

Binary Vector Dissimilarity Measures for Handwriting Identification

Bin Zhang and Sargur N. Srihari

CEDAR, Computer Science and Engineering Department
State University of New York at Buffalo, Buffalo, NY 14228

ABSTRACT

Several dissimilarity measures for binary vectors are formulated and examined for their recognition capability in handwriting identification for which the binary micro-features are used to characterize handwritten character shapes. Pertaining to eight dissimilarity measures, i.e., Jaccard-Needham, Dice, Correlation, Yule, Russell-Rao, Sokal-Michener, Rogers-Tanmoto and Kulzinsky, the discriminative power of ten individual characters and their combination is exhaustively studied. It is found that different character shapes and different positions of a character in a word bear different capability for identifying handwriting individuality. Moreover, a character alone usually is not informative enough to identify handwriting differences, but a combination of many characters presents good performance in handwriting identification. While five measures, Rogers-Tanmoto, Jaccard-Needham, Dice, Correlation, and Sokal-Michener, consistently perform very well in identification, Kulzinsky and Russell-Rao ones give low identification accuracy.

Keywords: Binary feature vector, Similarity measure, Handwriting identification

1. INTRODUCTION

Binary features have many applications in pattern recognition,¹ image retrieval² and statistical analysis.³ Tubbs summarized eight similarity measures for binary feature vectors and generalized these measures by weighing matches based on statistical and spatial information.⁴ However, up to now no comprehensive work has been done to evaluate the performance of these measures in pattern recognition.

The objectives of this paper are to examine the recognition capabilities of these similarity measures and to draw some general conclusions on how to use binary matches for constructing similarity or dissimilarity measures. Instead of using the similarity measures aforementioned, we formulate eight normalized binary vector dissimilarity measures corresponding to the eight similarity measures for binary vectors. Use of the normalized dissimilarities is necessary in order to match other distances in a system with multiple types of features.

We examine the recognition capability of each dissimilarity through a forensic document examination system, recently developed by CEDAR.⁵ Since this system can extract a very large number (5120) of binary features to characterize the handwriting shapes in a handwritten document and the number of classes in handwriting identification can be very big (we have 3000 handwritten samples by 1000 writers), it is an ideal testing-bed to evaluate these dissimilarity measures.

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 introduce the binary micro-features for individual handwritten characters and give several computation examples. In Section 4, we examine the discriminative power of micro-features and the recognition capability of each dissimilarity measure. We draw conclusions in Section 5.

2. DISSIMILARITY MEASURES FOR BINARY VECTORS

A binary vector Z with N dimensions is defined as:

$$Z = (z_1, z_2, \dots, z_N) \tag{1}$$

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

Further author information: (Send correspondence to Bin Zhang)
By E-mail: binzhang@cedar.buffalo.edu

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

$$\bar{Z} = I - Z \quad (2)$$

We consider eight similarity measures for binary feature vectors summarized by Tubbs,⁴ i.e., Jaccard-Needham, Dice, Correlation, Yule, Russell-Rao, Sokal-Michener, Rogers-Tanmoto and Kulzinsky.

Let S_{ij} ($i, j \in \{0, 1\}$) be the number of occurrences of matches with i in the first pattern and j in the second pattern at the corresponding positions. Given two binary feature vectors $X \in \Omega$ and $Y \in \Omega$, each similarity measure $S(X, Y)$ above uses all or some of the four possible values, i.e., S_{00}, S_{01}, S_{10} and S_{11} . Table 1 shows the definitions of eight similarity measures based on S_{00}, S_{01}, S_{10} and S_{11} .

Similarity	Definition	Range
Jaccard-Needham	$\frac{S_{11}}{S_{11} + S_{10} + S_{01}}$	$[0, 1]$
Dice	$\frac{S_{11}}{2S_{11} + S_{10} + S_{01}}$	$[0, \frac{1}{2}]$
Correlation	$\frac{S_{11}S_{00} - S_{10}S_{01}}{((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}}$	$[-1, 1]$
Yule	$\frac{S_{11}S_{00} - S_{10}S_{01}}{S_{11}S_{00} + S_{10}S_{01}}$	$[-1, 1]$
Russell-Rao	$\frac{S_{11}}{N}$	$[0, 1]$
Sokal-Michener	$\frac{S_{11} + S_{00}}{N}$	$[0, 1]$
Rogers-Tanmoto	$\frac{S_{11} + S_{00}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$	$[0, 1]$
Kulzinsky	$\frac{S_{11}}{(S_{10} + S_{01})}$	$[0, \infty)$

Table 1. Eight similarity measures for binary feature vectors.

