


A new type of distance metric and its use for clustering

Xiaowei Gu¹  · Plamen P. Angelov¹ · Dmitry Kangin¹ · Jose C. Principe²

Received: 4 May 2017 / Accepted: 27 June 2017 / Published online: 26 July 2017
© Springer-Verlag GmbH Germany 2017

Abstract In order to address high dimensional problems, a new ‘direction-aware’ metric is introduced in this paper. This new distance is a combination of two components: (1) the traditional Euclidean distance and (2) an angular/directional divergence, derived from the cosine similarity. The newly introduced metric combines the advantages of the Euclidean metric and cosine similarity, and is defined over the Euclidean space domain. Thus, it is able to take the advantage from both spaces, while preserving the Euclidean space domain. The direction-aware distance has wide range of applicability and can be used as an alternative distance measure for various traditional clustering approaches to enhance their ability of handling high dimensional problems. A new evolving clustering algorithm using the proposed distance is also proposed in this paper. Numerical examples with benchmark datasets reveal that the direction-aware distance can effectively improve the clustering quality of the k-means algorithm for high dimensional problems and demonstrate the proposed evolving clustering

algorithm to be an effective tool for high dimensional data streams processing.

Keywords Cosine similarity · Distance metric · Metric space · Clustering · High dimensional data streams processing

1 Introduction

The widely used clustering techniques may use different kind of distances to measure the separation between data samples. The well-known Euclidean distance is currently the most frequently used metric space for the established clustering algorithms (MacQueen 1967; Fukunaga and Hostetler 1975). Other metric spaces, using the Mahalanobis (McLachlan 1999), city block, Hamming, Minkowski types of distances, etc., are also widely used in different clustering algorithms for different purposes. It is often the case that clustering algorithms employing divergences, i.e. pairwise dissimilarity, which does not obey all the properties of distances (e.g. cosine similarity), could generate meaningless conclusions.

One problem the traditional distance metrics are facing is the so-called “curse of dimensionality” (Domingos 2012; Aggarwal et al. 2001). Many clustering techniques, which use the traditional distance metrics work well in low dimensional space, however, become intractable for high dimensional problems. Research results have shown that in high dimensional space, the concept of distance may not even be qualitatively meaningful (Aggarwal et al. 2001; Beyer et al. 1999). Under certain reasonable conditions, it has been found that the distances of the nearest and farthest neighbours to a given data sample are the same for a number of distance metrics in high dimensional space (Beyer et al. 1999). This phenomenon is frequently seen in

✉ Plamen P. Angelov
p.angelov@lancaster.ac.uk

Xiaowei Gu
x.gu3@lancaster.ac.uk

Dmitry Kangin
d.kangin@lancaster.ac.uk

Jose C. Principe
principe@cnel.ufl.edu

¹ School of Computing and Communications, Lancaster University Lancaster, B24, InfoLab21, Bailrigg, Lancaster LA1 4WA, UK

² Computational NeuroEngineering Laboratory, Department of Electrical and Computer Engineering, University of Florida, Gainesville, USA

the cases that some dimensions of the data are highly irrelevant. This is not hard to understand because our intuitions come from a three-dimensional world only, which may not be applicable to high dimensional ones.

Compared with the commonly used distance metrics including the Euclidean, Mahalanobis, Minkowski distances, etc., which measure the magnitude of vector difference, cosine similarity focuses much more on the directional similarity. Therefore, it is more often used in the natural language processing (NLP) problems (Allah et al. 2008; Dehak et al. 2010, 2011; Setlur and Stone 2016; Senoussaoui et al. 2013). In NLP problems, machine learning algorithms, for example, k-means (Allah et al. 2008; Setlur and Stone 2016), mean shift (Senoussaoui et al. 2013), etc., are used to cluster very high dimensional vectors representing the documents together based on the cosine similarity. Nonetheless, the cosine similarity is a pseudo metric because it does not obey the triangle inequality [it obeys the Cauchy–Schwarz inequality (Callebaut 1965)]. Consequently, the cosine similarity between two vectors can be misleading and hides information, especially in cases where the vectors are sparse or orthogonal.

In this paper, a new “direction-aware” distance is introduced. This new metric space is a combination of a distance (in this paper, we consider Euclidean), and an angular/directional component, which is based on the cosine similarity. Therefore, it takes the advantages of the both components while still obeys all the properties of a distance metric (McCune et al. 2002) as we will demonstrate.

The proposed distance in this paper is applicable to various machine learning algorithms including the recently published ones (Angelov et al. 2014; Rong et al. 2006, 2011; Precup et al. 2014; Lughofer et al. 2015) as an alternative distance measure and can enhance the ability of the algorithms to handle high dimensional problems. A new evolving clustering algorithm is also proposed for streaming data processing. This algorithm employs the new direction-aware distance only and is able to start “from scratch”. Therefore, it is very suitable for handling the high dimensional data streams.

Numerical examples using benchmark datasets demonstrate the potential of the direction-aware distance against many traditional metrics in high dimensional problems. It is also shown that the proposed clustering algorithm is able to produce top quality clustering results on various problems with high computational efficiency.

The remainder of this paper is organised as follows. Section 2 describes the newly proposed direction-aware distance and provides the proof for the proposed distance to be a full metric. Section 3 introduces the application of the newly proposed direction-aware distance to traditional clustering algorithms. The new evolving clustering algorithm based on the proposed distance is presented in Sect. 4. Section 5 presents numerical examples. The paper is concluded by Sect. 6.

2 Direction-aware distance and proof of metric axioms

2.1 The new direction-aware distance

In this section, we introduce the direction-aware distance, and prove that it is a distance over the space of real numbers. If no specific declaration is provided, all the derivations in this paper are conducted over the real numbers.

