. 26, pp.7–10.