

Validating a Biometric Authentication System: Sample Size Requirements

Sarat Dass*, Yongfang Zhu*, and Anil Jain*

Abstract—Authentication systems based on biometric features (e.g., fingerprint impressions, iris scans, human face images, etc.) are increasingly gaining widespread use and popularity. Often, vendors and owners of these commercial biometric systems claim impressive performance that is estimated based on some proprietary data. In such situations, there is a need to independently validate the claimed performance levels. System performance is typically evaluated by collecting biometric templates from n different subjects, and for convenience, acquiring multiple instances of the biometric for each of the n subjects. Very little work has been done in (i) constructing confidence regions based on the ROC curve for validating the claimed performance levels, and (ii) determining the required number of biometric samples needed to establish confidence regions of pre-specified width for the ROC curve. To simplify the analysis that address these two problems, several previous studies have assumed that multiple acquisitions of the biometric entity are statistically independent. This assumption is too restrictive and is generally not valid. We have developed a validation technique based on multivariate copula models for correlated biometric acquisitions. Based on the same model, we also determine the minimum number of samples required to achieve confidence bands of desired width for the ROC curve. We illustrate the estimation of the confidence bands as well as the required number of biometric samples using a fingerprint matching system that is applied on samples collected from a small population.

Index Terms—Biometric authentication, Error estimation, Gaussian copula models, bootstrap, ROC confidence bands.

I. INTRODUCTION

The purpose of a biometric authentication system is to validate the claimed identity of a user based on his/her physiological characteristics. In such a system operating in the verification mode, we are interested in accepting queries which are “close” or “similar” to the template of the claimed identity, and rejecting those that are “far” or “dissimilar”. Suppose a user with true identity I_t supplies a biometric query Q and a claimed identity I_c . We are interested in testing the hypothesis

$$H_0 : I_t = I_c \quad \text{vs.} \quad H_1 : I_t \neq I_c \quad (1)$$

Manuscript received September 3, 2004; revised April 1, 2006.

Sarat Dass and Yongfang Zhu are in the Department of Statistics & Probability at Michigan State University. Address: A-430 Wells Hall, E Lansing, MI 48824. E-mail: {sdass,zhuyongf}@msu.edu. Phone: 517-355-9589. Fax: 517-432-1405. Anil Jain is in the Department of Computer Science & Engineering at Michigan State University. Address: 3115 EB, E Lansing, MI 48824. E-mail: jain@cse.msu.edu. Phone: 517-355-9282. Fax: 517-432-1061.

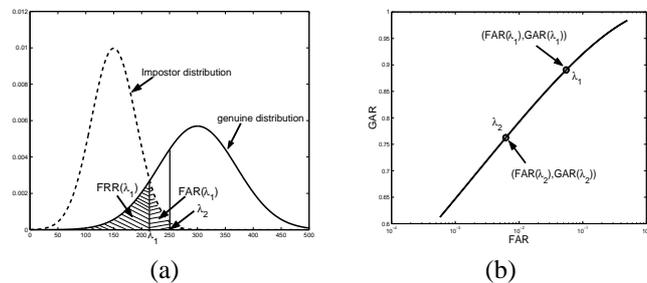


Fig. 1. Obtaining the ROC curve by varying the threshold λ . Panel (a) shows the FRR and FAR corresponding to a threshold λ_1 . λ_2 is another threshold different from λ_1 . Panel (b) shows the ROC curve obtained when λ varies. The values of (FAR, GAR) on the ROC curve corresponding to the thresholds λ_1 and λ_2 are shown.

based on the query Q and the template T of the claimed identity in the database; in Equation (1), H_0 (respectively, H_1) is the null (alternative) hypothesis that the user is genuine (impostor). The testing in (1) is carried out by computing a similarity measure, $S(Q, T)$ where large (respectively, small) values of S indicate that T and Q are close to (far from) each other. A threshold, λ , is specified so that all similarity values lower (respectively, greater) than λ lead to the rejection (acceptance) of H_0 . Thus, when a decision is made whether to accept or reject H_0 , the testing procedure (1) is prone to two types of errors: the false reject rate (FRR) is the probability of rejecting H_0 when in fact the user is genuine, and the false accept rate (FAR) is the probability of accepting H_0 when in fact the user is an impostor. The genuine accept rate (GAR) is $1 - FRR$, which is the probability that the user is accepted given that he/she is genuine. Both the FRR (and hence GAR) and the FAR are functions of the threshold value λ (see Figure 1 (a)). The Receiver Operating Curve (ROC) is a graph that expresses the relationship between the FAR versus GAR when λ varies, that is,

$$ROC(\lambda) = (FAR(\lambda), GAR(\lambda)), \quad (2)$$

and is commonly used to report the performance of a biometric authentication system (see Figures 1 (a) and (b)).

