

# A New Type of Distance Metric and Its Use for Clustering

Xiaowei Gu<sup>1</sup>, Plamen Angelov<sup>1</sup>, Dmitry Kangin<sup>1</sup> and Jose Principe<sup>2</sup>

<sup>1</sup> School of Computing and Communications, Lancaster University Lancaster, LA1 4WA, UK. e-mail: {x.gu3, p.angelov, d.kangin}@lancaster.ac.uk

<sup>2</sup> Computational NeuroEngineering Laboratory, Department of Electrical and Computer Engineering, University of Florida, USA. e-mail: principe@cnel.ufl.edu

**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: *i*) the traditional Euclidean distance and *ii*) 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.

**Index Terms-** 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 [1], [2]. Other metric spaces, using the Mahalanobis [3], 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” [4], [5]. 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 [5], [6]. This phenomenon is frequently seen in 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 [7]–[11]. In NLP problems, machine learning algorithms, for example, *k-means* [7], [10], *mean shift* [11], 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 [12]). Consequently, the cosine similarity between two vectors can be

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, where the weights of the Euclidean and angular components are under the user control. Therefore, it takes the advantages of the both components while still obeys all the properties of a distance metric [13] as we will demonstrate.

The proposed distance in this paper is applicable to various traditional clustering algorithms 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 section 4. Section 5 presents numerical examples. The paper is concluded by section 6.

## 2. Direction-Aware Distance and Proof of Metric Axioms

### A. 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:

- i) a Euclidean component,  $d_M(\mathbf{x}, \mathbf{y})$ , and
- ii) a direction-aware component,  $d_A(\mathbf{x}, \mathbf{y})$ ,

and is expressed as:

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

where  $\mathbf{x} = [x_1, x_2, \dots, x_m]^T$  and  $\mathbf{y} = [y_1, y_2, \dots, y_m]^T$ ,  $\mathbf{x}, \mathbf{y} \in \mathbf{R}^m$ ;  $\lambda_M, \lambda_A$  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:

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

where  $\theta_{xy}$  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:

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

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

### B. Proof of Metric Axioms

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

- i) non-negativity:  $d(\mathbf{x}, \mathbf{y}) \geq 0$ ;

$$ii) \text{ identity of indiscernibles: } d(\mathbf{x}, \mathbf{y}) = 0 \text{ iff } \mathbf{x} = \mathbf{y}; \quad (5)$$

$$iii) \text{ symmetry: } d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x}); \quad (6)$$

$$iv) \text{ triangle inequality: } d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y}) \geq d(\mathbf{x}, \mathbf{y}). \quad (7)$$

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

**Theorem:**  $d_{DA}(\mathbf{x}, \mathbf{y})$  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  $d_{DA}(\mathbf{x}, \mathbf{y})$  obeys the four distance axioms stated in equations (5)-(6) and inequalities (4) and (7) one by one.

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

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

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

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

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

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

**Proof:** For the Euclidean metric, it is true that:

$$d_{DA}(\mathbf{x}, \mathbf{y}) = \sqrt{\lambda_M^2 \|\mathbf{x} - \mathbf{y}\|^2 + \frac{\lambda_A^2}{2} \left\| \frac{\mathbf{x}}{\|\mathbf{x}\|} - \frac{\mathbf{y}}{\|\mathbf{y}\|} \right\|^2} = \sqrt{\lambda_M^2 \|\mathbf{y} - \mathbf{x}\|^2 + \frac{\lambda_A^2}{2} \left\| \frac{\mathbf{y}}{\|\mathbf{y}\|} - \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\|^2} = d_{DA}(\mathbf{y}, \mathbf{x}) \quad (8)$$

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

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

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

$$d_{DA}(\mathbf{x}, \mathbf{z}) > d_{DA}(\mathbf{x}, \mathbf{y}) + d_{DA}(\mathbf{y}, \mathbf{z}) \quad (9)$$

