

Discovery of Tri-Edge Inequality with Several Binary Vector Dissimilarity Measures

Bin Zhang

Departments of Human Genetics and Biostatistics
School of Medicine, UCLA
Los Angeles, CA 90095
binzhang@mednet.ucla.edu

Sargur Srihari

CEDAR, Department of Computer Science and Engineering
State University of New York at Buffalo
Buffalo, NY 14228
srihari@cedar.buffalo.edu

Abstract

In certain spaces using some distance measures, the sum of any two distances is always bigger than the third one. Such a special property is called as the tri-edge inequality (TEI). In this paper, the tri-edge inequality characterizing several binary distance measures is mathematically proven and experimentally verified, and the implications of TEI are discussed.

1 Introduction

Searching in feature spaces based on dissimilarity measures, either metric or non-metric, has many applications in pattern recognition, data mining, and multimedia information retrieval [2]. Efficient search algorithms based on metrics usually take advantage of the essential property of metric, i.e., the triangle inequality [1, 4, 5, 7]. Specifically, all samples to be searched are organized into clusters and the center and the radius of each cluster are recorded. The distance between each sample and its corresponding cluster center is pre-computed as well. In the search process, computation of distances between a query sample and most samples in a cluster can be avoided by applying some reduction rules based on the triangle inequality.

However, we discover that the reduction rules utilized by the aforementioned traditional fast algorithms are inapplicable to the dissimilarity measures, which are characterized by a special property, the tri-edge inequality (TEI). The tri-edge inequality means that in certain spaces using some

distance measures the sum of *any two distances* is always bigger than the third one. Notice that here three edges corresponding to three distances don't necessarily constitute a triangle. In this paper, the tri-edge inequality characterizing several binary distance measures is first mathematically proven and then experimentally verified. We also discuss in detail some important implications of TEI.

2 Binary Vector Dissimilarity Measures

In this section, we introduce the existing eight similarity measures for binary vectors summarized by Tubbs [6], then define the associated dissimilarity measures.

A binary vector Z with N dimensions is defined as:

$$Z = (z_1, z_2, \dots, z_N) \quad (1)$$

where, $z_i \in \{0, 1\}, \forall i \in \{1, 2, \dots, N\}$.

Let Ω be the set of all N -dimensional binary vectors, then the **unit binary vector** $I \in \Omega$ is a binary vector with every element equal to 1. The **complement** of a binary vector $Z \in \Omega$ is given by $\bar{Z} = I - Z$.

We define the **magnitude** of a binary vector Z as the number of 1's in Z , given by $|Z| = \sum_{i=1}^N z_i$.

Given two vectors $X \in \Omega$ and $Y \in \Omega$, we define S_{ij} ($i, j \in \{0, 1\}$) as the number of occurrences of different matches between the components in X and Y :

$$S_{00}(X, Y) = \bar{X} \cdot \bar{Y}, S_{01}(X, Y) = \bar{X} \cdot Y \quad (2)$$

and

$$S_{10}(X, Y) = X \cdot \bar{Y}, S_{11}(X, Y) = X \cdot Y \quad (3)$$

Hamming distance (here equivalent to Euclidean), the most commonly used distance measure for binary vectors, is computed as the sum of S_{10} and S_{01} . Based on S_{00} , S_{01} , S_{10} , and S_{11} , there exist eight measures to evaluate the similarity, $S(X, Y)$, between X and Y , as indicated by Tubbs [6]. The second column in Table 1 shows the definitions of eight similarity measures, i.e., Jaccard-Needham, Dice, Correlation, Yule, Russell-Rao, Sokal-Michener, Rogers-Tanmoto and Kulzinsky.

