

A Weighted Distance Metric on Categorical Data

Fuyuan Cao^a, Jie Wen^a

^a*Key Laboratory of Computational Intelligence and Chinese Information Processing of
Ministry of Education, School of Computer and Information Technology, Shanxi University,
Taiyuan 030006, China*

Abstract

Recently most studies of similarity analysis of categorical data assume that attributes are mutual independent. However, in real data sets, attributes are more or less associated with certain relationships. In this paper, we define a weighted dissimilarity measure on categorical objects. It involves the frequency-based weighted intra-coupled distance within an attribute and the weighted inter-coupled distance upon value co-occurrences between attributes. Then we prove the dissimilarity measure is a metric. The experimental results on categorical data have shown the good effect of the proposed distance measure, comparing to other measures.

Keywords: Clustering, Coupled attribute distance, Weighting metric

1. Introduction

The unsupervised learning being one of the machine learning task, which aims to reveal the inherent nature and regularities of the the unlabeled training objects, provides the basis for further data analysis. The most widely used approaches in this learning task is clustering, which useful in pattern-analysis, grouping, decision-making, and machining-learning situations, such as data mining. Data clustering is a data analysis technique and has been considered as a primary data mining method for knowledge discovery [7]. The basic problem involved in the process of data clustering is distance calculation, which is also

*Corresponding author

Email addresses: cfy@sxu.edu.cn (Fuyuan Cao), 1967688145@qq.com (Jie Wen)

called similarity calculation. As the distance between objects becomes larger, the similarity becomes smaller, and vice versa.

In the studies of the similarity between two numerical variables, the definition of distance among objects has Euclidean, Minkowski and so on. This works fine when the defining attributes of a data set are purely numeric in nature. However, these distance measure fails to capture the similarity of data elements when attributes are categorical [1]. Clustering categorical data sets into meaningful groups is a challenging task in which a good distance measure, which can adequately capture data similarities, has to be used in conjunction with an efficient clustering algorithm [1]. The current measure approaches on categorical data are divided into three categories.

- Attribute irrelevance: we assume that attributes are mutual independent in these measures, including simple matching similarity (SMS), which uses 0s and 1s to distinguish the similarity between distinct and identical categorical values, the occurrence frequency (OF) [4] and the information-theoretical similarity (Lin) [4], to discuss the similarity between categorical values.
- Attribute values co-occurrence: These measures compute dissimilarity between two categorical values of same attribute depending upon co-occurrence probabilities of these two values with respect to every other attributes of data set, such as ALGO DISTANCE measure [2] proposed by Ahmad.
- Intra-attribute coupling and Inter-attribute coupling: These methods take into consideration both the frequency-based intra-attribute dissimilarity within an attribute and the inter-attribute dissimilarity between attributes, such as Wang proposed the coupled attribute dissimilarity for objects (CADO) measure [13].

Below, we take a real example to illustrate the challenge of analyzing the categorical data similarity using above measures. A real data application exam-

ple is described in Table 1. As shown in Table 1, six movie objects with three

Table 1: Instances of The Movie Database

<i>Movie</i>	<i>Actor</i>	<i>Genre</i>	<i>Director</i>	<i>Class</i>
Godfaher II	De Niro	Crime	Scorsese	L1
Good Fellas	De Niro	Crime	Coppola	L1
Vertigo	Stewart	Thriller	Hitchcock	L2
Harvey	Stewart	Comedy	Koster	L2
N by NW	Grant	Thriller	Hitchcock	L2
Bishop’s Wife	Grant	Comedy	Koster	L2

categorical attributes are divided into two classes. The SMS measure between actors Stewart and Grant is 0, but Stewart and Grant are very similar. Another observation by SMS is that the similarity between Stewart and Grant is equal to that between De Niro and Stewart. In the same way, the ALGO DISTANCE measure between actors Stewart and Grant is 0.5, which is equivalent to the distance between De Niro and Stewart. However, the similarity of the former pair should be greater because both actors belong of the same class L2.

The above examples show that it is highly complex to analyse the similarity between categorical variables. The SMS and the ALGO DISTANCE fail to capture a global picture of the real relationship for categorical data.

Wang has put forward an effective algorithm [13] (CADO) to improve the shortcomings of the above measures. Wang makes a point about the data-driven intra-coupled similarity and inter-coupled similarity, as well as their global aggregation in unsupervised learning on categorical data.

However, the measure of intra-coupled dissimilarity in CADO does not greatly show the high intra-class similarity and the low inter-class similarity. Moreover, the compute method of weight parameter for attributes is not given in the CADO algorithm. For example, the CADO measure similarity between Godfather II’s De Niro and Good Fellas’s De Niro is 0.5 and so is the similarity between Good Fellas’s De Niro and Vertigo’s Stewart. However, since both actor of the former pair belong to the same class L1, so the similarity should be

greater.

The shortcomings of the above measures are summarized in three aspects.

- The difference of coupling between attributes is not taken into account. Even if the CADO measure considers the inter-attribute coupling, all of the weights for attributes are assumed to be the same.
- All of them don't consider the differences of weight for an attribute when the dissimilarity or similarity between objects is computed. In real data sets, different attributes have their own features, so the contribution of the dissimilarity or similarity among attributes is different.
- For some measures, the intra-attribute dissimilarity between value pairs is independent of the probability of occurrence. For example, SMS and ALGO DISTANCE assume that the value pairs are mutual independent within an attribute.