In marketing commercial biometric systems, it is often the case that error rates are either not reported or poorly reported (i.e., reported without giving details on how it was determined). In a controlled environment such as in laboratory experiments, one may achieve very high

accuracies when the underlying biometric templates are of very good quality. However, these accuracies may not reflect the true performance of the biometric system in real field applications where uncontrolled factors such as noise and distortions can significantly degrade the system's performance. Thus, the problem we address in this paper is the validation of a claimed ROC curve, $ROC_c(\lambda)$, by a biometric vendor. Of course, reporting just $ROC_c(\lambda)$ does not give the complete picture. One should also report as much information as one can about the underlying biometric samples, such as the quality, the sample acquisition process, sample size as well as a brief description of the subjects themselves. If the subjects used in the experiments for reporting $ROC_c(\lambda)$ are not representative of the target population, then $ROC_c(\lambda)$ is not very useful. But assuming that the underlying samples are representative and can be replicated by other experimenters under similar conditions, one can then proceed to give margins of errors for validating $ROC_c(\lambda)$.

The process of obtaining biometric samples usually involves selecting n individuals (or, subjects) and using c different biometric instances or entities¹ from each individual. Additional biometric samples can be obtained by sampling each biometric multiple times, d , over a period of time. It is well known that multiple acquisitions corresponding to each biometric exhibit a certain degree of dependence (or, correlation); see, for example, [1], [3], [10], [16]–[19]). There have been several earlier efforts to validate the performance of a biometric system based on multiple biometric acquisitions. Bolle et al. [4] first obtained confidence intervals for the FRR and FAR assuming that the multiple biometric acquisitions were independent of each other. To account for correlation, Bolle et al. [2], [3] introduced the subsets bootstrap approach to construct confidence intervals for the FAR, FRR and the ROC curve. Schuckers [16] proposed the beta-binomial family to model the correlation between the multiple biometric acquisitions as well as to account for varying FRR and FAR values for different subjects. He showed that the beta-binomial model gives rise to extra variability in the FRR and FAR estimates when correlation is present. However, a limitation of this approach is that it models correlation for a single threshold value. Thus, this method cannot be used to obtain a confidence region for the entire ROC curve. Further, Schucker's approach is strictly model-based; inference drawn from this model may be inappropriate when the true underlying model does not belong to the beta-binomial family.

To construct confidence bands for the ROC curve, Bolle et al. [3] select T threshold values, $\lambda_1, \lambda_2, \dots, \lambda_T$ and compute the 90% confidence intervals for the associated FARs and GARs. At each threshold value λ_i , combining these 90% confidence intervals results in a

confidence rectangle for $ROC(\lambda_i)$ (see (2)). Repeating this procedure for each $i = 1, 2, \dots, T$ and combining the confidence rectangles obtained gives rise to a confidence region for $ROC(\lambda)$. A major limitation of this approach is that the 90% confidence intervals for the FARs and GARs will neither automatically guarantee a 90% confidence rectangle at each λ_i nor a 90% confidence region for the ROC curve. In other words, ensuring a confidence level of 90% for each of the individual intervals cannot, in general, ensure a specific confidence level for the combined approach. This is the well-known problem of combining evidence from simultaneous hypothesis testing scenarios [9], [11], [12]: In essence, for each i , we are performing the tests

$$H_{0,i} : FAR(\lambda_i) = FAR_c(\lambda_i) \quad \text{vs.} \quad H_{1,i} : \text{not } H_{0,i}, \quad (3)$$

and

$$H_{0,i}^* : GAR(\lambda_i) = GAR_c(\lambda_i) \quad \text{vs.} \quad H_{1,i}^* : \text{not } H_{0,i}^*, \quad (4)$$

where $FAR(\lambda_i)$ (respectively, $FAR_c(\lambda_i)$) are the true but unknown (respectively, claimed) FAR at λ_i , and $GAR(\lambda_i)$ (respectively, $GAR_c(\lambda_i)$) are the true but unknown (respectively, claimed) GAR at λ_i . To test each $H_{0,i}$ (and $H_{0,i}^*$) individually, the 90% confidence interval for FAR (and GAR) can be used, and the resulting decision has a FRR of at most $100 - 90 = 10\%$. The confidence region for the ROC curve combines the $2T$ confidence intervals above and is used to test the hypothesis

$$H_0 : \cap_{i=1}^T \{ H_{0,i} \cap H_{0,i}^* \} \quad \text{versus} \quad H_1 : \text{not } H_0. \quad (5)$$

However, the combined confidence region is not guaranteed to have a confidence level of 90%. In other words, the decision of whether to accept or reject H_0 does not have an associated FRR of 10% as in the case of the individual hypotheses. In fact, for a number α where $0 < \alpha < 1$, combining $2T$ $100(1 - \alpha)\%$ level confidence intervals based on a-priori selected thresholds can only guarantee a lower bound of $100(1 - 2T\alpha)\%$ on the confidence level. This fact is based on Bonferroni's inequality, and is well-known in the statistics literature. Instead of trying to derive this inequality, we point the reader to the relevant literature in statistics on simultaneous hypotheses testing procedures; see, for example, the following references [9], [11], [12]. The lower bound $100(1 - 2T\alpha)\%$ on the confidence level is not useful when T is large; in this case, $100(1 - 2T\alpha)\%$ is negative, and we know that any confidence level should range between 0% and 100%. In Bolle et al.'s procedure, the value of T is large since the confidence rectangles are reported at various locations of the entire ROC curve.

In this paper, we present a new approach for constructing confidence regions for the ROC curve with a guaranteed pre-specified confidence level. In fact, we are able to construct confidence regions for a *continuum*

¹By entities we mean different fingers from each individual, or iris images from the left and right eyes from each individual, etc.