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

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

where,  $\boldsymbol{\chi} = \left[ \lambda_M \mathbf{x}^T, \frac{\lambda_A \mathbf{x}^T}{\sqrt{2} \|\mathbf{x}\|} \right]^T = \left[ \lambda_M x_1, \lambda_M x_2, \dots, \lambda_M x_m, \frac{\lambda_A x_1}{\sqrt{2} \|\mathbf{x}\|}, \frac{\lambda_A x_2}{\sqrt{2} \|\mathbf{x}\|}, \dots, \frac{\lambda_A x_m}{\sqrt{2} \|\mathbf{x}\|} \right]^T$  and

$$\boldsymbol{\psi} = \left[ \lambda_M \mathbf{y}^T, \frac{\lambda_A \mathbf{y}^T}{\sqrt{2} \|\mathbf{y}\|} \right]^T = \left[ \lambda_M y_1, \lambda_M y_2, \dots, \lambda_M y_m, \frac{\lambda_A y_1}{\sqrt{2} \|\mathbf{y}\|}, \frac{\lambda_A y_2}{\sqrt{2} \|\mathbf{y}\|}, \dots, \frac{\lambda_A y_m}{\sqrt{2} \|\mathbf{y}\|} \right]^T.$$

$$\text{Similarly, for } \boldsymbol{\zeta} = \left[ \lambda_M \mathbf{z}^T, \frac{\lambda_A \mathbf{z}^T}{\sqrt{2} \|\mathbf{z}\|} \right]^T, \text{ we can see that } d_{DA}(\mathbf{x}, \mathbf{z}) = d_M(\boldsymbol{\chi}, \boldsymbol{\zeta}), d_{DA}(\mathbf{y}, \mathbf{z}) = d_M(\boldsymbol{\psi}, \boldsymbol{\zeta}).$$

Considering an auxiliary algebraic data space  $\mathbf{R}^{2m}$ , for  $\boldsymbol{\chi}, \boldsymbol{\psi}, \boldsymbol{\zeta}$ , it follows that:

$$d_M(\boldsymbol{\chi}, \boldsymbol{\zeta}) \leq d_M(\boldsymbol{\chi}, \boldsymbol{\psi}) + d_M(\boldsymbol{\psi}, \boldsymbol{\zeta}) \quad (11)$$

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  $d_{DA}(\mathbf{x}, \mathbf{y})$  always satisfies the triangle inequality:  $d_{DA}(\mathbf{x}, \mathbf{z}) \leq d_{DA}(\mathbf{x}, \mathbf{y}) + d_{DA}(\mathbf{y}, \mathbf{z})$ .

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.

### C. The Property of the Proposed Distance

The proposed direction-aware distance metric is a combination of two components: *i)* the traditional Euclidean distance and *ii)* 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 [7], [10], mean-shift clustering [14], k nearest neighbour classification [15] 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,  $\lambda_M$  and  $\lambda_A$  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  $\lambda_A$ , and vice versa.

#### 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 evolve its system structure and update 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:

- i.  $k \leftarrow 1$ , the current time instance;
- ii.  $C \leftarrow 1$ , the number of existing clusters;
- iii.  $\boldsymbol{\mu}_M \leftarrow \mathbf{x}_1$ , the global mean of  $\{\mathbf{x}\}_k$ ;
- iv.  $X_M \leftarrow \|\mathbf{x}_1\|^2$ , the global average scalar product of  $\{\mathbf{x}\}_k$ ;
- v.  $\boldsymbol{\mu}_A \leftarrow \frac{\mathbf{x}_1}{\|\mathbf{x}_1\|}$ , the global mean of  $\left\{ \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\}_k$ , which is also the normalized global mean of  $\{\mathbf{x}\}_k$ .
- vi.  $X_A \leftarrow \left\| \frac{\mathbf{x}_1}{\|\mathbf{x}_1\|} \right\|^2 = 1$ , 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:

- i.  $\Xi^1 \leftarrow \{\mathbf{x}_1\}$ , the first cluster;
- ii.  $\mathbf{f}_M^1 \leftarrow \mathbf{x}_1$ , the centre of the first cluster, which is also the mean of  $\Xi^1$ ;
- iii.  $X_M^1 \leftarrow \|\mathbf{x}_1\|^2$ , the average scalar product of  $\Xi^1$ ;
- iv.  $\mathbf{f}_A^1 \leftarrow \frac{\mathbf{x}_1}{\|\mathbf{x}_1\|}$ , the normalized mean of  $\Xi^1$ ;
- v.  $X_A^1 \leftarrow 1$ , the normalized average scalar product of  $\Xi^1$ , which is always equal to 1 as well;
- vi.  $S^1 \leftarrow 1$ , 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,  $\boldsymbol{\mu}_M$ ,  $X_M$  and  $\boldsymbol{\mu}_A$  are updated using the following equations [16]:

$$i. \boldsymbol{\mu}_M \leftarrow \frac{k}{k+1} \boldsymbol{\mu}_M + \frac{1}{k+1} \mathbf{x}_{k+1} \quad (13a)$$

$$ii. X_M \leftarrow \frac{k}{k+1} X_M + \frac{1}{k+1} \|\mathbf{x}_{k+1}\|^2 \quad (13b)$$

$$iii. \boldsymbol{\mu}_A \leftarrow \frac{k}{k+1} \boldsymbol{\mu}_A + \frac{1}{k+1} \frac{\mathbf{x}_{k+1}}{\|\mathbf{x}_{k+1}\|} \quad (13c)$$

$$iv. k \leftarrow k+1 \quad (13d)$$

Then, the condition A is checked to see whether the new data sample denoted by  $\mathbf{x}_k$  is associated with a new cluster:

$$\text{Condition A: } IF \left( d_{DA}(\mathbf{x}_k, \boldsymbol{\mu}_M) > \max_{j=1,2,\dots,C} \left( d_{DA}(\mathbf{f}_M^j, \boldsymbol{\mu}_M) \right) \right) OR \left( d_{DA}(\mathbf{x}_k, \boldsymbol{\mu}_M) < \min_{j=1,2,\dots,C} \left( d_{DA}(\mathbf{f}_M^j, \boldsymbol{\mu}_M) \right) \right) \quad (14)$$

THEN ( $\mathbf{x}_k$  creates a new cluster)

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  $\lambda_M$  and the inverse of the average  $d_A$  as  $\lambda_A$ , correspondingly. However, for streaming data processing, it is less efficient to keep all the observed data samples in the memory and recalculate  $\lambda_M$  and  $\lambda_A$  every time when a new data sample is observed. Therefore we introduce the recursive forms for calculating the pair of scaling coefficients as follows [16]:

$$\lambda_M = \frac{1}{\sqrt{\frac{\sum_{i=1}^k \sum_{j=1}^k d_M^2(\mathbf{x}_i, \mathbf{x}_j)}{k^2}}} = \frac{1}{\sqrt{2(X_M - \|\boldsymbol{\mu}_M\|^2)}} \quad (15a)$$

$$\lambda_A = \frac{1}{\sqrt{\frac{\sum_{i=1}^k \sum_{j=1}^k d_A^2(\mathbf{x}_i, \mathbf{x}_j)}{k^2}}} = \frac{1}{\sqrt{1 - \|\boldsymbol{\mu}_A\|^2}} \quad (15b)$$

If condition A is satisfied, a new cluster is added with  $\mathbf{x}_k$  as its centre:

- i.  $C \leftarrow C+1$ , the number of existing clusters;
- ii.  $\Xi^C \leftarrow \{\mathbf{x}_k\}$ , the new cluster;
- iii.  $\mathbf{f}_M^C \leftarrow \mathbf{x}_k$ , the centre of the new cluster/ mean of  $\Xi^C$ ;
- iv.  $X_M^C \leftarrow \|\mathbf{x}_k\|^2$ , the average scalar product of  $\Xi^C$ ;
- v.  $\mathbf{f}_A^C \leftarrow \frac{\mathbf{x}_k}{\|\mathbf{x}_k\|}$ , the normalized centre of the new cluster/ normalized mean of  $\Xi^C$ ;
- vi.  $X_A^C \leftarrow 1$ , the normalized average scalar product of  $\Xi^C$ ,
- vii.  $S^C \leftarrow 1$ , the support of the new cluster.

In contrast, if condition A is not met,  $\mathbf{x}_k$  is assigned to the cluster with the nearest centre, denoted by  $\mathbf{f}_M^n$  as:

$$\mathbf{f}_M^n = \arg \min_{i=1,2,\dots,C} \left( d_{DA}(\mathbf{x}_k, \mathbf{f}_M^i) \right) \quad (16)$$

The meta-parameters of the cluster with the nearest centre are updated as follows [16]:

$$i. \Xi^n \leftarrow \Xi^n \cup \{\mathbf{x}_k\} \quad (17a)$$

$$ii. \mathbf{f}_M^n \leftarrow \frac{S^n}{S^n + 1} \mathbf{f}_M^n + \frac{1}{S^n + 1} \mathbf{x}_k \quad (17b)$$

$$iii. X_M^n \leftarrow \frac{S^n}{S^n + 1} X_M^n + \frac{1}{S^n + 1} \|\mathbf{x}_k\|^2 \quad (17c)$$

$$iv. \mathbf{f}_A^n \leftarrow \frac{S^n}{S^n + 1} \mathbf{f}_A^n + \frac{1}{S^n + 1} \frac{\mathbf{x}_k}{\|\mathbf{x}_k\|} \quad (17d)$$

$$v. S^n \leftarrow S^n + 1 \quad (17e)$$

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

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

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

For each cluster centre,  $\mathbf{f}_M^i$  ( $i = 1, 2, \dots, C$ ), we calculate its weighted unimodal density as [16]:

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

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,  $D^W$ :

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

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 equation (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:  $k, C, \boldsymbol{\mu}_M, X_M, \boldsymbol{\mu}_A, X_A$ ;

- 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  $\mu_M, X_M, \mu_A$  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:  $\Xi^C, f_M^C, X_M^C, f_A^C, X_A^C, S^C$ ;

- **Else**

1. Find the nearest cluster  $\Xi^n$  using equation (16);
2. Update the meta-parameter of this cluster using equation (17):  $\Xi^n, f_M^n, X_M^n, f_A^n, S^n$ .

- **End If**

\* **End If**

**ii. End While**

**iii.** Find the neighbouring clusters  $\{\Xi_n^i\}$  for each existing cluster  $\Xi^i$  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:

- i)* Dim256 dataset [17];
- ii)* Dim512 dataset [17];
- iii)* Dim1024 dataset [17];
- iv)* Dim15 dataset [17];
- v)* Steel plate faults dataset [18];
- vi)* Pen-based recognition of handwritten digits dataset [19];
- vii)* Optical recognition of handwritten digits dataset [20];
- viii)* Cardiotocography dataset [21];

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 5 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 I.

TABLE I. 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	5	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

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 II. The algorithms used in this paper were implemented within MATLAB 2015b; the performance was evaluated on a PC with dual core Intel i7 processor with clock frequency 3.4GHz 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 equation (12) and we consider the Calinski-Harabasz (CH) index [22] to evaluate the quality of the clustering results. Higher Calinski-Harabasz (CH) index indicates a better clustering quality.