First of all, let us define a metric space, \mathbf{R}^m , \mathbf{x} and \mathbf{y} are two data points within the space, m is the dimensionality of the metric space \mathbf{R}^m . The newly introduced direction-aware distance, $d_{DA}(\mathbf{x}, \mathbf{y})$ consists of two terms:

1. A Euclidean component, $d_M(\mathbf{x}, \mathbf{y})$, and
2. A direction-aware component, $d_A(\mathbf{x}, \mathbf{y})$, and is expressed as:

$$d_{DA}(\mathbf{x}, \mathbf{y}) = \sqrt{\lambda_M^2 (d_M(\mathbf{x}, \mathbf{y}))^2 + \lambda_A^2 (d_A(\mathbf{x}, \mathbf{y}))^2}, \tag{1}$$

where $\mathbf{x} = [x_1, x_2, \dots, x_m]^T$ and $\mathbf{y} = [y_1, y_2, \dots, y_m]^T$, $\mathbf{x}, \mathbf{y} \in \mathbf{R}^m$; λ_M, λ_A are a pair of scaling coefficients and $\lambda_M > 0, \lambda_A > 0$; $d_M(\mathbf{x}, \mathbf{y})$ denotes the Euclidean distance between \mathbf{x} and \mathbf{y} , $d_M(\mathbf{x}, \mathbf{y}) =$

$$\sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})} = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}.$$

The direction-aware component $d_A(\mathbf{x}, \mathbf{y})$ is derived based on the cosine similarity expressed by:

$$d_A(\mathbf{x}, \mathbf{y}) = \sqrt{1 - \cos(\Theta_{xy})}, \tag{2}$$

where Θ_{xy} is the angle between \mathbf{x} and \mathbf{y} . In the Euclidean space, since $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^m x_i y_i$ and $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$, thus $\cos(\Theta_{xy}) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}$. Therefore, the directional component $d_A(\mathbf{x}, \mathbf{y})$ can be rewritten as:

$$\begin{aligned} d_A(\mathbf{x}, \mathbf{y}) &= \sqrt{1 - \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}} = \sqrt{1 - \frac{\sum_{i=1}^m x_i y_i}{\|\mathbf{x}\| \|\mathbf{y}\|}} \\ &= \sqrt{\frac{\sum_{i=1}^m x_i^2}{2\|\mathbf{x}\|^2} + \frac{\sum_{i=1}^m y_i^2}{2\|\mathbf{y}\|^2} - \frac{\sum_{i=1}^m x_i y_i}{\|\mathbf{x}\| \|\mathbf{y}\|}} \\ &= \sqrt{\frac{1}{2} \sum_{i=1}^m \left(\frac{x_i}{\|\mathbf{x}\|} - \frac{y_i}{\|\mathbf{y}\|} \right)^2} \\ &= \frac{1}{\sqrt{2}} \left\| \frac{\mathbf{x}}{\|\mathbf{x}\|} - \frac{\mathbf{y}}{\|\mathbf{y}\|} \right\|. \end{aligned} \tag{3}$$

One can notice that, if x or y is equal to 0, $d_A(x, y) = 0$.

2.2 Proof of metric axioms

In this subsection, we will prove that the proposed distance is a full metric. For a distance $d(x, y)$ in the space to be a full metric, \mathbf{R}^m , it is required to satisfy the following properties for $\forall x, y$ (13):

1. Non-negativity:

$$d(x, y) \geq 0; \tag{4}$$

2. Identity of indiscernibles:

$$d(x, y) = 0 \text{ iff } x = y; \tag{5}$$

3. Symmetry:

$$d(x, y) = d(y, x); \tag{6}$$

4. Triangle inequality:

$$d(x, z) + d(z, y) \geq d(x, y). \tag{7}$$

In this paper, we propose a new theorem as follows:

Theorem $d_{DA}(x, y)$ is a distance within the metric space over the domain \mathbf{R}^m .

In the rest of this subsection, we will prove this theorem by proving that $d_{DA}(x, y)$ obeys the four distance axioms stated in Eqs. (5–6) and inequalities (4) and (7) one by one.

Lemma 1 $\forall x, y \in \mathbf{R}^m, d_{DA}(x, y) \geq 0$

Proof It can be seen directly from the Eq. (5) that $d_{DA}(x, y)$ is always non-negative.

Lemma 2 $\forall x, y \in \mathbf{R}^m, d_{DA}(x, y) = 0$ iff $x = y$.

Proof It is clear that if $x = y$, then $d_A(x, y) = \sqrt{1 - 1} = 0$, $d_M(x, y) = 0$ and $d_{DA}(x, y) = 0$.

The directional component $d_A(x, y)$ alone does not obey this property because as we can see from equations (2) and (3), if x and y are nonzero and orthogonal, $d_A(x, y) = 0$, so it is not true. However, in this case, due to the fact that if $x \neq y$, $d_M(x, y) \neq 0$, $d_{DA}(x, y)$ will still be nonzero as $\lambda_M, \lambda_A > 0$. Therefore, one can still conclude that $d_{DA}(x, y) = 0$ if and only if $x = y$.

Lemma 3 $\forall x, y \in \mathbf{R}^m, d_{DA}(x, y) = d_{DA}(y, x)$

Proof For the Euclidean metric, it is true that:

$$\begin{aligned} d_{DA}(x, y) &= \sqrt{\lambda_M^2 \|x - y\|^2 + \frac{\lambda_A^2}{2} \left\| \frac{x}{\|x\|} - \frac{y}{\|y\|} \right\|^2} \\ &= \sqrt{\lambda_M^2 \|y - x\|^2 + \frac{\lambda_A^2}{2} \left\| \frac{y}{\|y\|} - \frac{x}{\|x\|} \right\|^2} \\ &= d_{DA}(y, x). \end{aligned} \tag{8}$$

Therefore, $d_{DA}(x, y) = d_{DA}(y, x)$.

Lemma 4 $\forall x, y, z \in \mathbf{R}^m, d_{DA}(x, z) \leq d_{DA}(x, y) + d_{DA}(y, z)$

Proof Firstly, let us assume that there is a triplet data samples x, y, z , which make d_{DA} break the triangle rule, namely:

$$d_{DA}(x, z) > d_{DA}(x, y) + d_{DA}(y, z). \tag{9}$$

By including Eq. (3) in Eq. (2), the direction-aware distance $d_{DA}(x, y)$ can be rewritten as:

$$\begin{aligned} d_{DA}(x, y) &= \sqrt{\lambda_M^2 (d_M(x, y))^2 + \lambda_A^2 (d_A(x, y))^2} \\ &= \sqrt{\lambda_M^2 \|x - y\|^2 + \frac{\lambda_A^2}{2} \left\| \frac{x}{\|x\|} - \frac{y}{\|y\|} \right\|^2} \\ &= \sqrt{\lambda_M^2 \sum_{i=1}^m (x_i - y_i)^2 + \frac{\lambda_A^2}{2} \sum_{i=1}^m \left(\frac{x_i}{\|x\|} - \frac{y_i}{\|y\|} \right)^2} \\ &= \|\chi - \psi\| = d_M(\chi, \psi), \end{aligned} \tag{10}$$

where, $\chi = \left[\lambda_M x^T, \frac{\lambda_A x^T}{\sqrt{2}\|x\|} \right]^T = \left[\lambda_M x_1, \lambda_M x_2, \dots, \lambda_M x_m, \frac{\lambda_A x_1}{\sqrt{2}\|x\|}, \frac{\lambda_A x_2}{\sqrt{2}\|x\|}, \dots, \frac{\lambda_A x_m}{\sqrt{2}\|x\|} \right]^T$ and $\psi = \left[\lambda_M y^T, \frac{\lambda_A y^T}{\sqrt{2}\|y\|} \right]^T = \left[\lambda_M y_1, \lambda_M y_2, \dots, \lambda_M y_m, \frac{\lambda_A y_1}{\sqrt{2}\|y\|}, \frac{\lambda_A y_2}{\sqrt{2}\|y\|}, \dots, \frac{\lambda_A y_m}{\sqrt{2}\|y\|} \right]^T$.