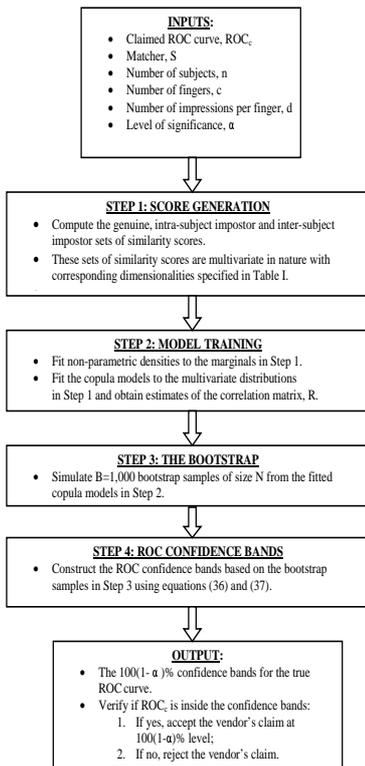


Fig. 2. The main steps involved in constructing the ROC confidence bands for validating the claim of a fingerprint vendor.

of threshold values, and not just for finite pre-selected threshold values. In contrast to the non-parametric bootstrap approach of [3], we develop a semi-parametric approach for constructing confidence regions for $\text{ROC}(\lambda)$. This is done by estimating the genuine and impostor distributions of similarity scores obtained from multiple biometric acquisitions of the n subjects where the marginals are first estimated non-parametrically (without any model assumptions), and then coupled together to form a multivariate joint distribution via a parametric family of Gaussian copula models [13]. The parametric form of the copula models enables us to investigate how correlation between the multiple biometric acquisitions affects the confidence regions. Confidence regions for the ROC are constructed using bootstrap re-samples from our estimated semi-parametric model. The main steps of our procedure are shown in Figure 2. Note that our approach based on modeling the distribution of similarity scores is fundamentally different from that of [16], where binary (0 and 1) observations are used to construct confidence intervals for the FRRs and FARs.

Our approach also varies from that of [1], [3], [10],

[16] in several respects. First, we explicitly model the correlation via a parametric copula model, and thus, are able to demonstrate the effects of varying the correlation on the width of the ROC confidence regions. We also obtain a confidence *band*, rather than confidence rectangles as in [3], consisting of upper and lower bounds for the ROC curve. Further, the confidence bands come with a guaranteed confidence level for the *entire* ROC in the region of interest. Thus, we are able to perform tests of significance for the ROC curve and report error rates corresponding to our decision of whether to accept or reject the claimed ROC curve.

Another important issue that we address is that of the test sample size: How many subjects and how many biometric acquisitions per subject should be considered in order to obtain a confidence band for the ROC with a pre-specified width? Based on the multivariate Gaussian copula model for correlated biometric acquisitions, we give the minimum number of subjects required to achieve the desired width. In presence of non-zero correlation, increasing the number of subjects is more effective in reducing the width of the confidence band compared to increasing the number of biometric acquisitions per subject. For achieving the desired confidence level, the required number of subjects based on our method is much smaller compared to the subset bootstrap. Rules of thumb such as the Rule of 3 [20] and the Rule of 30 [14] grossly underestimate the number of users required to obtain a specific width. The underestimation becomes more severe as the correlation between any two acquisitions of a subject increases.

The paper is organized as follows: Section II presents the problem formulation. Section III discusses the use of multivariate copula functions to model the correlation between multiple queries per subject for the genuine and impostor similarity score distributions. Section IV presents the construction of confidence bands for the ROC curve. Section V discusses the minimum number of biometric samples required for obtaining confidence bands of a pre-specified width for the ROC curve. Some of the more technical details and experimental results have been moved to the Appendix due to space restrictions; interested readers can also refer to the paper [6] which incorporates the relevant details into appropriate sections of the main text.

II. PRELIMINARIES

Suppose we have n subjects available for validating a biometric authentication system. Often, during the data collection stage, multiple biometric entities (e.g., different fingers) from the same subject are used. We denote the number of biometric entities used per subject by c . To obtain additional data, each biometric of a subject is usually sampled a multiple number of times, d , over a period of time. Thus, at the end of the data collection stage, we acquire a total of ncd biometric samples from the n subjects. This collection of ncd

biometric samples will be denoted by \mathcal{B} . To obtain similarity scores, a pair of biometric samples, B and B' with $B \neq B'$, are taken from \mathcal{B} and a matcher S is applied to them, resulting in the similarity score $S(B, B')$. We will consider asymmetric matchers for S in this paper: The matcher S is asymmetric if $S(B, B') \neq S(B', B)$ for the pair of biometric samples (B, B') (a symmetric matcher implies that $S(B, B') = S(B', B)$).

In the subsequent text, we will use a fingerprint authentication system as the generic biometric system that needs to be validated. Thus, the c different biometric entities will be represented as c different fingers from each subject, and the d acquisitions will be represented by d impressions of each finger. When B and B' are multiple impressions of the same finger from the same user, the similarity score $S(B, B')$ is termed as a genuine similarity score, whereas when B and B' are impressions from either (i) different fingers from the same subject, or (ii) different subjects, the similarity score $S(B, B')$ is termed as an impostor score. The impostor scores arising from (i) (respectively, (ii)) are termed as the intra-subject (respectively, inter-subject) impostor scores.

