

Enhancing Binary Feature Vector Similarity Measures

Sung-Hyuk Cha

Charles Tappert

*Computer Science Department, Pace University
861 Bedford Road, Pleasantville, New York, 10570 USA*

Sungsoo Yoon

*Department of Computer Science and Engineering, Ewha Womans University
Daehyun Dong, Seodaemun Gu, Seoul, Korea*

scha@pace.edu

ctappert@pace.edu

ssyoon@ewha.ac.kr

Received Jul. 25, 2006. Published Nov. 30, 2006.

Abstract

Similarity and dissimilarity measures play an important role in pattern classification and clustering. For a century, researchers have searched for a good measure. Here, we review, categorize, and evaluate various binary vector similarity / dissimilarity measures. One of the most contentious disputes in the similarity measure selection problem is whether the measure includes or excludes negative matches. While inner-product based similarity measures consider only positive matches, other conventional measures credit both positive and negative matches equally. Hence, we propose an enhanced similarity measure that gives variable credits and show that it is superior to conventional measures in IRIS biometric authentication and offline handwritten character recognition applications. Finally, the proposed similarity measure can be further boosted by applying weights and we demonstrate that it outperforms the weighted Hamming distance.

Keywords: Binary similarity, Distance metric, Nearest neighbor, Genetic algorithm, IRIS biometric verification, Handwriting recognition.

1. Introduction

A common method for classifying an unknown input vector involves finding the top k similar vectors in a reference set. The k -nearest neighbor, or simply k -NN, has wide acceptance in pattern classification problems (e.g. [1, 2] for extensive surveys). There are two important aspects in this approach. One is extracting important features from the pattern, and the other is selecting an appropriate similarity measure. Although there are many types of features, in this paper we consider only binary features and their similarity / dissimilarity measures.

Distance and similarity measures are encountered in various fields such as image retrieval [3], information retrieval, chemistry [4], ecology [5], psychology, and biological taxonomy [6]. There is a wealth of literature regarding the similarity measure selection problem dating as far back as 1908. Extensive lists of similarity measures can be found in [7, 8]. Conventional definitions of similarity or dissimilarity measures include Hamming [9], inner-product, Tanimoto distance, etc., and for a century, researchers have searched for a good measure. The most recent comparative works on similarity measures include Tubbs [10] who summarized seven conventional similarity measures for the template matching problem [10], and Zhang et al. who compared these seven measures for their recognition capability in handwriting identification [11]. Yet another distance measure that has been used for binary features extracted from offline character images can be found in [12, 13, 14]. Here, we categorize and review various similarity and distance measures.

The most favored distance measure is the Hamming distance when the features are binary. To

further improve the performance, there are two approaches. First, weights can be applied to features [14] and optimized using techniques such as genetic algorithms [15, 16]. Another approach is to use a similarity measure that gives full credit to features present in both patterns, less credit to those not present in either pattern, and no credit to those present in only one of the patterns to be matched [14]. Both approaches have been reported to perform better than the simple Hamming distance approach. In this paper, we create a new measure that combines these two approaches, and we present experimental results that demonstrate its superiority over the other measures.

One of the most contentious problems in binary feature vector similarity measures is whether the measure includes or excludes the number of negative matches. When features are absent in both patterns, should we consider it as important as those present in both patterns? This problem has been argued in [7, 17] and several binary vector coefficients include the negative matches as well as the positive matches equally. In this paper, the proposed measure includes both positive and negative matches, and has weights that can be optimized to control the degree to which the positive and negative matches are considered.

To evaluate similarity measures for binary features, we chose two binary feature databases: an IRIS code database and an offline handwritten character database. First, we consider the problem of IRIS biometric verification which often uses a distance or similarity between two samples of the same class and between samples of two different classes. Two patterns are categorized into one of only two classes – the patterns are either from the same class or from two different classes. Given two IRIS biometric samples, the feature distance between the two samples is classified as intra-person (identity) or inter-person (non-identity). The intra- and inter-person distances form two distributions having some overlap with each other. Two types of errors, False Accept Rate (FAR) and False Reject Rate (FRR), are used to evaluate the various similarity measures.

In the offline handwritten character image database, Gradient, Structural, and Concavity binary features, or simply GSC, have been developed and utilized in character recognition [12]. While GSC features, which are binary in type, are considered significant ones, relatively little study has been conducted on selecting and designing a good similarity measure for these features. In this paper, we evaluate numerous similarity measures and determine the optimal one.

The subsequent sections are organized as follows. Section 2 enumerates many similarity measures and their weighted variations. Sections 3 and 4 evaluate similarity measures on IRIS code and offline handwritten character databases, respectively. Finally, section 5 concludes the paper.

2. Overview of Similarity Measures

In this section, we give the definitions of conventional binary vector similarity and dissimilarity (distance) measures and then show how some of these measures can be refined with weights that can be optimized to enhance their discrimination capability.

2.1 Basic Binary Similarity Measures

Let x , y , and z be binary feature vectors of fixed length d , and let x_i denote the i -th feature value, which is equal to either 0 or 1. One of the most popular measures in comparing two fixed-length bit patterns is the Hamming distance given in (1), which is the count of the bits that differ in the two patterns [9]. It is a simple geometrical L_1 distance, also known as Manhattan or city block distance, applied to d -dimensional binary space.

$$\begin{aligned}
 D_H(x, y) &= x^T \bar{y} + \bar{x}^T y \\
 D_H(x, y) &= \sum_{i=1}^d |x_i - y_i| \\
 S_H(x, y) &= d - D_H(x, y) = x^T y + \bar{x}^T \bar{y} \\
 S_H(x, y) &= \sum_{i=1}^d s_i, \text{ where } s_i = 1 \text{ if } x_i = y_i \text{ and } s_i = 0 \text{ otherwise.}
 \end{aligned} \tag{1}$$

The term $x^T y$ denotes the positive matches, i.e. the number of “1” bits that match between x and y . The term $\bar{x}^T \bar{y}$ is the negative matches, i.e. the number of “0” matching bits. The terms $x^T \bar{y}$ and $\bar{x}^T y$ denote the number of bit mismatches – the first where pattern x has a “1” and pattern y has a “0”, and the second where pattern x has a “0” and pattern y has a “1”.