To sum up, these methods have certain challenges to accurately reflect the dissimilarity or similarity on categorical data.

In this paper, we explicitly define a measure on categorical data. The metric takes into account the characteristics of the categorical values, including the frequency probability within an attribute, the coupling inter-attribute and the feature of every attribute. The key contributions are as follows.

- An intra-attribute weighting scheme for categorical attributes is presented, which assigns different weight according to the different distribution of each attribute's values. The intra-attribute weighting not only takes into account the weight between different attributes, but also takes into consideration the occurrence frequency of the value pairs within an attribute.
- A weighted coupled attribute distance metric between objects (W-CADO) is proposed, which is based on CADO algorithm [13]. By using the intra-attribute weight and the inter-attribute weight [6] in the distance calculation, the comprehensive characteristics between objects are revealed.

This paper is organized as follows. In Section 2, we review some related definitions. The intra-attribute weight and inter-attribute weight are given in Section 3. Section 4 define the weighted intra-coupled distance, the weighted inter-coupled distance, and their integration. We describe the W-CADO algorithm in Section 5. The effectiveness of W-CADO is empirically studied in Section 6. Finally, we conclude this paper in Section 7.

2. Related Definitions

Given a data set $X = \{x_1, x_2, \dots, x_n\}$ with n objects represented by d categorical attributes $\{A_1, A_2, \dots, A_d\}$. Suppose that V_r is a set of attribute values from the attribute A_r ($1 \leq r \leq d$) with m_r possible values and x_{ir} is the value of the object x_i in the attribute A_r and v_r represent the value of any object in the attribute A_r . Obviously, $x_{ir}, v_r \in V_r$.

Definition 1. [11] ($p_{A_r}(x_{ir})$): For any $x_i \in X$, the probability of x_{ir} in the attribute A_r is defined as

$$p_{A_r}(x_{ir}) = \frac{\sigma_{A_r=x_{ir}}(X)}{\sigma_{A_r \neq NULL}(X)}. \quad (1)$$

Here, the operation $\sigma_{A_r=x_{ir}}(X)$ counts the number of objects in X that have the value x_{ir} for the attribute A_r and the symbol $NULL$ refers to the empty.

Definition 2. [11] ($p_{A_r}^-(x_{ir})$): The estimated probability of attribute A_r in representing a value equal to x_{ir} in X is defined as

$$p_{A_r}^-(x_{ir}) = \frac{\sigma_{A_r=x_{ir}}(X) - 1}{\sigma_{A_r \neq NULL}(X) - 1}. \quad (2)$$

Obviously, we can obtain $p_{Actor}(Stewart) = \frac{1}{3}$, $p_{Actor}^-(Stewart) = \frac{1}{5}$ from Table 1.

Definition 3. [13] (ICP): Given a subset $V_r' \subseteq V_r$ having m_r' possible values in the attribute A_r , and a value $v_l \in V_l$ in the attribute A_l , then the information conditional probability (ICP) of V_r' with respect to v_l is defined as

$$P_{A_r|A_l}(V_r'|v_l) = \frac{\sum_{v_r \in V_r'}^{m_r'} \sigma_{A_r=v_r \wedge A_l=v_l}(X)}{\sigma_{A_l=v_l}(X)}. \quad (3)$$

Intuitively, when given all the objects with the value v_l in the attribute A_l , ICP is the percentage of common objects whose values of attribute A_r fall in subset V_r' and whose values of attribute A_l are exactly v_l as well. Hence, ICP quantifies the relative overlapping ratio of attribute values in terms of objects. for example, $P_{Actor|Genre}(\{Grant\}|Thriller) = 0.5$.

Definition 4. [13] *The inter-coupled relative similarity based on intersection set between values x_{ir} and x_{jr} of attribute A_r based on another attribute A_l is defined as*

$$\lambda_{A_r|A_l}(x_{ir}, x_{jr}, V_l) = \sum_{v_l \in \mathcal{X}} \min\{P_{A_l|A_r}(\{v_l\}|x_{ir}), P_{A_l|A_r}(\{v_l\}|x_{jr})\}, \quad (4)$$

where $v_l \in \mathcal{X}$ denote $v_l \in$ the intersection of the set of value in attribute A_l corresponds with x_{ir} in the attribute A_r and the set of value in attribute A_l corresponds with x_{jr} in the attribute A_r . In [13], Wang consider four measures for the inter-coupled similarity to calculate the similarity between two categorical values by considering their relationships with other attributes in terms of power set, universal set, joint set, and intersection set. Wang reveals the equivalent accuracy and superior efficiency of the measure based on the intersection set by theoretical analysis.

3. New weight-computing method for Categorical data

In [13], the attribute couplings include intra-attribute coupling and inter-attribute coupling. Since Wang doesn't consider the difference among attributes and the coupling between attributes, we propose a new weight-computing method. The weight for intra-attribute and inter-attribute are formalized and exemplified below.

3.1. Intra-attribute Weight

Recently most studies of similarity analysis of categorical data treat each attribute equally in data sets. However, it's not always reasonable in real data sets. As we know, unusual features generally can provide more information for

the comparison between objects, so when we compare two objects, we usually pay more attention to the special features they have [6]. In other words, different attribute features should have different contributions in the distance calculation, furthermore, different value pairs within an attribute should have different weights. Considering this phenomenon, we can further adjust the weight according to the following criterion.