In Table 1, the ranges of the eight measures are given. Notice that some similarity measures like Correlation, Yule, and Kulzinsky are not normalized. 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, $S(\bar{I}, \bar{I}) = \infty$.

For each similarity $S(X, Y)$ defined above, we define an associated **dissimilarity measure** $D(X, Y)$. The corresponding eight dissimilarities are defined in Table 2. Notice that the range of each dissimilarity measure has been normalized as $[0, 1]$.

Dissimilarity	Definition
Jaccard-Needham	$\frac{S_{10} + S_{01}}{S_{11} + S_{10} + S_{01}}$
Dice	$\frac{S_{10} + S_{01}}{2S_{11} + S_{10} + S_{01}}$
Correlation	$\frac{1}{2} - \frac{S_{11}S_{00} - S_{10}S_{01}}{2((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}}$
Yule	$\frac{S_{10}S_{01}}{S_{11}S_{00} + S_{10}S_{01}}$
Russell-Rao	$\frac{N - S_{11}}{N}$
Sokal-Michener	$\frac{N - S_{11} - S_{00}}{N}$
Rogers-Tanmoto	$\frac{2S_{10} + 2S_{01}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$
Kulzinsky	$\frac{S_{10} + S_{01} - S_{11} + N}{(S_{10} + S_{01} + N)}$

Table 2. Eight dissimilarity measures for binary feature vectors.

Similarly, some assumptions have to be made for the functions in Table 2 to ensure that they are valid over the whole domain: (i) for Jaccard-Needham and Dice dissimilarity measures, $D(\bar{I}, \bar{I}) = 0$, (ii) for Correlation

and Yule dissimilarity measures, $D(I, I) = D(\bar{I}, \bar{I}) = 0$ and $D(I, \bar{I}) = D(\bar{I}, I) = 1$, and (iii) for Kulzinsky dissimilarity measure, $D(\bar{I}, \bar{I}) = 0$.

Among eight dissimilarity measures defined above, some like Sokal-Michener are metric, but some others like Russell-Rao are non-metric. The metric/non-metric properties of each dissimilarity measure will be detailed in another paper.

Considering the unequal importance of “zero” matches and “one” matches, we define a weighted version for each dissimilarity measure by weighing S_{00} with β ($0 \leq \beta \leq 1$). The proposition of weighing “zero” matches resulted from the observation that a bit '0' provides less information of separability than a bit '1'. Here, we don't follow Tubbs' rule to weight individual bit matches because in some cases, like handwriting identification, it is almost impossible to have sufficient handwritten samples from individual writers to yield useful statistical information. In the rest of the paper, any dissimilarity measure mentioned refers to the corresponding weighted dissimilarity measure.

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 weighing S_{00} makes no difference.
- The Yule dissimilarity measure becomes a constant when the coefficient β equals 0.
- The Jaccard-Needham measure is similar to the Dice measure. Given n points in Ω , there are C_n^2 pairs of points. If we sort these point pairs based on their distances in either ascendant or descendant order, the Jaccard-Needham measure and the Dice result in the same order. Thus, they should behave in the same way in k -nearest neighbor classification.

3. BINARY MICRO-FEATURES

Binary micro-features have been used in handwritten character recognition.¹ For an individual character, the micro-features consist of 512 bits corresponding to gradient (192 bits), structural (192 bits), and concavity (128 bits) features. Each of these three sets of features rely on dividing the scanned image of the allograph (character or combination of characters) into a 4 x 4 region. The gradient features capture the frequency of the direction of the gradient, as obtained by convolving the image with a Sobel edge operator, in each of 12 directions and then thresholding the resultant values to yield a 192-bit vector. The structural features capture, in the gradient image, the presence of corners, diagonal lines, and vertical and horizontal lines, as determined by 12 rules. The concavity features capture, in the binary image, major topological and geometrical features including direction of bays, presence of holes, and large vertical and horizontal strokes.

Three examples of handwritten digits and their corresponding binary micro-feature vectors are shown in Figure 1. For each digit image, a 512-dimensional micro-feature vector is extracted by the character recognizer.¹

Table 3 shows the dissimilarities between any two handwritten digit images in Figure 1 based on four measures, i.e., Jaccard-Needham, Correlation, Sokal-Michener, and Rogers-Tanmoto. From Table 3, we may find out, the dissimilarity between a vector and itself is 0 and Figure 1(a) has a smaller dissimilarity from Figure 1(c) than that from Figure 1(b). However, the four measures results in different proximal degrees among the handwritten images.

