

Received September 16, 2021, accepted October 28, 2021, date of publication November 30, 2021, date of current version December 27, 2021.

Digital Object Identifier 10.1109/ACCESS.2021.3131368

Dynamic Metric Accelerated Method for Fuzzy Clustering

SHENGBING XU^{1,2}, WEI CAI¹, HONGXI XIA¹, BO LIU³, AND JIE XU¹

¹School of Mathematics and Statistics, Guangdong University of Technology, Guangzhou 510520, China

²School of Computers, Guangdong University of Technology, Guangzhou 510006, China

³School of Automation, Guangdong University of Technology, Guangzhou 510006, China

Corresponding author: Wei Cai (caiwei-email@qq.com)

This work was supported in part by the Natural Science Foundation of China under Grant 61876044 (Algorithm research and application of positive sample and unmarked sample learning).

ABSTRACT Features in data samples usually need a unified dimension by a standardization process before clustering. However, there still exists a non-standardized metric in which the distance between samples is greater than 1 after features are standardized. It is difficult to find the optimal search path if the data sample metrics are not standardized. To address this problem, we develop a dynamic-metric accelerated method for fuzzy clustering by introducing a metric matrix, whose diagonal elements consist of infinite norms of the metric matrix into the Fuzzy C-Means (FCM) clustering algorithm and its derived algorithms. More specifically, we focus on constructing a dynamic metric matrix that is used to unify the metric between data samples and updating cluster centers to optimize the search path of the cluster center. In addition, we propose a new evaluation index named the Coefficient of Variation Metric (CVM) to evaluate metric effectiveness. The dynamic metric accelerated method, whose complexity remains unchanged, can effectively accelerate the iteration speed of fuzzy clustering. The comparisons between the algorithm using the dynamic metric accelerated method and the corresponding algorithm on UCI, business district and COVID-19 CT image datasets show the superiority of the dynamic metric accelerated method in accelerating effect and clustering performance.

INDEX TERMS Dynamic metric, fuzzy clustering, iteration acceleration, coefficient of variation metric.

I. INTRODUCTION

Generally, clustering depends on a metric that describes the similarity between data samples [1]–[4]. Metrics can express rich information, and metric-based clustering approaches in which the metrics of data samples are the main component of objective functions have been widely used in industrial applications [5]. For the above reasons, research on clustering algorithms focuses on algorithm design and metric research [6].

What's more, Metric plays a more important part in fuzzy clustering. There are many metrics in fuzzy clustering, such as Euclidean distance [7], Minkowski distance [8], Manhattan distance [9], Chebyshev distance [10], Mahalanobis distance [11]–[13], angular cosine [14]–[16], correlation coefficient [17], entropy [18]–[20], and Hamming distance [21]. In addition, in recent years, many studies have been done on metric learning of the kernel method [22]–[26]. The learning

process of the fuzzy clustering algorithm depends on the metric that describes associations between different objects. Fuzzy c-means (FCM) uses Euclidean distance to measure relationship between samples and distance centers [27]; SFCM measures relationship between data samples, distance centers and supervised information by using Euclidean distance [28]; eSFCM measures data samples and cluster centers and entropy to measure supervised information by using Euclidean distance [29]; SMUC measures relationship between samples and distance centers and entropy to measure supervised information by using Mahalanobis distance [30]; In addition, there are some kernel methods [31], [32]. However, the learning process of the relationship between data samples will be affected due to the limitations of metrics. The learning process of the relationship between data samples will be affected due to the limitations of metrics. The clustering algorithms based on Euclidean distance have the advantages of fast iteration convergence and stable results [33]. However, in the face of complex data in practical applications, these clustering algorithms are sensitive to data sample

The associate editor coordinating the review of this manuscript and approving it for publication was Md. Abdur Razzaque¹.

dimensions and cannot apply to nonlinear data [34]. The introduction of the Mahalanobis distance can solve these problems, but it will also lead to the exponential growth of algorithm computation requirements [30].

In the past, research mainly focuses on three kinds of metrics in fuzzy clustering: Entropy [17], [25], [35], [36], Euclidean distance [37]–[40], and Mahalanobis distance [11], [41], [42]. Relative entropy measures the difference of data from the perspective of probability distribution and uncertainty [43]–[45]; it has the advantage of being insensitive to common noises [46], [47]. Euclidean distance is the most popular metric in objective functions of fuzzy clustering because it can reflect the real distance in the sample space [48]. However, there exists a case that Euclidean distance ignores the differences among different features in the sample, which makes it difficult to reflect the associations between data samples [49]. The Mahalanobis distance can better reflect the correlation between data samples, so the process of fuzzy clustering is not affected by the feature dimensions, but the calculation of the inverse matrix of the covariance matrix in Mahalanobis distance greatly increases the complexity of the calculation [50]–[52]. In addition, the three metrics mentioned above don't change with iterations. We call this kind of metrics the static metrics. In our research, the standardization of static metrics can achieve normalization at the feature level, but cannot guarantee the normalization of distance. This paper analyzes this phenomenon and proposes a dynamic metric accelerated method based on dynamic measurement.

Our work is summarized as follows:

- 1) We introduce distance standardization and how it affects the clustering process using static metric algorithms
- 2) We propose a dynamic metric accelerated method for fuzzy clustering that has better adaptability than the traditional static metric method for fuzzy clustering. In addition, we analyze the time complexity and optimality of metric
- 3) For the dynamic metric, we define a Coefficient of Variation Metric (CVM) to evaluate the effectiveness of the metric.
- 4) We use the dynamic metric accelerated method to improve the metric effect in clustering and set experiments on the classic UCI dataset, classic image dataset, business circle dataset [53] and COVID-19-CT dataset [54] to verify the effect of the dynamic metric accelerated method.

II. RELATED WORK

A. FUZZY C-MEANS CLUSTERING (FCM)

FCM is a widely used unsupervised fuzzy clustering algorithm. The Euclidean distance which is used in FCM is widely used in clustering algorithms.

¹https://github.com/ChoiNgai/paper_DynamicMetricClustering/tree/main/data/4.18

Let us first assume that the sample set to be clustered is $X = \{x_1, x_2, \dots, x_n\}$, where $x_j \in R^d$ ($1 \leq j \leq n$) in the d -dimensional Euclidean space, and c is the number of clusters. The objective function of FCM can be expressed as [27]:

$$J(U, V) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m (x_j - v_i)^T I (x_j - v_i) \quad (1)$$

For the convenience of introduction, the objective function is expanded as follows:

$$J(U, V) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m (x_j - v_i)^T A (x_j - v_i) \quad (2)$$

where I is the identity matrix, A is a metric matrix, only when $A = I$, Eq.(2) is equivalent to Eq.(1). and the corresponding metric is Euclidean metric. m is any real number ($m > 1$) which denotes the degree of fuzziness, u_{ij} is the membership degree of the j -th sample x_j belonging to the i -th cluster whose centroid is v_i , $U = (u_{ij})$, $V = [v_1, v_2, \dots, v_c]$, $1 \leq i \leq c$, $1 \leq j \leq n$, $2 \leq c < n$ and u_{ij} satisfies the following constraint condition:

$$\sum_{i=1}^c u_{ij} = 1, \quad u_{ij} \geq 0 \quad (3)$$

By minimizing (1) and using the Lagrange optimization, we obtain the following alternative update equations for the cluster center v_i and the membership degree u_{ij} :

$$v_i = \frac{\sum_{j=1}^n u_{ij}^m x_j}{\sum_{j=1}^n u_{ij}^m} \quad (4)$$

$$u_{ij} = \frac{((x_j - v_i)^T A (x_j - v_i))^{\frac{2}{m-1}}}{\left(\sum_{h=1}^c (x_j - v_h)^T A (x_j - v_h) \right)^{\frac{2}{m-1}}} \quad (5)$$

B. KERNEL-BASED FUZZY C-MEANS CLUSTERING OF METRIC ACCELERATED (KFCM)

KFCM is an unsupervised fuzzy clustering algorithm that projects the original data into the kernel space and takes the inner product of the kernel space as the metric of the algorithm.

Let $V = \{v_1, v_2, \dots, v_c\}$ be the V cluster centers in the kernel space. c is the number of clusters. This minimizes the following objective function subject to conditions as considered in FCM:

$$\begin{aligned} J_{KFCM}(U, V) &= \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m (\varphi(x_j) - \varphi(v_i))^T I (\varphi(x_j) - \varphi(v_i)) \\ &= \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m (2 - 2 \cdot K(x_j, v_i)) \end{aligned} \quad (6)$$

where $K(x_j, v_i)$ is the Gaussian radial basis function, and its form is as follows:

$$K(x_j, v_i) = e^{-\frac{\|x_j - v_i\|^2}{2\sigma^2}} \quad (7)$$

In the iterative process of KFCM, the memberships and the cluster centers are updated as follows:

$$u_{ij} = \frac{(1 - K(x_j, v_i))^{\frac{-1}{m-1}}}{\sum_{k=1}^c (1 - K(x_j, v_k))^{\frac{-1}{m-1}}} \quad (8)$$

$$v_i = \frac{\sum_{j=1}^n \mu_{ij}^m K(x_j, v_i) x_j}{\sum_{j=1}^n \mu_{ij}^m K(x_j, v_i)} \quad (9)$$

C. SEMI-SUPERVISED FUZZY C-MEANS CLUSTERING (SFCM)

SFCM is a semi-supervised fuzzy clustering algorithm that utilizes Euclidean distance to express the relationship of the membership matrix and the prior membership matrix. Its objective function is as follows:

$$J_{SFCM}(U, V) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m (x_j - v_i)^T I(x_j - v_i) + \alpha \sum_{i=1}^c \sum_{j=1}^n (u_{ij} - \tilde{u}_{ij}^2)(x_j - v_i)^T I(x_j - v_i) \quad (10)$$

where $m(m > 1)$ denotes the degree of fuzziness. As m tends toward 1, SFCM approaches HCM. $u_{ij}(0 \leq u_{ij} \leq 1)$ is the membership degree of the j -th sample x_j belonging to the i -th cluster and v_i is its centroid. $V = [v_1, v_2, \dots, v_c]$, $1 \leq i \leq c$, $1 \leq j \leq n$, $2 \leq c \leq n$, $u_j = (u_{1j}, \dots, u_{cj})$, and u_{ij} satisfies the following constraint condition:

$$s.t. u_{ij} \in [0, 1]; \quad \sum_{i=1}^c u_{ij} = 1, \quad (j = 1, 2, \dots, n) \quad (11)$$

m is the weighted index, which is an empirical value, and the value is usually 2. Then, we obtain the iterative formula of the membership degree and clustering center:

$$u_{ij} = \frac{1}{1 + \alpha} \left[\frac{1 + \alpha(1 - \sum_{i=1}^c \tilde{u}_{ij}^{m-1})}{\frac{(x_j - v_i)^T I(x_j - v_i)}{\sum_{h=1}^c (x_j - v_h)^T I(x_j - v_h)}} + \alpha \tilde{u}_{ij} \right] \quad (12)$$

$$v_i = \frac{\sum_{j=1}^n \mu_{ij}^m x_j}{\sum_{j=1}^n \mu_{ij}^m} \quad (13)$$

\tilde{u}_{ij} is an a priori membership matrix, which is transformed from label information; α is a predetermined suppression coefficient.