The contribution of the distance between two attribute values to the whole object distance is inverse to the probability of these two values' situation in the whole data set. That is, if two data objects have different values along one attribute, the greater the probability that two data objects have different values along this attribute in the data set, the less contribution of the distance between these two values to the entire data distance, and vice versa. What's more, distance metric should assign different weights according to different attribute features. According to the situation that two objects have the same or different value on an attribute, we regard the probability that have the same values or different values from this attribute as the weight of the attribute.

For an attribute A_r , the probability that two data objects from X have the same values along A_r is calculated by

$$p_s(A_r) = \sum_{x_{*r} \in V_r}^{m_r} p_r(x_{*r})p_r^-(x_{*r}). \quad (5)$$

For example, $p_s(Actor) = \frac{1}{5}$, $p_s(Director) = \frac{2}{15}$.

Correspondingly, the probability that two data objects from X have different values along A_r is given by

$$p_f(A_r) = 1 - p_s(A_r). \quad (6)$$

Subsequently, following the proposed criterion, the weight of attribute A_r should be:

$$\eta(A_r) = \begin{cases} p_s(A_r), & \text{if } x_{ir} = x_{jr} \\ p_f(A_r), & \text{otherwise} \end{cases} \quad (7)$$

In a word, $\eta(A_r)$ shows the special features of attribute A_r among attributes.

For the two values x_{ir} and x_{jr} of the attribute A_r , the distance weight between these two values is calculated by

$$\theta(x_{ir}, x_{jr}) = \frac{\sigma_{A_r=x_{ir}}(X)}{\sigma_{A_r \neq NULL}(X)} * \frac{\sigma_{A_r=x_{jr}}(X)}{\sigma_{A_r \neq NULL}(X)}. \quad (8)$$

Alternatively, we could consider other forms of distance weight between two values of attribute A_r according to the data structure, such as $\theta(x_{ir}, x_{jr}) = \alpha \cdot \frac{\sigma_{A_r=x_{ir}}(X)}{\sigma_{A_r \neq NULL}(X)} + \gamma \cdot \frac{\sigma_{A_r=x_{jr}}(X)}{\sigma_{A_r \neq NULL}(X)}$ or $\theta(x_{ir}, x_{jr}) = -\alpha \cdot \log \frac{\sigma_{A_r=x_{ir}}(X)}{\sigma_{A_r \neq NULL}(X)} - \gamma \cdot \log \frac{\sigma_{A_r=x_{jr}}(X)}{\sigma_{A_r \neq NULL}(X)}$, where $0 \leq \alpha, \gamma \leq 1 (\alpha + \gamma = 1)$ are the corresponding weights. $\theta(x_{ir}, x_{jr})$ show the difference between two values of attribute A_r .

Subsequently, the intra-attribute weight is between value x_{ir} and x_{jr} for attribute A_r is

$$\omega(A_r, x_{ir}, x_{jr}) = \eta(A_r) * \theta(x_{ir}, x_{jr}). \quad (9)$$

Alternatively, we could consider other combination forms of η and θ according to the data structure, such as $\omega(A_r, x_{ir}, x_{jr}) = \alpha \cdot \eta(A_r) + \gamma \cdot \theta(x_{ir}, x_{jr})$, where $0 \leq \alpha, \gamma \leq 1 (\alpha + \gamma = 1)$ are the corresponding weights. Thus, η and θ can be controlled flexibly to display in which cases the former is more significant than the latter, and vice versa.

The intra-attribute weighting can adjust the contribution of distance along each attribute to the whole object distance. Moreover, the weight reflects the distance between different value pairs within an attribute on the basis of the occurrence frequency.

3.2. Inter-attribute Weight

Most existing distance or similarity measures for categorical data assume that each attribute is independent. However, in real data, we often have some attributes that are highly dependent on each other. Therefore, the computation of similarity or distance for categorical attribute should be considered based on frequently co-occurring items [5]. That is, the similarity between two values from one attribute should be calculated by considering the other attributes that are highly correlated with this one. In order to utilize the useful relationship information accompanying with each pair of attributes well, the interdependence

redundancy measure [3] has been introduced to evaluate the dependence degree between different attributes. Subsequently, the distance between two values from one attribute is measured not only by their own frequency probabilities but also by the values of other attributes that are highly relevant to this one. In particular, given the data set X , the dependence degree between each pair of attributes A_r and A_l ($r, l \in \{1, 2, \dots, d\}$) can be quantified based on the mutual information [10] between them, which is defined as

$$I(A_r; A_l) = \sum_{x_{ir} \in V_r}^{m_r} \sum_{x_{jl} \in V_l}^{m_l} p(x_{ir}, x_{jl}) \log\left(\frac{p(x_{ir}, x_{jl})}{p_{A_r}(x_{ir})p_{A_l}(x_{jl})}\right). \quad (10)$$

Here, the items $p_{A_r}(x_{ir})$ and $p_{A_l}(x_{jl})$ stand for the frequency probability of the two attribute values in the while data set.