4. USE OF THE DISSIMILARITY MEASURES IN HANDWRITING IDENTIFICATION

The recognition capability of each dissimilarity is evaluated through a forensic document examination system, recently developed by CEDAR.⁵ This forensic document examination system characterizes a handwritten document with 12 floating macro-features for parameterizing against document, paragraph and word levels, and

5120 binary micro-features for individual character level. The identification process is as following. The macro-features are used as a filter that reduces the number of documents from 3,000 to 100, among which micro-features were then used to identify the writer.⁶ When the Sokal-Michener dissimilarity measure with $\beta = 0.5$ was used to measure similarity between micro-feature sequences, the experiments with 3000 documents written by 1000 representatives of US population showed that combination of the two sets of features presents a writer identification accuracy of 87.3% for 900 writers. However, we will see soon that the Sokal-Michener measure is not necessarily the best choice for this task.

We examine the performance variations of the eight dissimilarity measures in handwriting identification under the following scenarios: (i) using individual characters, (ii) using all characters, and (iii) combining the 12 macro-features and the 5120 dimensional binary micro-features.

We first address the experimental settings, then present experimental results with regard to the aforementioned scenarios, and finally give some general remarks on the recognition capabilities of the eight dissimilarity measures and the discriminative power of micro-features.

4.1. Experimental Settings

We conducted experiments with 2826 documents written by 942 writers. Each writer copied three times a source document specially designed by CEDAR.⁵ We randomly chose 900 writers out of 942, and for each writer one document out of three were randomly selected as testing sample. Finally, 900 documents by 900 writers were taken as the testing set and the other 1926 documents by 942 writers as the training set.

For each document, twelve macro-features were extracted to parameterize against document, paragraph and word levels, and the eight characters in "referred" as well as the letters "b" and "h" from the words "been" and "Cohen" were used as sample allographs in all of the tests for inter and intra-writer variability. These eleven allographs are used for extracting micro-features. Therefore, a 5120-dimensional binary micro-feature vector corresponds to each document in both the training and the testing sets.

The experiments were designed to test the following subjects: (i) the discriminatory powers of individual characters and their combination, (ii) the discriminatory power of the combination of the macro-features and the micro-features, (iii) the recognition capabilities of the eight dissimilarity measures, (iv) the effect of the weighing coefficient β for each dissimilarity measure, and (v) the effect of combination of the macro-features and micro-features.

4.2. Identification with Micro-Features Using Binary Vector Dissimilarities

When only a certain individual character in one word "referred" is considered for writership identification, the performance variations of the eight dissimilarity measures are shown in Figure 2, Figure 3, and Figure 4. It is easy to find that for any individual character different measures do have different performance in identification except that Jaccard-Needham and Dice measures behave in the same way. Moreover, the value of β has a very big impact on the performance when certain measures are used for the nearest neighbor classification. From these figures, we can also compare the discriminative power of individual characters.

For convenience, the discriminative power of individual characters using the Sokal-Michener dissimilarity measure, as a function of β , is assembled in Figure 5. Obviously, the first “r” in the word “referred” achieves the highest identification accuracy, thus exhibits more individuality of handwriting than others. As shown in Figure 5 (b), when $\beta = 0.5$, the first “r” in the word “referred” results in an accuracy of 20.67%, but the identification accuracy with the second “r” in the same word is only 10.89%.

When the micro-features from all ten characters were used for identification, Figure 6 shows the performance variation of the dissimilarity measures, which are the functions of the coefficient β . It is observed that: (i) three measures, i.e., Correlation, Sokal-Michener, and Roger-Tanmoto, can achieve very high accuracy (over 82%) when the optimal value of β is assigned, (ii) Jaccard-Needham and Dice measures have the same performance, which is comparable to the best one, (iii) Russell-Rao presents the lowest accuracy (less than %60), and (iv) the performance of the Sokal-Michener measure is more sensitive to the coefficient β than that of any other one.