Fact 1. *The Hamming distance has been shown to be metric [6].*

While the Hamming distance is the number of bits that are different in the two patterns, the Hamming similarity is the number of identical bits in the two patterns. Sokal and Michener (SM) normalized the Hamming similarity [18]

$$S_{SM}(x, y) = \frac{x^T y + \bar{x}^T \bar{y}}{d}. \tag{2}$$

An alternative normalized Hamming similarity is given by Rogers and Tanimoto (RT) [19]

$$S_{RT}(x, y) = \frac{x^T y + \bar{x}^T \bar{y}}{x^T y + \bar{x}^T \bar{y} + 2x^T \bar{y} + 2\bar{x}^T y}. \tag{3}$$

The term $x^T y$ is the inner product (IP) of two vectors, which yields a scalar, and it is sometimes called the scalar product or dot product. It can be converted to a distance by subtracting it from d , and this distance is clearly non-metric because of the reflexivity violation; $D_{IP}(x, y) = 0$ iff $x = y$ and $|x| = |y| = d$ and $D_{IP}(x, y) > 0$, otherwise.

Fact 2. *Nonnegativity, symmetry, and triangle inequality are trivial and preserved in the inner product [20].*

$$\begin{aligned}
 S_{IP}(x, y) &= x^T y \\
 S_{IP}(x, y) &= \sum_{i=1}^d s_i \\
 &\text{where } s_i = 1 \text{ if } x_i = y_i = 1, \text{ and } s_i = 0 \text{ otherwise}
 \end{aligned} \tag{4}$$

A normalized inner product (NIP) is given in (5) [6]. Alternative normalizations are Russell-Rao (RR) given in (6) [7], Jaccard-Needham (JN) given in (7) [21], Dice (D) given in (8) [22], and Kulzinsky (K) given in (9) [23]:

$$S_{NIP}(x, y) = \frac{x^T y}{\|x\| \|y\|} = \frac{x^T y}{\sqrt{x^T x y^T y}} \tag{5}$$

$$S_{RR}(x, y) = \frac{x^T y}{d} \tag{6}$$

$$S_{JN}(x, y) = \frac{x^T y}{x^T y + x^T \bar{y} + \bar{x}^T y} \tag{7}$$

$$S_D(x, y) = \frac{x^T y}{2x^T y + x^T \bar{y} + \bar{x}^T y} \quad (8)$$

$$S_K(x, y) = \frac{x^T y}{x^T \bar{y} + \bar{x}^T y} . \quad (9)$$

The Jaccard, Dice, and Kulzinsky similarity measures differ in their ranges: the Jaccard measure ranges from 0 to 1, the Dice measure from 0 to 0.5, and the Kulzinsky measure from 0 to ∞ . Equations 7-9 can be generalized to generalized Jaccard (GJ) similarity measure given in (10), which for $\sigma=0$ becomes the Kulzinsky coefficient, for $\sigma=1$ the Jaccard coefficient, and for $\sigma=2$ the Dice coefficient:

$$S_{GJ}(x, y) = \frac{x^T y}{\sigma x^T y + x^T \bar{y} + \bar{x}^T y} . \quad (10)$$

Another popular distance measure between binary feature vectors is the Tanimoto (T) metric defined in (11) [6] where n_x and n_y are the numbers of “1” bits in x and y , respectively, and $n_{x,y}$ is $x^T y$:

$$D_T(x, y) = \frac{n_x + n_y - 2n_{x,y}}{n_x + n_y - n_{x,y}} = \frac{x^T \bar{y} + \bar{x}^T y}{x^T \bar{y} + \bar{x}^T y + x^T y} . \quad (11)$$

The Tanimoto coefficient [6], defined in (12), is another variation of the normalized inner product which is frequently encountered in the fields of information retrieval and biological taxonomy:

$$S_T(x, y) = \frac{x^T y}{x^T x + y^T y - x^T y} \quad (12)$$

Most similarity measures are variations either of Hamming or of the inner product measures. Generally, the former ones treat the presence, $x^T y$, and the absence, $\bar{x}^T \bar{y}$, of features equally while the later take only the presence, $x^T y$, into account and exclude $\bar{x}^T \bar{y}$. The decision to include or exclude the $\bar{x}^T \bar{y}$ term is a difficult and contentious one [7, 17]. Prior to 1950 when the Hamming distance was introduced, the use of inner product based similarity coefficients flourished. Sokal and Michener made a good argument to include the negative matches [7, 17, 18], however, both positive and negative matches were equally weighted. Hence, we propose a new measure with variable credit for the $\bar{x}^T \bar{y}$ term, defined in (13), where σ is the contribution factor, and $0 \leq \sigma < \infty$. We call it the Alter Zero Zero One One (AZZOO) similarity measure because we can alter the credit for the zero-zero matches relative to that for the one-one matches.

$$S_{AZZOO}(x, y) = \bar{x}^T \bar{y} + \sigma \bar{x}^T \bar{y} = \sum_{i=1}^d x_i y_i + \sigma \sum_{i=1}^d (1-x_i)(1-y_i) \quad (13)$$

$$S_{AZZOO}(x, y) = \sum_{i=1}^d s_i, \text{ where } s_i = 1 \text{ if } x_i = y_i = 1, s_i = \sigma \text{ if } x_i = y_i = 0, \text{ and } s_i = 0 \text{ otherwise}$$

Note that for $\sigma=0$, S_{AZZOO} becomes the inner product, and for $\sigma=1$, S_{AZZOO} becomes the Hamming similarity measure. Although S_{AZZOO} requires finding the optimal σ factor, the experimental results in later sections show that it outperforms both the Hamming and inner product similarity measures.

Originally, the half-credit similarity, S_{00-11} was used in an offline handwriting recognition

system [12], and it is the same as S_{AZZOO} with $\sigma=0.5$. It gives full credit to features present in both patterns, $x^T y$, half credit to those not present in either pattern, $\bar{x}^T \bar{y}$, and no credit to those present in only one of the patterns, $x^T \bar{y}$ and $\bar{x}^T y$, as defined in (14) [12, 13, 14]. In this paper we generalized the half-credit similarity to S_{AZZOO} .

$$S_{00-11}(x, y) = x^T y + 0.5 \bar{x}^T \bar{y} \quad (14)$$

The range of S_{AZZOO} is $[0, d]$ if $0 \leq \sigma \leq 1$ and $[0, \sigma d]$ if $\sigma > 1$. Assuming $0 \leq \sigma \leq 1$, S_{AZZOO} can be converted to a distance measure for metric property testing:

$$D_{AZZOO}(x, y) = d - S_{AZZOO}(x, y) = d - (x^T y + \sigma \bar{x}^T \bar{y}). \quad (15)$$

Nonnegativity and symmetry are trivial and preserved. Reflexivity is violated, however, because $D_{AZZOO}(x, y) = 0$ iff $x = y$ and $|x| = |y| = d$ and $D_{AZZOO}(x, y) \neq 0$ otherwise. Similarly, $S_{AZZOO}(x, y) = d$ iff $x = y$ and $|x| = |y| = d$ and $\sigma d \leq S_{00-11}(x, y) < d$ if $x = y$ and $|x| < d$.