Measure	$S(X, Y)$	$D(X, Y)$
Jaccard-Needham	$\frac{S_{11}}{S_{11}+S_{10}+S_{01}}$	$\frac{S_{10}+S_{01}}{S_{11}+S_{10}+S_{01}}$
Dice	$\frac{S_{11}}{2S_{11}+S_{10}+S_{01}}$	$\frac{S_{10}+S_{01}}{2S_{11}+S_{10}+S_{01}}$
Correlation	$\frac{S_{11}S_{00}-S_{10}S_{01}}{\sigma}$	$\frac{1}{2} - \frac{S_{11}S_{00}-S_{10}S_{01}}{2\sigma}$
Yule	$\frac{S_{11}S_{00}-S_{10}S_{01}}{S_{11}S_{00}+S_{10}S_{01}}$	$\frac{S_{10}S_{01}}{S_{11}S_{00}+S_{10}S_{01}}$
Russell-Rao	$\frac{S_{11}}{N}$	$\frac{N-S_{11}}{N}$
Sokal-Michener	$\frac{S_{11}+S_{00}}{N}$	$1 - \frac{S_{11}+S_{00}}{N}$
Rogers-Tanmoto	$\frac{S_{11}+S_{00}}{S_{11}+S_{00}+2S_{10}+2S_{01}}$	$\frac{2S_{10}+2S_{01}}{S_{11}+S_{00}+2S_{10}+2S_{01}}$
Rogers-Tanmoto-a	$\frac{S_{11}+S_{00}}{S_{11}+S_{00}+2S_{10}+2S_{01}}$	$\frac{2(N-S_{11}-S_{00})}{2N-S_{11}-S_{00}}$
Kulzinsky	$\frac{S_{11}}{(S_{10}+S_{01})}$	$\frac{S_{10}+S_{01}-S_{11}+N}{(S_{10}+S_{01}+N)}$

Table 1. The binary vector similarity measures and their associated dissimilarity measures.

$$\sigma = ((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}.$$

For each similarity measure $S(X, Y)$ given in Table 1, we define its associated **dissimilarity measure** $D(X, Y)$. The corresponding nine dissimilarities are defined in the third column of Table 1. Notice that the range of each dissimilarity measure has been normalized as $[0, 1]$.

Considering the unequal importance of “zero” matches and “one” matches, we define a weighted version for each dissimilarity measure by weighting S_{00} with β ($0 \leq \beta \leq 1$) [3]. The proposition of weighting “zero” matches resulted from the observation that a bit ‘0’ provides less information of separability than a bit ‘1’ in pattern recognition. In case a bit ‘0’ is more important than a bit ‘1’, the difference between two vectors can be computed by any weighted dissimilarity measure based on their corresponding complement vectors.

3 Tri-Edge Inequality (TEI)

Given two constants λ and M , Θ as a subset of Ω is defined as follows: the magnitude of each vector is less than

or equal to λ and the inner product of a vector with any other vector in the set is less than or equal to M .

$\forall X, Y, U, V, P, Q \in \Theta$, the tri-edge inequality for a dissimilarity measure $D(\cdot, \cdot)$ is expressed as:

$$D(X, Y) + D(U, V) \geq D(P, Q) \quad (4)$$

In the following, we examine the tri-edge inequality of the Sokal-Michener, Rogers-Tanmoto and Russell-Rao measures by checking $2D_{min}(\cdot, \cdot) \geq D_{max}(\cdot, \cdot)$, which is a much stronger condition.

3.1 Sokal-Michener Dissimilarity Measure

When $\beta = 1$, the dissimilarity $D^{sm}(X, Y)$ is the Manhattan metric normalized by the feature-vector dimensionality. In the case of $\beta \neq 1$, the Sokal-Michener measure $D^{sm}(X, Y)$ is not a metric because it only satisfies three properties of a metric, i.e., non-negativity, communitativity, and triangle inequality, and violates reflexivity in general. For convenience, $D^{sm}(X, Y)$ is expressed as

$$D^{sm}(X, Y) = 1 - \beta + \frac{\beta(|X| + |Y|)}{N} - \frac{(1 + \beta)X \cdot Y}{N} \quad (5)$$

$\forall X, Y \in \Theta$, by (5) the following relationships hold:

$$\begin{aligned} D^{sm}(X, Y) &= (1 - \beta) + \frac{\beta(|X| + |Y|)}{N} - \frac{(1 + \beta)X \cdot Y}{N} \\ &\geq (1 - \beta) + \frac{\beta(X \cdot Y + X \cdot Y)}{N} - \frac{(1 + \beta)X \cdot Y}{N} \\ &= (1 - \beta) \left(1 - \frac{X \cdot Y}{N}\right) \\ &\geq (1 - \beta) \left(1 - \frac{M}{N}\right). \end{aligned}$$