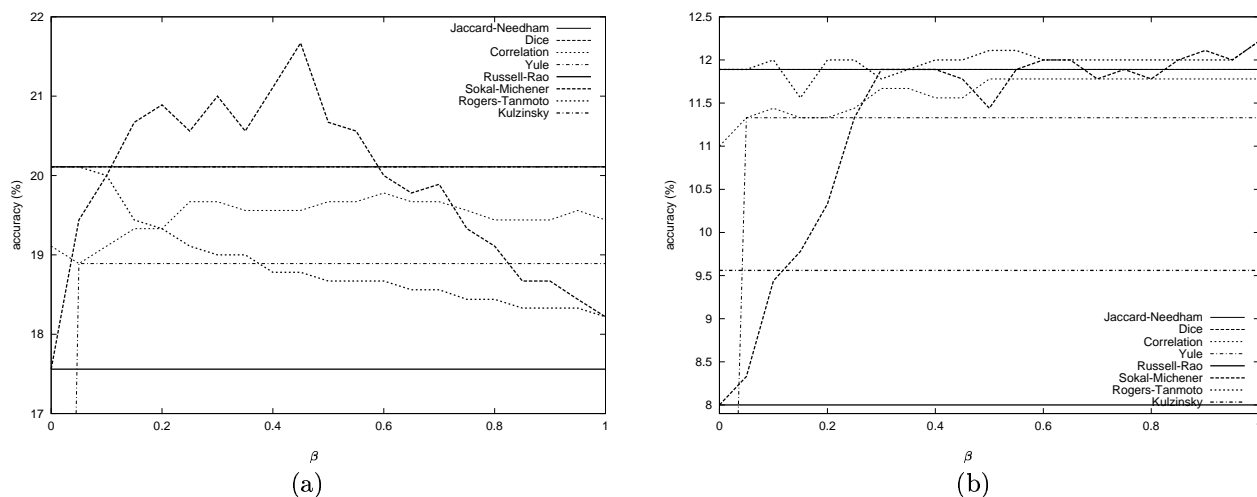


Figure 2. The identification performance of individual characters using the eight dissimilarity measures. (a) the first “r” in the word “referred”. (b) the second “r” in the word ‘referred’.

4.3. Identification with Macro-Features and Micro-Features

Instead of using the two-stage classification method developed by Srihari et al.,⁵ here we use an alternative strategy to combine the macro-features and the micro-features; that is, we take all 5120 binary micro-features as the 13th macro-feature and use the nearest neighbor classification based on a hybrid distance measure defined as following. Given two handwriting samples, each with 13 macro-features, the hybrid distance between them is computed in a way similar to Euclidean distance except that the difference between the 13th component is calculated by the dissimilarity measures defined in Table 2.

Figure 7 shows the identification performance of the combination with regard to the eight binary vector dissimilarity measures. We can find that the results here are inconsistent with those in Figure 6. A few

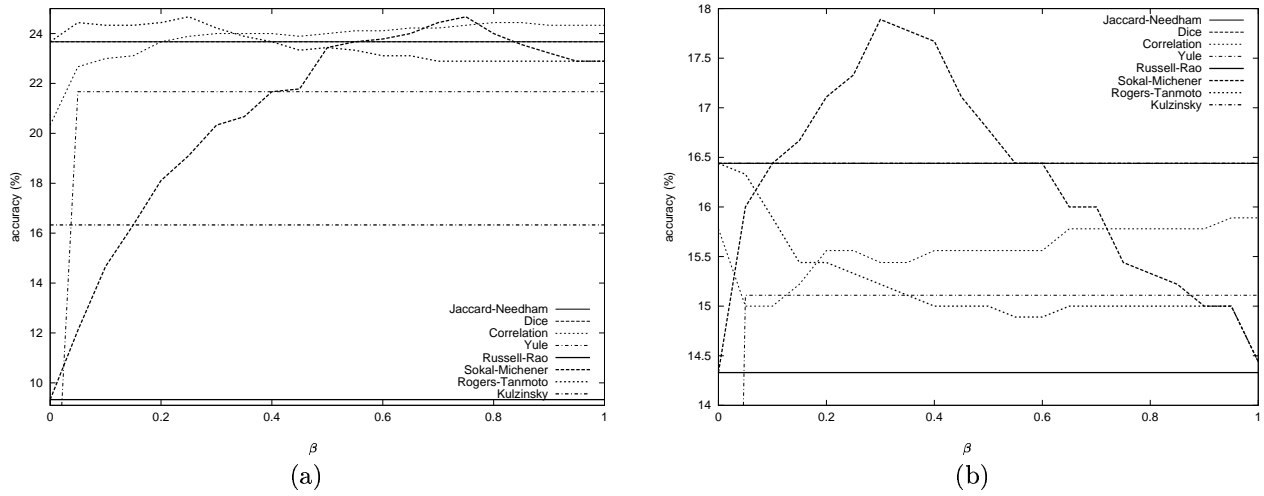


Figure 3. The identification performance of individual characters using the eight dissimilarity measures. (a) the first “e” in the word “referred”. (b) the third “e” in the word “referred”.