D. ENTROPY SEMI-SUPERVISED FUZZY C-MEANS CLUSTERING (eSFCM)

eSFCM is a semi-supervised fuzzy clustering algorithm that utilizes information entropy to express the relationship of the membership matrix and the prior membership matrix. Its objective function is as follows:

$$J(U, V) = \sum_{i=1}^c \sum_{j=1}^n \mu_{ij}(x_j - v_i)^T I(x_j - v_i) + \lambda^{-1} \sum_{i=1}^c \sum_{j=1}^n (|\mu_{ij} - f_{ij} b_j| \ln |\mu_{ij} - f_{ij} b_j|) \quad (14)$$

$$v_i = \frac{\sum_{j=1}^n \mu_{ij} x_j}{\sum_{j=1}^n \mu_{ij}} \quad (15)$$

$$u_{ij} = \tilde{u}_{ij} + \frac{e^{-\lambda(x_j - v_i)^T I(x_j - v_i)}}{\sum_{h=1}^c e^{-\lambda(x_j - v_h)^T I(x_j - v_h)}} (1 - \sum_{h=1}^c \tilde{u}_{ij}) \quad (16)$$

Among them, λ is an empirical parameter.

E. SEMI-SUPERVISED METRIC-BASED FUZZY CLUSTERING (SMUC)

SMUC is a semi-supervised fuzzy clustering algorithm using Mahalanobis distance and entropy. By introducing the Mahalanobis distance on the basis of eSFCM [29], semi-supervised metric-based fuzzy clustering (SMUC) algorithm [30] was proposed, and its objective function is as follows:

$$J_{SMUC}(U, V) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}(x_j - v_i)^T M(x_j - v_i) + \lambda \sum_{i=1}^c \sum_{j=1}^n (u_{ij} - \tilde{u}_{ij}) \ln(u_{ij} - \tilde{u}_{ij}) \quad (17)$$

where is subject to conditions (2) and (5) and the following conditions:

$$M = \left[\frac{1}{n} \sum_{i=1}^c \sum_{j=1}^n \tilde{u}_{ij} (x_j - \frac{\sum_{j=1}^n \tilde{u}_{ij}^2 x_j}{\sum_{j=1}^n \tilde{u}_{ij}^2}) (x_j - \frac{\sum_{j=1}^n \tilde{u}_{ij}^2 x_j}{\sum_{j=1}^n \tilde{u}_{ij}^2})^T \right]^{-1} \quad (18)$$

There is the following optimal solution [30]:

$$v_i = \frac{\sum_{j=1}^n u_{ij} x_j}{\sum_{j=1}^n u_{ij}} \quad (19)$$

$$u_{ij} = \tilde{u}_{ij} + \frac{e^{-\lambda(x_j - v_i)^T M(x_j - v_i)}}{\sum_{h=1}^c e^{-\lambda(x_j - v_h)^T M(x_j - v_h)}} (1 - \sum_{h=1}^c \tilde{u}_{ij}) \quad (20)$$

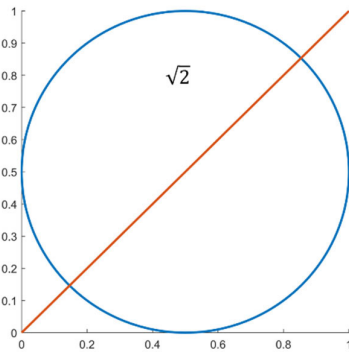


FIGURE 1. The numerical range is beyond [0,1] in the iterative process.

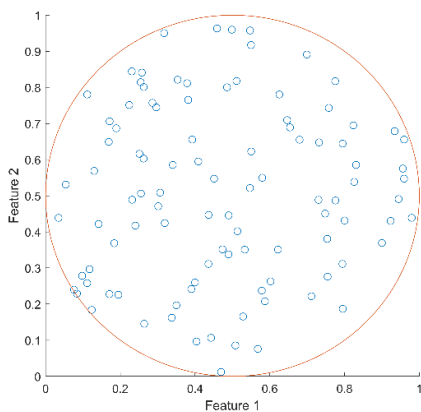


FIGURE 2. Synthetic dataset¹ in which the distance between samples is less than 1.

III. ACCELERATED ITERATIVE METHOD FOR FUZZY CLUSTERING

A. MOTIVATION

At present, clustering algorithms usually need to standardize sample features to eliminate the influence of different feature dimensions before conducting iterative calculations. After data normalization, the value of each sample feature is limited to [0,1] for iterative calculations, but the value of the sample distance will exceed [0,1] after calculations, as shown in Fig 1.

It can be seen from Fig 1 that the standardized feature does not mean that distance (i.e., metric) between samples is also standardized.

To verify the necessity of standardizing metrics, we built artificial datasets. The distance between samples in dataset (a) was less than 1, and some of the distances in the dataset (b) were more than 1 (as shown in Fig 2 and Fig 3).

Fig 2 and Fig 3 represent the simulation dataset of dataset 1 and dataset 2, respectively, whose feature number is 2 (two-dimensional space) and sample number is 100.

Let us assume that the number of cluster centers is 2, the initial cluster center is [0.3, 0.5; 7, 0.5], the maximum number of iterations of clustering in dataset (a) is 50, and the solution path-length of the two clustering centers is

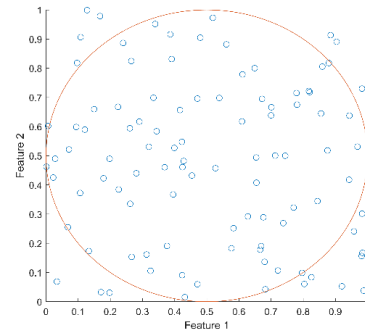
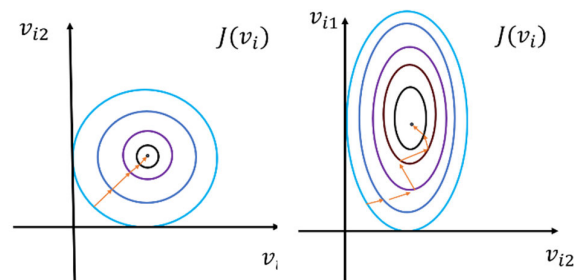


FIGURE 3. Synthetic dataset in which the distance between samples exceeds 1.



(a)standardized metric (b)non-standardized metric

FIGURE 4. Two possible cases caused by traditional methods (longer search distance).

0.0029. The maximum number of iterations of clustering in dataset (b) is 52, and the solution path length is 0.0097.

If the sample distribution exhibits the phenomenon in Fig 3, then after the usual iterative process, the search path may become longer and indirect, as shown in Fig 4b:

We are committed to making the solving path closer to Fig 4a instead of Fig 4b.

B. METRIC COMPLEXITY

The Mahalanobis distance, which has a good effect on eliminating the dimensional effect, is commonly used to metric the continuous type of feature. At present, the Mahalanobis distance has a good influence on data of different dimensions, so the Mahalanobis distance is widely used for metric data of different dimensions. However, it is difficult for the Mahalanobis distance to apply to image recognition and other fields because it needs to calculate the covariance matrix and its inverse [30], resulting in a large amount of calculations (as shown in Fig 5).

In Fig 5, the number of randomly generated samples is fixed at 1000, and the number of features is from 1 to 1000. The time consumption of calculating the sample and distance and the sample data and the randomly generated $3 \times d$ size matrix is gradually calculated, where the number of features is d . With the increase of features number, the amount of computation increases exponentially.

C. DYNAMIC STANDARDIZATION OF EUCLIDEAN METRIC

In the iterative process of the algorithm, there exists a case in which the distance between the sample and the cluster center

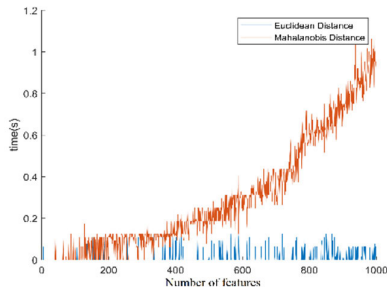


FIGURE 5. Comparison of computational complexity (time-consuming) between the Mahalanobis metric and Euclidean metric.

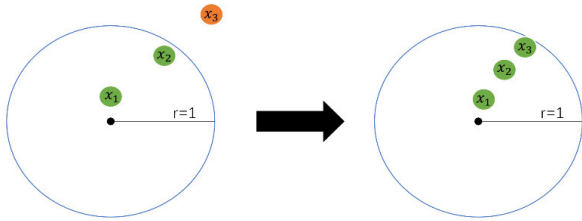


FIGURE 6. Distance standardization process (compress).

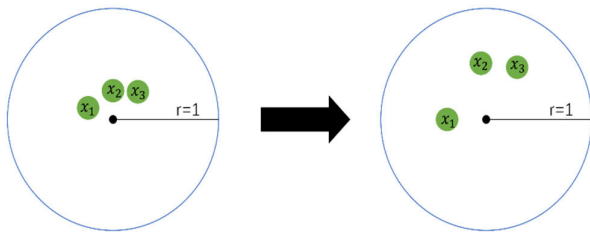


FIGURE 7. Distance standardization process (expand).

is greater than 1, so it is necessary to compress the distances so that all the distances after the transformation are in $[0,1]$ (as shown in Fig 6).