Theorem 1. *The triangle inequality property is valid for $D_{AZZOO}(x, y)$, i.e. $D_{AZZOO}(x, y) + D_{AZZOO}(y, z) \geq D_{AZZOO}(x, z)$.*

Proof

$$\bar{x}^T \bar{y} + \bar{y}^T \bar{z} \geq \bar{x}^T \bar{z} \text{ by Fact 1} \quad \text{line 1}$$

$$d - (x^T y + \bar{x}^T \bar{y}) + d - (y^T z + \bar{y}^T \bar{z}) \geq d - (x^T z + \bar{x}^T \bar{z}) \text{ by Fact 2} \quad \text{line 2}$$

Now, evaluate

$$\begin{aligned} D_{AZZOO}(x, y) + D_{AZZOO}(y, z) &\geq D_{AZZOO}(x, z) \\ d - (x^T y + \sigma \bar{x}^T \bar{y}) + d - (y^T z + \sigma \bar{y}^T \bar{z}) &\geq d - (x^T z + \sigma \bar{x}^T \bar{z}) \\ d - (x^T y + \bar{x}^T \bar{y}) + d - (y^T z + \bar{y}^T \bar{z}) + (1 - \sigma) \bar{x}^T \bar{y} &\geq d - (x^T z + \bar{x}^T \bar{z}) + (1 - \sigma) \bar{x}^T \bar{z} \end{aligned}$$

Hence, the theorem is true by line 1 and line 2 of proof. \square

Similarly, since $S_{AZZOO}(x, y) = x^T y + \sigma \bar{x}^T \bar{y} = S_{IP}(x, y) + \sigma S_{IP}(\bar{x}, \bar{y})$, the properties of the AZZOO similarity measure are similar to those of the inner product measure.

Other popular similarity measures utilize correlation coefficients (CC) and have been used frequently in both psychology and ecology studies [7]. The correlation similarity measure is given in (16) and Yule and Kendall (YK) [24] suggested a similar coefficient given in (17).

$$S_{CC}(x, y) = \frac{x^T y \times \bar{x}^T \bar{y} - x^T \bar{y} \times \bar{x}^T y}{\sqrt{(x^T \bar{y} + x^T y)(\bar{x}^T y + \bar{x}^T \bar{y})(x^T y + \bar{x}^T y)(\bar{x}^T \bar{y} + x^T \bar{y})}} \quad (16)$$

$$S_{YK}(x, y) = \frac{x^T y \times \bar{x}^T \bar{y} - x^T \bar{y} \times \bar{x}^T y}{x^T y \times \bar{x}^T \bar{y} + x^T \bar{y} \times \bar{x}^T y} \quad (17)$$

While Hamming based similarity measures are additive forms of the positive and negative matches, the correlation based measures are multiplicative forms. Nonetheless, contribution factors of positive and negative matches are considered equally important in correlation based similarity measures as well as Hamming based ones.

Historically, all the measures enumerated above have had great value in their respective fields. In the following sections 3 and 4, we evaluate these measures in the applications of IRIS biometric authentication and offline handwriting character recognition. Fig. 1 shows a

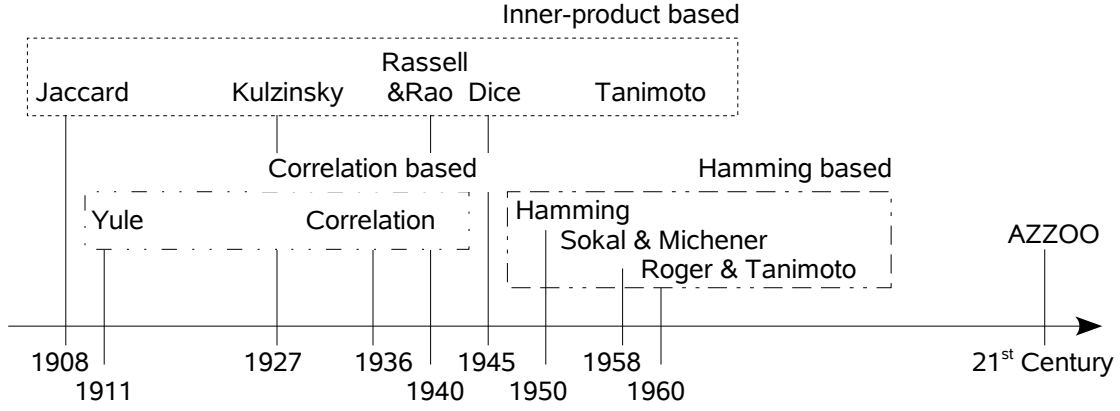


Fig. 1: A chronological table for binary vector similarity measures.

chronological table for binary feature vector similarity measures in which these conventional measures are categorized into three major groups: inner product, Hamming, and correlation based groups.

2.2 Binary Similarity Measures with Weights

To further improve discrimination capability, weights can be applied to distance or similarity measures [14] and optimized using techniques such as genetic algorithms [15, 16]. When features have numeric values, a scaling problem arises. In order to mitigate this problem, one can combine the nonlinear accuracy weighting with the Minkowski distance (WM) concept as shown in (18), where $P(C|i)$ is the probability of being correct when only feature i is used [25, 26].

$$D_{WM}(x, y) = \sum_{i=1}^d P(C|i)^a |x_i - y_i|^r \quad (18)$$

When features are binary, one can still generalize (18) to weighted Hamming (WH), given in (19), by setting $r=1$ and $P(C|i)^a = w_i$.

$$D_{WH}(x, y) = \sum_{i=1}^d w_i |x_i - y_i| \quad \text{and} \quad S_{WH}(x, y) = \sum_{i=1}^d w_i (x_i y_i + \bar{x}_i \bar{y}_i) \quad (19)$$

The weighted Hamming distance has been applied to numerous applications such as image template matching [27, 28] and object recognition [28]. The weighted Hamming distance provides an improvement over the simple Hamming distance for discriminating between similar images [27, 28]. This distance measure gives greater importance to error pixels which appear in close proximity to other error pixels. Error pixels that appear close together tend to correspond to structurally meaningful features. In [29], a slightly different weighted Hamming distance was introduced to optimize the distance measure for object detection by adding a null weight, w_0 . Similarly, the weighted inner product (WIP) similarity measure is given by (20).

$$S_{WIP}(x, y) = \sum_{i=1}^d w_i x_i y_i \quad (20)$$

Here, we claim that the performance can be further improved by optimizing the similarity measure rather than distance measure. Since the Hamming distance is the number of mismatches, the weights are applied to the mismatched bits, whereas in a similarity measure the weights are applied to the matching bits. As discussed in the earlier section, there are two kinds of matches:

positive and negative matches. Although the Hamming similarity can be improved by applying the equal weights are applied to both positive and negative matches, we claim that if different weights are applied, the performance is further improved, and the proposed weighted AZZOO similarity measure is given in (21).

$$S_{WAZZOO}(x, y) = \sum_{i=1}^d w_i^+ x_i y_i + \sum_{i=1}^d w_i^- \bar{x}_i \bar{y}_i \quad (21)$$

Note that if w^+ and w^- are identical, $S_{WAZZOO} = S_{WH}$, and if $w^- = 0$, $S_{WAZZOO} = S_{WIP}$.