Similarly, for $\zeta = \left[\lambda_M z^T, \frac{\lambda_A z^T}{\sqrt{2}\|z\|} \right]^T$, we can see that

$$d_{DA}(x, z) = d_M(\chi, \zeta), d_{DA}(y, z) = d_M(\psi, \zeta).$$

Considering an auxiliary algebraic data space \mathbf{R}^{2m} , for χ, ψ, ζ , it follows that:

$$d_M(\chi, \zeta) \leq d_M(\chi, \psi) + d_M(\psi, \zeta). \tag{11}$$

As we can see from inequalities (9) and (11), the two equations have the same algebraic form, but there are different signs ($>$ and \leq). For Euclidean distance in \mathbf{R}^{2m} , the triangle rule is always conformed, therefore, we can conclude that $d_{DA}(x, y)$ always satisfies the triangle inequality: $d_{DA}(x, z) \leq d_{DA}(x, y) + d_{DA}(y, z)$.

Based on the proofs of the four Lemmas, the proposed Theorem is proven. Therefore, we can conclude that the proposed *direction-aware* distance, d_{DA} is a full distance in the Euclidean space.

2.3 The property of the proposed distance

The proposed direction-aware distance metric is a combination of two components: (1) the traditional Euclidean distance and (2) an angular/directional divergence, derived from the cosine similarity. It defines a metric space as a combination of Euclidean metric space and cosine similarity pseudo-metric space, and consequently, can effectively combine information extracted from both spaces and takes into account both spatial and angular divergences. Therefore, the direction aware distance can serve as a more representative distance metric than the traditional distance metric.

3 The application of the proposed distance to traditional clustering approaches

In this section, we will describe the applications of the proposed distance to the traditional offline clustering approaches. First of all, let us define the dataset in the metric space as $\{\mathbf{x}\}_N = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \in \mathbf{R}^m$, $\mathbf{x}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,d}]^T \in \mathbf{R}^d$, $i = 1, 2, \dots, N$, where N is the number of data samples in the dataset.

The newly proposed direction-aware distance can be used in various clustering, classification as well as regression approaches. For example, the k-means (Allah et al. 2008; Setlur and Stone 2016), mean-shift clustering (Comaniciu and Meer 2002), k nearest neighbour classification (Keller and Gray 1985) algorithms may use the newly introduced direction-aware distance to enhance the ability in dealing with high dimensional data.

Since the traditional offline algorithms have been studied well for many years, in this paper, we will not focus on the algorithm themselves. Instead, we will look at the direction-aware distance and introduce the strategy of using the proposed distance in the algorithms for different purposes.

The direction-aware distance has a pair of scaling factors, the values of which can be adjusted for various problems. For example, if without losing generality, we want to allocate the same importance to the Euclidean and directional components, λ_M and λ_A can be set as the inverse of average d_M and d_A , respectively (the data is taken without pre-processing):

$$\lambda_M = \frac{1}{\sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N d_M^2(\mathbf{x}_i, \mathbf{x}_j)}{N^2}}}, \quad (12a)$$

$$\lambda_A = \frac{1}{\sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N d_A^2(\mathbf{x}_i, \mathbf{x}_j)}{N^2}}}. \quad (12b)$$

Alternatively, if the data has been re-scaled to the range $[0, 1]$ in advance, the values of d_M and d_A are within the ranges $[0, \sqrt{m}]$ and $[0, 1]$, respectively, thus, the pair of the scaling coefficients within the proposed distance can be set to $\lambda_M = \frac{1}{\sqrt{m}}$ and $\lambda_A = 1$ if we aim to allocate the same importance to each component.

While for some problems like NLP, where the directional similarity plays a more important role compared with magnitude differences, we can enhance the importance of the directional component in the distance measures by increasing the value of λ_A , and vice versa. The scaling factors λ_M and λ_A that allow the clustering approaches to achieve the best performance with the proposed direction-aware distance are always problem-specific, which incorporates the *prior* knowledge of the problem. We believe that this choice is out of the scope of this paper.

4 The applications of the proposed distance to evolving clustering

Similarly, the direction-aware distance can also be employed in the evolving clustering approaches. In this section, we propose a new evolving clustering approach with the direction-aware distance. This algorithm is able to “start from scratch” and consistently evolves its system structure and updates the meta-parameters based on the newly arrived data samples.

The main procedure of the proposed algorithm is described as follows. In this section, we consider $\{\mathbf{x}\}_k = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\} \in \mathbf{R}^m$ as a data stream and the subscript indicates the time instance that the data sample arrives.

Stage 1 Initialization

The first data sample \mathbf{x}_1 in the data stream is used for initializing the system and its meta-parameters. In the proposed algorithm, the system has the following initialized global meta-parameters:

1. $k \leftarrow 1$, the current time instance;
2. $C \leftarrow 1$, the number of exiting clusters;
3. $\boldsymbol{\mu}_M \leftarrow \mathbf{x}_1$, the global mean of $\{\mathbf{x}\}_k$;
4. $X_M \leftarrow \|\mathbf{x}_1\|^2$, the global average scalar product of $\{\mathbf{x}\}_k$;
5. $\boldsymbol{\mu}_A \leftarrow \frac{\mathbf{x}_1}{\|\mathbf{x}_1\|}$, the global mean of $\left\{ \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\}_k$, which is also the normalized global mean of $\{\mathbf{x}\}_k$.
6. $X_A \leftarrow \left\| \frac{\mathbf{x}_1}{\|\mathbf{x}_1\|} \right\|^2 = 1$, the global average scalar product of $\left\{ \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\}_k$, which is always equal to 1.

The local meta-parameters of the first cluster are initialized as follows:

1. $\Xi^1 \leftarrow \{x_1\}$, the first cluster;
2. $f_M^1 \leftarrow x_1$, the centre of the first cluster, which is also the mean of Ξ^1 ;
3. $X_M^1 \leftarrow \|x_1\|^2$, the average scalar product of Ξ^1 ;
4. $f_A^1 \leftarrow \frac{x_1}{\|x_1\|}$, the normalized mean of Ξ^1 ;
5. $X_A^1 \leftarrow 1$, the normalized average scalar product of Ξ^1 , which is always equal to 1 as well;
6. $S^1 \leftarrow 1$, the support (population) of the first cluster.