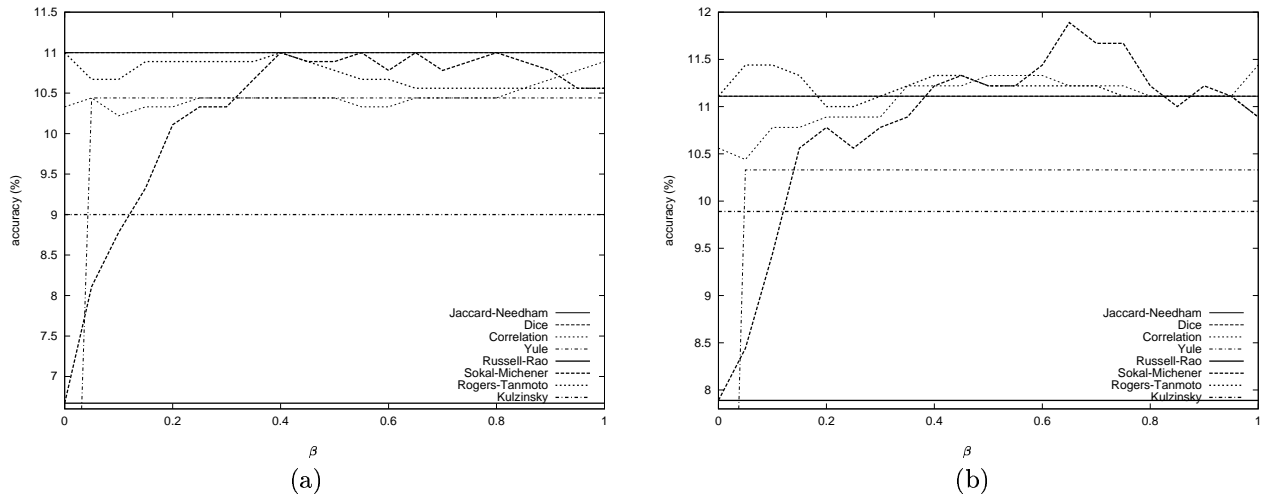


Figure 4. The identification performance of individual characters using the eight dissimilarity measures. (a) “d” in the word “referred”. (b) “f” in “referred”.

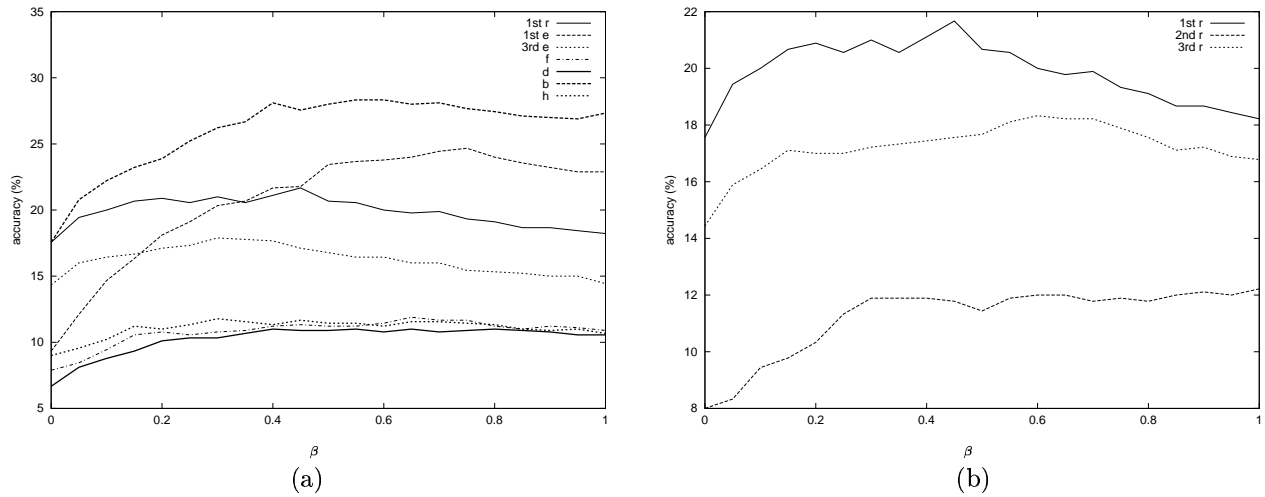


Figure 5. Comparison of the discriminative power of individual characters using the Sokal-Michener dissimilarity measure. (a) seven characters. (b) three “r”s in the word “referred”.