We give some intuitive understanding of why similarity scores arising from certain pairs of fingerprint impressions in \mathcal{B} are correlated (or, dependent). During the fingerprint acquisition process, multiple impressions of a finger are obtained by successive placement of the finger onto the sensor. Therefore, given the first impression, B , and two subsequent impressions B_1 and B_2 , the similarity scores $S(B, B_1)$ and $S(B, B_2)$ are most likely going to be correlated. Further, the fingerprint acquisition process is prone to many different types of uncontrollable factors such as fingertip pressure, fingertip moisture and skin elasticity factor. These factors cause some level of dependence between fingerprint impressions of two different fingers of the same user. If this is the case, then we expect to see some level of correlation between the similarity scores $S(B_1, B_2)$ where B_1 and B_2 are impressions from different fingers. Also, as noted in [3], even the scores $S(B_1, B_2)$ from different fingers of different subjects could be correlated. All these facts lead us to statistically model the correlation for similarity scores in the three major categories, namely the genuine, intra-user impostor and inter-user impostor similarity scores.

In order to develop the framework that incorporates correlation, we need to introduce some notation. We denote the set consisting of the d impressions of finger f , $f = 1, 2, \dots, c$, from subject i by $\mathcal{M}_{i,f}$. The notation $\mathcal{S}(i, j, f, f') =$

$$\{S(B_u, B_v); B_u \in \mathcal{M}_{i,f}, B_v \in \mathcal{M}_{j,f'}, B_u \neq B_v\} \quad (6)$$

represents the set of all similarity scores available from matching the fingerprint impressions of finger f from subject i and those of finger f' from subject j . Three disjoint sets of (6) are of importance, namely, the set

Entities	\mathcal{G}_i	\mathcal{I}_i	\mathcal{I}_{ij}
Dimension, K	$cd(d-1)$	$c(c-1)d^2$	c^2d^2

TABLE I

VALUES OF K FOR THE DIFFERENT SETS $\mathcal{G}_i, \mathcal{I}_i$ AND \mathcal{I}_{ij} . HERE c IS THE NUMBER OF FINGERS AND d IS THE NUMBER OF IMPRESSIONS PER FINGER.

of genuine similarity scores (taking $i = j$ and $f = f'$ in (6)), the set of intra-subject impostor scores ($i = j$ and $f \neq f'$), and the set of inter-subject impostor scores ($i \neq j$). We denote the genuine, intra-subject impostor and inter-subject impostor score sets by

$$\begin{aligned} \mathcal{G}_i &\equiv \bigcup_{f=1}^c \mathcal{S}(i, i, f, f), & \mathcal{I}_i &\equiv \bigcup_{f=1}^c \bigcup_{\substack{f'=1 \\ f' \neq f}}^c \mathcal{S}(i, i, f, f'), \\ \text{and } \mathcal{I}_{ij} &\equiv \bigcup_{f=1}^c \bigcup_{f'=1}^c \mathcal{S}(i, j, f, f') \end{aligned} \quad (7)$$

where $i \neq j$, respectively.

We give the cardinality or dimension (the number of possibly distinct similarity scores) of each of the sets discussed above. The dimensions of $\mathcal{G}_i, \mathcal{I}_i$ and \mathcal{I}_{ij} are $cd(d-1)$, $c(c-1)d^2$ and c^2d^2 , respectively, when the matcher S is asymmetric. In all of these scenarios, we will denote the dimension corresponding to each set by K (see Table I). The total number of sets of similarity scores arising from the genuine, intra- and inter-impostor cases will be denoted by N ; we have that $N = n$, $N = n$ and $N = n(n-1)$, respectively, for the total number of sets of genuine, intra-subject impostor and inter-subject scores.

When the matcher S is symmetric, the dimension associated with each of the genuine, intra-subject impostor and inter-subject impostor sets of similarity scores gets reduced since many of the similarity scores in each of the three sets will be identical to each other. In the subsequent text, we outline the methodology for validating a vendor's claim for an asymmetric matcher. Our methodology for constructing the ROC confidence bands for a symmetric matcher can be handled in a similar fashion, keeping in mind the reduction in dimensions of each of the three sets of similarity scores discussed above.

Subsequently, N will denote the total number of independent sets of similarity scores, and K will denote the dimension of each of these N sets. For $i = 1, 2, \dots, N$, the i -th set of similarity scores will be denoted by the K -dimensional vector

$$\underline{S}_i = (s(i, 1), s(i, 2), \dots, s(i, K))^T, \quad (8)$$

where $s(i, k)$ is the generic score corresponding to the k -th component of \underline{S}_i , for $k = 1, 2, \dots, K$.