TABLE II. Experimental Results

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

As we can see from Table II, 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 I 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:

- i) Subtractive clustering algorithm [23];
- ii) Mean-shift clustering algorithm [14];
- iii) DBScan clustering algorithm [24];
- iv) Mode identification based clustering algorithm [25];
- v) Random swap algorithm [26];
- vi) Density peak algorithm [27].

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 III.

TABLE III. Experimental Settings of the Algorithms

Algorithm	Free Parameter(s)	Experimental setting
Subtractive	initial cluster radius, $r$	$r = 0.3$ [23]
Mean-shift	<i>i</i> ) bandwidth, $p$	<i>i</i> ) $r = 0.15$ [28]
	<i>ii</i> ) kernel function type	<i>ii</i> ) Gaussian kernel
DBScan	<i>i</i> ) cluster radius, $r$	<i>i</i> ) the value of the knee point of the sorted m-dist graph
	<i>ii</i> ) minimum number of data samples within the radius, $m$	<i>ii</i> ) $m=4$ [24]
Mode identification	grid size	Default [25]
Random swap	number of class	number of class [26]
Density peak	<i>i</i> ) minimum distance, $\rho$	<i>i</i> ) relatively high, $\rho$
	<i>ii</i> ) local density, $\delta$	<i>ii</i> ) high, $\delta$ [27]

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

*i*) 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.

*ii*) Calinski Harabasz index ( $CH$ ) [22], the higher the Calinski Harabasz index is, the better the clustering result is;

*iii*) Purity ( $P$ ) [28], 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  $i^{\text{th}}$  cluster. The higher purity the clustering result has, the stronger separation ability the clustering algorithm exhibits.

*iv*) Davies-Bouldin ( $DB$ ) index [29], the lower Davies-Bouldin index is, the better the clustering result is.

*v*) 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 IV. 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.

TABLE IV. Experimental Results

Dataset	Algorithm	$C$	$CH$	$P$	$DB$	Time	Validity <sup>a</sup>
D256	<b>The proposed</b>	16	<b>203865.1622</b>	<b>1.0000</b>	<b>0.0248</b>	1.61	O
	Subtractive	16	<b>203865.1622</b>	<b>1.0000</b>	<b>0.0248</b>	2.86	O
	Mean-shift	103	44374.6685	1.0000	0.3728	<b>0.19</b>	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	<b>The proposed</b>	16	<b>330337.8605</b>	<b>1.0000</b>	<b>0.0204</b>	2.15	O
	Subtractive	16	<b>330337.8605</b>	<b>1.0000</b>	<b>0.0204</b>	4.22	O
	Mean-shift	149	56283.7373	1.0000	0.3974	0.52	×
	DBScan	16	203.2336	0.7891	1.0046	<b>0.32</b>	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	<b>The proposed</b>	16	<b>718469.7967</b>	<b>1.0000</b>	<b>0.0132</b>	3.66	O

	Subtractive	16	<b>718469.7967</b>	<b>1.0000</b>	<b>0.0132</b>	11.37	O
	Mean-shift	120	126798.4888	1.0000	0.4496	0.88	×
	DBScan	16	381.3919	0.8721	0.9975	<b>0.57</b>	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	×
	<b>The proposed</b>	9	<b>302436.3684</b>	<b>1.0000</b>	<b>0.1177</b>	13.18	O
D15	Subtractive	9	<b>302436.3684</b>	<b>1.0000</b>	<b>0.1177</b>	11.28	O
	Mean-shift	9	<b>302436.3684</b>	<b>1.0000</b>	<b>0.1177</b>	<b>0.04</b>	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	×
	<b>The proposed</b>	23	<b>2784.0320</b>	<b>0.5064</b>	1.8149	1.62	O
ST	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	<b>0.3034</b>	69.05	O
	Random swap	5	1.6124	0.4086	13.2612	2.15	O
	Density peak	3	1224.2338	0.3478	0.4226	2.40	×
	<b>The proposed</b>	161	<b>572.8011</b>	<b>0.9446</b>	<b>1.3937</b>	10.09	O
PE	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	<b>9.24</b>	O
	Density peak	7	2559.6071	0.5993	1.3044	12.65	×
	<b>The proposed</b>	139	<b>80.4085</b>	<b>0.9247</b>	<b>2.0033</b>	17.46	O
OP	Subtractive	5620	NaN	1.0000	0.0000	42.07	×
	Mean-shift		No result after 10 hours				×
	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	<b>14.35</b>	O
	Density peak	8	71.5796	0.2962	1.4627	6.16	×
<b>The proposed</b>	113	<b>231.0072</b>	<b>0.8758</b>	1.0824	1.93	O	
CA	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	<b>0.43</b>	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	<b>0.5081</b>	2.71	O