There are twice as many coefficients to optimize in this new similarity measure than in the weighted Hamming or inner product similarity measures. This is a multi-dimensional, space optimization problem, and one can use a genetic algorithm to determine the weights from training data. A genetic algorithm can be a general optimization method that searches a large space of candidate objects to find one that performs near optimal according to the fitness function [15, 16]. Genetic algorithms offer a number of advantages: they search from a set of solutions rather than from a single one, they are not derivative-based, and they explore and exploit the parameter space. For the weight adaptive model, we create a numerical optimization model that depends on a set of weights.

3. Similarity Measure Evaluation on IRIS Biometric Verification

In order to evaluate the binary vector similarity measures, we consider an IRIS biometric database. Daugman proposed the degrees of freedom of IRIS mismatch score distribution as a measure of the individuality or uniqueness of an IRIS pattern [30]. The biometric verification problem is a simple dichotomy classification problem that places the input into one of only two categories – that is, given two randomly selected biometric samples, the problem is to determine whether the two samples belong to the same person or two to different people. Fig. 2 depicts the biometric verification model.

First, features are extracted from IRIS biometric image data x and $y: \{x_1, x_2, \dots, x_d\}$ and $\{y_1, y_2, \dots, y_d\}$. Let $c(x)$ denote the class (the person) to which x belongs. The IRIS code,

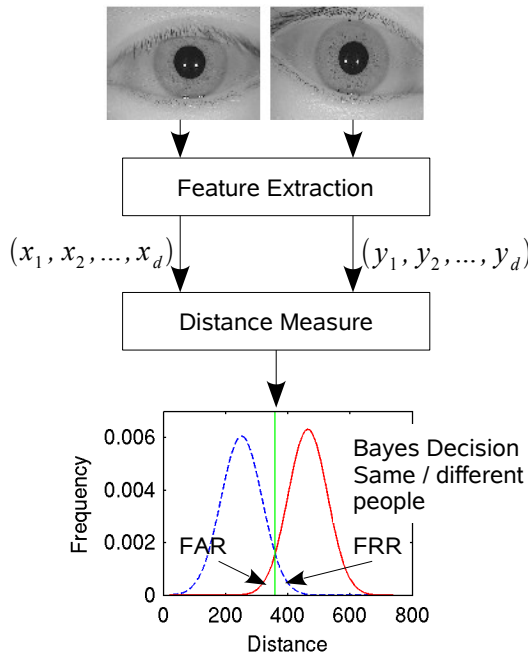


Fig. 2: IRIS verification process.

proposed by Daugman [30], is a 8×256 binary feature extracted from an IRIS image by applying a 2-D Gabor wavelet filter, and Daugman used the Hamming distance. When a distance measure is applied, two distributions are generated. One distribution, called the intra-distance (or within person) distribution, occurs when $c(x) = c(y)$. The other distribution, called the inter-distance (or between two different people) distribution, occurs when $c(x) \neq c(y)$. By assuming the distributions are normal one can easily find the decision threshold to minimize the false accept rate (FAR) and the false reject rate (FRR).

3.1 Performance Evaluation Method

In Fig. 2, FAR is the probability of error that one classifies two biometric data as coming from the same person even though they belong to two different people, and FRR is the probability of error that one classifies two biometric samples as coming

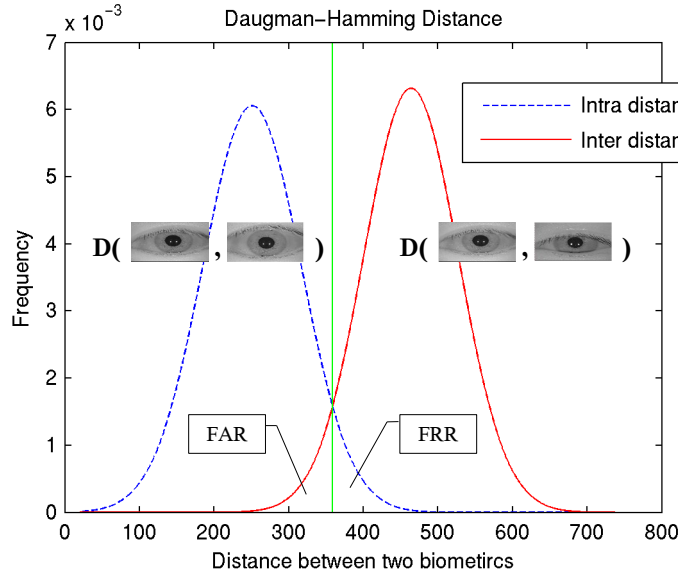


Fig. 3: Intra and inter distance distribution using Hamming distance.

from different people even though they belong to a same person. Note that when a distance measure is used, the intra distance distribution tends to be close to 0 whereas the inter distance distribution tends to be far from 0. Thus, FAR is usually the left side area of the decision boundary. When a similarity measure is used, on the other hand, FAR is the right-side area of the decision boundary because the larger value means that the two biometric samples are similar, as defined in (22) and (23) where T is the threshold value for the Bayesian decision.

$$\text{FAR} = \Pr(S(x, y) \geq T | c(x) \neq c(y)) \text{ and } \text{FRR} = \Pr(S(x, y) < T | c(x) = c(y)) \quad (22, 23)$$

The overall performance is the number of correctly classified instances divided by the total testing database size.

3.2 Experimental Results

In this section, we compare the experimental results obtained by using several similarity measures. From the IRIS biometric image database [31], we selected 10 left bare eye samples of 52 subjects. In order to test the described models, two sets of samples are required: intra-class

Table 1: Performance evaluation of the similarity measures on the iris database.