The ordered indices $1, 2, \dots, K$ are associated to the elements of each of the sets $\mathcal{G}_i, \mathcal{I}_i$ and \mathcal{I}_{ij} defined

in (7) in the following way: Let $s(B_{f,u}, B_{f',v})$ denote the similarity score obtained when matching impression u of finger f , $B_{f,u}$, with impression v of finger f' , $B_{f',v}$. In the case of a genuine set (that is, $\underline{S}_i = \mathcal{G}_i$), the order of the genuine scores is taken as $\underline{s}(f) \equiv (s(B_{f,u}, B_{f,v}), v = 1, 2, \dots, (u-1), (u+1), \dots, d, u = 1, 2, \dots, d)$ and $\underline{S}_i = (\underline{s}(1), \underline{s}(2), \dots, \underline{s}(c))$. In the case when $\underline{S}_i = \mathcal{I}_i$, the order of the scores is taken as $\underline{s}(f, f') \equiv (s(B_{f,u}, B_{f',v}), v = 1, 2, \dots, d, u = 1, 2, \dots, d)$ and $\underline{S}_i = (\underline{s}(f, f'), f' = 1, 2, \dots, (f-1), (f+1), \dots, c, f = 1, 2, \dots, c)$. Finally, in the case when \underline{S}_i is an inter-subject impostor set (one of \mathcal{I}_{ij}), the order of the scores are taken as $\underline{s}(f, f') \equiv (s(B_{f,u}, B_{f',v}), v = 1, 2, \dots, d, u = 1, 2, \dots, d)$ and $\underline{S}_i = (\underline{s}(f, f'), f' = 1, 2, \dots, c, f = 1, 2, \dots, c)$.

If the scores $s(i, k)$ are bounded between two numbers a and b , the order preserving transformation

$$\mathcal{T}(s(i, k)) = \log \left(\frac{s(i, k) - a}{b - s(i, k)} \right) \quad (9)$$

converts each score onto the entire real line. This transformation yields better non-parametric density estimates for the marginal distribution of similarity scores. The transformed scores will be represented by the same notation $s(i, k)$. The distribution function for each \underline{S}_i will be denoted by F , that is,

$$P\{s(i, k) \leq s_k, 1 \leq k \leq K\} = F(s_1, s_2, \dots, s_K), \quad (10)$$

for real numbers s_1, s_2, \dots, s_K . Note that (i) F is a multivariate joint distribution function on R^K , and (ii) we assume that F is the common distribution function for every $i = 1, 2, \dots, N$. The distribution function F has K associated marginals; we denote the marginals by $F_k, k = 1, 2, \dots, K$, where

$$P\{s(i, k) \leq s_k\} = F_k(s_k). \quad (11)$$

III. COPULA MODELS FOR F

We propose a semi-parametric family of Gaussian copula models as models for F . Let H_1, H_2, \dots, H_K be K continuous distribution functions on the real line. Suppose that H is a K -dimensional distribution function with the k -th marginal given by H_k for $k = 1, 2, \dots, K$. According to Sklar's Theorem [13], there exists a unique function $C(u_1, u_2, \dots, u_K)$ from $[0, 1]^K$ to $[0, 1]$ satisfying

$$H(s_1, s_2, \dots, s_K) = C(H_1(s_1), H_2(s_2), \dots, H_K(s_K)), \quad (12)$$

where s_1, s_2, \dots, s_K are K real numbers. The function C is known as a K -copula function that ‘‘couples’’ the one-dimensional distribution functions $H_k, k = 1, 2, \dots, K$ to obtain H . Basically, K -copula functions are K -dimensional distribution functions on $[0, 1]^K$ whose marginals are uniform. Equation (12) can also be used to construct K -dimensional distribution function H whose marginals are the pre-specified distributions

$H_k, k = 1, 2, \dots, K$: choose a copula function C and define the function H as in (12). It follows that H is a K -dimensional distribution function with marginals $H_k, k = 1, 2, \dots, K$.

The choice of C we consider in this paper is the K -dimensional Gaussian copulas [5] given by

$$C_R(u_1, u_2, \dots, u_K) = \Phi_R^K(\Phi^{-1}(u_1), \Phi^{-1}(u_2), \dots, \Phi^{-1}(u_K)) \quad (13)$$

where each $u_k \in [0, 1]$ for $k = 1, 2, \dots, K$, $\Phi(\cdot)$ is the distribution function of the standard normal, $\Phi^{-1}(\cdot)$ is its inverse, and Φ_R^K is the K -dimensional distribution function of a normal random vector with component means and variances given by 0 and 1, respectively, and with correlation matrix R . Note that R is a positive definite matrix with diagonal entries equal to unity. The distribution function F will be assumed to be of the form (12) with $H_k = F_k$ for $k = 1, 2, \dots, K$, and $C = C_R$; thus, we have

$$F(s_1, s_2, \dots, s_K) = C_R(F_1(s_1), F_2(s_2), \dots, F_K(s_K)). \quad (14)$$