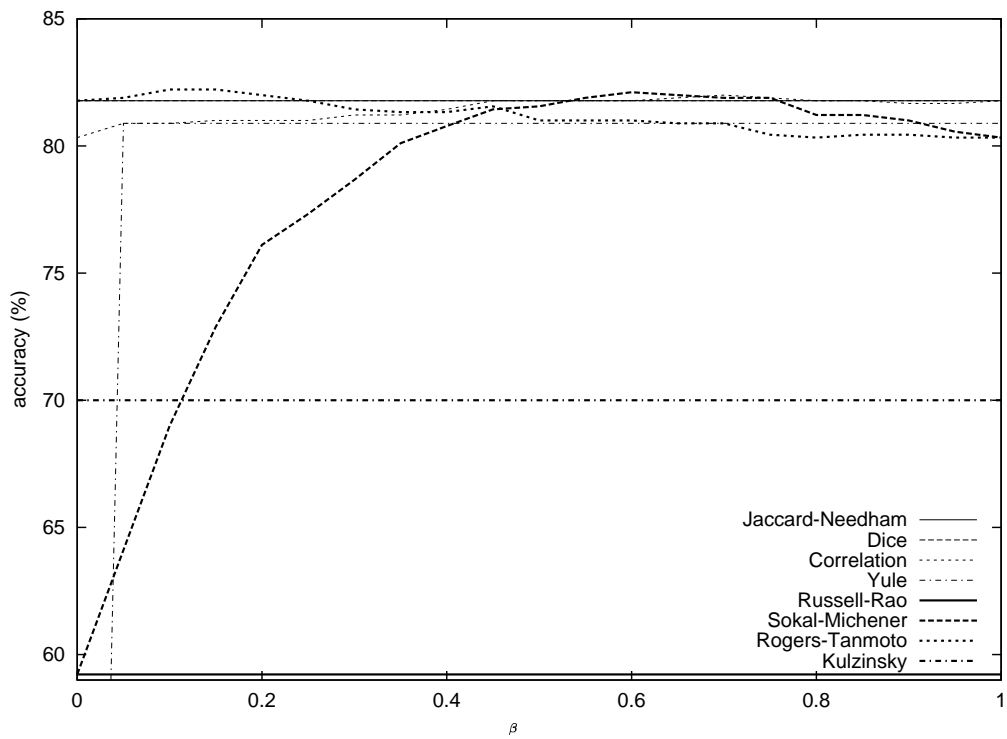


Figure 6. The identification performance of the 5120-dimensional micro-features using eight dissimilarity measures.

observations are: (i) Rogers-Tanmoto measure exhibits the highest accuracy of 92.11% with $\beta = 0$, (ii) Jaccard-Needham, Dice, Correlation, and Yule measures all give very good performance (with accuracy between 91% and 91.5%), (iii) With $\beta = 0.45$ Sokal-Michener gives its best performance of only 89.22%, about 2% less than that of Rogers-Tanmoto measure, (iv) Kulzinsky measure presents an accuracy of 88.33% and Russell-Rao measure performs worst.