After the initialization of the system, the proposed algorithm updates the system structure and meta-parameters with the arrival of each new data samples.

Stage 2 System structure and meta-parameters update

With each newly arrived data sample, the system’s global meta-parameters, μ_M , X_M and μ_A are updated using the following equations (Angelov et al. 2017):

$$\mu_M \leftarrow \frac{k}{k+1} \mu_M + \frac{1}{k+1} x_{k+1}, \tag{13a}$$

$$X_M \leftarrow \frac{k}{k+1} X_M + \frac{1}{k+1} \|x_{k+1}\|^2, \tag{13b}$$

$$\mu_A \leftarrow \frac{k}{k+1} \mu_A + \frac{1}{k+1} \frac{x_{k+1}}{\|x_{k+1}\|}, \tag{13c}$$

$$k \leftarrow k + 1. \tag{13d}$$

Then, the condition A is checked to see whether the new data sample denoted by x_k is associated with a new cluster:

$$\text{Condition A: } IF \left(d_{DA}(x_k, \mu_M) > \max_{j=1,2,\dots,C} \left(d_{DA}(f_M^j, \mu_M) \right) \right) \text{ OR } \left(d_{DA}(x_k, \mu_M) < \min_{j=1,2,\dots,C} \left(d_{DA}(f_M^j, \mu_M) \right) \right) \text{ THEN } (x_k \text{ creates a new cluster}) \tag{14}$$

Based on the previous subsection, without a loss of generality, we use the inverse of the average d_M between the existing data samples as λ_M and the inverse of the average d_A as λ_A , correspondingly. However, for streaming data processing, it is less efficient to keep all the observed data samples in the memory and recalculate λ_M and λ_A , every time when a new data sample is observed. Therefore we introduce the recursive forms for calculating the pair of scaling coefficients as follows (Angelov et al. 2017):

$$\lambda_M = \frac{1}{\sqrt{\frac{\sum_{i=1}^k \sum_{j=1}^k d_M^2(x_i, x_j)}{k^2}}} = \frac{1}{\sqrt{2(X_M - \|\mu_M\|^2)}}, \tag{15a}$$

$$\lambda_A = \frac{1}{\sqrt{\frac{\sum_{i=1}^k \sum_{j=1}^k d_A^2(x_i, x_j)}{k^2}}} = \frac{1}{\sqrt{1 - \|\mu_A\|^2}}. \tag{15b}$$

If condition A is satisfied, a new cluster is added with x_k as its centre:

1. $C \leftarrow C + 1$, the number of existing clusters;
2. $\Xi^C \leftarrow \{x_k\}$, the new cluster;
3. $f_M^C \leftarrow x_k$, the centre of the new cluster/ mean of Ξ^C ;
4. $X_M^C \leftarrow \|x_k\|^2$, the average scalar product of Ξ^C ;
5. $f_A^C \leftarrow \frac{x_k}{\|x_k\|}$, the normalized centre of the new cluster/ normalized mean of Ξ^C ;
6. $X_A^C \leftarrow 1$, the normalized average scalar product of Ξ^C ;
7. $S^C \leftarrow 1$, the support of the new cluster.

In contrast, if condition A is not met, x_k is assigned to the cluster with the nearest centre, denoted by f_M^n as:

$$f_M^n = \arg \min_{i=1,2,\dots,C} (d_{DA}(x_k, f_M^i)). \tag{16}$$

The meta-parameters of the cluster with the nearest centre are updated as follows (Angelov et al. 2017):

$$1. \ \Xi^n \leftarrow \Xi^n \cup \{x_k\}, \tag{17a}$$

$$2. \ f_M^n \leftarrow \frac{S^n}{S^n + 1} f_M^n + \frac{1}{S^n + 1} x_k, \tag{17b}$$

$$3. \ X_M^n \leftarrow \frac{S^n}{S^n + 1} X_M^n + \frac{1}{S^n + 1} \|x_k\|^2, \tag{17c}$$

$$4. \ f_A^n \leftarrow \frac{S^n}{S^n + 1} f_A^n + \frac{1}{S^n + 1} \frac{x_k}{\|x_k\|}, \tag{17d}$$

$$5. \ S^n \leftarrow S^n + 1. \tag{17e}$$

After the update of the global and local meta-parameters, the system is ready for the arrival of the next data sample and begins a new processing cycle.

Stage 3 Clusters adjusting

In this stage, all the existing clusters will be examined and adjusted to avoid the possible overlap. For each existing

cluster Ξ^i ($i = 1, 2, \dots, C$), firstly, we find its neighbouring clusters, denoted by $\{\Xi\}_{neighbour}^i$ based on the following condition:

$$\text{Condition B: IF } \left(d_{DA}(\mathbf{f}_M^i, \mathbf{f}_M^j) > \frac{\sum_{p=1}^C \sigma_{DA}^p}{C} \right) \text{ THEN } \left(\{\Xi\}_{neighbour}^i \leftarrow \Xi^i \cup \{\Xi\}_{neighbour}^i \right), \quad (18)$$

where $(\sigma_{DA}^p)^2 = \frac{\sum_{x \in \Xi^p} \sum_{y \in \Xi^p} d_{DA}^2(x, y)}{S_p^2} = 2 \left(X_M^p - \|\mathbf{f}_M^p\|^2 \right) + \left(1 - \|\mathbf{f}_A^p\|^2 \right)$ is the average square direction-aware distance between all the members within the p^{th} cluster.

For each cluster centre, \mathbf{f}_M^i ($i = 1, 2, \dots, C$), we calculate its weighted unimodal density as (Angelov et al. 2017):

$$D^W(\mathbf{f}_M^i) = S^i \frac{\sum_{l=1}^C \sum_{j=1}^C d_{DA}^2(\mathbf{f}_M^l, \mathbf{f}_M^j)}{2C \sum_{j=1}^C d_{DA}^2(\mathbf{f}_M^i, \mathbf{f}_M^j)}, \quad (19)$$

and we also compare $D^W(\mathbf{f}_M^i)$ with the D^W of its neighbouring clusters denoted by $\{D^W(\mathbf{f}_M)\}_{neighbour}^i$, to identify

the local maxima of the weighted unimodal density, D^W :

$$\text{Condition C: IF } \left(D^W(\mathbf{f}_M^i) > \max \left(\{D^W(\mathbf{f}_M)\}_{neighbour}^i \right) \right) \text{ THEN } \left(\mathbf{f}_M^i \text{ is one of the local maxima of } D^W \right). \quad (20)$$

By identifying all the local maxima, denoted by $\{\mathbf{f}_M\}_o$ and assigning each data sample to the cluster with the nearest centre using Eq. (16), the whole clustering processing is finished. The parameters of the clusters can be extracted post factum.