When the distance between the sample and cluster center is small, the rounding error produced in the calculation process will have a great impact on the clustering process. Therefore, it is necessary to expand the distance by transformation. The effect of stretching is shown in Fig 7:

In Fig 6 and Fig 7, the left side is a non-standardized distance, and the right side is a standardized distance. To ensure that the mathematical relationship between distances is unchanged and that the algorithm is fast iterative, we satisfy the distance with the following relationship between the standardized distance and the non-standardized distance:

$$\frac{d_{i,1}}{d_{i,2}} = \frac{d_{i,1}^*}{d_{i,2}^*}, \frac{d_{i,2}}{d_{i,3}} = \frac{d_{i,2}^*}{d_{i,3}^*}, \dots, \frac{d_{i,j-1}}{d_{i,j}} = \frac{d_{i,j-1}^*}{d_{i,j}^*} \quad (21)$$

where $d_{i,j}$ is the non-standardized distance between x_j and v_i ; $d_{i,j}^*$ is the standardized distance between x_j and v_i .

D. DYNAMIC METRIC ACCELERATION

The Euclidean distance between the sample and the cluster center expression in clustering expands from $\|x_j - v_i\|^2$ to $(x_j - v_i)^T A (x_j - v_i)$.

- 1) When $A = I$, where I is the identity matrix, the corresponding metric is Euclidean distance.
- 2) When $A = M^{-1}$, where M is from (18), the corresponding metric is the Mahalanobis distance.
- 3) When A is a diagonal matrix whose elements consist of the positive infinite norm of the distance matrix, the corresponding metric is a dynamic metric whose matrix changes with each iteration. We call the method using dynamic metrics for fuzzy clustering as the dynamic metric accelerated method for fuzzy clustering.

Calculation steps of the metric matrix in the dynamic metric accelerated method:

- 1) Compute the Euclidean distance between the sample and the cluster center to get the Euclidean distance matrix;
- 2) Calculate the infinite norm according to the distance matrix;
- 3) Calculate the dynamic measure matrix A based on the infinite norm.

The Euclidean distance between x_j and v_i is denoted as $D = \{d_{ij}\}$; the infinite norm of the matrix is:

$$\|D\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |d_{ij}| \quad (22)$$

Then, the metric matrix A composed of the infinite norm of the metric matrix is as follows:

$$C = \begin{bmatrix} \|D\|_\infty & 0 & \dots & 0 \\ 0 & \|D\|_\infty & \dots & 0 \\ 0 & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \|D\|_\infty \end{bmatrix}_{d \times d} \quad (23)$$

where $\|D\|_\infty$ is the infinite norm of the metric matrix:

$$A = C^{-1} \quad (24)$$

The metric matrix $A = \text{diag}(m_1, \dots, m_d)$, and the expression of element A is:

$$m_k = \frac{1}{\|D\|_\infty} \quad (25)$$

where d is the number of features, $k = \{1, 2, \dots, d\}$.

Therefore, the distance calculation expression of the dynamic metric is:

$$(x_j - v_i)^T A (x_j - v_i) \quad (26)$$

IV. ANALYSIS OF WEIGHTING EXPONENT M IN FCM

The dynamic metric will inevitably affect the membership matrix when the iterative process changes, and the weighting exponent M of the FCM model plays a role in determining the validity of FCM partitions. Therefore, we explore the impact of the weighting exponent M in FCM when using the dynamic metric acceleration method. In this part, the effect of the value of M is analyzed through experiments using the Wisconsin breast cancer dataset.² The internal evaluation index data

²<http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>

reconstruction error rate [56], the external evaluation index purity [57] and the number of iterations are used to analyze the influence of the m value on the dynamic metric and conventional Euclidean metric in FCM.

A. EXTERNAL EVALUATION INDEX

$$Purity(\Omega, C) = \frac{1}{n} \sum_c \max_j |\omega_c \cap y_j| \quad (27)$$

where n is the number of samples, $\Omega = \{\omega_1, \dots, \omega_c\}$ is the clustering result, and $Y = \{y_1, \dots, y_j\}$ is the real classification. Purity reflects the similarity between clustering results and reality. We propose an internal evaluation index $Purity(\Omega, C)/iter$ to evaluate the improvement level of the clustering accuracy of the algorithm:

$$Purity(\Omega, C)/iter = \frac{1}{n} \sum_k \max_j |\omega_c \cap y_j|/t \quad (28)$$

where $iter$ represents the number of iterations of the algorithm

B. INTERNAL EVALUATION INDEX

V_{RE} is the error rate of data reconstruction, which is defined as follows:

$$V_{RE} = \frac{1}{n} \sum_{i=1}^n \|I'_g(i) - I_g(i)\|^2 \quad (29)$$

It analyzes the difference between reconstructed data and the original data after clustering.

$I'_g(i)$ is the gray level of the i -th sample of the reconstructed data:

$$I'_g(i) = \frac{\sum_{k=1}^c u_{ki}^2 I_g(i)}{\sum_{k=1}^c u_{ki}^2} \quad (30)$$

Table 1 shows the influence of m value on the iteration speed of the algorithm. In terms of cluster purity, dynamic metric does not affect the role of the m value in FCM, as shown in Table 2. Table 3 shows that with the increase in the m value, the clustering error rate of the FCM algorithm increases. At the same time, combined with Nikhil R. pal's conclusion [55], we can conclude that $M = 2$ is a better choice in the dynamic acceleration matrix.

C. EVALUATING INDICATOR AND VERIFICATION

To measure the effectiveness of the standardization of metrics, we analyzed the maximum and minimum distances between data sample and cluster centers.

FCM experiments were carried out on the Wisconsin breast cancer dataset³ with 2 clusters.

The maximum and minimum distance curves are shown in Fig 8. In order to better describe this process, we define some metric related concepts

³<http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>

TABLE 1. Influence of m value on iterations.

M value	Euclidean metric	Dynamic metric
m=2	15	11
m=3	14	12
m=4	18	11
m=5	19	11
m=6	21	13
m=7	19	10

TABLE 2. Influence of m value on purity.

M value	Traditional metric	Dynamic metric
m=2	0.958	0.958
m=3	0.955	0.955
m=4	0.955	0.955
m=5	0.955	0.955
m=6	0.955	0.955
m=7	0.958	0.958

³ <http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>

TABLE 3. Influence of m value on V_{RE} .

M value	Traditional metric	Dynamic metric
m=2	5.07×10^{-33}	4.76×10^{-33}
m=3	6.04×10^{-33}	5.68×10^{-33}
m=4	6.37×10^{-33}	5.98×10^{-33}
m=5	5.07×10^{-33}	6.33×10^{-33}
m=6	7.06×10^{-33}	6.82×10^{-33}
m=7	7.19×10^{-33}	6.78×10^{-33}

Definition 1 (Iterative Max/Min Distance): The maximum/minimum distance from the sample's t -th iteration to v_i is:

$$d_i^{(t)} = \|x_j - v_i^{(t)}\|^2 \quad (31)$$

where T is the maximum number of iterations, $1 < t < T$, $1 < i < c$.

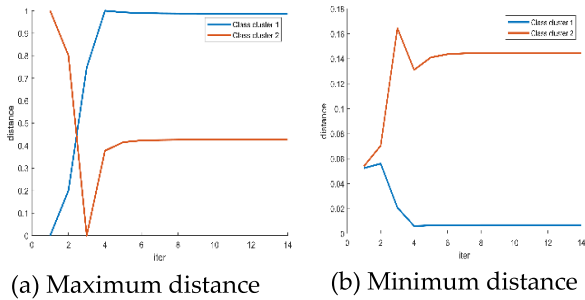


FIGURE 8. The relationship between the number of iterations and the distance of the samples and each cluster center.

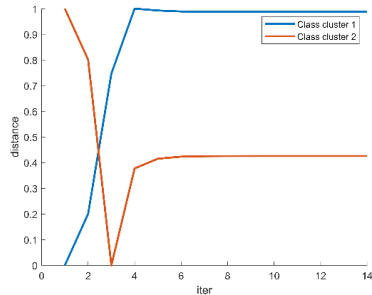


FIGURE 9. Difference in the distance extreme deviations between different clusters.

Definition 2 (The Extreme Deviations of Distance): The extreme deviations of the distance from the sample to $v_i^{(t)}$ are calculated as follows:

$$\begin{aligned} \Gamma_i^{(t)} &= \max_{1 < t < T} d_i^{(t)} - \min_{1 < t < T} d_i^{(t)} \\ &= \max_{1 < t < T} \|x_j - v_i^{(t)}\|^2 - \min_{1 < t < T} \|x_j - v_i^{(t)}\|^2 \end{aligned} \quad (32)$$

However, due to $\text{Max_}d_i^{(t)} \gg \text{Min_}d_i^{(t)}$, the range of distance is mainly controlled by $\text{Max_}d_i^{(t)}$. To avoid this case and express the intensity of the transformation between distances better, we normalize it to obtain the improved range. To express the normalization of the distance range conveniently, we define the following function:

$$\Phi(d_i^{(t)}) = \frac{d_i^{(t)} - \min |d_i^{(t)}|}{\max |d_i^{(t)}| - \min |d_i^{(t)}|} \quad (33)$$

The improved extreme deviations of the distance are:

$$R_i^{(t)} = \Phi(\text{Max_}d_i^{(t)}) - \Phi(\text{Min_}d_i^{(t)}) \quad (34)$$

The improved distance extreme deviations diagram is drawn, as shown in Fig 9:

The difference between the maximum and the minimum of the distance between the sample and each cluster center is increasing. The change process of the distance is shown in Fig 9.

We call the above distance as the **Coefficient of Variation Metric (CVM)**. The CVM is used to evaluate this case.

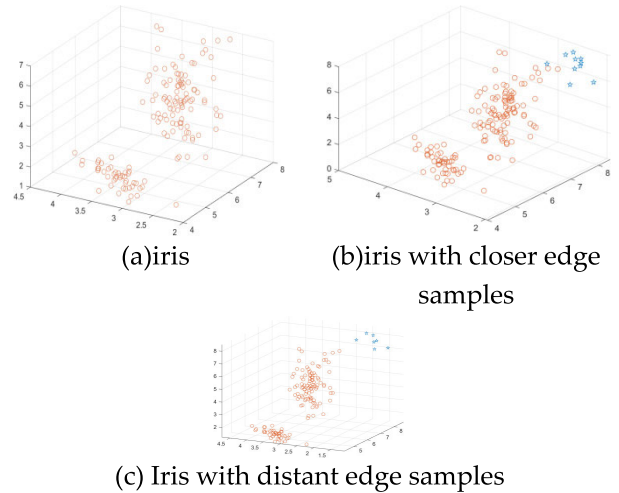


FIGURE 10. Add data points to simulate the change of edge points.