<sup>a</sup> “×” stands for invalid results, “O” stands for valid result

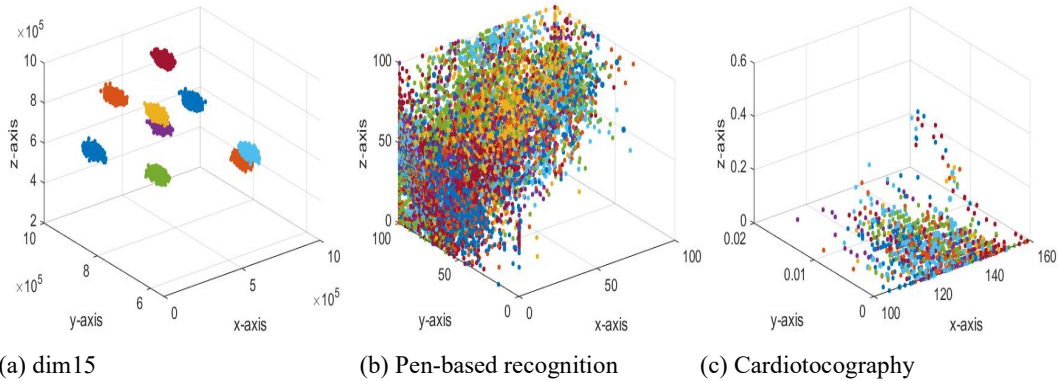


Fig. 1. Visualization of clustering results

From Table IV 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 3 clustering quality measures. Mode identification based clustering algorithm is a so-called “non-parametric” 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. 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.

## 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: *i)* the traditional Euclidean distance and *ii)* 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.

## References

- [1] J. B. MacQueen, “Some methods for classification and analysis of multivariate observations,” *5th Berkeley Symp. Math. Stat. Probab. 1967*, vol. 1, no. 233, pp. 281–297, 1967.
- [2] K. Fukunaga and L. Hostetler, “The estimation of the gradient of a density function, with applications in pattern recognition,” *IEEE Trans. Inf. Theory*, vol. 21, no. 1, pp. 32–40, 1975.
- [3] G. J. McLachlan, “Mahalanobis distance,” *Resonance*, vol. 4, no. 6, pp. 20–26, 1999.
- [4] P. Domingos, “A few useful things to know about machine learning,” *Commun. ACM*, vol. 55, no. 10, pp. 78–87, 2012.
- [5] C. C. Aggarwal, A. Hinneburg, and D. A. Keim, “On the surprising behavior of distance metrics in high dimensional space,” in *International Conference on Database Theory*, 2001, pp. 420–434.
- [6] K. Beyer, J. Goldstein, R. Ramakrishnan, and U. Shaft, “When is ‘nearest neighbors’ meaningful?,” in *International Conference on Database Theory*, 1999, pp. 217–235.
- [7] F. A. Allah, W. I. Grosky, and D. Aboutajdine, “Document clustering based on diffusion maps and a comparison of the k-means performances in various spaces,” in *IEEE Symposium on Computers and Communications*, 2008, pp. 579–584.
- [8] N. Dehak, R. Dehak, J. Glass, D. Reynolds, and P. Kenny, “Cosine Similarity Scoring without Score Normalization Techniques,” *Proc. Odyssey 2010 - Speak. Lang. Recognit. Work. (Odyssey 2010)*, pp. 71–75, 2010.
- [9] N. Dehak, P. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, “Front end factor analysis for speaker verification,” *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 19, no. 4, pp. 788–798, 2011.
- [10] V. Setlur and M. C. Stone, “A linguistic approach to categorical color assignment for data visualization,” *IEEE Trans. Vis. Comput. Graph.*, vol. 22, no. 1, pp. 698–707, 2016.
- [11] M. Senoussaoui, P. Kenny, P. Dumouchel, and T. Stafylakis, “Efficient iterative mean shift based cosine dissimilarity for multi-recording speaker clustering,” *IEEE Int. Conf. Acoust. Speech Signal Process.*, pp. 7712–7715, 2013.
- [12] D. K. Callebaut, “Generalization of the Cauchy-Schwarz inequality,” *J. Math. Anal. Appl.*, vol. 12, no. 3, pp. 491–494, 1965.
- [13] B. McCune, J. B. Grace, and D. L. Urban, *Analysis of Ecological Communities*. 2002.