The expression $p(x_{ir}, x_{jl})$ is to calculate the joint probability of these two attribute values, i.e., the frequency probability of objects in X having $A_i = x_{ir}$ and $A_j = x_{jl}$, which is given by

$$p(x_{ir}, x_{jl}) = p(A_r = x_{ir} \wedge A_l = x_{jl} | X) = \frac{\sigma_{A_r=x_{ir} \wedge A_l=x_{jl}}(X)}{\sigma_{A_r \neq NULL \wedge A_l \neq NULL}(X)}. \quad (11)$$

The mutual information between the two attributes actually measures the average reduction in the uncertainty of an attribute by learning the value of another attribute [10]. A larger value of mutual information usually indicates a greater dependency. However, the disadvantage of using this index is that its value increase with the number of possible values that can be chosen by each attribute. Therefore, Au et al. [3] proposed to normalize the mutual information with a joint entropy, which yields the interdependence redundancy measure denoted as

$$R(A_r; A_l) = \frac{I(A_r; A_l)}{H(A_r; A_l)}. \quad (12)$$

where the joint entropy $H(A_r, A_l)$ is calculated by

$$H(A_r; A_l) = - \sum_{x_{ir} \in V_r}^{m_r} \sum_{x_{jl} \in V_l}^{m_l} p(x_{ir}, x_{jl}) \log(p(x_{ir}, x_{jl})). \quad (13)$$

This interdependence redundancy measure evaluates the degree of deviation from independence between two attributes [3]. In particular, $R(A_r; A_l) = 1$

means that the attributes A_r and A_l are strictly dependent on each other while $R(A_r; A_l) = 0$ indicates that they are statistically independent. If the value of $R(A_r; A_l)$ is between 0 and 1, we can say that these two attributes are partially dependent. Since the number of attribute values has no effect on the result of independence redundancy measure, it is perceived as a more ideal index to measure the dependence degree between different categorical attributes.

In the process of experiments, we maintain a $d * d$ relationship matrix ξ to store the dependence degree of each pair of attributes [6]. Each element $\xi(r, l)$ of this matrix is given by $\xi(r, l) = R(A_r; A_l)$. It is obvious that ξ is a symmetric matrix with all diagonal elements equal to 1. To consider the independent attributes simultaneously in distance measure, for each attribute A_r , we find out all the attributes that have obvious interdependence with it and store them in a set denoted as S_r [6]. In particular, the set S_r is constructed by

$$S_r = \{A_l | R(A_r; A_l) > \beta, 1 \leq l \leq d\}. \quad (14)$$

where β is a specific threshold.

4. Coupled Attribute Distance

4.1. The Weighted Intra-Coupled Distance

According to CADO algorithm [13], the intra-coupled attribute similarity for values (IaASV) between values x_{ir} and x_{jr} for attribute A_r is

$$\delta_{A_r}^{IaASV}(x_{ir}, x_{jr}) = \frac{\sigma_{A_r=x_{ir}}(X) \cdot \sigma_{A_r=x_{jr}}(X)}{\sigma_{A_r=x_{ir}}(X) + \sigma_{A_r=x_{jr}}(X) + \sigma_{A_r=x_{ir}}(X) \cdot \sigma_{A_r=x_{jr}}(X)}. \quad (15)$$

For example, in Table 1, we have $\delta_{Actor}^{IaASV}(Stewart, DeNiro) = \delta_{Actor}^{IaASV}(DeNiro, DeNiro) = 0.5$ since both De Niro and Stewart appear twice.

However, the measure of intra-coupled similarity in CADO algorithm does not show the similarity in the same class and the dissimilarity between different classes. For instance, the similarity of the Godfather II's De Niro and Good Fellas's De Niro should be greater than the Good Fellas's De Niro and Harvey's

Stewart because Godfather's Actor and Good Fellas's Actor belong to the same class L1.

Here, Wang considers $h_1(t) = 1/t - 1$ to reflect the complementarity between similarity and dissimilarity measures. We adopt the same complementarity in the algorithm proposed in this paper. In order to overcome the above shortcomings of CADO algorithm, we use the intra-attribute weight that is described in Section 3 to calculate the inter-coupled distance. Subsequently, the weight intra-coupled attribute distance for values (W-IaADV) between values x_{ir} and x_{jr} for attribute A_r is

$$\delta_{A_r}^{W-IaADV}(x_{ir}, x_{jr}) = \omega(A_r, x_{ir}, x_{jr}) * \left(\frac{1}{\delta_{A_r}^{IaASV}(x_{ir}, x_{jr})} - 1 \right). \quad (16)$$

For example, $\delta_{Actor}^{W-IaADV}(DeNiro, DeNiro) = \frac{1}{45}$, $\delta_{Actor}^{W-IaADV}(DeNiro, Stewart) = \frac{4}{45}$. They correspond to the fact that the distance between Good Fellas's De Niro and Vertigo's Stewart is larger than the distance between Godfather II's De Niro and Good Fellas's De Niro.

4.2. The Weighted Inter-Coupled Distance

According to CADO algorithm, the inter-coupled attribute similarity for values (IeASV) between attribute value x_{ir} and x_{jr} of attribute A_r is

$$\delta_{A_r}^{IeASV}(x_{ir}, x_{jr}, \{V_l\}_{l \neq r}) = \sum_{l=1, l \neq r}^d \alpha_l \lambda_{A_r|A_l}(x_{ir}, x_{jr}, V_l). \quad (17)$$

where α_l is the weight parameter for attribute A_l . In CADO algorithm, Wang assign $\alpha_l = \frac{1}{d-1}$. Here, Wang consider $h_2(t) = 1 - t$ to reflect the complementarity between similarity and dissimilarity measures.

