

# Properties of Binary Vector Dissimilarity Measures

Bin Zhang and Sargur N. Srihari

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

## Abstract

*This study is to examine the metric or non-metric properties of binary vector dissimilarity measures, on which no comprehensive research has been conducted so far. Especially, the triangle-inequality invariance characterizing several dissimilarity measures is revealed. Moreover, an entropy-based measurement is proposed to estimate the discriminative capability of a binary vector dissimilarity measure, the effectiveness of the measurement is experimentally justified.*

## 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, 4, 7]. Since efficient search algorithms using metrics are usually not applicable to search a space based on non-metrics, metric/non-metric verification for a dissimilarity measure is critical for developing and choosing the proper fast search algorithms with regard to the measure [1, 3, 5, 6].

Tubbs summarized eight similarity measures for binary feature vectors and generalized these measures by weighting matches based on statistical and spatial information [8]. However, up to now no comprehensive work has been done to verify the properties of binary vector similarity measures. This paper aims at formulating a set of binary vector dissimilarity measures, examining their metric/non-metric properties, and further exploring their discriminative capability in pattern recognition.

The rest of the paper is organized as follows. In Section 2, we formulate eight normalized dissimilarity measures for binary feature vectors based on the existing binary vector similarity measures. In Section 3, we summarize the metric/non-metric properties for each dissimilarity measure. In Section 4, we compare the recognition capability of these measures. Section 5 concludes the paper.

## 2 Binary Vector Dissimilarity Measures

In this section, we give the existing eight similarity measures for binary vectors, then define the associated dissimilarity measures.

### 2.1 Binary Vector Similarity 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  $\Omega$  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  $Y \in \Omega$ , let  $S_{ij}$  ( $i, j \in \{0, 1\}$ ) be the number of occurrences of matches with  $i$  in  $X$  and  $j$  in  $Y$  at the corresponding positions. It is easy to find that  $S_{11}(X, Y) = X \cdot Y$  and  $S_{00}(X, Y) = \bar{X} \cdot \bar{Y}$ .

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 [8]. 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.

From Table 1 we may see that four similarity measures are not normalized, i.e., Dice with a range  $[0, \frac{1}{2}]$ , Correlation with a range  $[-1, 1]$ , Yule with a range  $[-1, 1]$ , and Kulzinsky  $[0, \infty)$ . Moreover, there are some assumptions underlying the ranges given in Table 1: (i) for Jaccard-Needham similarity measure,  $S(\bar{I}, \bar{I}) = 1$ , (ii) for Dice similarity measure,  $S(\bar{I}, \bar{I}) = \frac{1}{2}$ , (iii) for Correlation and Yule similarity measures,  $S(I, I) = S(\bar{I}, \bar{I}) = 1$ , and  $S(I, \bar{I}) = S(\bar{I}, I) = -1$ , and (iv) for Kulzinsky similarity measure,  $S(\bar{I}, \bar{I}) = \infty$ .

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	$\frac{S_{11}+S_{00}}{N}$	$\frac{2S_{10}+2S_{01}}{S_{11}+S_{00}+2S_{10}+2S_{01}}$
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.**

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

## 2.2 Binary Vector Dissimilarity Measures

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  $[0, 1]$ .

To ensure that these dissimilarity measures are valid over the whole domain, some assumptions are made for some dissimilarity functions: (i) for Jaccard-Needham and Dice dissimilarity measures,  $D(\bar{I}, \bar{I}) = 0$ , and (ii) for Correlation and Yule dissimilarity measures,  $D(I, I) = D(\bar{I}, \bar{I}) = 0$  and  $D(I, \bar{I}) = D(\bar{I}, I) = 1$ .

Considering the unequal importance of “zero” matches and “one” matches, we define a weighted version for each dissimilarity measure by weighting  $S_{00}$  with  $\beta$  ( $0 \leq \beta \leq 1$ ). 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 that a bit ‘0’ is more important than bit ‘1’, the difference between two vectors can be computed by any weighted dissimilarity measure based on their corresponding complement vectors.