The main procedure of the algorithm is summarised in the form of pseudo code as follows.

i. While a new data sample \mathbf{x}_k of the data stream is available (or until interrupted)

- * **If** (it is the first data sample) **Then**
 - Initialise global meta-parameters: $k, C, \mu_M, X_M, \mu_A, X_A$;
 - Initialise local meta-parameters of the first cluster: $\Xi^1, \mathbf{f}_M^1, X_M^1, \mathbf{f}_A^1, X_A^1, S^1$;
- * **Else**
 - Update μ_M, X_M, μ_A and k using equation (13);
 - **If** (Condition A is met) **Then**
 1. $C \leftarrow C + 1$;
 2. Initialise local meta-parameters of the new cluster: $\Xi^C, \mathbf{f}_M^C, X_M^C, \mathbf{f}_A^C, X_A^C, S^C$;
 - **Else**
 1. Find the nearest cluster Ξ^n using equation (16);
 2. Update the meta-parameter of this cluster using equation (17): $\Xi^n, \mathbf{f}_M^n, X_M^n, \mathbf{f}_A^n, S^n$.
 - **End If**
- * **End If**

ii. End While

iii. Find the neighbouring clusters Ξ_n^i for each existing cluster Ξ^i using equation (18) ($i = 1, 2, \dots, C$).

iv. Calculate the weighted unimodal densities at the centres of the clusters using equation (19);

v. Identify the local maxima of the weighted unimodal density using equation (20);

vi. Assign each data sample to the cluster with the nearest centre using equation (16).

5 Numerical examples and analysis

In this section, a number of numerical experiments are conducted to demonstrate the performance of the newly proposed direction-aware distance for high dimensional problems. Analysis based on the numerical examples will be provided.

Firstly, we use the standard k-means algorithm as a benchmark. We consider the following problems to test the performance of the k-means algorithm with different type of distance/similarity including Euclidean distance, cosine similarity, cityblock distance and the proposed direction-aware distance:

1. Dim256 dataset (22);
2. Dim512 dataset (22);
3. Dim1024 dataset (22);
4. Dim15 dataset (22);
5. Steel plate faults dataset (23);
6. Pen-based recognition of handwritten digits dataset (24);
7. Optical recognition of handwritten digits dataset (25);
8. Cardiotocography dataset (26);

The dim256, dim512, dim1024 and dim15 datasets are sampled from Gaussian distributions, and, thus, the four datasets are ideal for testing the ability of the algorithms in separating high dimensional data samples from different classes. The other five datasets are real benchmark problems and we use them to evaluate the performance of the algorithms on real, non-Gaussian problems. The details of the benchmark datasets are given in Table 1.

Because of the complexity of the high-dimensional problems, the clustering results of the k-means algorithm may exhibit some degree of randomness, for each dataset and each type of distance/similarity, we did 100 Monte Carlo experiments and tabulated the average values of the five different measures in Table 2. The algorithms used in this paper were implemented within MATLAB 2015b; the

Table 1 Details of the datasets

| Abbreviation | Dataset | Samples | Classes | Attributes |
|--------------|-----------------------|---------|---------|----------------|
| D256 | dim256 | 1024 | 16 | 256 + 1 label |
| D512 | dim512 | 1024 | 16 | 512 + 1 label |
| D1024 | dim1024 | 1024 | 16 | 1024 + 1 label |
| D15 | dim15 | 10125 | 9 | 15 + 1 label |
| ST | Steel plates faults | 1941 | 7 | 27 + 1 label |
| PE | Pen-based recognition | 10992 | 10 | 16 + 1 label |
| OP | Optical recognition | 5620 | 64 | 64 + 1 label |
| CA | Cardiotocography | 2126 | 3 | 22 + 1 label |

performance was evaluated on a PC with dual core Intel i7 processor with clock frequency 3.4 GHz each and 16 GB RAM. In the experiment, without loss of generality, the pair of the scaling parameters of the direction-aware distance is set by Eq. (12) and we consider the Calinski-Harabasz (*CH*) index (Caliński and Harabasz 1974) to evaluate the quality of the clustering results. Higher Calinski-Harabasz (*CH*) index indicates a better clustering quality.

As we can see from Table 2, in the previous section, the performance of the k-means algorithm is largely influenced by the choice of the type of distance/similarity. Based on the Calinski-Harabasz (*CH*) indexes of the clustering results, one can see that the k-means algorithm with the proposed direction-aware distance can produce higher quality clusters compared with the one with traditional distances/dissimilarities.

Then, numerical experiments for the same benchmark problems as tabulated in Table 1 are conducted to evaluate the performance of the evolving algorithm employing the direction-aware distance. To better demonstrate the performance of the evolving algorithm using the direction-aware distance, we involve the following algorithms for comparison:

1. Subtractive clustering algorithm (Chiu 1994);
2. Mean-shift clustering algorithm (Comaniciu and Meer 2002);
3. DBScan clustering algorithm (Ester et al. 1996);
4. Mode identification based clustering algorithm (Li et al. 2007);
5. Random swap algorithm (Franti et al. 2008);
6. Density peak algorithm (Rodriguez and Laio 2014).

As the k-means algorithm exhibits certain degree of randomness, we exclude it from the comparison. In the experiments, due to the insufficient prior knowledge, we use the recommended settings of the free parameters from the published literature. The experimental setting of the free parameters of the algorithms are presented in Table 3.

To objectively compare the performance of different algorithms, we consider the following measures:

1. Number of clusters (*C*), which should be equal or larger than the number of classes in the dataset. However, if *C* is too large (in our paper, we consider $C > 0.1 \times \text{Number of Samples}$ as too large) or is smaller than the number of classes in the dataset, the clustering result should be considered as an invalid one. The former case indicates that there are too many trivial clusters generated which are hard for users to understand. The latter case implies that the clustering algorithm fails to separate the data samples from different classes.