However, this assignment method does not take into account the degree of correlation between the different attributes. To overcome the shortcomings of CADO algorithm, we use the relationship matrix that is described in Section 3. Subsequently, the weighted inter-coupled attribute similarity for values (W-IeASV) between values x_{ir} and x_{jr} for attribute A_r is

$$\delta_{A_r}^{W-IeASV}(x_{ir}, x_{jr}, \{V_l\}_{l \neq r}) = \sum_{l=1, l \neq r}^d \xi(r, l) \lambda_{A_r|A_l}(x_{ir}, x_{jr}, V_l). \quad (18)$$

In order to meet the reflexivity, the weighted inter-coupled attribute distance for values (W-IeADV) between values x_{ir} and x_{jr} for attribute A_r , that is, the convert between similarity and dissimilarity measure, is

$$\delta_{A_r}^{W-IeADV}(x_{ir}, x_{jr}, \{V_l\}_{l \neq r}) = \sum_{l=1, l \neq r}^d \xi(r, l) - \delta_{A_r}^{W-IeASV}. \quad (19)$$

4.3. Coupling Integration

So far, we have built formal definitions for both W-IaADV and W-IeADV measures. The W-IaADV emphasizes the attribute value occurrence frequency, while W-IeADV focuses on the co-occurrence comparison of ICP with inter-coupled relative dissimilarity options. Then, the W-CADV is naturally derived by simultaneously considering both measures.

The W-CADV between attribute values x_{ir} and x_{jr} of attribute A_r is

$$\delta_{A_r}^{W-CADV}(x_{ir}, x_{jr}, \{V_l\}_{l=1}^d) = \delta_{A_r}^{W-IaADV}(x_{ir}, x_{jr}) \cdot \delta_{A_r}^{W-IeADV}(x_{ir}, x_{jr}, \{V_l\}_{l \neq r}). \quad (20)$$

where $V_l(l \neq r)$ is a value set of attribute A_l different from A_r to enable the weight inter-coupled interaction. $\delta_{A_r}^{W-IaADV}$ and $\delta_{A_r}^{W-IeADV}$ are W-IaADV and W-IeADV.

As indicated in Eq.(20), we choose the multiplication of these two components. W-IaADV is associated with the frequency of the value, while W-IeADV reflects the extent of the value difference brought by other attributes, hence intuitively, the multiplication of them indicates the total amount of attribute value difference. Alternatively, we could consider other combination forms of W-IaADV and W-IeADV according to the data structure, such as $\delta_{A_r}^{W-CADV}(x_{ir}, x_{jr}, \{V_k\}_{k=1}^d) = \alpha \cdot \delta_{A_r}^{W-IaADV}(x_{ir}, x_{jr}) + \gamma \cdot \delta_{A_r}^{W-IeADV}(x_{ir}, x_{jr}, \{V_k\}_{k \neq j})$, where $0 \leq \alpha, \gamma \leq 1 (\alpha + \gamma = 1)$ are the corresponding weights. Thus, W-IaADV and W-IeADV can be controlled flexibly to display in which cases the intra-coupled interaction is more significant than the inter-coupled interaction, and vice versa.

5. Weighted Coupled Attribute Distance Algorithm

In previous sections, we have discussed the construction of W-CADV. In this section, a weighted coupled attribute distance between objects (W-CADO) is built based on W-CADV.

Given the data set X , the W-CADO between object x_i and x_j is

$$W - CADO(x_i, x_j) = \sum_{r=1}^d \delta_{A_r}^{W-CADV}(x_{ir}, x_{jr}, \{V_k\}_{k=1}^d). \quad (21)$$

We can prove that the dissimilarity measure $W - CADO(\cdot, \cdot)$ is a distance metric satisfying three properties as follows.

- 1) Nonnegativity: $W - CADO(x_i, x_j) \geq 0$ and $W - CADO(x_i, x_i) = 0$;
- 2) Symmetry: $W - CADO(x_i, x_j) = W - CADO(x_j, x_i)$;
- 3) Triangle inequality: $W - CADO(x_i, x_j) + W - CADO(x_j, x_k) \geq W - CADO(x_i, x_k)$.

Obviously, we can easily prove the first two properties according to the previous description. The triangle inequality as the third property is verified as follows.

Proof 1. *To prove the inequality*

$$W - CADO(x_i, x_j) + W - CADO(x_j, x_k) \geq W - CADO(x_i, x_k),$$

we only need to demonstrate

$$\sum_{r=1}^d \delta_{A_r}^{W-CADV}(x_{ir}, x_{jr}, \{V_l\}_{l=1}^d) + \sum_{r=1}^d \delta_{A_r}^{W-CADV}(x_{jr}, x_{kr}, \{V_l\}_{l=1}^d) \geq \sum_{r=1}^d \delta_{A_r}^{W-CADV}(x_{ir}, x_{kr}, \{V_l\}_{l=1}^d).$$