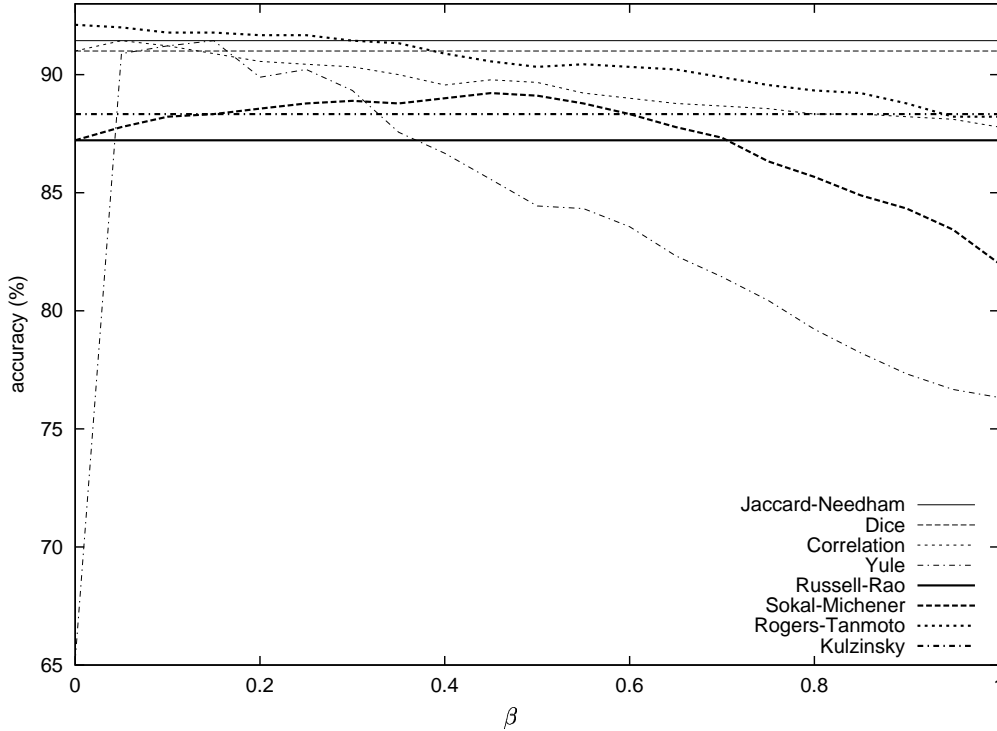


Figure 7. Identification performance of the combination of macro- and micro-features with regard to the eight dissimilarity measures, which are functions of β .

The inconsistency between Figure 7 and Figure 6 indicates that some mismatches may occur in a hybrid distance when certain binary dissimilarity measures are used for computing the difference between two micro-feature vectors, i.e., the improper scale of the difference between two 13th components will cause micro-features to either dominate or have small contribution to the distance value, resulting in biased classification.

To examine the effect of the mismatch between two types of features, real-valued macro and binary micro, another experiment was designed by weighing (with the coefficient α) the difference between two 13th components in the hybrid distance (called the weighted hybrid distance). Figure 8 shows the effect of the weighted hybrid distance with regard to the different binary vector dissimilarity measures.

Some interesting results are presented in Figure 8. Considering the best performance with each measure, we come to the following conclusions: (i) Rogers-Tanmoto measure exhibits the highest accuracy of 92.11%, (ii) each of the four measures, Jaccard-Needham, Dice, Correlation, and Sokal-Michener, gives an accuracy over 91.78%