Table 2 Experimental results

| Dataset | Distance/dissimilarity | <i>CH</i> | Dataset | Distance/dissimilarity | <i>CH</i> |
|---------|------------------------|-------------------|---------|------------------------|-----------------|
| D256 | Euclidean | 405.2386 | ST | Euclidean | 20.2314 |
| | Cosine | 448.0036 | | Cosine | 21.769 |
| | Cityblock | 424.2804 | | Cityblock | 17.4560 |
| | Direction-aware | 509.2634 | | Direction-aware | 25.8675 |
| D512 | Euclidean | 373.8111 | PE | Euclidean | 575.0739 |
| | Cosine | 405.8308 | | Cosine | 609.6965 |
| | Cityblock | 410.8807 | | Cityblock | 487.6149 |
| | Direction-aware | 802.3132 | | Direction-aware | 633.2244 |
| D1024 | Euclidean | 368.2901 | OP | Euclidean | 406.5342 |
| | Cosine | 514.7207 | | Cosine | 418.5355 |
| | Cityblock | 721.6852 | | Cityblock | 361.4222 |
| | Direction-aware | 838.6839 | | Direction-aware | 434.6537 |
| D15 | Euclidean | 30834.3331 | CA | Euclidean | 81.8571 |
| | Cosine | 27464.4951 | | Cosine | 109.6599 |
| | Cityblock | 19788.1358 | | Cityblock | 84.0488 |
| | Direction-aware | 36783.2175 | | Direction-aware | 115.3565 |

The best results in each example are bolded

2. Calinski Harabasz index (*CH*) (Caliński and Harabasz 1974), the higher the Calinski Harabasz index is, the better the clustering result is;
3. Purity (*P*) (Dutta Baruah and Angelov 2012), which is calculated based on the result and the ground truth:

$$P = \frac{\sum_{i=1}^N S_D^i}{K}, \quad (21)$$

where S_D^i is the number of data samples with the dominant class label in the i th cluster. The higher purity the clustering result has, the stronger separation ability the clustering algorithm exhibits.

4. Davies–Bouldin (*DB*) index (Davies and Bouldin 1979), the lower Davies–Bouldin index is, the better the clustering result is.
5. Time: the execution time (in seconds) should be as small as possible.

The experiment results obtained by the proposed evolving algorithm as well as other clustering algorithms are given in Table 4. The clustering results of the dim15, Pen-based recognition and Cardiotocography datasets obtained by the proposed algorithm are depicted in Fig. 1, where dots in different colours represent data samples in different clusters.

From Table 4 one can see that the subtractive clustering algorithm is able to produce high quality clustering results on the datasets with Gaussian distribution. However, for the more complex benchmark datasets, it fails to give valid results. The mean-shift clustering algorithm is one of the most efficient algorithms, but it can only perform high-quality clustering with low dimensional datasets. The DBScan algorithm is very efficient as well, but the quality of its clustering results is very limited in terms of the three clustering quality measures. Mode identification based clustering algorithm is a so-called “non-parametric”

Table 3 Experimental settings of the algorithms

| Algorithm | Free parameter(s) | Experimental setting |
|---------------------|--|--|
| Subtractive | Initial cluster radius, r | $r = 0.3$ (Chiu 1994) |
| Mean-shift | 1. Bandwidth, p 2. Kernel function type | 1. $r = 0.15$ (Dutta Baruah and Angelov 2012) 2. Gaussian kernel |
| DBScan | 1. Cluster radius, r 2. Minimum number of data samples within the radius, m | 1. The value of the knee point of the sorted m -dist graph 2. $m = 4$ (Ester et al. 1996) |
| Mode identification | Grid size | Default (Li et al. 2007) |
| Random swap | Number of class | Number of class (Franti et al. 2008) |
| Density peak | 1. Minimum distance, ρ 2. Local density, δ | 1. Relatively high, ρ 2. High, δ (Rodriguez and Laio 2014) |