Before defining *CVM*, we need to explain the **Variation Metric (VM)**.

Definition 3 (Variation Metric): As appendix Fig 1 shows, we define the D-value as $VM^{(t)}$, and its mathematical expression is:

$$VM^{(t)} = \max |R^{(t)}| - \min |R^{(t)}| \quad (35)$$

where t is the number of iterations. The smaller $VM^{(t)}$ is, the smoother the iteration process and the better the normalization effect.

Definition 4 (Coefficient of Variation Metric):

$$CVM = \frac{\sum_{t=1}^{T-1} \left[\left(\frac{VM^{(t)}}{VM^{(t)} + VM^{(t+1)}} \right) \cdot \frac{t}{T} \right]}{\frac{1}{T} \cdot \sum_{i=1}^{T-1} i} \quad (36)$$

The coefficient of variation metric satisfies the following properties:

If the mean value of *VM* is 0.5 and the standard deviation is 0, then $CVM = 0.5$.

If *VM* and $VM^{(0)} > 0.5$ is monotonically increasing, then $0 < CVM < 0.5$.

If *VM* and $VM^{(0)} < 0.5$ is monotonically decreasing, then $0.5 < CVM < 1$.

In addition, the value of *CVM* is independent of the initial value

One example of dynamic metric acceleration is the iteration process, and the variation of metric is shown in appendix Fig 2.

A larger *CVM* represents a more stable iteration for the metric.

The *CVM* is calculated and shown in Fig 9 as: $CVM = 0.5013 + 0.0010$.

After acceleration by the dynamic metric in Fig 11, $CVM = 0.8132 \pm 0.0110$.

From the comparison of appendix Fig 1 and Fig 2, we can see that the measured acceleration method effectively

TABLE 4. (FCM, FCM-M) algorithm performance.

Datasets	iterations	Purity	V _{RE}	Max-distance
Iris	(30,19)	(0.89,0.89)	$(1.07 \times 10^{-32}, 1.07 \times 10^{-32})$	(3.08,0.03)
Wine	(13,11)	(0.96,0.96)	$(2.83 \times 10^{-32}, 2.60 \times 10^{-32})$	(3.79,0.04)
Seed	(16,9)	(0.90,0.90)	$(1.46 \times 10^{-32}, 1.45 \times 10^{-32})$	(3.74,0.04)
Breast	(14,10)	(0.96,0.96)	$(4.93 \times 10^{-32}, 4.76 \times 10^{-32})$	(6.56,0.03)
Thyroid	(23,16)	(0.93,0.93)	$(5.57 \times 10^{-32}, 5.56 \times 10^{-32})$	(6.16,0.003)

TABLE 5. Time complexity analysis of different metrics.

Euclidean	Mahalanobis	Dynamic
$O(n \cdot d \cdot c)$	$O(d^4 \cdot n \cdot c + n^2 \cdot d)$	$O(n \cdot d \cdot c)$

improves the optimization of the iteration. In order to explore the optimization effect of using dynamic indexes, we use FCM algorithm to experiment and try to verify its optimality. The experiment reflects the optimal dynamic metric by using the change of internal and external evaluation indices before and after dynamic metric acceleration in multiple datasets. In parentheses, the left side is the algorithm without dynamic metric, and the right side is the algorithm with dynamic metric. We let $\alpha = 5$, $e = 10^{-5}$, the maximum number of iterations is 100, $\lambda = 10$, and randomly initialize cluster centers in experiments. Five labeled samples were used as supervised information in all experiments.

We compare the traditional fuzzy clustering algorithm with the dynamic metric accelerated fuzzy clustering algorithm, FCM-M, for five datasets with the same m value.

The results are shown in Table 4. The results are the mean of 10 repeated experiments. The dataset used in the experiment is shown on first column.

V. OPTIMIZATION OF DYNAMIC METRICS

The experimental results show that the dynamic metric not only maintains the FCM clustering effect, but also achieves a better iterative speed for most data sets and ensures the standardization of the metric. Combined with the above results, we can conclude that the dynamic accelerated method can achieve iterative acceleration without affecting the optimal solution, and can stably achieve the purpose of normalized distance.

VI. DISCUSSION ON TIME COMPLEXITY

In the previous section, we present a dynamic accelerated method based on infinite norms. This section discusses the time complexity of various metric calculations.

Here, n is the number of samples, d is the number of features, and c is the number of cluster centers.

The time complexity of calculating Euclidean metric is $O(n \cdot d \cdot c)$. The time complexity of computing the covariance matrix by the Mahalanobis metric is $O(n^2 \cdot d)$, and the time complexity of finding all elements in the Mahalanobis metric matrix is $O(d^4 \cdot n \cdot c)$. Thus, the time complexity of the Mahalanobis metric is $O(d^4 \cdot n \cdot c + n^2 \cdot d)$.

The time complexity of the three kinds of metrics are shown in Table 5.

VII. INFLUENCE ANALYSIS OF EDGE POINTS

The Euclidean distance of the sample points at the edge of the cluster is farther from other points, so change of the edge points will affect the infinite norm in the metric matrix, and our metric accelerated method depends on the norm in the dynamic metric. Therefore, this section discusses the impact of edge points in clusters on clustering speed and accuracy.

In this part, the sample pairs far away from the main points are analyzed by using iris experiment. Readers can obtain the preprocessed image datasets here.⁴

Fig 10a is one of the projection of iris dataset in a three-dimensional space. Fig 10b and Fig 10c add ten edge samples which far away from the main point in one cluster of the iris dataset, and the edge sample points added in Fig 10c are farther away from the main point.

Table 6 shows the influence of edge samples on the acceleration algorithm. In terms of the maximum distance, the change of edge samples changes the maximum distance, while in terms of clustering purity and iteration speed, the change of edge samples do not affect the accuracy and speed of the algorithm.

VIII. EXPERIMENT

This section investigates the acceleration effect and performance of the dynamic-metric acceleration method applied to different algorithms in real datasets to verify its effectiveness.

⁴[https://github.com/ChoiNgai/paper_DynamicMetricClustering/tree/main/data/edge points/](https://github.com/ChoiNgai/paper_DynamicMetricClustering/tree/main/data/edge%20points/)

TABLE 6. FCM-M algorithm performance.

Datasets	iterations	Purity	Max- distance
Iris	4	0.0022	4.0588
iris with closer edge samples	4	0.0022	4.0867
Iris with distant edge samples	4	0.0022	4.1232

TABLE 7. The details of the datasets.

Dataset	samples	features	Cluste:
Iris	150	4	3
Wine	178	13	3
Seed	569	30	2
Thyroid	7200	21	3
UKM	257	5	4
MNIST (test)	10000	784	10
Business circle	431	5	3
COVID-19-CT	544	36	2

A. EXPERIMENTAL SETUP

The experiment was conducted on a series of real datasets. These datasets include data with clusters that are not linearly separable, with noise, with numerous features and with partial redundancy [55]. In addition, image data classification [56] is also included. The statistical data of all datasets are summarized in Table 7. The image datasets have been preprocessed in the experiment, which can obtain the preprocessed image dataset here.⁵

To evaluate the effectiveness of the methods presented in this paper, the dynamic metric method is used to improve the fuzzy clustering algorithm based on Euclidean distance, Mahalanobis distance and entropy. In the experiment, we analyze the effect of the dynamic metric method by comparing the performance before and after the improvement.

B. ALGORITHM IMPROVEMENT

We choose the fuzzy clustering algorithm based on the Euclidean distance, Mahalanobis distance and entropy to

⁵https://github.com/ChoiNgai/paper_DynamicMetricClustering/tree/main/data/

verify the acceleration effect of the dynamic metric method and we choose FCM as the unsupervised fuzzy clustering method based on the Euclidean distance, SFCM as the semi-supervised fuzzy clustering method based on the Euclidean distance, eSFCM as the unsupervised fuzzy clustering method based on the Mahalanobis distance, and SMUC as the semi-supervised clustering algorithm based on the Mahalanobis distance and entropy. the FCM, KFCM, SFCM, eSFCM and SMUC algorithms to replace the static metric of the original algorithm. The algorithms that introduced dynamic metrics are called FCM-M, KFCM-M, SFCM-M, eSFCM-M and SMUC-M.

The calculation flow of the improved algorithm is as follows:

1) FUZZY C-MEANS CLUSTERING OF METRIC ACCELERATED (FCM-M)

The membership matrix and the iterative formula of the clustering center of FCM are as follows (three iterative formulas):

$$v_i = \frac{\sum_{j=1}^n (u_{ij})^m x_j}{\sum_{j=1}^n (u_{ij})^m} \tag{37}$$

$$A = \text{diag}(m_1, \dots, m_d) \tag{38}$$

$$u_{ij} = \frac{[(x_j - v_i)^T A (x_j - v_i)]^{\frac{2}{m-1}}}{\sum_{h=1}^c [(x_j - v_h)^T A (x_j - v_h)]^{\frac{2}{m-1}}} \tag{39}$$

Algorithm 1. FCM-M

Input: n: the number of samples; d: samples dimensions; c: cluster number; ϵ : a small enough error; X: dataset; T: maximum iteration times. V: random initialized clustering centroid. \tilde{U} : prior membership degree. α : trade-off parameter

Output: V: clustering centroid matrix; U: membership degree matrix

Initialize metric matrix A by (38);
Initialize membership degree matrix $U = [u_{ij}]$ by (39).

Repeat

Update the cluster prototype matrix $V = [v_i]$ by (37);

Update the Metric matrix A by (38);

Update the membership degree U matrix by (38);

Calculate the objective function value $J(t)$ by (2);

Iteration time $t + +$;

Until

$$J(t + 1) - J(t) < \epsilon \text{ or } t = T$$

2) KERNEL-BASED FUZZY C-MEANS CLUSTERING OF METRIC-ACCELERATED (KFCM-M)

The membership matrix and the iterative formula of the clustering center of KFCM are as follows (three

iterative formulas):

$$v_i = \frac{\sum_{j=1}^n (u_{ij})^m \kappa(x_j, v_i) x_j}{\sum_{j=1}^n (u_{ij})^m \kappa(x_j, v_i)} \quad (40)$$

$$A = \text{diag}(m_1, \dots, m_d) \quad (41)$$

$$u_{ij} = \frac{(1 - e^{-\frac{|x_j - v_i|^T A(x_j - v_i)}{-2\sigma^2}})^{\frac{-1}{m-1}}}{\sum_{k=1}^c (1 - e^{-\frac{(x_j - v_k)^T A(x_j - v_k)}{-2\sigma^2}})^{\frac{-1}{m-1}}} \quad (42)$$

Algorithm 2. KFCM-M

Input: n : the number of samples; d : samples dimensions; c : cluster number; ε : a small enough error; X : dataset; T : maximum iteration times. V : random initialized clustering centroid. \tilde{U} : prior membership degree. α : trade-off parameter

Output: V : clustering centroid matrix; U : Initialize Metric matrix by (41); Initialize membership degree matrix by (42).