Therefore, the lower bound of $D^{sm}(X, Y)$ is given by

$$D_{min}^{sm}(X, Y) = (1 - \beta) \left(1 - \frac{M}{N}\right) \quad (6)$$

The upper bound of $D^{sm}(X, Y)$ is estimated by

$$\begin{aligned} D^{sm}(X, Y) &= (1 - \beta) + \frac{\beta(|X| + |Y|)}{N} - \frac{(1 + \beta)X \cdot Y}{N} \\ &\leq (1 - \beta) + \frac{\beta(|X| + |Y|)}{N} \\ &\leq (1 - \beta) + \frac{2\lambda\beta}{N} \end{aligned}$$

Hereby,

$$D_{max}^{sm}(X, Y) = (1 - \beta) + \frac{2\lambda\beta}{N} \quad (7)$$

By solving the equation, $2D_{min}^{sm}(X, Y) \geq D_{max}^{sm}(X, Y)$, for β , we have,

$$\beta \leq \frac{1}{1 + \frac{2\lambda}{N-2M}} = \beta_t \quad (8)$$

where, $\lambda = \max_{X \in \Theta} \{|X|\}$ and $M = \max_{X, Y \in \Theta} \{X \cdot Y\}$.

Thus, given $M < \frac{N}{2}$, $\forall 0 \leq \beta \leq \beta_t$, the Sokal-Michener dissimilarity measure satisfies the tri-edge inequality (4).

3.2 Rogers-Tanmoto Dissimilarity Measure

Considering the relationship between $D^{sm}(X, Y)$ and $D^{rta}(X, Y)$ given by (??), we can easily obtain the bound of $D^{rta}(X, Y)$, $D_{max}^{rta}(X, Y)$ and $D_{min}^{rta}(X, Y)$, from (7) and (6):

$$D_{min}^{rta}(X, Y) = \frac{2D_{min}^{sm}(X, Y)}{1 + D_{min}^{sm}(X, Y)} \quad (9)$$

$$D_{max}^{rta}(X, Y) = \frac{2D_{max}^{sm}(X, Y)}{1 + D_{max}^{sm}(X, Y)} \quad (10)$$

Similarly, we solve the equation, $2D_{min}^{rta}(X, Y) \geq D_{max}^{rta}(X, Y)$, for β . The detailed solving process is omitted here because of the page limitation. The solutions are presented as follows.

If $\lambda = \frac{N}{2}$, then, the tri-edge inequality is satisfied with $\beta \leq \beta_{t0}$ given by

$$\beta_{t0} = 1 - \frac{N}{3(N - M)} \quad (11)$$

If $\lambda \neq \frac{N}{2}$, the equation has two roots, β_{t1} and β_{t2} :

$$\beta_{t1} = \frac{a + \sqrt{\Delta}}{b}, \beta_{t2} = \frac{a - \sqrt{\Delta}}{b} \quad (12)$$

where, $a = 3N^2 + 2\lambda M - 4NM$, $b = 2(N - 2\lambda)(N - M)$, and $\Delta = (2M(N + \lambda) - N^2)^2 + 8\lambda(2N - 3M)$.

When $\lambda < \frac{N}{2}$, the tri-edge inequality is satisfied with $0 \leq \beta \leq \beta_{t2}$; when $\lambda > \frac{N}{2}$, $\beta \in [0, 1] \cap [\beta_{t2}, \beta_{t1}]$.

3.3 Russell-Rao Dissimilarity Measure

We can easily estimate range of $D_{min}^{rr}(X, Y)$ as $[1 - \frac{M}{N}, 1]$. To satisfy the tri-edge inequality, it is required that $2(1 - \frac{M}{N}) \geq 1$, resulting in $M \leq \frac{N}{2}$.

3.4 Other Measures

As $D^{yu}(X, Y)$, $D^{di}(X, Y)$, and $D^{co}(X, Y)$ violate the triangle inequality, the tri-edge inequality is not applicable to them.

Some simple numeral computations have shown that the tri-edge inequality is not applicable for $D^{rt}(X, Y) |_{\beta \neq 1}$, $D^{jn}(X, Y)$, or $D^{ku}(X, Y)$.

4 Verification of TEI

TEI was verified with the Sokal-Michener dissimilarity measure. Experiments were conducted on two data sets, a training set \mathcal{G} of 14,678 handwritten numerals and a testing set \mathcal{F} of 13,565 handwritten numerals. For each numeral image, the binary 512-dimensional gradient, structure, and concavity features were extracted.