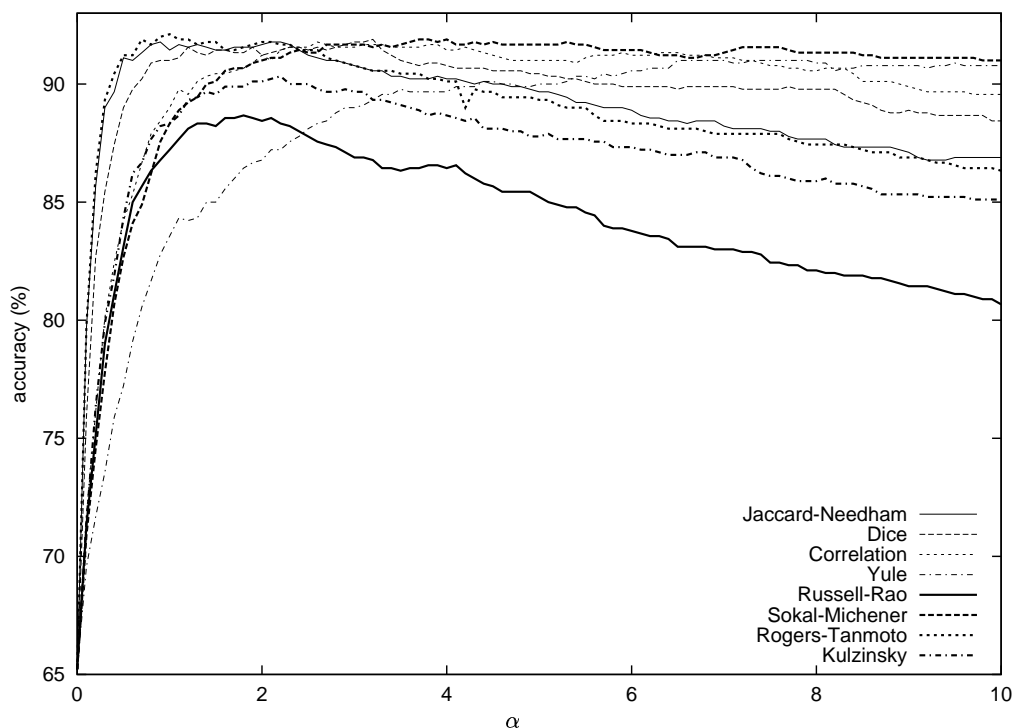


Figure 8. Identification performance of the weighted combination of macro- and micro-features with regard to the eight dissimilarity measures, where for Rogers-Tanmoto measure β equals to 0 and for others β is set as 0.6.

but under 92%, (iii) Yule measure comes as the next one with an accuracy of 91%, (iv) Kulzinsky measure presents an accuracy of 90.33% and Russell-Rao measure gives the lowest accuracy of 88.67%. Obviously these results are consistent with those in Figure 6.

4.4. Some General Remarks

Following the observations and analysis in the subsections 4.2 and 4.3, we can make some general remarks upon the effects of micro-features and the dissimilarity measures on handwriting identification.

- Different character shapes and different positions of a character in a word bear different capability for identifying handwriting individuality.
- A character alone usually is not informative enough to identify handwriting differences, but a combination of many characters presents good performance in handwriting identification.
- A combination of macro-features and micro-features gives a high identification accuracy. However, the problem of mismatching will occur in computing the hybrid distance. This brings out a question about how to weight feature components for optimal performance.

- The eight measures are ranked in descending order of identification accuracy as follows. The first group consists of five measures, e.g., Rogers-Tanmoto, Jaccard-Needham, Dice, Correlation, and Sokal-Michener. The second group includes Yule measure, Kulzinsky comes as the next tier, and Russell-Rao belongs to the fourth group.
- When constructing a similarity or dissimilarity measure for two binary vectors, all possible matches, S_{ij} , $i, j \in \{0, 1\}$, should be considered for better classification performance.
- Weighing S_{00} with β in three measures, Correlation, Sokal-Michener, and Rogers-Tanmoto, is very effective for boosting classification performance. At the moment, the optimal value of β can only be estimated by experimentation.

5. CONCLUSION

The eight dissimilarity measures for binary vectors were formulated and examined for their recognition capability in a forensic document examination system, where the binary micro-features are used to characterize handwritten character shapes. For each dissimilarity measure, a large number of experiments were conducted to evaluate the writership identification effects of individual characters and their combination. According to the identification performance of the eight dissimilarity measures in the various scenarios, these measures are ranked as follows. The first tier consists of five measures, e.g., Rogers-Tanmoto, Jaccard-Needham, Dice, Correlation, and Sokal-Michener. The second tier includes Yule measure, Kulzinsky measure comes as the next tier, and Russell-Rao measure belongs to the fourth tier.

In addition, a new strategy was developed to combine the micro- and macro-features for handwriting identification, and a high identification accuracy of 92.11% was obtained by using the optimal dissimilarity measure for computing the difference between two micro-feature vectors when 900 writers are considered.

The work in this paper has led to several interesting topics which are worth further exploration: (i) how to compute the optimal value of the coefficient β in three measures, Correlation, Sokal-Michener, and Rogers-Tanmoto, and (ii) how to weight feature components in the hybrid distance for the optimal performance.

REFERENCES

1. J. T. Favata and G. Srikantan, "A multiple feature/resolution approach to handprinted character/digit recognition," *International Journal of Imaging Systems and Technology* **7**, pp. 304–311, 1996.
2. 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.
3. R. R. Sokal and P. H. Sneath, *Principles of Numerical Taxonomy*, W. H. Freeman, 1963.
4. J. D. Tubbs, "A note on binary template matching," *Pattern Recognition* **22**(4), pp. 359–365, 1989.
5. S. N. Srihari, S.-H. Cha, H. Arora, and S. Lee, "Individuality of handwriting," *Journal of Forensic Sciences* **47**, July 2002.
6. S. Lee, S.-H. Cha, and S. N. Srihari, "Combining macro and micro features for writer identification," in *Proceedings SPIE, Document Recognition and Retrieval IX*, pp. 155–166, (San Jose, CA), January 2002.