Repeat

Update the cluster prototype matrix by (40);
Update the Metric matrix by (41);
Update the membership degree matrix by (42);
Calculate the objective function value by (6);
Iteration time;

Until

$$J(t+1) - J(t) < \varepsilon \text{ or } t = T$$

3) SEMI-SUPERVISED FUZZY C-MEANS CLUSTERING OF METRIC-ACCELERATED (SFCM-M)

The membership matrix and the iterative formula of the clustering center of SFCM are as follows (three iterative formulas):

$$v_i = \frac{\sum_{j=1}^n (u_{ij})^m x_j}{\sum_{j=1}^n (u_{ij})^m} \quad (43)$$

$$A = \text{diag}(m_1, \dots, m_d) \quad (44)$$

$$u_{ij} = \frac{1}{1 + \alpha} \left[\frac{1 + \alpha (1 - \sum_{i=1}^c \tilde{u}_{ij})^{m-1}}{\frac{[(x_j - v_i)^T A(x_j - v_i)]^2}{[\sum_{h=1}^c (x_j - v_h)^T A(x_j - v_h)]^2}} + \alpha \tilde{u}_{ij} \right] \quad (45)$$

4) ENTROPY SEMI-SUPERVISED FUZZY C-MEANS CLUSTERING OF METRIC-ACCELERATED (eSFCM-M)

The membership matrix and the iterative formula of the clustering center of KFCM are as follows (three

Algorithm 3. SFCM-M

Input: n : the number of samples; d : samples dimensions; c : cluster number; ε : a small enough error; X : dataset; T : maximum iteration times. V : random initialized clustering centroid. \tilde{U} : prior membership degree. α : trade-off parameter

Output: V : clustering centroid matrix; U : membership degree matrix

Initialize Metric matrix A by (44);

Initialize membership degree matrix $U = [u_{ij}]$ by (45).

Repeat

Update the cluster prototype matrix $V = [v_i]$ by (43);
Update the Metric matrix A by (44);
Update the membership degree U matrix by (45);
Calculate the objective function value $J(t)$ by (10);

Iteration time $t + +$;

Until

$$J(t+1) - J(t) < \varepsilon \text{ or } t = T$$

iterative formulas):

$$v_i = \frac{\sum_{j=1}^n (u_{ij})^m x_j}{\sum_{j=1}^n (u_{ij})^m} \quad (46)$$

$$A = \text{diag}(m_1, \dots, m_d) \quad (47)$$

$$u_{ij} = \tilde{u}_{ij} + \frac{e^{-\lambda(x_j - v_i)^T A(x_j - v_i)}}{\sum_{h=1}^c e^{-\lambda(x_j - v_h)^T A(x_j - v_h)}} (1 - \sum_{i=1}^c \tilde{u}_{ij}) \quad (48)$$

Algorithm 4. eSFCM-M

Input: n : the number of samples; d : samples dimensions; c : cluster number; ε : a small enough error; X : dataset; T : maximum iteration times. V : random initialized clustering centroid. \tilde{U} : prior membership degree. α : trade-off parameter

Output: V : clustering centroid matrix; U : membership degree matrix

Initialize metric matrix A by (47);

Initialize membership degree matrix $U = [u_{ij}]$ by (48).

Repeat

Update the cluster prototype matrix $V = [v_i]$ by (46);
Update the Metric matrix A by (47);
Iteration time $t + +$;

Until

$$J(t+1) - J(t) < \varepsilon \text{ or } t = T$$

5) SEMI-SUPERVISED METRIC-BASED FUZZY CLUSTERING OF METRIC-ACCELERATED (SMUC-M)

The membership matrix and the iterative formula of the clustering center of SMUC-M are as follows (three

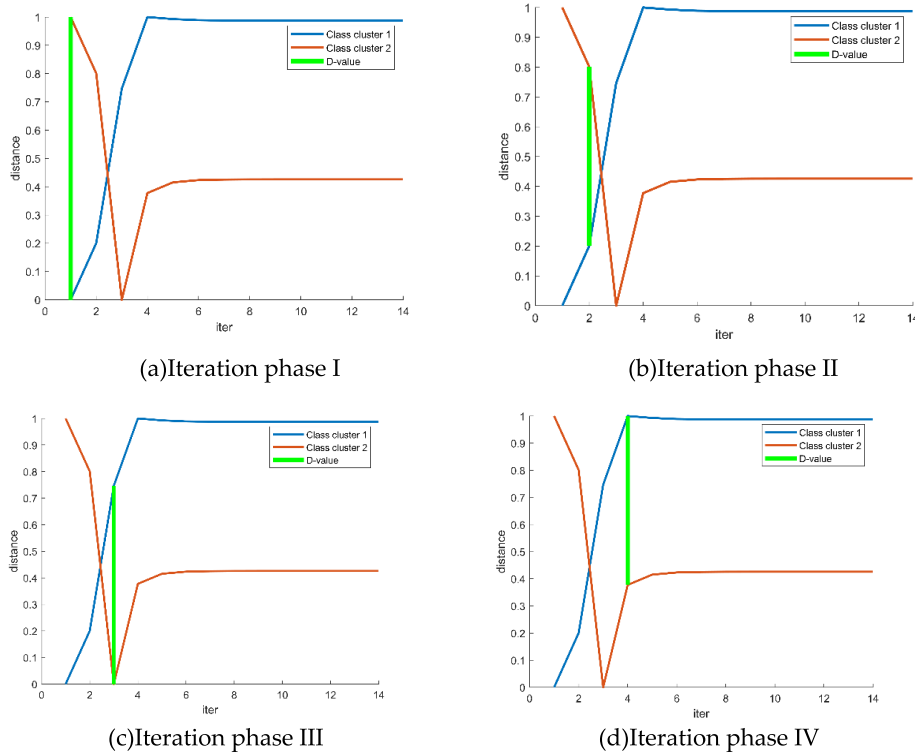


FIGURE 11. Schematic diagram of the change process of sample distance difference.

iterative formulas):

$$v_i = \frac{\sum_{j=1}^n (u_{ij})^m x_j}{\sum_{j=1}^n (u_{ij})^m} \quad (49)$$

$$A = \text{diag}(m_1, \dots, m_d) \quad (50)$$

$$u_{ij} = \tilde{u}_{ij} + \frac{e^{-\lambda(x_j - v_i)^T M A (x_j - v_i)}}{\sum_{h=1}^c e^{-\lambda(x_j - v_h)^T M A (x_j - v_h)}} \left(1 - \sum_{i=1}^c \tilde{u}_{ij}\right) \quad (51)$$

C. EVALUATION INDEX

Experiments compare the clustering accuracy and running time before and after the improvement of the algorithm. The internal evaluation index data reconstruction error rate, the external evaluation index purity, and the number of iterations are used to evaluate the optimization effect of the iterative process.

In the following experimental appendix Fig 3, Fig 4 and Fig 5, the dynamic metric acceleration method is represented by the L_∞ metric speedup. Moreover, we make $\alpha = 5$, $e = 10^{-5}$, the maximum number of iterations is 100, $\lambda = 10$, and randomly initialize cluster centers in experiments.

D. EXTERNAL EVALUATION INDEX

Purity reflects the similarity between clustering results and reality. We propose an internal evaluation index

Algorithm 5. SMUC-M

Input: n : the number of samples; d : samples dimensions; c : cluster number; ε : a small enough error; X : dataset; T : maximum iteration times. $V(0)$: random initialized clustering centroid. \tilde{U} : prior membership degree. α : trade-off parameter

Output: V : clustering centroid matrix; U : membership degree matrix

Initialize Metric matrix A by (50);

Initialize membership degree matrix $U = [u_{ij}]$ by (51).

Repeat

Update the cluster prototype matrix $V = [v_i]$ by (49);

Update the Metric matrix A by (50);

Update the membership degree U matrix by (51);

Calculate the objective function value $J(t)$ by (17);

Iteration time $t + +$;

Until

$$J(t + 1) - J(t) < \varepsilon \text{ or } t = T$$

$Purity(\Omega, C)/iter$ to evaluate the improvement level of the clustering accuracy of the algorithm:

$$Purity(\Omega, C)/iter = \frac{1}{n} \sum_k \max_j |\omega_c \cap y_j|/t \quad (52)$$

where $iter$ represents the number of iterations of the algorithm

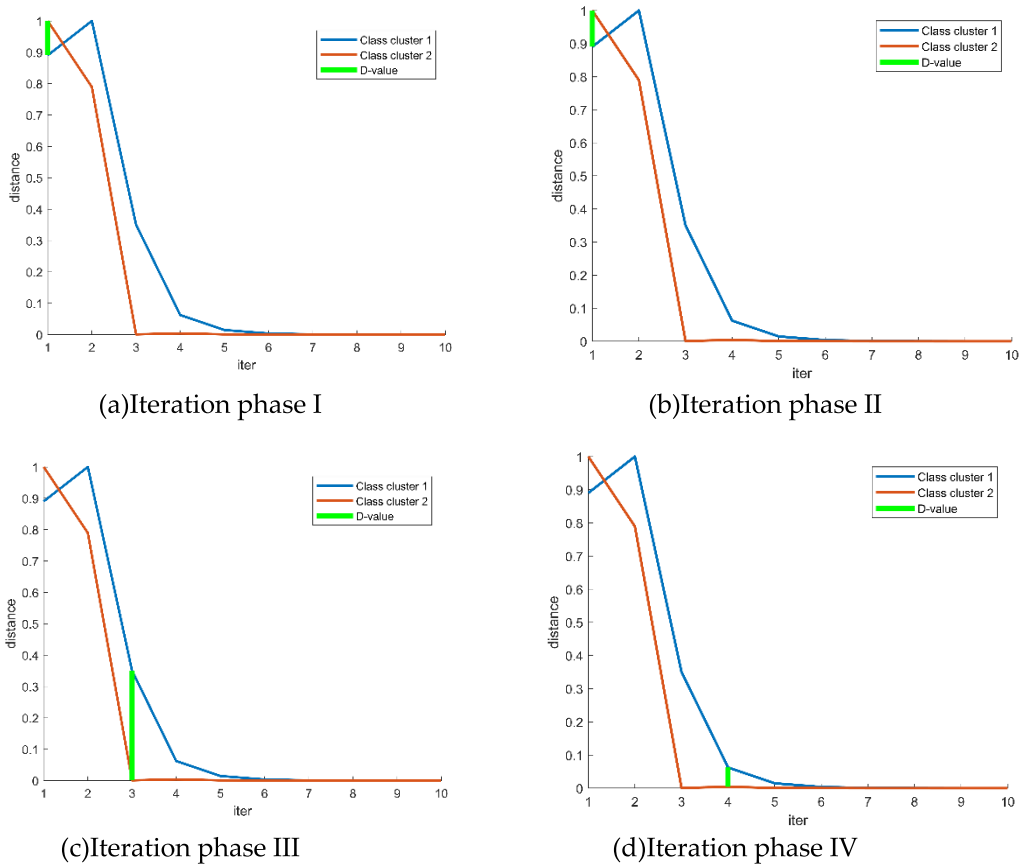


FIGURE 12. Schematic diagram of the change process of the sample distance difference.