With Eq.(21), the inequality above can be rewritten as

$$\begin{aligned} & \sum_{r=1}^d (\delta_{A_r}^{W-IaADV}(x_{ir}, x_{jr}) \cdot \delta_{A_r}^{W-IeADV}(x_{ir}, x_{jr}, \{V_l\}_{l \neq r})) \\ & + \sum_{r=1}^d (\delta_{A_r}^{W-IaADV}(x_{jr}, x_{kr}) \cdot \delta_{A_r}^{W-IeADV}(x_{jr}, x_{kr}, \{V_l\}_{l \neq r})) \\ & = \sum_{r=1}^d ((\frac{1}{\sigma_{A_r=x_{ir}}(X)} + \frac{1}{\sigma_{A_r=x_{jr}}(X)}) \cdot \omega(A_r, x_{ir}, x_{jr}) \cdot \delta_{A_r}^{W-IeADV}(x_{ir}, x_{jr}, \{V_l\}_{l \neq r})) \\ & + \sum_{r=1}^d ((\frac{1}{\sigma_{A_r=x_{jr}}(X)} + \frac{1}{\sigma_{A_r=x_{kr}}(X)}) \cdot \omega(A_r, x_{jr}, x_{kr}) \cdot \delta_{A_r}^{W-IeADV}(x_{jr}, x_{kr}, \{V_l\}_{l \neq r})) \\ & \geq \sum_{r=1}^d ((\frac{1}{\sigma_{A_r=x_{ir}}(X)} + \frac{1}{\sigma_{A_r=x_{kr}}(X)}) \cdot \omega(A_r, x_{ir}, x_{kr}) \cdot \delta_{A_r}^{W-IeADV}(x_{ir}, x_{kr}, \{V_l\}_{l \neq r})) \\ & = \sum_{r=1}^d (\delta_{A_r}^{W-IaADV}(x_{ir}, x_{kr}) \cdot \delta_{A_r}^{W-IeADV}(x_{ir}, x_{kr}, \{V_l\}_{l \neq r})) \\ & = \sum_{r=1}^d \delta_{A_r}^{W-CADV}(x_{ir}, x_{kr}, \{V_l\}_{l=1}^d) \end{aligned}$$

The above proof verifies that the triangle inequality property holds on all attribute. It follows that we have $W - CADO(x_i, x_j) + W - CADO(x_j, x_k) \geq W - CADO(x_i, x_k)$. Therefore, the dissimilarity measure $W - CADO(\cdot, \cdot)$ is a distance metric.

Algorithm 1 Weighted Coupled Attribute Distance for Objects (W-CADO)

- 1: **Input:** data set $X = \{x_1, x_2, \dots, x_n\}$.
 - 2: **Output:** $D(x_i, x_j)$ for $i, j \in \{1, 2, \dots, n\}$.
 - 3: Calculate $p_s(A_r)$ and $p_f(A_r)$ for each attribute A_r according to Eq.(5) and Eq.(6).
 - 4: For each pair of attributes $(A_r, A_l) (r, l \in \{1, 2, \dots, d\})$ calculate $R(A_r; A_l)$ according to Eq.(14).
 - 5: Construct the relationship matrix ξ .
 - 6: Get the index set S_r for each attribute A_r by $S_r = \{l | \xi(r, l) > \beta, 1 \leq l \leq d\}$.
 - 7: Choose two objects x_i and x_j from X .
 - 8: Let $D(x_i, x_j) = 0$.
 - 9: **for** attribute $A_r, r = 1$ to d **do**
 - 10: // Compute the weight intra-coupled distance for two attribute values x_{ir} and x_{jr}
 - 11: $W\text{-IaADV} = \delta_{A_r}^{W\text{-IaADV}}(x_{ir}, x_{jr});$
 - 12: //Compute the weight inter-coupled distance for two attribute values x_{ir} and x_{jr}
 - 13: $W\text{-IeADV} = \delta_{A_r}^{W\text{-IeADV}}(x_{ir}, x_{jr}, \{V_l\}_{l \neq r});$
 - 14: //Compute coupled distance between two attribute values x_{ir} and x_{jr}
 - 15: $W\text{-CADV} = W\text{-IaADV} \cdot W\text{-IeADV};$
 - 16: //Compute coupled distance between two objects x_i and x_j
 - 17: $W\text{-CADO} = \text{sum}(W\text{-CADV});$
 - 18: **end for**
 - 19: $D(x_i, x_j) = W\text{-CADO};$
 - 20: return $D(x_i, x_j);$
-

We then design the W-CADO algorithm, given in Algorithm 1, to compute

the coupled object distance.

6. Experiments

To investigate the effectiveness of the distance metric for the categorical data proposed in this paper, we mainly make some experiments on the five UCI data sets, Balloons data set, Soybean-small data set, Zoo data set, Congressional Voting Records data set and Breast Cancer data set. We firstly describe the information of the five data sets. Then five evaluation indexes are introduced. Finally, we show the comparison results of the W-CADO algorithm with the CADO algorithm.

In our experiments, the value of the threshold parameter β in the proposed metric is set equal to the average interdependence redundancy of all attribute pairs [6]. That is, β is calculated by

$$\beta = \frac{1}{d^2} \sum_{r=1}^d \sum_{l=1}^d R(A_r; A_l). \quad (22)$$

6.1. Data Description

The information of the data sets we utilized is as follows.

<i>Data Set</i>	<i>Instance</i>	<i>Attribute</i>	<i>Class</i>
Balloons Data Set	20	4	2
Soybean-small Data Set	47	35	4
Zoo Data Set	101	16	7
Congressional Voting Records Data Set	435	16	2
Breast Cancer Data Set	699	10	2

6.2. Evaluation Indexes

To evaluate the effectiveness of the W-CADO algorithm, we used the following five external criterions: (1) adjusted rand index (ARI) [8], (2) normalized mutual information (NMI) [12], (3) accuracy (AC), (4) precision (PR) and (5)

recall (RE) to compare the obtained cluster of each object with that provided by data label.