Table 4 Experimental results

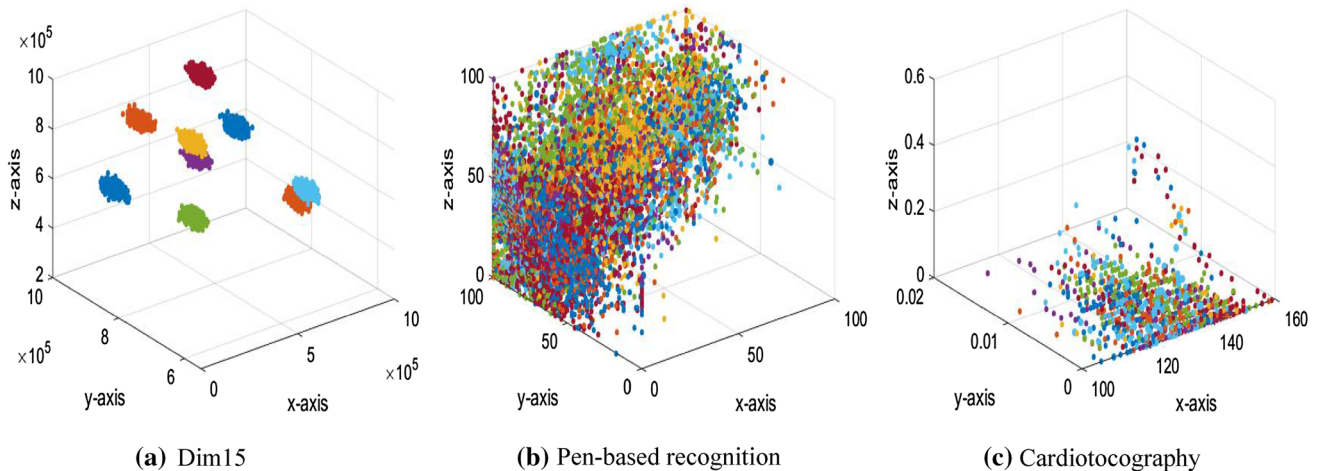
| Dataset | Algorithm | <i>C</i> | <i>CH</i> | <i>P</i> | <i>DB</i> | Time | Validity ^a |
|---------|---------------------|----------------------|--------------------|---------------|---------------|--------------|-----------------------|
| D256 | The proposed | 16 | 203865.1622 | 1.0000 | 0.0248 | 1.61 | O |
| | Subtractive | 16 | 203865.1622 | 1.0000 | 0.0248 | 2.86 | O |
| | Mean-shift | 103 | 44374.6685 | 1.0000 | 0.3728 | 0.19 | O |
| | DBScan | 16 | 173.1715 | 0.7598 | 1.0104 | 0.21 | O |
| | Mode identification | 112 | 41989.1015 | 1.0000 | 0.3736 | 66.68 | × |
| | Random swap | 16 | 1.0259 | 0.1221 | 15.2841 | 16.03 | O |
| | Density peak | 14 | 597.5327 | 0.8750 | 0.6610 | 1.52 | × |
| D512 | The proposed | 16 | 330337.8605 | 1.0000 | 0.0204 | 2.15 | O |
| | Subtractive | 16 | 330337.8605 | 1.0000 | 0.0204 | 4.22 | O |
| | Mean-shift | 149 | 56283.7373 | 1.0000 | 0.3974 | 0.52 | × |
| | DBScan | 16 | 203.2336 | 0.7891 | 1.0046 | 0.32 | O |
| | Mode identification | 1024 | NaN | 1.0000 | 0.0000 | 724.09 | × |
| | Random swap | 16 | 1.1962 | 0.1260 | 15.0519 | 30.76 | O |
| | Density peak | 12 | 291.1243 | 0.7500 | 0.8889 | 1.66 | × |
| D1024 | The proposed | 16 | 718469.7967 | 1.0000 | 0.0132 | 3.66 | O |
| | Subtractive | 16 | 718469.7967 | 1.0000 | 0.0132 | 11.37 | O |
| | Mean-shift | 120 | 126798.4888 | 1.0000 | 0.4496 | 0.88 | × |
| | DBScan | 16 | 381.3919 | 0.8721 | 0.9975 | 0.57 | O |
| | Mode identification | 1024 | NaN | 1.0000 | 0.0000 | 2080.58 | × |
| | Random swap | 16 | 0.9093 | 0.1152 | 16.3316 | 71.11 | O |
| | Density peak | 14 | 529.5497 | 0.8750 | 0.6965 | 3.29 | × |
| D15 | The proposed | 9 | 302436.3684 | 1.0000 | 0.1177 | 13.18 | O |
| | Subtractive | 9 | 302436.3684 | 1.0000 | 0.1177 | 11.28 | O |
| | Mean-shift | 9 | 302436.3684 | 1.0000 | 0.1177 | 0.04 | O |
| | DBScan | 9 | 20602.0570 | 0.9586 | 1.2317 | 10.82 | O |
| | Mode identification | 3 | 4327.2420 | 0.3333 | 0.5837 | 141.34 | O |
| | Random swap | 9 | 126.0758 | 0.2575 | 10.8063 | 7.54 | O |
| | Density peak | 4 | 4533.2627 | 0.4444 | 0.6696 | 12.23 | × |
| ST | The proposed | 23 | 2784.0320 | 0.5064 | 1.8149 | 1.62 | O |
| | Subtractive | 4 | 494.1967 | 0.3988 | 0.9100 | 0.66 | × |
| | Mean-shift | 1555 | 24.7451 | 0.9948 | 9.8535 | 2.92 | × |
| | DBScan | 18 | 57.8279 | 0.48583 | 1.7112 | 0.42 | O |
| | Mode identification | 9 | 690.3357 | 0.3653 | 0.3034 | 69.05 | O |
| | Random swap | 7 | 1.1539 | 0.4096 | 24.1123 | 2.15 | O |
| | Density peak | 3 | 1224.2338 | 0.3478 | 0.4226 | 2.40 | × |
| PE | The proposed | 161 | 572.8011 | 0.9446 | 1.3937 | 10.09 | O |
| | Subtractive | 187 | 382.6055 | 0.8454 | 1.9995 | 12.38 | O |
| | Mean-shift | 8501 | 154.0923 | 0.9999 | 0.3652 | 169.14 | × |
| | DBScan | 38 | 312.9177 | 0.6209 | 1.4997 | 14.04 | O |
| | Mode identification | 4316 | 46.6194 | 0.9968 | 0.4969 | 4243.31 | × |
| | Random swap | 10 | 1.1696 | 0.1160 | 77.2047 | 9.24 | O |
| | Density peak | 7 | 2559.6071 | 0.5993 | 1.3044 | 12.65 | × |
| OP | The proposed | 139 | 80.4085 | 0.9247 | 2.0033 | 17.46 | O |
| | Subtractive | 5620 | NaN | 1.0000 | 0.0000 | 42.07 | × |
| | Mean-shift | No result after 10 h | | | | | × |
| | DBScan | 5 | 80.5137 | 0.2190 | 5.5459 | 3.88 | × |
| | Mode identification | 5620 | NaN | 1.0000 | 0.0000 | 27368.18 | × |
| | Random swap | 10 | 1.7029 | 0.1142 | 31.2458 | 14.35 | O |
| | Density peak | 8 | 71.5796 | 0.2962 | 1.4627 | 6.16 | × |

Table 4 (continued)

| Dataset | Algorithm | <i>C</i> | <i>CH</i> | <i>P</i> | <i>DB</i> | Time | Validity ^a |
|---------|---------------------|----------|-----------------|---------------|---------------|-------------|-----------------------|
| CA | The proposed | 113 | 231.0072 | 0.8758 | 1.0824 | 1.93 | O |
| | Subtractive | 254 | 140.7584 | 0.9147 | 1.3239 | 0.65 | × |
| | Mean-shift | 1594 | 181.2899 | 0.9962 | 0.4175 | 2.91 | × |
| | DBScan | 13 | 35.8486 | 0.8053 | 1.5204 | 0.43 | O |
| | Mode identification | 328 | 63.5207 | 0.9008 | 0.6740 | 40.26 | × |
| | Random swap | 3 | 47.2156 | 0.7785 | 5.2548 | 1.42 | O |
| | Density peak | 3 | 63.5735 | 0.7813 | 0.5081 | 2.71 | O |

The best results in each example are bolded

^a“×” stands for invalid results, “O” stands for valid result

**Fig. 1** Visualization of clustering results

clustering algorithm. Nonetheless, its performance is very limited on high dimensional problems; its computational efficiency is also not very good. The quality of the clustering results obtained by the random swap algorithm is also very limited. In addition, this algorithm requires the number of classes to be known in advance in order to perform valid clustering results; its computational efficiency is also relatively lower. The density peak clustering algorithm is highly efficient, however, based on the recommended input selection, the algorithm failed to separate data samples from different classes in many cases. In addition, with the growth of the number of data samples, the difficulty of deciding the input selection for the users is also increasing.

In contrast, the proposed evolving clustering algorithm consistently produces the top quality clustering results on various problems. In the core of the method is the idea to incorporate into the direction-aware distance both the spatial divergence and the angular similarity. Its computational efficiency does not deteriorate with the increase of dimensionality. Therefore, one can conclude that the proposed evolving clustering algorithm is the top one in the comparison. Nonetheless, we have to admit that the

computational complexity of some clustering algorithms using the direction-aware distance will be inevitably higher compared with the same ones using the traditional distance/dissimilarity.