Method	Data 1			Data 2			Data 3			Data 4			Total
	FAR	FRR	Rate	FAR	FRR	Rate	FAR	FRR	Rate	FAR	FRR	Rate	Rate
AZZOO	5.0	4.4	95.3	6.4	3.2	95.2	7.6	3.2	94.6	4.4	3.8	95.9	95.3
Norm. I.P.	4.8	4.8	95.2	6.0	4.4	94.8	7.4	3.6	94.5	5.0	4.4	95.3	95.0
Sokal	5.0	4.8	95.1	6.4	3.6	95.0	7.8	3.2	94.5	4.6	3.8	95.8	95.1
Rogers	5.0	4.8	95.1	6.4	3.6	95.0	7.8	3.2	94.5	4.6	3.8	95.8	95.1
Russell	12.8	12.2	87.5	11.8	10.4	88.9	11.4	8.6	90.0	11.2	10.4	89.2	88.9
Jaccard	4.8	4.8	95.2	6.2	4.2	94.8	7.4	3.6	94.5	5.0	4.0	95.5	95.0
Dice	4.8	4.8	95.2	6.0	4.2	94.9	7.4	3.6	94.5	5.0	4.4	95.3	95.0
Kulzinsky	6.8	3.8	94.7	7.4	3.0	94.8	9.0	2.6	94.2	6.4	3.4	95.1	94.7
Tanimoto	4.8	4.8	95.2	6.2	4.2	94.8	7.4	3.6	94.5	5.0	4.0	95.5	95.0
Correlation	5.4	4.6	95.0	6.4	3.8	94.9	7.8	3.2	94.5	4.2	3.6	96.1	95.1
Yule	4.8	4.8	95.2	5.6	4.8	94.8	7.8	3.0	94.6	4.0	3.8	96.1	95.2

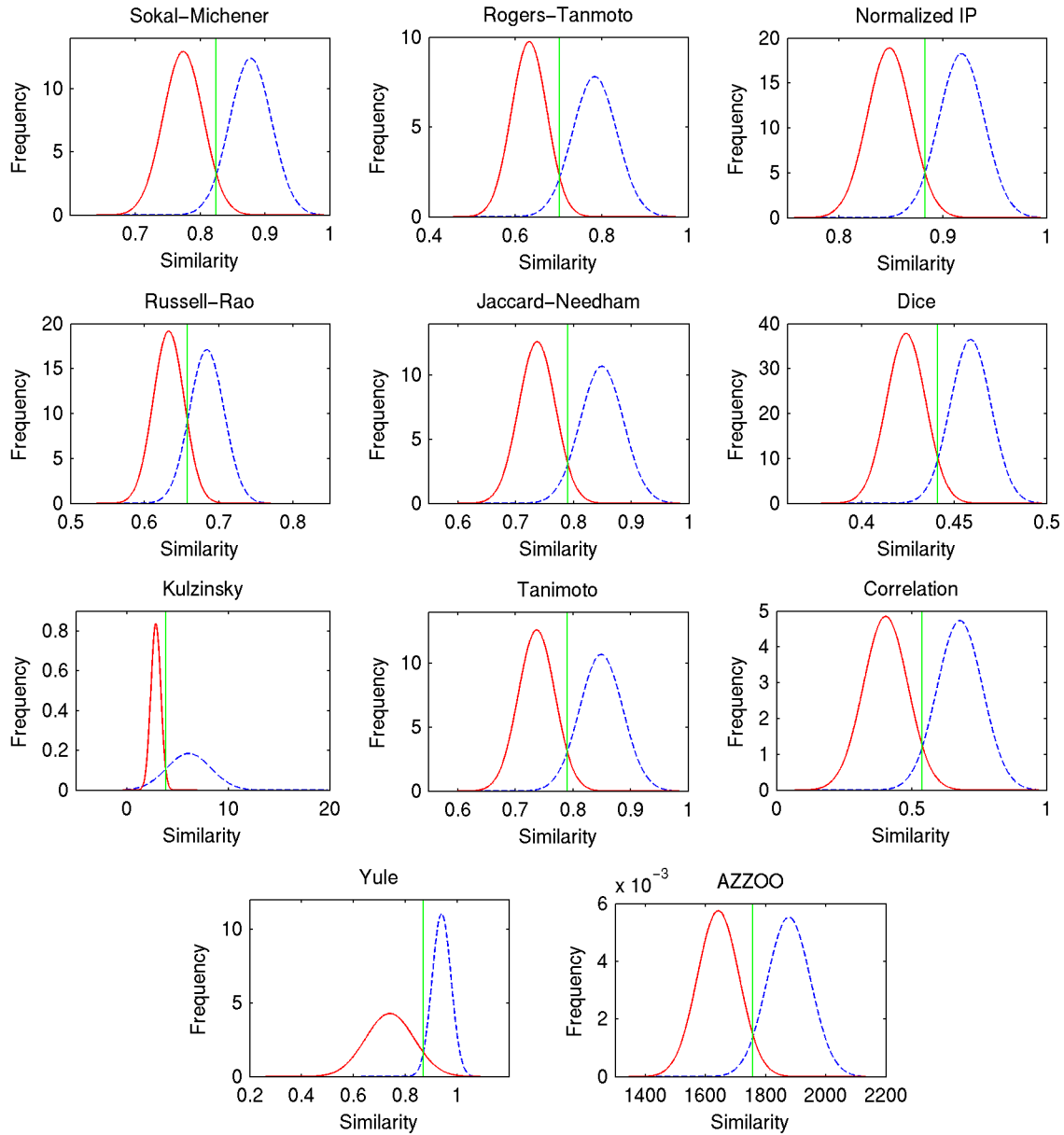


Fig. 4: Intra and inter distance distributions for the various similarity measures.

distance and inter-class distance sets. The intra-class distance sample is acquired by randomly selecting two IRIS data from the same subject while the inter-class distance sample is obtained by randomly selecting two IRIS data from two different subjects. We prepared three sets of inter and intra distance data for training and three independent ones for testing, each of size 1000 (500 intra-class and 500 inter-class pairs).

The IRIS biometric verification models were trained on 500 distance or similarity values obtained from the intra- and inter-class sets. These scalar values form distributions and the mean and variance can be computed for each distribution. Assuming normal distributions, one can easily find the Bayes decision threshold. For testing, each scalar distance value is classified into the intra or inter person class by comparing to the threshold value. First, we used the Hamming distance as Daugman originally proposed [30] and Fig. 3 shows the results on the database used in this experiment [31].