E. INTERNAL EVALUATION INDEX

V_{RE} analyzes the difference between reconstructed data and the original data after clustering. It reconstructs data by using the membership matrix and clustering center obtained by clustering. The reconstructed data should be as similar to the original data as possible. A small V_{RE} value means that the algorithm created a well-clustered effect.

Fewer iterations mean a faster search speed. The experiments on real data sets show that the variance of V_{RE} and iteration times is small and stable.

By combining the characteristics of these two evaluation indicators, we propose an internal evaluation index V_{RE*} iter to evaluate the improvement effect of the algorithm on iterative performance:

$$V_{RE} \cdot iter = \frac{1}{n} \sum_{i=1}^n \|I'(i) - I(i)\|^2 \cdot t \quad (53)$$

where t is the number of iterations of the algorithm. Iter represents the number of current iterations

F. EXPERIMENTAL RESULTS ANALYSIS

From the experimental results in appendix Fig 3, except for the KFCM algorithm, the algorithms have better effects after using the metric accelerated iteration method, which

shows that this method is not suitable for the fuzzy clustering algorithm of the kernel function.

From the experimental results in appendix Fig 4, the metric acceleration method still has a poor effect on the KFCM algorithm. Some datasets will not be improved on eSFCM and SMUC, and the effect is the same as that of the unimproved method, indicating that this method still has room for improvement in the nonlinear metrics.

From the experimental results in appendix Fig 5, the metric acceleration method performs well in algorithms other than KFCM. From the comparison of SMUC algorithm with SFCM and eSFCM, it can be seen that our acceleration method is weaker than the algorithm based on Mahalanobis distance in the optimization effect of the algorithm based on Euclidean distance. The experiments of FCM and SFCM show that the optimization effect of unsupervised algorithm is lower than that of semi supervised algorithm.

From the above experimental results, it can be seen that the metric acceleration method in this paper achieves high-quality results in the Euclidean distance and Mahalanobis distance metrics. However, for some nonlinear metric methods, the effect is not improved (such as KFCM and metrics with the kernel method), which indicates that this method is not suitable for nonlinear metrics, such as the kernel

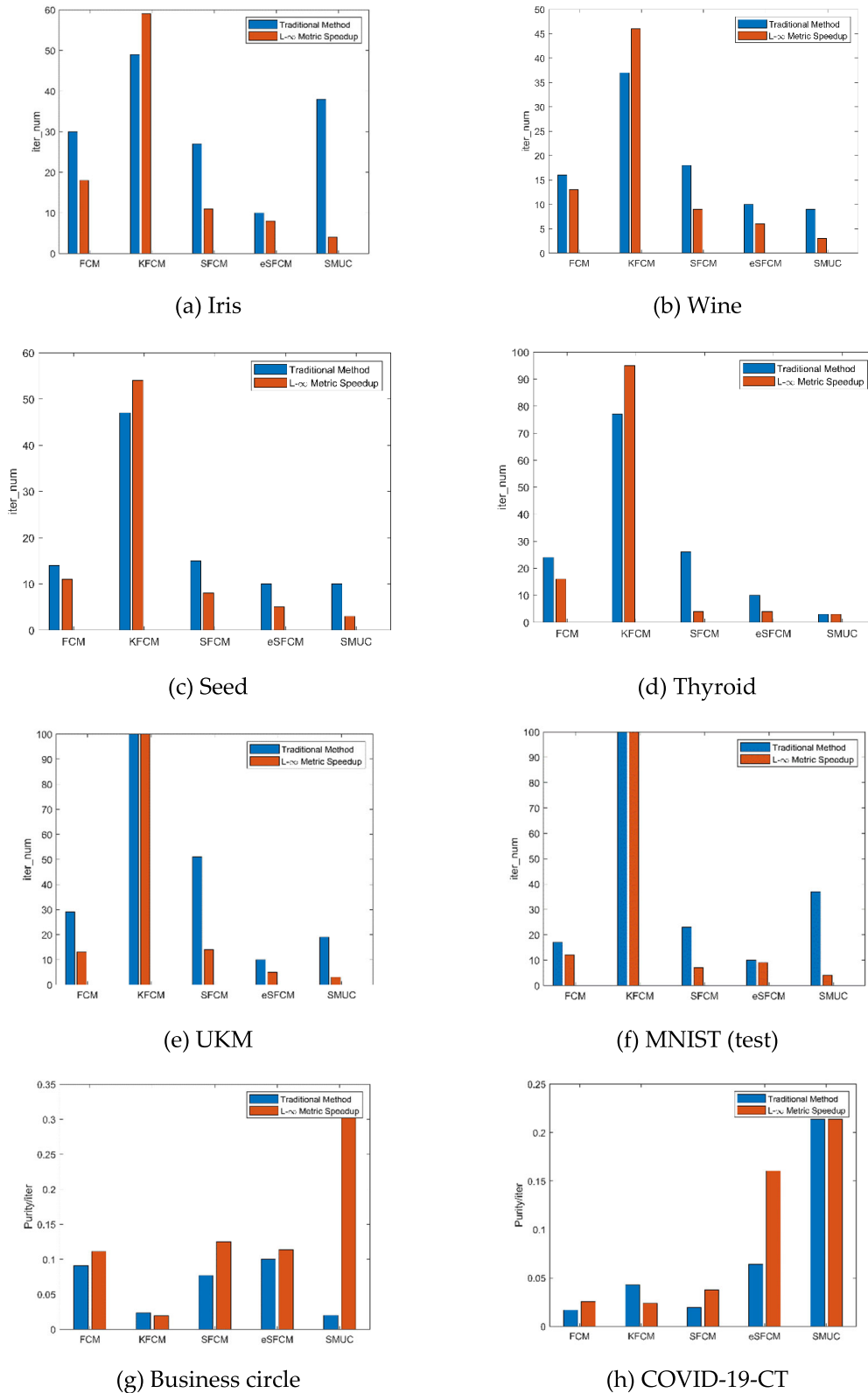


FIGURE 13. Comparison of the number of iterations required to run various algorithms.

function. In the experiment, while combining the external evaluation index and iteration speed, the acceleration effect of

dynamic metric is not substantial when dealing with linearly inseparable data. The optimization effect of our acceleration

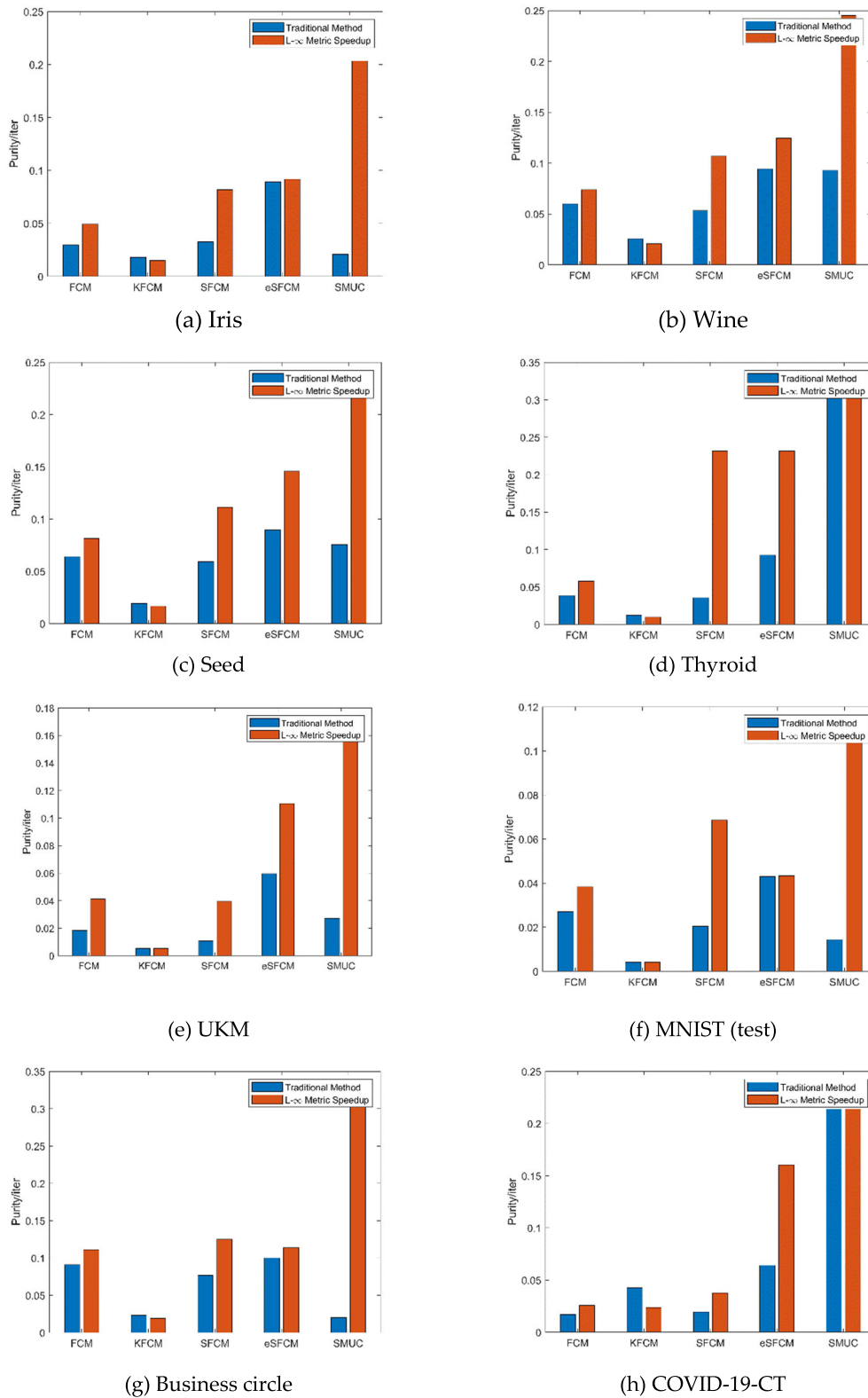


FIGURE 14. Iterative/purity level complete with the dynamic metric accelerated method and traditional method.

method based on Euclidean distance is weaker than that based on Markov distance, and the optimization effect of

unsupervised algorithm is lower than that of semi supervised algorithm.