Notice the following facts about the weighted dissimilarity measures:

- Four measures, i.e., Jaccard-Needham, Dice, Russell-Rao, and Kulzinsky, are independent of  $S_{00}$ . Thus, for these measures weighting  $S_{00}$  makes no difference.
- The Yule dissimilarity measure becomes a constant when the coefficient  $\beta$  equals 0.
- Rogers-Tanmoto-a dissimilarity measure is in fact a

variation of Sokal-Michener dissimilarity measure, the two will behave in the same way.

- Jaccard-Needham measure is similar to the Dice measure.

In the rest of the paper, any dissimilarity measure mentioned refers to the corresponding weighted dissimilarity measure.

## 3 Examination of Metric/Non-metric Properties

A *metric* defined in a set  $S$  is a *distance function*  $d$  which assigns to each pair of points of  $S$  a distance between them, and which satisfies the following four properties: (i) non-negativity:  $d(x, y) \geq 0, \forall x, y \in S$ , (ii) commutativity:  $d(x, y) = d(y, x), \forall x, y \in S$ , (iii) triangle Inequality:  $d(x, z) \leq d(x, y) + d(y, z), \forall x, y, z \in S$ , and (iv) reflexivity:  $d(x, y) = 0 \iff x = y$ .

Obviously all nine dissimilarity measures are negative and commutative. A fact about  $S_{11}$  is that  $S_{11}(X, Y) = X \cdot Y \leq |X|$  and  $S_{11}(X, Y) \leq |Y|, \forall X, Y \in \Omega$ .

Let  $\Theta$  be a subset of  $\Omega$ , defined as:  
 $\Theta = \{X, Y \in \Omega | 0 \leq |X|, |Y| \leq \lambda, 0 \leq X \cdot Y \leq M\}$ .

$\forall X, Y, U, V, P, Q \in \Theta$ , the Triangle-Inequality Invariance (TII) for a dissimilarity measure  $D(\cdot, \cdot)$  is expressed as:

$$D(X, Y) + D(U, V) \geq D(P, Q) \quad (2)$$

By examining the inequality  $2D_{min}(\cdot, \cdot) \geq D_{max}(\cdot, \cdot)$ , three measures, Russell-Rao, Sokal-Michener and Rogers-Tanmoto-a, have been proved to comply with triangle-inequality invariance. Because of the length limitation, the detailed justification process is not included here.

For the Rogers-Tanmoto-a measure, TII holds under the following conditions: (1)  $\beta \leq 1 - \frac{N}{3(N-M)}$  if  $\lambda = \frac{N}{2}$ ; (2)  $\beta \leq \beta_{t2}$  if  $\lambda < \frac{N}{2}$ ; (3)  $\beta \in [0, 1] \cap [\beta_{t2}, \beta_{t1}]$  if  $\lambda > \frac{N}{2}$ . Here  $\lambda = \max\{|X|\}$ ,  $M = \max\{X \cdot Y\}$ , and  $\beta_{t1}$  and  $\beta_{t2}$  are given by

$$\beta_{t1} = \frac{a + \sqrt{\Delta}}{b}, \beta_{t2} = \frac{a - \sqrt{\Delta}}{b} \quad (3)$$

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

Table 2 summarizes the metric, nonmetric and TII properties of the nine measures. Three measures, Sokal-Michener, Rogers-Tanmoto and Rogers-Tanmoto-a, are metric only when  $\beta = 1$ . Dice, Correlation, Yule and Russell-Rao are nonmetric.

Measure	metric	TH
Jaccard-Needham	unknown	no
Dice	no	no
Correlation	no	no
Yule	no	no
Russell-Rao	no	$M \leq N/2$
Sokal-Michener	$\beta = 1$	$\beta \leq \frac{1}{1 + \frac{2\lambda}{N-2M}}$
Rogers-Tanmoto	$\beta = 1$	no
Rogers-Tanmoto-a	$\beta = 1$	yes
Kulzinsky	unknown	no

**Table 2. Metric/nonmetric properties of the binary vector similarity measures. Here,  $\lambda = \max\{|X|\}$ ,  $M = \max\{X \cdot Y\}$  and  $\sigma = ((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}$ .**

## 4 Discriminative Power of Binary Dissimilarity Measures

In this section, we propose an entropy-based method for evaluating dissimilarity measure discriminability and justify it by a computation example and a recognition experiment.