The first experiment was designed to generate the dissimilarity distribution for a particular β . Each template in \mathcal{G} was paired with each in \mathcal{F} and the dissimilarity between two templates in each pair was calculated by the Sokal-Michener dissimilarity measure $D^{sm}(\cdot, \cdot)$ with $\beta = 0.5$. In total, there are $13565 * 14678 = 1.991070e + 08$ dissimilarities. These dissimilarities cluster around a mean of 311. The minimum dissimilarity is 207 and the maximum is 373. Since $2 * 207 > 373$, the TEI property holds for $D^{sm}(\cdot, \cdot)$ with $\beta = 0.5$.

To further examine the TEI property of $D^{sm}(\cdot, \cdot)$, the experiment above was repeated for a number of different values of β , minimum and maximum dissimilarities were recorded each time. Figure 1(a) shows the curves of minimum dissimilarity ($\min\{D^{sm}(\cdot, \cdot)\}$) vs β and maximum dissimilarity ($\max\{D^{sm}(\cdot, \cdot)\}$) vs β . In Figure 1, the curve of $2 * \min\{D^{sm}(\cdot, \cdot)\}$ vs β is also drawn for comparison. It is observed that $2 * \min\{D^{sm}(\cdot, \cdot)\}$ is bigger than $\max\{D^{sm}(\cdot, \cdot)\}$ when $\beta \leq 0.58$. Therefore, under the current experimental settings, $D^{sm}(\cdot, \cdot)$ demonstrates the property of TEI if $\beta \leq 0.58$.

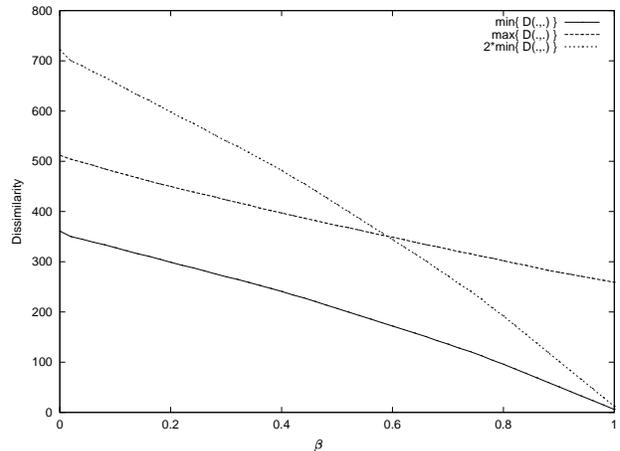


Figure 1. Verification of the TEI property with the Sokal-Michener dissimilarity: $\min\{D^{sm}(\cdot, \cdot)\}$, $\max\{D^{sm}(\cdot, \cdot)\}$, and $2 * \min\{D^{sm}(\cdot, \cdot)\}$ as functions of β .

The experiments also indicate: $M = \max(S_{11}) = 93$ and $\lambda = 216$. By (8), β_t is 0.43. This number is different from the threshold 0.58 shown in Figure ??(a) since the estimated upper bound of $D^{sm}(\cdot, \cdot)$ is bigger than the actual maximum value in the experiment where $\mathcal{F} \cup \mathcal{G}$ is a subset of Θ .

Another experiment was performed to evaluate how the discriminative power (recognition accuracy) of $D^{sm}(\cdot, \cdot)$ varies with β . Each template in \mathcal{F} was recognized by k-nearest neighbor classification over the training set \mathcal{G} . The best performance tops 99.42% when $\beta = 0.42$. The recognition rate drops down to 98.85% when $\beta = 1.0$. When $\beta = 0.5$, the rate reaches 99.38%.

5 Implications of TEI

In this section, we reveal the limitations of the existing fast k -nn search algorithms: the reduction rules utilized by many traditional fast algorithms are inapplicable to the dissimilarity measures, which are characterized by TEI.

The branch and bound algorithm [4] is representative of the traditional fast nearest neighbor search methods. This algorithm applies two rules to reduce the computations, but we show that the reduction rules are not applicable to the nonmetric measures such as $D(\cdot, \cdot)$.