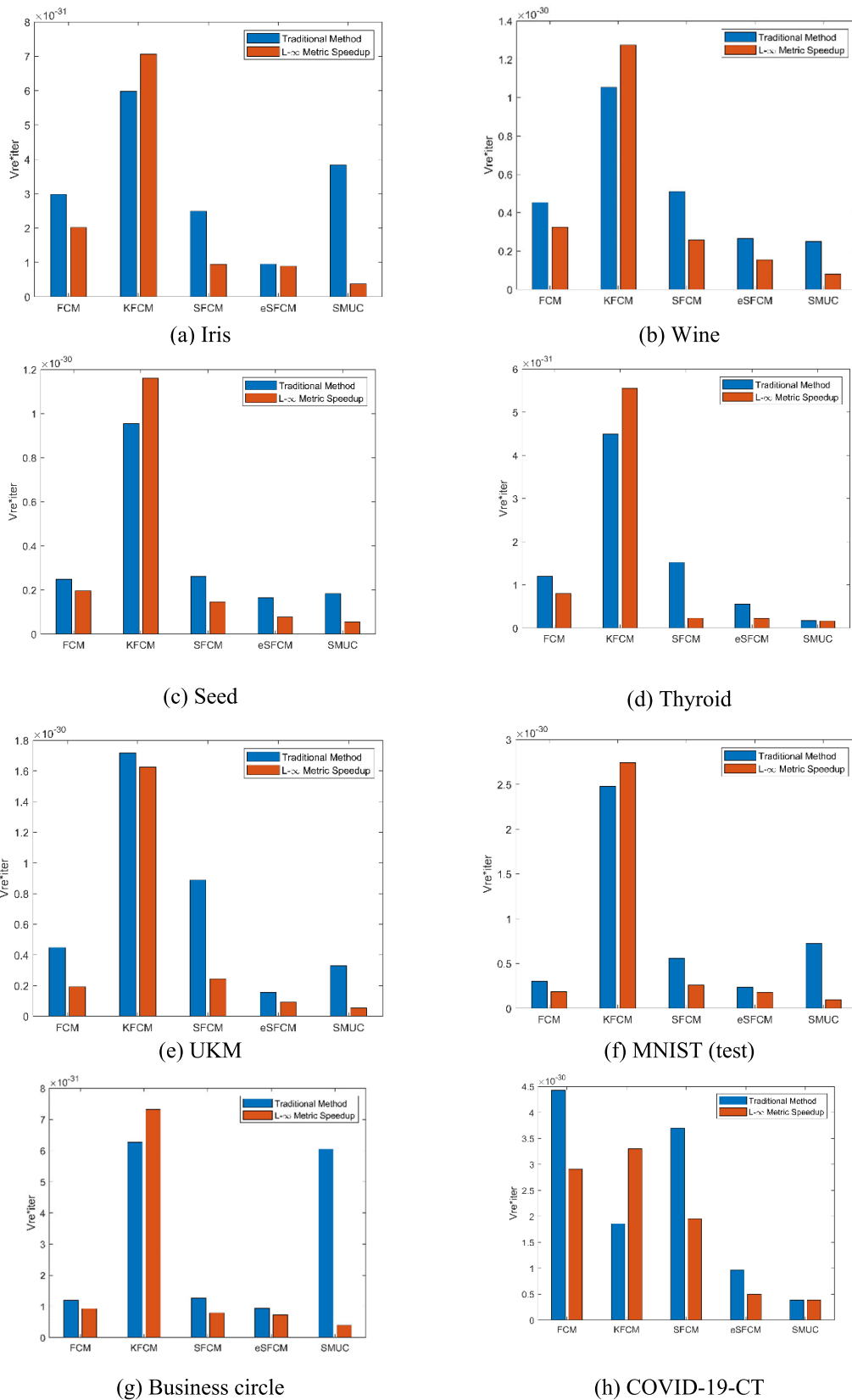


FIGURE 15. Metric accelerates the iterative $V_{RE} \cdot \text{iter}$ promotion level.

IX. CONCLUSION

While focusing on the problem of non-standardized distance in clustering algorithms, we develop a fuzzy clustering acceleration method based on dynamic metrics. This method standardizes the distance between instances to improve the search efficiency of the iteration process. We prove its stability in metric normalization, and verify that its optimization performance is not affected by the change of edge points in clusters. The method is applied to five common fuzzy clustering algorithms, and the experimental results on eight real datasets show that the method has an effective acceleration, while maintaining the stability of clustering accuracy. Furthermore, the algorithm uses the Chebyshev distance to standardize between features, which keeps the time complexity of the algorithm constant during metric iterations. In addition, the experimental results show that the dynamic metric acceleration method is less effective in the algorithm with the Gauss kernel function.

APPENDIX

See Figs. 11–15.