6 Conclusion

In this paper, a new type of distance, named “direction-aware”, is proposed and proved to be a full metric. The proposed distance is defined as a combination of two components: (1) the traditional Euclidean distance and (2) a cosine similarity based angular/directional divergence. Therefore, it is able to consider both spatial and angular divergences. It is using the advantages of one of them to compensate for the disadvantages of the other. The proposed distance is applicable to various traditional machine learning algorithms as an alternative distance measure. A new direction-aware distance based evolving clustering algorithm is also proposed for streaming data processing. Numerical examples demonstrate that the proposed distance can improve the clustering quality of the k-means algorithm for high

dimensional problems. They also show the validity and effectiveness of the proposed evolving algorithm for handling high dimensional streaming data.

As future work, we will apply the proposed distance to various high dimensional problems including, but not limited to, the NLP, image processing problems, etc. We will also study the convergence of the evolving clustering algorithm.

References

- Aggarwal CC, Hinneburg A, Keim DA (2001) On the surprising behavior of distance metrics in high dimensional space. In: International conference on database theory, pp 420–434
- Allah FA, Grosky WI, Aboutajdine D (2008) Document clustering based on diffusion maps and a comparison of the k-means performances in various spaces. In: IEEE symposium on computers and communications, pp 579–584
- Angelov P, Sadeghi-Tehran P, Ramezani R (2014) An approach to automatic real-time novelty detection, object identification, and tracking in video streams based on recursive density estimation and evolving Takagi–Sugeno fuzzy systems. *Int J Intell Syst* 29(2):1–23
- Angelov P, Gu X, Kangin D (2017) Empirical data analytics. *Int J Intell Syst*. doi:10.1002/int.21899
- Beyer K, Goldstein J, Ramakrishnan R, Shaft U (1999) When is ‘nearest neighbors’ meaningful? In: International conference on database theory, pp 217–235
- Caliński T, Harabasz J (1974) A dendrite method for cluster analysis. *Commun Stat Methods* 3(1):1–27
- Callebaut DK (1965) Generalization of the Cauchy–Schwarz inequality. *J Math Anal Appl* 12(3):491–494
- Cardiotocography Dataset. <https://archive.ics.uci.edu/ml/datasets/Cardiotocography>. Accessed 19 July 2017
- Chiu SL (1994) Fuzzy model identification based on cluster estimation. *J Intell Fuzzy Syst* 2(3):267–278
- Clustering datasets. <http://cs.joensuu.fi/sipu/datasets/>. Accessed 19 July 2017
- Comaniciu D, Meer P (2002) Mean shift: a robust approach toward feature space analysis. *IEEE Trans Pattern Anal Mach Intell* 24(5):603–619
- Davies DL, Bouldin DW (1979) A cluster separation measure. *IEEE Trans Pattern Anal Mach Intell* 2:224–227
- Dehak N, Dehak R, Glass J, Reynolds D, Kenny P (2010) Cosine similarity scoring without score normalization techniques. In: *Proceeding Odyssey 2010—Speaker Language Recognition Work (Odyssey 2010)*, pp 71–75
- Dehak N, Kenny P, Dehak R, Dumouchel P, Ouellet P (2011) Front end factor analysis for speaker verification. *IEEE Trans Audio Speech Lang Process* 19(4):788–798
- Domingos P (2012) A few useful things to know about machine learning. *Commun ACM* 55(10):78–87
- Dutta Baruah R, Angelov P (2012) Evolving local means method for clustering of streaming data. In: *IEEE international conference fuzzy system*, pp 10–15
- Ester M, Kriegel HP, Sander J, Xu X (1996) A density-based algorithm for discovering clusters in large spatial databases with noise. *Int Conf Knowl Discov Data Min* 96:226–231
- Franti P, Virtajoki O, Hautamaki V (2008) Probabilistic clustering by random swap algorithm. In: *IEEE international conference on pattern recognition*, pp 1–4
- Fukunaga K, Hostetler L (1975) The estimation of the gradient of a density function, with applications in pattern recognition. *IEEE Trans Inf Theory* 21(1):32–40
- Keller JM, Gray MR (1985) A fuzzy k-nearest neighbor algorithm. *IEEE Trans Syst Man Cybern* 15(4):580–585
- Li J, Ray S, Lindsay BG (2007) A nonparametric statistical approach to clustering via mode identification. *J Mach Learn Res* 8(8):1687–1723
- Lughofer E, Cernuda C, Kindermann S, Pratama M (2015) Generalized smart evolving fuzzy systems. *Evol Syst* 6(4):269–292
- MacQueen JB (1967) Some methods for classification and analysis of multivariate observations. In: *5th Berkeley symposium mathematical statistics and probability 1967*, vol 1, no 233, pp 281–297
- McCune B, Grace JB, Urban DL (2002) Analysis of ecological communities, vol 28. *MJM Software Design*, Gleneden Beach
- McLachlan GJ (1999) Mahalanobis distance. *Resonance* 4(6):20–26
- Optical Recognition of Handwritten Digits Dataset. <https://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits>. Accessed 19 July 2017
- Pen-Based Recognition of Handwritten Digits Dataset. <http://archive.ics.uci.edu/ml/datasets/Pen-Based+Recognition+of+Handwritten+Digits>. Accessed 19 July 2017
- Precup RE, Filip HI, Radac MB, Petriu EM, Preitl S, Dragoş CA (2014) Online identification of evolving Takagi–Sugeno–Kang fuzzy models for crane systems. *Appl Soft Comput J* 24:1155–1163
- Rodriguez A, Laio A (2014) Clustering by fast search and find of density peaks. *Science* (80-) 344(6191):1493–1496
- Rong HJ, Sundararajan N, Bin Huang G, Saratchandran P (2006) Sequential adaptive fuzzy inference system (SAFIS) for nonlinear system identification and prediction. *Fuzzy Sets Syst* 157(9):1260–1275
- Rong HJ, Sundararajan N, Bin Huang G, Zhao GS (2011) Extended sequential adaptive fuzzy inference system for classification problems. *Evol Syst* 2(2):71–82
- Senoussaoui M, Kenny P, Dumouchel P, Stafylakis T (2013) Efficient iterative mean shift based cosine dissimilarity for multi-recording speaker clustering. In: *IEEE international conference acoustics speech and signal processing*, pp 7712–7715
- Setlur V, Stone MC (2016) A linguistic approach to categorical color assignment for data visualization. *IEEE Trans Vis Comput Graph* 22(1):698–707
- Steel Plates Faults Dataset. <https://archive.ics.uci.edu/ml/datasets/Steel+Plates+Faults>. Accessed 19 July 2017