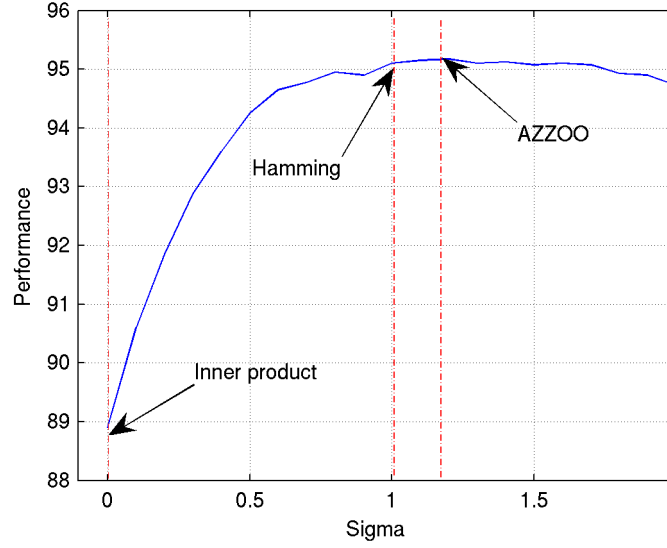


Fig. 5: Performance vs. the contribution factor σ .

We then obtained results on the other similarity measures. Fig. 4 depicts the intra and inter similarity distributions using various similarity measures and Table 1 shows the comparative results of the overall performances. Finally, Fig. 5 shows the performance as a function of the contribution factor, σ , and highlights the relative performance of the inner product, Hamming, and AZZOO measures. The S_{AZZOO} with $\sigma=1.175$ yields the best performance.

4. Similarity Measure Evaluation on Character Recognition

To further evaluate the binary vector similarity measures, we consider an offline handwritten character image database.

4.1 Binary Feature Extraction

Among many features, the Gradient, Structural, and Concavity (GSC) feature set has been shown to have high accuracy in offline character recognition problems [12] based on the philosophy that feature sets can be designed to extract certain types of information from the image. These types are gradient, structural, and concavity information. Gradient features use the stroke shapes on a small scale, structural features use stroke trajectories on an intermediate scale, and concavity features use stroke relationships at long distances.

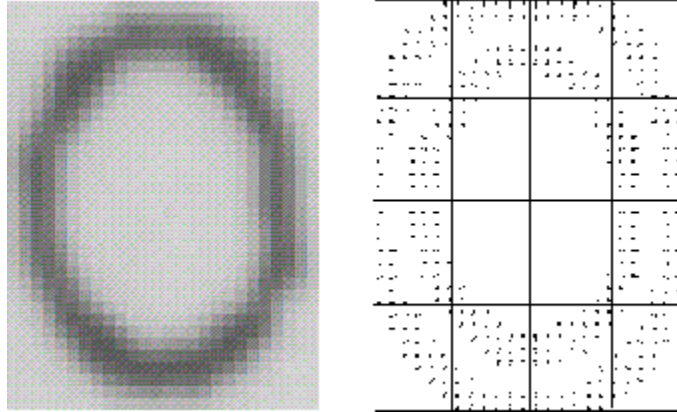


Fig. 6: Character recognition 4x4 grid.



```

Gradient      000000000011000000001100001110000000111000000011000000110
(192bits)    001000000001100000000000001110011000111110000111100000000
              100101000001000111001111100111110000010000010000000000000
              000000001000001001000

Structural    000000000000000000001100001110001000010000100000010000000
(192bits)    000000100101000000000011000010100110000110000000000000100
              100011001100000000000000110010100000000000011000000000000
              000000000000000010000

Concavity     111101101001111101100110000001101111011010011001000001100
(128bits)    000111000000000000000000000000000000000000000000000111111000
              000000000000000
    
```

Fig. 7: A sample character and its GSC feature vector.

The input character image is a binarized and slant-normalized image. A bounding box is placed around the image and divided into a 4 x 4 grids, which is known as a quasi-multiresolution approach shown in Fig. 6. For each grid region, all directional rules and various concavity features are checked, resulting in 192 Gradient, 192 Structural, and 128 Concavity features, for a total of 512 features as listed in Table 2. A sample vector for a character “A” is given in Fig. 7. See [12] for a detailed description of the rules.

4.2 Experimental Results

The problem of offline handwritten character recognition is to classify an unknown handwritten character image as one of the 26 letters of the alphabet. There are 800 samples per letter of 512 binary feature vectors in the database: 400 samples per letter are used for the reference set and

Table 2: GSC features where $x=0...3$ and $y=0...3$.

Grid Pos.	Gradient		Structural		Concavity Features	
	ID	Directional	ID	Rule	ID	Concavity
(0,0)	G01-00	1°~30°	S01-00	r_1	C-CP-00	Coarse pixel density
...
(3,3)	G01-33	1°~30°	S01-33	r_1	C-CP-33	Coarse pixel density
(x,y)	G02-xy	31°~60°	S02-xy	r_2	C-HR-xy	Horizontal run length
(x,y)	G03-xy	61°~90°	S03-xy	r_3	C-VR-xy	Vertical run length
(x,y)	G04-xy	91°~120°	S04-xy	r_4	C-UC-xy	Upward concavity
(x,y)	G05-xy	121°~150°	S05-xy	r_5	C-DC-xy	Downward concavity
(x,y)	G06-xy	151°~180°	S06-xy	r_6	C-LC-xy	Left concavity
(x,y)	G07-xy	181°~210°	S07-xy	r_7	C-RC-xy	Right concavity
(x,y)	G08-xy	211°~240°	S08-xy	r_8	C-HC-xy	Hole concavity
(x,y)	G09-xy	241°~270°	S09-xy	r_9		
(x,y)	G10-xy	271°~300°	S10-xy	r_{10}		
(x,y)	G11-xy	301°~330°	S11-xy	r_{11}		
(x,y)	G12-xy	331°~360°	S12-xy	r_{12}		

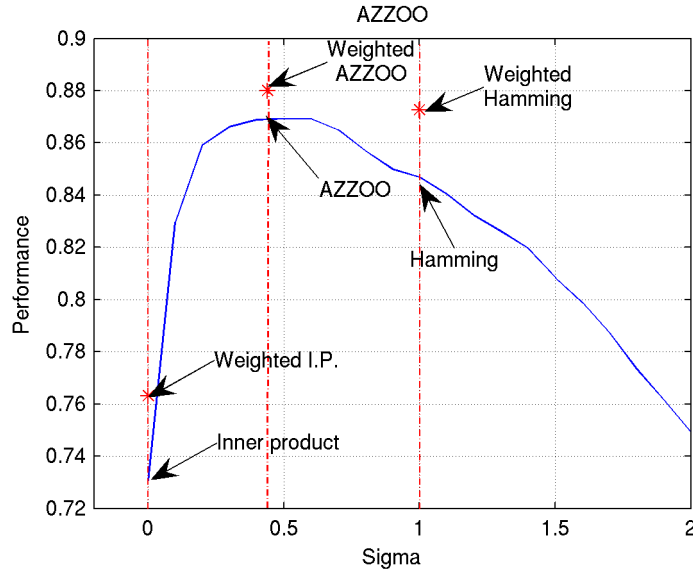


Fig. 8: Performance vs. the contribution factor σ .