As described in the Section 2, X represents a data set, $C = \{C_1, C_2, \dots, C_k\}$ be a clustering result of X , $P = \{P_1, P_2, \dots, P_k\}$ be a real partition in X . The overlap between C and P can be summarized in a contingency table shown in Table 3, where n_{ij} denotes the number of objects in common between P_i and C_j , $n_{ij} = |P_i \cap C_j|$. p_i and c_j are the number of objects in P_i and C_j , respectively.

Table 3: The contingency table.

	C_1	C_2	\dots	$C_{k'}$	$Sums$
P_1	n_{11}	n_{12}	\dots	$n_{1k'}$	p_1
P_2	n_{21}	n_{22}	\dots	$n_{2k'}$	p_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
P_k	n_{k1}	n_{k2}	\dots	$n_{kk'}$	p_k
$Sums$	c_1	c_2	\dots	$c_{k'}$	n

The five evaluation indexes are defined as follows:

$$ARI = \frac{\sum_{i,j} C_{n_{ij}}^2 - [\sum_i C_{p_i}^2 \sum_j C_{c_j}^2] / C_n^2}{\frac{1}{2} [\sum_i C_{p_i}^2 + \sum_j C_{c_j}^2] - [\sum_i C_{p_i}^2 \sum_j C_{c_j}^2] / C_n^2},$$

$$NMI = \frac{\sum_{i=1}^k \sum_{j=1}^{k'} n_{ij} \log(\frac{n_{ij}n}{p_i c_j})}{\sqrt{\sum_{i=1}^k p_i \log(\frac{p_i}{n}) \sum_{j=1}^{k'} c_j \log(\frac{c_j}{n})}},$$

$$AC = \frac{1}{n} \max_{j_1 j_2 \dots j_k \in S} \sum_{i=1}^k n_{ij_i},$$

$$PE = \frac{1}{k} \sum_{i=1}^k \frac{n_{ij_i^*}}{p_i},$$

$$RE = \frac{1}{k'} \sum_{i=1}^{k'} \frac{n_{ij_i^*}}{c_i},$$

where $n_{1j_1^*} + n_{2j_2^*} + \dots + n_{kj_k^*} = \max_{j_1j_2\dots j_k \in S} \sum_{i=1}^k n_{ij_i}$ ($j_1^*j_2^*\dots j_k^* \in S$) and $S = \{j_1j_2\dots j_k : j_1, j_2, \dots, j_k \in \{1, 2, \dots, k\}, j_i \neq j_t \text{ for } i \neq t\}$ is a set of all permutations of $1, 2, \dots, k$. For AC, PE, RE , k is equal to k' in general case. In addition, we consider that the higher the values of ARI , NMI , AC , PE and RE are, the better the clustering solution is.

6.3. Comparisons between CADO Alogrithm and W-CADO Alogrithm

One of the clustering approaches is the k -Mode algorithm, designed to cluster categorical data sets. The main idea of k -Mode is to specify the number of clusters k and then to select k initial modes, followed by allocating every object to the nearest mode. The other is a branch of graph-based clustering, i.e., Spectral Clustering (SC), which makes use of Laplacian Eigenmaps on a distance matrix to perform dimensionality reduction for clustering before the k -means algorithm. Below, we aim to compare the performance of W-CADO against CADO as used in data cluster analysis for further clustering evaluation.

The following tables report the results on five data sets with different scale, ranging from 20 to 699 in the increasing order. For each data, the average performance is computed over 50 tests for k -Mode and SC with distinct start points. Note that the highest measure score of each experimental setting is highlighted in boldface.

Table 4: The Comparison of result on Balloons Data Set

	<i>Algorithm</i>	<i>AC</i>	<i>NMI</i>	<i>ARI</i>	<i>PR</i>	<i>RE</i>
<i>k</i> -Mode	CADO	0.7300	0.3283	0.2280	0.7783	0.8417
	W-CADO	0.7600	0.3999	0.2943	0.8100	0.8333
SC	CADO	0.9200	0.8404	0.7986	0.9500	0.9333
	W-CADO	0.9600	0.9202	0.8993	0.9750	0.9667

As table listed above indicates, the clustering methods with W-CADO, whether KM or SC, outperform those with CADO on both AC, NMI, PR, RE and ARI. The reason is that the weight of the attribute added in our algorithm improves the similarity between similar objects and the differences between dif-

Table 5: The Comparison of result on Soybean-small Data Set

	<i>Algorithm</i>	<i>AC</i>	<i>NMI</i>	<i>ARI</i>	<i>PR</i>	<i>RE</i>
<i>k</i> -Mode	CADO	0.7000	0.6325	0.4422	0.8086	0.6784
	W-CADO	0.7993	0.7509	0.6465	0.88030	0.7842
SC	CADO	0.9894	0.9895	0.9797	0.9954	0.9875
	W-CADO	1.0000	1.0000	1.0000	1.0000	1.0000