REFERENCES

- [1] E. H. Ruspini, "A new approach to clustering," *Inf. Control*, vol. 15, no. 1, pp. 22–32, 1969.
- [2] S. Tamura, S. Higuchi, and K. Tanaka, "Pattern classification based on fuzzy relations," *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-1, no. 1, pp. 61–66, Jan. 1971.
- [3] K. Le, "Fuzzy relation compositions and pattern recognition," *Inf. Sci.*, vol. 89, nos. 1–2, pp. 107–130, 1996.
- [4] Z. Wu and R. Leahy, "An optimal graph theoretic approach to data clustering: Theory and its application to image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 15, no. 11, pp. 1101–1113, Nov. 1993.
- [5] G. Xinbo and X. Weixin, "Research progress in the development and application of fuzzy clustering theory," *Chin. Sci. Bull.*, vol. 44, no. 21, pp. 2241–2250, 1999.
- [6] D. Tianchen, "Metric learning and clustering algorithm base on transfer distance," Yangzhou Univ., Tech. Rep., 2019.
- [7] I. Dokmanic, R. Parhizkar, J. Ranieri, and M. Vetterli, "Euclidean distance matrices: A short walk through theory, algorithms and applications," *IEEE Signal Process. Mag.*, vol. 32, no. 6, pp. 12–30, Jan. 2015.
- [8] P. Groenen, U. Kaymak, and J. Rosmalen, "Fuzzy clustering with Minkowski distance functions," in *Advances in Fuzzy Clustering and its Applications*. Amsterdam, The Netherlands: Elsevier, 2001.
- [9] Y. Kanzawa, J. Shibaura, Y. Endo, and S. Miyamoto, "KL-divergence-based and Manhattan distance-based semisupervised entropy-regularized fuzzy C-means," *J. Adv. Comput. Intell. Inform.*, vol. 15, no. 8, pp. 1057–1064, Oct. 2011.
- [10] L. Jiawei, Y. Zhiqiang, and L. Shuyan, "Novel method for phylogenetic tree construction based on correlation feature and fuzzy clustering," *Appl. Res. Comput.*, vol. 28, no. 8, pp. 2844–2847, 2011.
- [11] X. Zhao, Y. Li, and Q. Zhao, "Mahalanobis distance based on fuzzy clustering algorithm for image segmentation," *Digit. Signal Process.*, vol. 43, pp. 8–16, Aug. 2015.
- [12] S. F. Huang, Y. H. Lin, and J. M. Yih, "Fuzzy clustering algorithm based on Mahalanobis distances with recursive process," *Int. J. Intell. Technol. Appl. Statist.*, vol. 11, no. 14, p. 10, 2017.
- [13] U. Zhiwen and L. Qin, "Mahalanobis distance fuzzy clustering algorithm based on particle swarm optimization," *J. Chongqing Univ. Posts Telecommun.*, vol. 31, no. 2, pp. 275–284, 2019.
- [14] M. Friedman, M. Last, Y. Makover, and A. Kandel, "Anomaly detection in web documents using crisp and fuzzy-based cosine clustering methodology," *Inf. Sci.*, vol. 177, no. 2, pp. 467–475, Jan. 2007.
- [15] K. Jaferzadeh, K. Kiani, and S. Mozaffari, "Acceleration of fractal image compression using fuzzy clustering and discrete-cosine-transform-based metric," *IET Image Process.*, vol. 6, no. 7, pp. 1024–1030, 2012.
- [16] H. He, Y. Tan, and J. Huang, "Unsupervised classification of smartphone activities signals using wavelet packet transform and half-cosine fuzzy clustering," in *Proc. IEEE Int. Conf. Fuzzy Syst.*, Jul. 2017, pp. 1–6.
- [17] N. X. Thao, M. Ali, and F. Smarandache, "An intuitionistic fuzzy clustering algorithm based on a new correlation coefficient with application in medical diagnosis," *J. Intell. Fuzzy Syst.*, vol. 36, no. 1, pp. 189–198, Feb. 2019.
- [18] N. B. Karayiannis, "MECA: Maximum entropy clustering algorithm," in *Proc. IEEE 3rd Int. Fuzzy Syst. Conf.*, Jun. 1994, pp. 630–635.
- [19] E. Yasunori, H. Yukihiro, and Y. Makito, "On semi-supervised fuzzy c-means clustering," in *Proc. IEEE Int. Conf. Fuzzy Syst.*, Dec. 2009, pp. 1119–1124.
- [20] M. Zarinbal, M. H. Fazel Zarandi, and I. B. Turksen, "Relative entropy fuzzy C-means clustering," *Inf. Sci.*, vol. 260, pp. 74–97, Mar. 2014.
- [21] M. Ionescu and A. Ralescu, "Image clustering for a fuzzy Hamming distance based CBIR system," in *Proc. Midwest Artif. Intell. Cogn. Sci. Conf.*, 2005, pp. 102–108.
- [22] J.-H. Chiang and P.-Y. Hao, "A new kernel-based fuzzy clustering approach: Support vector clustering with cell growing," *IEEE Trans. Fuzzy Syst.*, vol. 11, no. 4, pp. 518–527, Aug. 2003.
- [23] Y. Li and Y. Shen, "Robust image segmentation algorithm using fuzzy clustering based on kernel-induced distance measure," in *Proc. Int. Conf. Comput. Sci. Softw. Eng.*, Dec. 2008, pp. 1065–1068.
- [24] M. Gong, Y. Liang, J. Shi, W. Ma, and J. Ma, "Fuzzy C-means clustering with local information and kernel metric for image segmentation," *IEEE Trans. Image Process.*, vol. 22, no. 2, pp. 573–584, Feb. 2013.
- [25] F. Salehi, M. R. Keyvanpour, and A. Sharifi, "SMKFC-ER: Semi-supervised multiple kernel fuzzy clustering based on entropy and relative entropy," *Inf. Sci.*, vol. 547, pp. 667–688, Feb. 2021.
- [26] F. Zhao, Z. Zeng, H. Liu, R. Lan, and J. Fan, "Semisupervised approach to surrogate-assisted multiobjective kernel intuitionistic fuzzy clustering algorithm for color image segmentation," *IEEE Trans. Fuzzy Syst.*, vol. 28, no. 6, pp. 1023–1034, Jun. 2020.
- [27] J. C. Bezdek, W. Full, and R. Ehrlich, "FCM: The fuzzy C-means clustering algorithm," *Comput. Geosci.*, vol. 10, nos. 2–3, pp. 191–203, 1984.
- [28] Q.-T. Bui, B. Vo, V. Snael, W. Pedrycz, T.-P. Hong, N.-T. Nguyen, and M.-Y. Chen, "SFCM: A fuzzy clustering algorithm of extracting the shape information of data," *IEEE Trans. Fuzzy Syst.*, vol. 29, no. 1, pp. 75–89, Jan. 2021.
- [29] N. Grira, M. Crucianu, and N. Boujemaa, "Semi-supervised fuzzy clustering with pairwise-constrained competitive agglomeration," in *Proc. 14th IEEE Int. Conf. Fuzzy Syst.*, May 2005, pp. 867–872.
- [30] X. Yin, T. Shu, and Q. Huang, "Semi-supervised fuzzy clustering with metric learning and entropy regularization," *Knowl.-Based Syst.*, vol. 35, pp. 304–311, Nov. 2012.
- [31] C. Wu and Z. Cao, "Entropy-like divergence based kernel fuzzy clustering for robust image segmentation," *Expert Syst. Appl.*, vol. 169, May 2021, Art. no. 114327.
- [32] D.-Q. Zhang and S.-C. Chen, "Clustering incomplete data using kernel-based fuzzy C-means algorithm," *Neural Process. Lett.*, vol. 18, pp. 155–162, Dec. 2003.
- [33] S. Kapil and M. Chawla, "Performance evaluation of K-means clustering algorithm with various distance metrics," in *Proc. IEEE 1st Int. Conf. Power Electron., Intell. Control Energy Syst. (ICPEICES)*, Jul. 2016, pp. 1–4.
- [34] C. Shen, J. Kim, and L. Wang, "Scalable large-margin Mahalanobis distance metric learning," *IEEE Trans. Neural Netw.*, vol. 21, no. 9, pp. 1524–1530, Sep. 2010.
- [35] J. Yao, M. Dash, S. T. Tan, and H. Liu, "Entropy-based fuzzy clustering and fuzzy modeling," *Fuzzy Sets Syst.*, vol. 113, no. 3, pp. 381–388, 2000.
- [36] M.-S. Yang and Y. Nataliani, "A feature-reduction fuzzy clustering algorithm based on feature-weighted entropy," *IEEE Trans. Fuzzy Syst.*, vol. 26, no. 2, pp. 817–835, Apr. 2018.
- [37] J. Hartigan and M. Wong, "Algorithm AS 136: A K-means clustering algorithm," *Appl. Statist.*, vol. 28, no. 1, pp. 100–108, 1979.
- [38] E. Xing, M. Jordan, S. J. Russell, and A. Ng, "Distance metric learning with application to clustering with side-information," in *International Conference on Neural Information Processing Systems*. Cambridge, MA, USA: MIT Press, 2002.
- [39] F. Klawonn and A. Keller, "Fuzzy clustering based on modified distance measures," in *Proc. Int. Symp. Intell. Data Anal.* Berlin, Germany: Springer, 1999, pp. 291–301.

- [40] A. Flores-Sintas, J. Cadenas, and F. Martin, "A local geometrical properties application to fuzzy clustering," *Fuzzy Sets Syst.*, vol. 100, nos. 1–3, pp. 245–256, Nov. 1998.
- [41] K. Li and Y. Zhou, "An improved semi-supervised fuzzy clustering algorithm," *Int. J. Comput. Inf. Technol.*, vol. 2, pp. 115–120, Dec. 2013.
- [42] D. T. C. Lai, J. M. Garibaldi, and J. Reps, "Investigating distance metric learning in semi-supervised fuzzy C-means clustering," in *Proc. IEEE Int. Conf. Fuzzy Syst.*, Jul. 2014, pp. 1817–1824.
- [43] R. L. Carter, R. Morris, and R. K. Blashfield, "On the partitioning of squared Euclidean distance and its applications in cluster analysis," *Psychometrika*, vol. 54, no. 1, pp. 9–23, Mar. 1989.
- [44] J. C. Bezdek, "Pattern recognition with fuzzy objective function algorithms," *Adv. Appl. Pattern Recognit.*, vol. 22, no. 1171, pp. 203–239, 1981.
- [45] W. Pedrycz and J. Waletzky, "Fuzzy clustering with partial supervision," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 27, no. 5, pp. 787–795, Sep. 1997.
- [46] H. Suresh and G. Raj, "An unsupervised fuzzy clustering method for Twitter sentiment analysis," in *Proc. Int. Conf. Comput. Syst. Inf. Technol. for Sustain. Solutions (CSITSS)*, Oct. 2016, pp. 80–85.
- [47] L. Yang, X. Geng, and H. Liao, "A web sentiment analysis method on fuzzy clustering for mobile social media users," *EURASIP J. Wireless Commun. Netw.*, vol. 2016, no. 1, p. 128, Dec. 2016.
- [48] G. Cosma and G. Acampora, "A computational intelligence approach to efficiently predicting review ratings in E-commerce," *Appl. Soft Comput.*, vol. 44, pp. 153–162, Jul. 2016.
- [49] M. D. Malkauthekar, "Analysis of Euclidean distance and Manhattan distance measure in face recognition," in *Proc. 3rd Int. Conf. Comput. Intell. Inf. Technol. (CIIT)*, 2013, pp. 503–507.
- [50] H. Ichihashi, M. Ohue, and T. Miyoshi, "Fuzzy C-means clustering algorithm with pseudo Mahalanobis distances," in *Proc. Korean Inst. Intell. Syst. Conf.*, 1998, pp. 148–152.
- [51] A. M. Aziz, "A new nearest-neighbor association approach based on fuzzy clustering," *Aerosp. Sci. Technol.*, vol. 26, no. 1, pp. 87–97, Apr. 2013.
- [52] S. Zeng, X. Wang, X. Duan, S. Zeng, Z. Xiao, and D. Feng, "Kernelized Mahalanobis distance for fuzzy clustering," *IEEE Trans. Fuzzy Syst.*, vol. 29, no. 10, pp. 3103–3117, Oct. 2021.
- [53] *MATLAB Data Analysis and Data Mining*, China Machine Press, Beijing, China, 2015.
- [54] X. Yang, X. He, J. Zhao, Y. Zhang, S. Zhang, and P. Xie, "COVID-CT-dataset: A CT scan dataset about COVID-19," 2020, *arXiv:2003.13865*.
- [55] N. R. Pal and J. C. Bezdek, "On cluster validity for the fuzzy c-means model," *IEEE Trans. Fuzzy Syst.*, vol. 3, no. 3, pp. 370–379, Aug. 1995.
- [56] S. C. Sripada and M. S. Rao, "Comparison of purity and entropy of K-means clustering and fuzzy C means clustering," *Indian J. Comput. Sci. Eng.*, vol. 2, no. 3, pp. 343–346, 2011.
- [57] W. Min, L. Yuan, and Z. L. Cai, "Blind area segmentation algorithm based on texture features," *Inf. Commun.*, no. 7, pp. 23–26, Jul. 2017.



SHENGBING XU received the B.S. degree in computational mathematics and applied software and the M.Sc. degree in applied mathematics from Xiangtan University, in 1997 and 2001, respectively. He is currently pursuing the Ph.D. degree with the Guangdong University of Technology. His research interests include fuzzy set, mathematical modeling, machine learning, and their applications. He is also a Peer Reviewer of IEEE ACCESS.



WEI CAI studied applied statistics at the Guangdong University of Technology. His research areas include neural networks, fuzzy clustering, and optimization method. His research interests also include machine learning, data mining, and computer vision.



HONGXI XIA is currently pursuing the bachelor's degree in applied mathematics with the Guangdong University of Technology. His research interests include machine learning, fuzzy clustering, and optimization.



BO LIU is currently with the Faculty of Automation, Guangdong University of Technology. His research interests include machine learning and data mining. He has published papers on IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS, IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, *Knowledge and Information Systems*, IEEE International Conference on Data Mining (ICDM), SIAM International Conference on Data Mining (SDM), and ACM International Conference on Information and Knowledge Management (CIKM).



JIE XU received the M.S. and Ph.D. degrees in mathematics from East China Normal University, Shanghai, in 2003 and 2006, respectively. He is currently a Lecturer with the Department of Mathematics and Applied Mathematics, Guangdong University of Technology. His research area includes cyclotomic Brauer algebras. His research interests include algebraic representation theory and optimization methods.

...