We denote the observed genuine scores by $\mathcal{S}_0 \equiv \{s_0(i, k), k = 1, 2, \dots, K_0, i = 1, 2, \dots, N_0\}$ with $K_0 = cd(d-1)$ and $N_0 = n$. Each vector $(s_0(i, 1), s_0(i, 2), \dots, s_0(i, K_0))$ is assumed to be independently distributed according to (14) with correlation matrix R_0 and marginals $F_{k,0}, k = 1, 2, \dots, K_0$. Both R_0 and the K_0 marginals are unknown and have to be estimated from the observed scores. In Section V, we show how this is done based on similarity scores obtained from a fingerprint matching system. The observed intra-subject and inter-subject impostor similarity scores are denoted by $\mathcal{S}_{11} \equiv \{s_{11}(i, k), k = 1, 2, \dots, K_{11}, i = 1, 2, \dots, N_{11}\}$ with $K_{11} = c(c-1)d^2$ and $N_{11} = n$, and $\mathcal{S}_{12} \equiv \{s_{12}(i, k), k = 1, 2, \dots, K_{12}, i = 1, 2, \dots, N_{12}\}$ with $K_{12} = c^2d^2$ and $N_{12} = n(n-1)$, respectively. Each vector $(s_{11}(i, 1), s_{11}(i, 2), \dots, s_{11}(i, K_{11}))$ (respectively, $(s_{12}(i, 1), s_{12}(i, 2), \dots, s_{12}(i, K_{12}))$) is assumed to be independently distributed according to (14) with correlation matrix R_{11} (R_{12}) and marginals $F_{k,11}, k = 1, 2, \dots, K_{11}$ ($F_{k,12}, k = 1, 2, \dots, K_{12}$). The correlation matrices R_{11}, R_{12} and the associated marginals are estimated from the observed impostor scores in the same way as is done for the genuine case. Details of the estimation procedure for the impostor case are presented in the Appendix and [6].

IV. CONFIDENCE BANDS FOR THE ROC CURVE

The Receiver Operating Curve (ROC) is a graph that expresses the relationship between the Genuine Accept Rate (GAR) and the False Accept Rate (FAR), and is used to report the performance of a biometric authentication system. For the threshold λ , the empirical GAR and FAR can be computed using the formulas

$$GAR_e(\lambda) = \frac{1}{N_0 K_0} \sum_{i=1}^{N_0} \sum_{k=1}^{K_0} I\{s_0(i, k) > \lambda\}, \quad (15)$$

and

$$FAR_e(\lambda) = \frac{1}{N_1} \left\{ \sum_{i=1}^{N_{11}} \sum_{k=1}^{K_{11}} I\{s_{11}(i, k) > \lambda\} + \sum_{i=1}^{N_{12}} \sum_{k=1}^{K_{12}} I\{s_{12}(i, k) > \lambda\} \right\}, \quad (16)$$

where $I(A) = 1$ if property A is satisfied, and 0, otherwise, and $N_1 = N_{11}K_{11} + N_{12}K_{12}$ denotes the total number of impostor scores. The true but unknown values of $GAR(\lambda)$ and $FAR(\lambda)$ are the population versions of (15) and (16); the expression for the population $GAR(\lambda)$ is given by

$$\begin{aligned} E(GAR_e(\lambda)) &= \frac{1}{N_0 K_0} \sum_{i=1}^{N_0} \sum_{k=1}^{K_0} P\{s_0(i, k) > \lambda\} \\ &= \frac{1}{K_0} \sum_{k=1}^{K_0} P\{s_0(1, k) > \lambda\} \\ &\equiv G_0(\lambda), \end{aligned} \quad (17)$$

where each set $\{s_0(i, k), k = 1, 2, \dots, K_0\}$ for $i = 1, 2, \dots, N_0$ is independent and identically distributed according to the copula model (14). Subsequently, the probabilities in (17) are functions of the unknown genuine marginal distributions, $F_{k,0}$, $k = 1, 2, \dots, K_0$, and the genuine correlation matrix, R_0 . Also, the second equality in (17) is a consequence of the identically distributed assumption. In a similar fashion, the population $FAR(\lambda)$ is given by

$$\begin{aligned} E(FAR_e(\lambda)) &= \frac{1}{N_1} \left\{ \sum_{i=1}^{N_{11}} \sum_{k=1}^{K_{11}} P\{s_{11}(i, k) > \lambda\} + \sum_{i=1}^{N_{12}} \sum_{k=1}^{K_{12}} P\{s_{12}(i, k) > \lambda\} \right\} \\ &= \frac{N_{11}}{N_1} \sum_{k=1}^{K_{11}} P\{s_{11}(i, k) > \lambda\} \\ &\quad + \frac{N_{12}}{N_1} \sum_{k=1}^{K_{12}} P\{s_{12}(i, k) > \lambda\} \\ &\equiv G_1(\lambda), \end{aligned} \quad (18)$$

where now, elements within each of the sets $\{s_{11}(i, k), k = 1, 2, \dots, K_{11}\}$ for $i = 1, 2, \dots, N_{11}$, and $\{s_{12}(i, k), k = 1, 2, \dots, K_{12}\}$ for $i = 1, 2, \dots, N_{12}$ are independent and identically distributed according to the copula model (14) with corresponding correlation matrices and marginals. The probabilities in (18) are functions of the unknown marginal distributions, $F_{k,11}$ for $k = 1, 2, \dots, K_{11}$ and $F_{k,12}$ for $k = 1, 2, \dots, K_{12}$, and the correlation matrices, R_{11} and R_{12} , for the intra-subject and inter-subject impostor scores, respectively.