the remaining samples are divided into four sets of 100 samples per letter for testing and tuning purposes. The k -nearest neighbor (k -NN) approach is used, and Table 3 shows the number of errors by each of the similarity measures tested. The numbers in parentheses are the errors after optimizing the weights in the weighted variations of the measures. Among the numerous similarity measures without weights, AZZOO performed the best. Fig. 8 shows the performance as a function of the contribution factor, σ , and stresses the relative performance of the inner product, Hamming, and AZZOO measures, clearly showing the superiority of the AZZOO measure for $\sigma=0.44$.

Applying different weights further improves the performance. Note that the number of weights is $d=512$ in the weighted Hamming distance and $2 \times d=1024$ in the weighted AZZOO. We use a genetic algorithm to determine the weights and the weighted AZZOO significantly outperforms the weighted Hamming distance.

We previously conducted similar optimizing experiments with fewer weights used [14]. Since the number of weights is enormous when each feature is given its own weight, in those

Table 3: Similarity measures and their errors on handwritten character recognition.

Category	Measure	Error
Hamming	Azzoo	330 (312)
	Hamming	398 (331)
	Sokal-Michener	398
	Roger-Tanimoto	398
Inner Product	Inner-product	700 (616)
	Russell-Rao	700
	Normalized I.P.	343
	Jaccard-Needham	341
	Dice	341
	Kulzinsky	341
	Tanimoto	341
	Correlation	347
Correlation	Yule	474

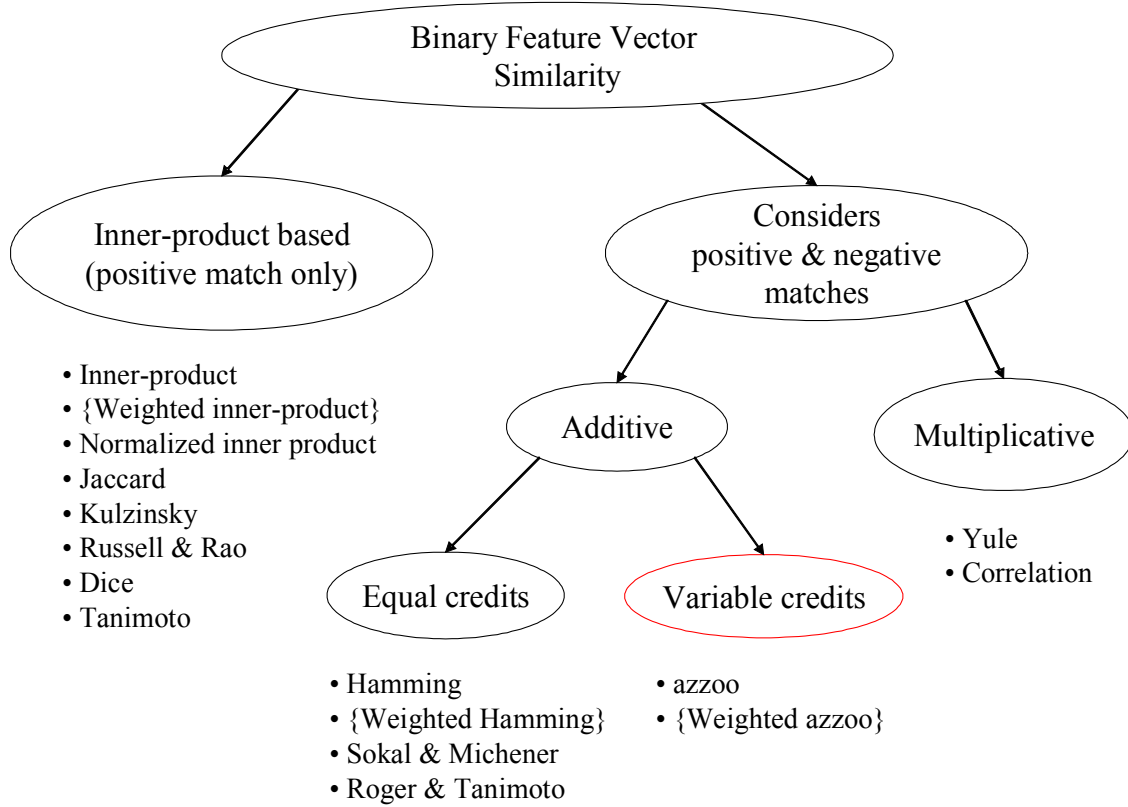


Fig. 9: Taxonomy of binary feature vector similarity measures.

experiments we simplified the features into three feature groups (gradient, structural, and concavity feature sets), and then used the additive model to combine these measures as follows:

$$S(x, y) = S^g(x_g, y_g) + S^s(x_s, y_s) + S^c(x_c, y_c) \quad (24)$$

where x_g , x_s , and x_c are the gradient, structural, and concavity subsets of x , so that $x = x_g \cup x_s \cup x_c$. First, we consider each individual set of GSC features. For $S^g(x_g, y_g)$, $S^s(x_s, y_s)$ and $S^c(x_c, y_c)$ we use the AZZOO similarity measure.

Next we associate weights with each feature group as shown in (25).

$$S_{WAZZOO}(x, y) = w_{g11} x_g y_g + w_{g00} \bar{x}_g \bar{y}_g + w_{s11} x_s y_s + w_{s00} \bar{x}_s \bar{y}_s + w_{c11} x_c y_c + w_{c00} \bar{x}_c \bar{y}_c \quad (25)$$

where w_{g11} , w_{g00} , w_{s11} , w_{s00} , w_{c11} , and w_{c00} are the weights for the Gradient, Structural, and Concavity feature groups. Optimizing these six coefficients, we found the number of errors to be 319 which, although an improvement over the AZZOO measure without weights, is not as good as the 312 errors when the full set of weights is used.

Conclusions

To conclude, we emphasize that selecting and designing a similarity measure is extremely important. First, we reviewed and categorized ten different binary feature vector similarity measures. Conventional similarity coefficients were categorized into three groups depending on their type: inner-product, Hamming, and correlation based similarity measures as depicted in Fig. 9. The first major division is between the inner-product based similarity measures that consider

positive matches only and those that credit both positive and negative matches. Next, those that consider both positive and negative matches are further categorized into additive forms and multiplicative forms (correlation based measures).