- [14] D. Comaniciu and P. Meer, "Mean shift: A robust approach toward feature space analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 5, pp. 603–619, 2002.
- [15] J. M. Keller and M. R. Gray, "A fuzzy k-nearest neighbor algorithm," *IEEE Trans. Syst. Man Cybern.*, vol. 15, no. 4, pp. 580–585, 1985.
- [16] P. Angelov, X. Gu, and D. Kangin, "Empirical data analytics," *Int. J. Intell. Syst.*, DOI 10.1002/int.21899, 2017.
- [17] "Clustering datasets," <http://cs.joensuu.fi/sipu/datasets/>.
- [18] "Steel Plates Faults Dataset," <https://archive.ics.uci.edu/ml/datasets/Steel+Plates+Faults>.
- [19] "Pen-Based Recognition of Handwritten Digits Dataset," <http://archive.ics.uci.edu/ml/datasets/Pen-Based+Recognition+of+Handwritten+Digits>.
- [20] "Optical Recognition of Handwritten Digits Dataset," <https://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits>.
- [21] "Cardiotocography Dataset," <https://archive.ics.uci.edu/ml/datasets/Cardiotocography>.
- [22] T. Caliński and J. Harabasz, "A dendrite method for cluster analysis," *Commun. Stat. Methods*, vol. 3, no. 1, pp. 1–27, 1974.
- [23] S. L. Chiu, "Fuzzy model identification based on cluster estimation.," *Journal of intelligent and Fuzzy systems*, vol. 2, no. 3, pp. 267–278, 1994.
- [24] M. Ester, H. P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *International Conference on Knowledge Discovery and Data Mining*, 1996, vol. 96, pp. 226–231.
- [25] J. Li, S. Ray, and B. G. Lindsay, "A nonparametric statistical approach to clustering via mode identification," *J. Mach. Learn. Res.*, vol. 8, no. 8, pp. 1687–1723, 2007.
- [26] P. Franti, O. Virtajoki, and V. Hautamaki, "Probabilistic clustering by random swap algorithm," in *IEEE International Conference on Pattern Recognition*, 2008, pp. 1–4.
- [27] A. Rodriguez and A. Laio, "Clustering by fast search and find of density peaks," *Science*, vol. 344, no. 6191, pp. 1493–1496, 2014.
- [28] R. Dutta Baruah and P. Angelov, "Evolving local means method for clustering of streaming data," *IEEE Int. Conf. Fuzzy Syst.*, pp. 10–15, 2012.
- [29] D. L. Davies and D. W. Bouldin, "A cluster separation measure," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 2, pp. 224–227, 1979.