### 4.1 Estimation of Discriminative Power by Entropy

Let  $p(x)$ ,  $x \in [0, \infty)$ , be the distribution of all dissimilarities in a  $N$ -dimensional binary space  $\Omega$  based on a certain dissimilarity measure  $D(\cdot, \cdot)$ , the entropy with the distribution is given by  $H(D) = -\int_0^\infty p(x) \ln p(x) dx$ .

Consider  $p(x) \sim N(\mu, \sigma^2)$ , a Gaussian distribution with mean  $\mu$  and standard deviation  $\sigma^2$ , its entropy is  $\ln[\sqrt{2\pi e}\sigma]$ . As the entropy is directly proportional to  $\sigma$ , the more scattered the dissimilarities, the more informative the distribution. Assume a pattern recognition problem with 2 classes, we can generally associate  $H(D)$  with the discriminative power of  $D(\cdot, \cdot)$  as following: if the dissimilarities cluster more densely around the mean, it would be more difficult for the measure  $D(\cdot, \cdot)$  to differentiate the inter-class dissimilarities from the intra-class dissimilarities, resulting in lower discriminative power; when the dissimilarities tend to spread evenly about the mean, differentiating the inter-class dissimilarities from the intra-class dissimilarities will be possibly more easy, resulting in high discriminative power of  $D(\cdot, \cdot)$ .

In a word, the entropy of the dissimilarity distribution characterizes in some way the discriminative power of a dissimilarity measure. We admit that an even distribution of dissimilarities more likely lead to the even intra- and inter-class distributions so that recognizing patterns from different classes may be difficult too. But here we are not at-

tempting to precisely measure the recognition accuracy but roughly predict the recognition power of similarity measures. It is believed that generally a measure resulting in a very concentrated distribution will not be good at differentiating a large number of pattern belonging to multiple classes.

In the next, we will justify this measurement by an exemplar computation and a real recognition experiment.

### 4.2 Experiments in 10-dimensional Binary Space

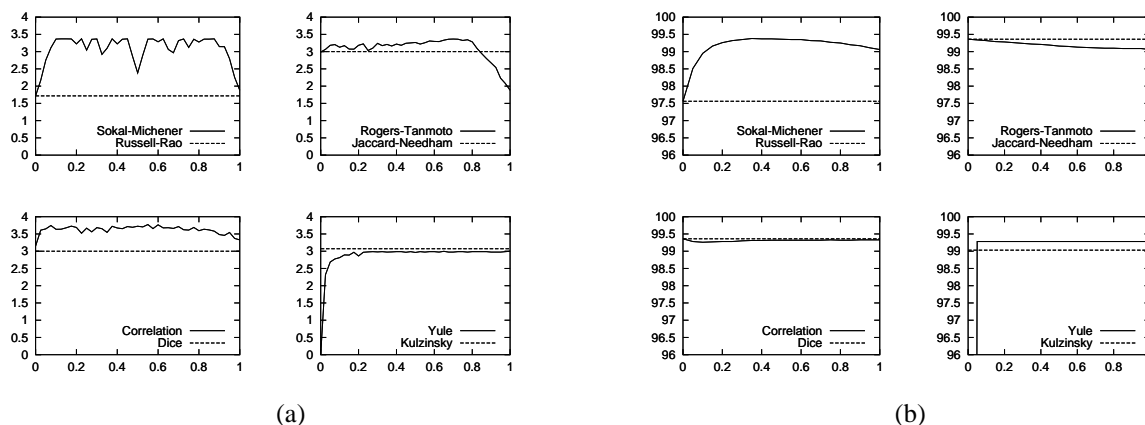
In the 10-dimensional binary space, we compare the entropies of the distributions of all possible dissimilarities using eight measures with  $\beta = 0.5$ . As expected, the distribution from Russell-Rao measure is very concentrated and has the lowest entropy 1.72. Correlation measure presents the high discriminative power of 3.72. Rogers-Tanmoto measures presents a high power of 3.26. Unexpectedly, Sokal-Michener measure has a low discriminative power of 2.38. The other measures present the similar discriminative power around 3.0.