Table 6: The Comparison of result on Zoo Data Set

	<i>Algorithm</i>	<i>AC</i>	<i>NMI</i>	<i>ARI</i>	<i>PR</i>	<i>RE</i>
<i>k</i> -Mode	CADO	0.7743	0.5113	0.4820	0.7963	0.5764
	W-CADO	0.8158	0.5623	0.6570	0.8423	0.5764
SC	CADO	0.8574	0.8158	0.7495	0.8335	0.7333
	W-CADO	0.8693	0.7890	0.7334	0.8745	0.7446

Table 7: The Comparison of result on Congressional Voting Records Data Set

	<i>Algorithm</i>	<i>AC</i>	<i>NMI</i>	<i>ARI</i>	<i>PR</i>	<i>RE</i>
<i>k</i> -Mode	CADO	0.7621	0.2675	0.3011	0.7703	0.7375
	W-CADO	0.8336	0.3869	0.4526	0.8387	0.8369
SC	CADO	0.8782	0.4895	0.5710	0.8717	0.8897
	W-CADO	0.8805	0.4994	0.5780	0.8743	0.8927

Table 8: The Comparison of result on Breast Cancer Data Set

	<i>Algorithm</i>	<i>AC</i>	<i>NMI</i>	<i>ARI</i>	<i>PR</i>	<i>RE</i>
<i>k</i> -Mode	CADO	0.7497	0.2010	0.2191	0.8032	0.6516
	W-CADO	0.8550	0.4879	0.5351	0.8514	0.8570
SC	CADO	0.9399	0.6956	0.7729	0.9260	0.9512
	W-CADO	0.9456	0.7126	0.7907	0.9276	0.9667

ferent classes of objects. Moreover, the consideration of a complete inter-coupled interaction leads to the largest improvement on clustering accuracy.

For *k*-Mode, the AC improving rate ranges from 4.0% (Balloons) to 14.2% (Soybean-small). With regard to SC, the AC rate takes the minimal and maximal ratios as 0.61% (Breast Cancer) and 4.3% (Balloons). In short, it can be

seen that the W-CADO algorithm is exactly better than the CADO algorithm. There is a significant observation that SC mostly outperforms k -Mode whenever it has the same distance metric. This is consistent with the finding in [9], indicating that SC very often outperforms k -means for numerical data.

7. Conclusion

We have proposed W-CADO, a weighted coupled attribute distance metric for objects incorporating both weighted intra-coupled attribute distance for values and weighted inter-coupled attribute distance for values based on CADO algorithm. By using the intra-attribute weight, the measure increases the intra-class aggregation and inter-class dissimilarity. Furthermore, the dependence degree between each pair of attribute is showed by the inter-attribute weight. Since considering inter-coupled interaction, W-CADO algorithm has improved the clustering accuracy largely. Experimental results on the five real data sets have shown that the W-CADO algorithm is better than the CADO algorithms in clustering categorical data.

Acknowledgements

This work is supported by the National Natural Science Foundation of China (under grants 61573229, 61473194, 61432011 and U1435212), the Natural Science Foundation of Shanxi Province (under grant 2015011048), the Shanxi Scholarship Council of China (under grant 2016-003) and the National Key Basic Research and Development Program of China (973) (under grant 2013CB329404).

References

- [1] Amir Ahmad and Lipika Dey. A k -mean clustering algorithm for mixed numeric and categorical data. *Data and Knowledge Engineering*, 63(2):503–527, 2007.

- [2] Amir Ahmad and Lipika Dey. *A method to compute distance between two categorical values of same attribute in unsupervised learning for categorical data set*. Elsevier Science Inc., 2007.
- [3] W. H. Au, K. C. Chan, A. K. Wong, and Y. Wang. Attribute clustering for grouping, selection, and classification of gene expression data. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 2(2):83–101, 2005.
- [4] Shyam Boriah, Varun Chandola, and Vipin Kumar. Similarity measures for categorical data: A comparative evaluation. pages 243–254, 2008.
- [5] Venkatesh Ganti, Johannes Gehrke, and Raghu Ramakrishnan. Cactus-clustering categorical data using summaries. pages 73–83, 1999.
- [6] Hong Jia and Yiu Ming Cheung. A new distance metric for unsupervised learning of categorical data. In *International Joint Conference on Neural Networks*, pages 1893–1899, 2014.
- [7] Lazhar Labiod and Youns Bennani. A spectral based clustering algorithm for categorical data with maximum modularity. In *Esann 2011, European Symposium on Artificial Neural Networks, Bruges, Belgium, April 27-29, 2011, Proceedings*, 2011.
- [8] Jiye Liang, Liang Bai, Chuangyin Dang, and Fuyuan Cao. The α -means-type algorithms versus imbalanced data distributions. *IEEE Transactions on Fuzzy Systems*, 20(4):728–745, 2012.
- [9] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.
- [10] MacKay and C Davidj. *Information theory, inference, and learning algorithms*. Cambridge University Press, 2003.
- [11] Tiago R. L. Dos Santos and Luis E. Zrate. Categorical data clustering: What similarity measure to recommend? *Expert Systems with Applications*, 42(3):1247–1260, 2015.

- [12] Alexander Strehl and Joydeep Ghosh. *Cluster ensembles — a knowledge reuse framework for combining multiple partitions*. JMLR.org, 2003.
- [13] Can Wang, Xiangjun Dong, Fei Zhou, Longbing Cao, and Chi Hung Chi. Coupled attribute similarity learning on categorical data. *IEEE Transactions on Neural Networks and Learning Systems*, 26(4):781, 2015.