Patterns can be analyzed by either distance or similarity. Pattern classification or clustering using Hamming distance will have the identical results as those using Sokal and Michener's normalized Hamming similarity measure. From the point of view of distance, positive and negative matches are treated equally. By first converting the Hamming distance into a similarity measure, we derived another similarity measure that distinguishes the positive and negative matches, and we called it the AZZOO similarity measure. In our version of a taxonomy, the AZZOO similarity measure is under the additive form of similarity measures that take both positive and negative matches into accounts.

We showed that the AZZOO measure outperforms all conventional measures in the applications of IRIS biometric verification and offline handwritten character recognition. While the AZZOO measure is superior to the other measures, it is interesting to note that the value of the contributing factor σ can vary considerably depending on the application data – in this case the optimal value was 1.175 on the IRIS data and 0.44 on the handwriting data.

Moreover, we explored enhancing the similarity measures by applying weights that can be optimized to specific application data. While the weighted Hamming similarity measure gives identical weights to both positive and negative matches, we demonstrated that the weighted AZZOO similarity measure that gives different weights to positive and negative matches can further improve the discrimination performance.

References

- [1] B. V. Dasarathy, "Visiting nearest neighbors - a survey of nearest neighbor pattern classification techniques," in *Proc. IEEE Int. Conf. on Cybernetics and Society*, Sep. 1977, pp. 630–636.
- [2] B. V. Dasarathy, *Nearest Neighbor Pattern Classification Techniques*. IEEE Computer Society Press, 1991.
- [3] J. R. Smith and S. -F. Chang, "Automated binary texture feature sets for image retrieval," *Int. Conf. Acoustics, Speech, Signal Processing*, Atlantic, GA, May 1996.
- [4] P. Willett, J. M. Barnard, and G. M. Downs, "Chemical similarity searching," *Journal of Chemical Information and Computer Sciences*, vol. 38, pp. 983-996, 1998.
- [5] L. C. Cole, "The measurement of partial interspecific association," *Ecology*, vol. 38, pp. 226-233, 1957.
- [6] R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification*. 2nd ed., John Wiley & Sons, Inc., 2000.
- [7] P. H. A. Sneath and R. R. Sokal, *Numerical Taxonomy*. London: Freeman, 1973.
- [8] D. H. T. Clifford and W. Stephenson, *An Introduction to Numerical Classification*. New York: Academic, 1975.
- [9] R. V. Hamming, "Error detecting and error correcting codes," *Bell Sys. Tech. Journal*, vol. 29, pp. 147-160, 1950.
- [10] J. D. Tubbs, "A note on binary template matching," *Pattern Recognition*, vol. 22, no. 4, pp. 359-365, 1989.
- [11] B. Zhang and S. N. Srihari, "Binary vector dissimilarities for handwriting identification," in *Proc. of SPIE, Document Recognition and Retrieval X*, 2003, pp. 150-166.
- [12] J. T. Favata and G. Srikantan, "A multiple feature / resolution approach to handprinted digit and character recognition," *International Journal of Imaging Systems and Technology*, vol. 7, pp. 304–311, 1996.
- [13] S. -H. Cha and S. N. Srihari, "A fast nearest neighbour search algorithm by filtration," *Pattern Recognition*, vol. 35, pp. 515-525, 2000.
- [14] S. -H. Cha and C. C. Tappert, "Optimizing binary feature vector similarity measure using genetic algorithm," *ICDAR*, Edinburgh, Scotland, 2003.

- [15] M. Mitchell, *An introduction to Genetic Algorithms*. Cambridge, MA: MIT Press, 1996.
- [16] L. Davis, *Handbook of Genetic Algorithms*. New York: Van Nostrand Reinhold, 1991.
- [17] G. Dunn and B. S. Everitt, *An Introduction to Mathematical Taxonomy*. Cambridge University Press, 1982.
- [18] R. R. Sokal and C. D. Michener, "A statistical method for evaluating systematic relationships," *University of Kansas Scientific Bulletin* 38, pp. 1409-1438, 1958.
- [19] D. J. Rogers and T. T. Tanimoto, "A computer program for classifying plants," *Science*, vol. 132, pp. 1115-1118, 1960.
- [20] P. R. Halmos, *An Introduction to Hilbert Spaces and the Theory of Spectral Multiplicity*. Chelsea Publishing, 2nd ed, 1957.
- [21] P. F. Russell and T. R. Rao, "On habitat and association of species of anopheline larvae in south-eastern Madras," *J. Malar. Inst. India*, vol. 3, pp. 153-178, 1940.
- [22] P. Jaccard, "Nouvelles recherches sur la distribution florale," *Bulletin de la Societe Vaudoise de Science Naturelle*, vol. 44, pp. 223-270, 1908.
- [23] L. R. Dice, "Measures of the amount of ecologic association between species," *Ecology*, vol. 26, pp. 297-302, 1945.
- [24] G. U. Yule and M. G. Kendall, *An Introduction to the Theory of Statistics*. 14th ed. Hafner, New York, pp. 701, 1950.
- [25] C. Gose, R. Johnsonbaugh, and S. Jost, *Pattern Recognition and Image Analysis*. Prentice Hall, Inc., p 172, 1996.
- [26] B. F. Wu, "Comparative performance evaluation of some techniques for ranking pattern recognition features," Ph.D. Dissertation, Bioengineering Program, University of Illinois at Chicago, 1977.
- [27] W. Pratt, P. Capitani, W. Chen, E. Hamilton, and R. Willis, "Combining symbol matching facsimile data compression system," in *Proc. IEEE*, vol. 68, pp. 786-796, 1980.
- [28] I. Witten, A. Moffat, and T. Bell, *Managing Gigabytes: Compressing and Indexing Documents and Images*. Van Nostrand Reinhold, 1994.
- [29] S. Mahamud and M. Hebert, "The optimal distance measure for object detection," *IEEE Computer Vision and Pattern Recognition*, Wisconsin, pp. 248-256, 2003.
- [30] J. G. Daugman, "High confidence visual recognition of persons by a test of statistical independence," *IEEE Tran. Pattern Analysis and Machine Intelligence*, vol. 15, no. 11, pp. 1148-1161, 1993.
- [31] G. Kee, Y. Byun, K. Lee and Y. Lee, "Improved techniques for an IRIS recognition system with high performance," *Advances in Artificial Intelligence, LNCS*, vol. 2256, pp. 177-184, 2001.