To have a complete view of the discriminative power of each measure in the 10-dimensional binary space, Figure 1(a) shows the discriminative power as a function of the weighting coefficient  $\beta$ . Interestingly, the discriminative power of Sokal-Michener measure reaches one of its minima when  $\beta = 0.5$ .

### 4.3 Handwriting Recognition with 512-dimensional Binary Feature Vectors

To testify further the effectiveness of the entropy-based measurement of discriminative power, another experiment was conducted on handwritten numeral recognition. Gradient-based binary features have been successfully used for recognizing handwritten characters [4]. For an individual character, the gradient-based features consist of 512 bits corresponding to gradient (192 bits), structural (192 bits), and concavity (128 bits) features, known as GSC features. In this experiment, we employed the  $k$ -nearest neighbor classification by using the aforementioned dissimilarity measures. The training set consists of 53449 512-dimensional binary feature vectors extracted from the images in NIST hsf0 data set, and the testing set contains 53313 512-dimensional binary feature vectors from the images in NIST hsf1 data set.

Figure 1(b) shows the recognition performance of each dissimilarity as function of  $\beta$ . As expected, Russell-Rao measure presents the worst recognition accuracy of 97.56%. Sokal-Michener one with  $\beta = 0.35$  achieves the highest rate of 99.38%, the recognition rates based on Rogers-Tanmoto, Correlation, Jaccard-Needham and Dice measures can top 99.36%. Yule measure ( $\beta \neq 0$ ) presents a ac-



**Figure 1. Discriminative Power of eight dissimilarity measures: (a) entropy variations with  $\beta$  in the 10-dimensional binary vector space, (b) variations of recognition accuracy (%) as functions of  $\beta$ .**

accuracy of 99.28% and Kulzinsky measure achieves 99.03% accuracy.

#### 4.4 Analysis

Based on the computation and experiments above, we can rank the eight measurements by the discriminative power as following: (i) Rogers-Tanmoto, Correlation and Sokal-Michener measures are the best, (ii) Russell-Rao measure has the worst power, and (iii) all the other four measure have moderate discriminative power. The similar ranking order has also been obtained based on the handwriting identification task using 5120-dimensional feature vectors[9]. It is expected that in high-dimensional binary space, this conclusion will hold.

#### 5 Conclusion

We systematically examine metric or non-metric properties of several dissimilarity measures for binary vectors and reveal the triangle-inequality invariance characterizing three measures. We also quantitatively analyze the recognition capability of these measures. The theoretic analysis provides solid ground for developing or selecting proper fast search algorithms for nonmetric binary vector similarity measures since most fast search algorithms using metrics usually fail in nonmetric spaces.

#### Acknowledgments

This work was supported in part by the Department of Justice, National Institute of Justice grant 2002-LT-BX-

K007.

#### References

- [1] C. L. Chang. Finding prototypes for nearest neighbor decision rule. *IEEE Trans. Computers*, C-23(11):1179–1184, Nov 1974.
- [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] A. Farago, T. Linder, and G. Lugosi. Fast nearest-neighbor search in dissimilarity spaces. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 15(9):957–962, Sept 1993.
- [4] 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.
- [5] 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.
- [6] B. S. Kim and S. B. Park. A fast k nearest neighbor finding algorithm based on the ordered partition. *IEEE Trans. Pattern Analysis and Machine Intelligence*, pami-8(6):761–766, Nov 1986.
- [7] J. R. Smith and S.-F. Chang. Automated binary texture feature sets for image retrieval. In *Intl Conf Acoust., Speech, Signal Processing*, Atlantic,GA, May 1996.
- [8] J. D. Tubbs. A note on binary template matching. *Pattern Recognition*, 22(4):359–365, 1989.
- [9] B. Zhang and S. N. Srihari. Binary vector dissimilarities for handwriting identification. In *Proceedings SPIE, Document Recognition and Retrieval X*, pages 155–166, Santa Clara, CA, January 2003.