In light of the notations used for the population versions of FAR and GAR, equations (15) and (16) are sample versions of $G_0(\lambda)$ and $G_1(\lambda)$. Thus, we define

$$\hat{G}_0(\lambda) \equiv GAR_e(\lambda) \quad \text{and} \quad \hat{G}_1(\lambda) \equiv FAR_e(\lambda). \quad (19)$$

The empirical ROC curve can be obtained by evaluating the expressions for GAR and FAR in (15) and (16) at various values λ based on the observed similarity scores, and plotting the resulting curve $(\hat{G}_1(\lambda), \hat{G}_0(\lambda))$. However, there is an alternative way in which an ROC curve can be constructed. Note that the ROC expresses the relationship between the FAR and GAR, and the threshold values are necessary only at the intermediate step for linking the FAR and GAR values. Thus, another representation of the ROC curve can be obtained by the following re-parameterization: we fix p as a value of FAR and obtain the threshold λ_* such that $\hat{G}_1(\lambda_*) = p$ or, $\lambda_* \equiv \hat{G}_1^{-1}(p)$. Substituting λ_* in (15) gives the ROC curve in the form $(p, \hat{W}(p))$, where

$$\hat{W}(p) = \hat{G}_0(\lambda_*) \equiv \hat{G}_0(\hat{G}_1^{-1}(p)). \quad (20)$$

Note that in the case when there is no λ_* such that $\hat{G}_1(\lambda_*) = p$, one can re-define the inverse, $\hat{G}_1^{-1}(p) \equiv \lambda_*$, where λ_* is the smallest λ satisfying $\hat{G}_1(\lambda) \leq p$. This definition of the inverse of \hat{G}_1 is more general and always yields a unique λ_* . The true but unknown ROC curve can be obtained in the same way as above by replacing the empirical versions with the corresponding population version; thus, we have

$$W(p) = G_0(G_1^{-1}(p)), \quad (21)$$

where $G_1^{-1}(p) \equiv \lambda_*$, where λ_* is the smallest λ satisfying $G_1(\lambda) \leq p$. The two representations of the ROC curves $(\hat{G}_1(\lambda), \hat{G}_0(\lambda))$ and $(p, \hat{W}(p))$, are close approximations of one another for large N_0 , and therefore we use the latter representation for deriving the confidence bands. For fixed numbers C_0 and C_1 satisfying $0 \leq C_0 < C_1 \leq 1$, let us consider all $p = FAR$ values that fall in $[C_0, C_1]$. A confidence band for the true (claimed) ROC curve of a biometric system at confidence level $100(1 - \alpha)\%$ gives two envelope functions, $e_L(p)$ and $e_U(p)$, so that for all p in $[C_0, C_1]$, the true ROC curve lies inside the interval $(e_L(p), e_U(p))$ with probability of at least $100(1 - \alpha)\%$. The numbers C_0 and C_1 form the lower and upper bounds of the range of FAR, and will be chosen to cover typical reported values of FAR in biometric applications. If $C_0 = 0$ and $C_1 = 1$, the resulting ROC confidence band is constructed for the true ROC curve for all p in $(0, 1)$.

For a specific $p = FAR$, the corresponding value of GAR, $W(p)$, is a proportion which takes values in $[0, 1]$. For proportions, the transformation

$$\sqrt{N_0}(\sin^{-1} \sqrt{\hat{W}(p)} - \sin^{-1} \sqrt{W(p)}) \quad (22)$$

is a variance stabilizing transformation [15]; the quantity in (22) is asymptotically distributed as a normal with zero mean and constant variance (independent of p and $W(p)$) for large N_0 . To obtain the envelopes, we first consider a continuum version of the absolute values of (22) for FAR values, p , in $[C_0, C_1]$, and take the

maximum over $p \in [C_0, C_1]$. This gives the statistic

$$z \equiv \max_{p: C_0 \leq p \leq C_1} \sqrt{N_0} |\sin^{-1} \sqrt{\hat{W}(p)} - \sin^{-1} \sqrt{W(p)}|. \quad (23)$$

Assume for the moment that the distribution of z is known. If $z_{1-\alpha}$ denotes the $100(1-\alpha)\%$ percentile of z , the envelopes are given by

$$e_L(p) = (\sin(\sin^{-1} \sqrt{\hat{W}(p)} - z_{1-\alpha}/\sqrt{N_0}))^2$$

and

$$e_U(p) = (\sin(\sin^{-1} \sqrt{\hat{W}(p)} + z_{1-\alpha}/\sqrt{N_0}))^2. \quad (24)$$

However, the distribution of z is difficult to obtain analytically, and thus, we present two approaches to approximate the distribution of z in (23) based on (i) the bootstrap methodology, and (ii) an asymptotic representation of the distribution of z for large N_0 .

A. The semi- and non-parametric bootstrap approaches

The value $z_{1-\alpha}$ will be found based on bootstrap samples from the fitted semi-parametric Gaussian copula models described in Section III. This bootstrap procedure requires the simulation of scores from the estimated distribution functions in (14) and is described in detail in the Appendix. Thus, we denote by $\mathcal{S}_0^* \equiv \{s_0^*(i, k), k = 1, 2, \dots, K_0, i = 1, 2, \dots, N_0\}$, $\mathcal{S}_{11}^* \equiv \{s_{11}^*(i, k), k = 1, 2, \dots, K_{11}, i = 1, 2, \dots, N_{11}\}$ and $\mathcal{S}_{12}^* \equiv \{s_{12}^*(i, k), k = 1, 2, \dots, K_{12}, i = 1, 2, \dots, N_{12}\}$ to be the sets of genuine, intra-impostor and inter-impostor similarity scores obtained by one simulation from the fitted copula models. Also let