Let n be a node of the search tree representing a set of patterns, T_n , let H_n be the hyper center of the set T_n , let R_n be the radius of the set T_n , and let Y be the current nearest pattern to Z with a distance L .

Then, the first rule shown as Figure 2(a) claims that no $X_i \in T_n$ can be the nearest neighbor to Z if $L + R_n < d(Z, H_n)$. The second rule shown as Figure 2(b) claims that $X_i \in T_n$ cannot be the nearest neighbor of Z if $L + d(X_i, H_n) < d(Z, H_n)$.

For a feature vector set $\Theta \in \Omega$ with β_t , k -nn search is performed based on $D(\cdot, \cdot)$ with β ($0 \leq \beta \leq \beta_t$). Any three of the four dissimilarities, L , R_n , $d(X_i, H_n)$, and $d(Z, H_n)$ in Figure 2, conform to the inequality (4) since they are computed from the vectors in Θ , therefore, the two conditions in the rules, $L + R_n < d(Z, H_n)$ and $L + d(X_i, H_n) < d(Z, H_n)$, can never be satisfied. Thus, no computation of dissimilarities can be pruned by the two rules.

6 Conclusion

An important characteristic, the tri-edge inequality, was discovered and proven with three dissimilarity measures. We demonstrated that the reduction rules, utilized by many traditional fast searching algorithms, are inapplicable to the dissimilarity measures which are characterized by the tri-edge inequality. For dissimilarity measures possessing the

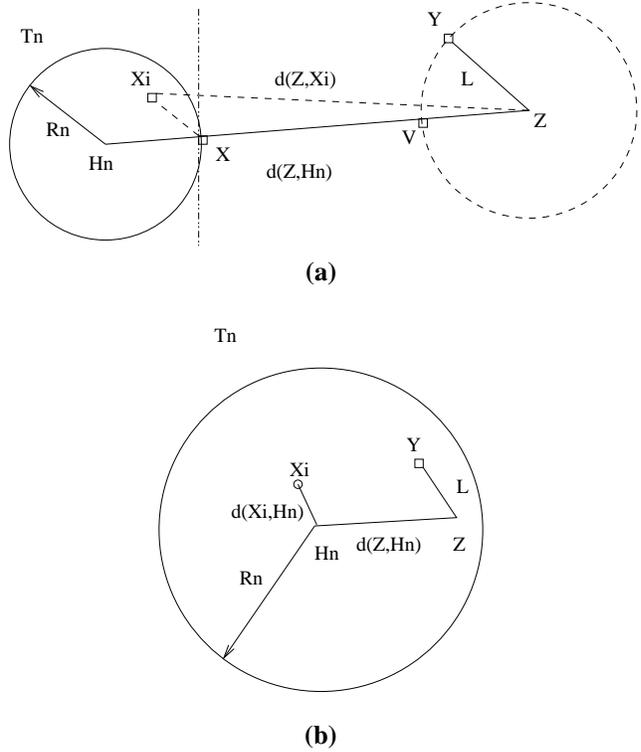


Figure 2. Illustration of the reduction rules

tri-edge inequality, novel fast search algorithms have to be developed without taking advantage of the triangle inequality.

References

- [1] A. J. Broder. Strategies for efficient incremental nearest neighbor search. *Pattern Recognition*, 23(1/2):171–178, Nov 1986.
- [2] E. Chavez, G. Navarro, R. Baeza-yates, and J. L. Marroquin. Searching in metric space. *ACM Computing Surveys*, 33(3):273–321, Sept. 2001.
- [3] J. T. Favata and G. Srikantan. A multiple feature/resolution approach to handprinted character/digit recognition. *International Journal of Imaging Systems and Technology*, 7:304–311, 1996.
- [4] K. Fukunaga and P. Narendra. A branch and bound algorithm for computing k-nearest neighbors. *IEEE Trans. Computers*, C-24(7):750–753, July 1975.
- [5] S. A. Nene and S. K. Nayar. A simple algorithm for nearest-neighbor search in high dimension. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 19(9):989–1003, Sept 1997.
- [6] J. D. Tubbs. A note on binary template matching. *Pattern Recognition*, 22(4):359–365, 1989.
- [7] E. Vidal. An algorithm for finding nearest neighbors in (approximately) constant average time. *Pattern Recognition Letters*, 4(3):145–157, July 1986.