$$W^*(p) = G_0^*(G_1^{*-1}(p)), \quad (25)$$

where $G_0^*(\lambda)$ (respectively, $G_1^*(\lambda)$) is obtained from equation (15) (respectively, (16)) with the bootstrap samples $s^*(i, k)$ used in place of the $s(i, k)$ s. We form the quantity

$$z^* \equiv \max_{C_0 \leq p \leq C_1} \sqrt{N_0} |\sin^{-1} \sqrt{W^*(p)} - \sin^{-1} \sqrt{\hat{W}(p)}|, \quad (26)$$

with $\hat{W}(p)$ and $W^*(p)$ defined as in equations (20) and (25), respectively. By repeating the above procedure a large number of times, $B^* = 1,000$, we obtain 1,000 values of z^* , $z_1^*, z_2^*, \dots, z_{1,000}^*$. The $100(1-\alpha)\%$ percentile of the distribution of z^* can be approximated by $z_{[1000(1-\alpha)]}^*$, which is the $[B^*(1-\alpha)]$ -th element in the ordered list of $z_1^*, z_2^*, \dots, z_{1000}^*$. Thus, we approximate $z_{1-\alpha}$ by $z_{[1000(1-\alpha)]}^*$.

In the non-parametric bootstrap approach, the set \mathcal{S}_0^* is obtained as follows: Sample with replacement one K_0 dimensional vector from the N_0 sets in \mathcal{S}_0 , and repeat this sampling N_0 times. The sets \mathcal{S}_{11}^* and \mathcal{S}_{12}^* , respectively, are obtained from the sets \mathcal{S}_{11} and \mathcal{S}_{12} in a similar fashion. The non-parametric bootstrap confidence bands are then constructed using the methodology outlined in the preceding paragraph.

B. An asymptotic representation of z

We approximate the distribution of z asymptotically when N_0 is large. Let $C_0 \equiv p_1 < p_2 < \dots < p_m < p_{m+1} < \dots < p_M \equiv C_1$ be a partition of the interval $[C_0, C_1]$. In the Appendix, we show that

$$z \equiv \max_{C_0 < p < C_1} \sqrt{N_0} |\sin^{-1} \sqrt{\hat{W}(p)} - \sin^{-1} \sqrt{W(p)}| \approx \max_{1 \leq m \leq M} |D_M \cdot \hat{G}_{0,M} + D_M \cdot \hat{G}_{1,M}|, \quad (27)$$

where D_M is a diagonal matrix with the (m, m) -th entry given by $1/\sqrt{4W(p_m)(1-W(p_m))}$, $D_M \cdot \hat{G}_{0,M}$ and $D_M \cdot \hat{G}_{1,M}$ are independent of each other, the distribution of $D_M \cdot \hat{G}_{0,M}$ (respectively, $D_M \cdot \hat{G}_{1,M}$) is approximately a M -dimensional multivariate normal with mean 0 (respectively, 0) and covariance matrix given by Γ_0 (respectively, Γ_1) given in equation (58) in the Appendix. The maximum in $[C_0, C_1]$ is approximated by the component of the multivariate normal that takes on the maximum absolute value. We define

$$\max_{1 \leq m \leq M} |D_M \cdot \hat{G}_{0,M} + D_M \cdot \hat{G}_{1,M}| \equiv z_M. \quad (28)$$

The distribution of z is approximated by the distribution of z_M for large M . Denoting the $100(1-\alpha)\%$ percentile of z_M by $z_{1-\alpha, M}$, the $100(1-\alpha)\%$ confidence interval for $W(p)$ is given by $(e_L(p), e_U(p))$ where

$$e_L(p) = (\sin(\sin^{-1} \sqrt{\hat{W}(p)} - z_{1-\alpha, M}/\sqrt{N_0}))^2$$

and

$$e_U(p) = (\sin(\sin^{-1} \sqrt{\hat{W}(p)} + z_{1-\alpha, M}/\sqrt{N_0}))^2. \quad (29)$$

C. Testing the claim of a biometric vendor

Suppose that a vendor of a biometric authentication system claims that his/her biometric authentication system has a ROC curve given by $ROC_c = (p, W_c(p))$, for p in some interval $[C_0, C_1]$. Based on acquisitions from n subjects, we can test the validity of this claim by generating our own genuine and impostor similarity scores, and obtaining the $100(1-\alpha)\%$ confidence band for the true ROC curve, $(p, W(p))$, for $p \in [C_0, C_1]$. We assume that the subjects as well as the scores generated from the subjects in the vendor's database are a representative sample from the underlying population of subjects and the corresponding distributions of genuine and impostor scores derived from this population. If this assumption is true, then the confidence bands constructed from the previous section can be used for validating the vendor's claim. We perform the test

$$H_0 : W(p) = W_c(p) \quad \text{vs.} \quad H_1 : W(p) \neq W_c(p), \quad (30)$$

for some p , and will accept H_0 (the claimed ROC curve) if

$$e_L(p) \leq W_c(p) \leq e_U(p) \quad (31)$$

for all $p \in (C_0, C_1)$; otherwise, we will reject it. We can also perform a test for claims of specific values of FRR

and FAR , FRR_c and FAR_c . At $p_c = FAR_c$, we obtain the upper and lower limits of $GAR(p_c)$, $GAR_L(p_c)$ and $GAR_U(p_c)$. We will accept the claimed error rates if

$$GAR_L(p_c) \leq GAR_c \leq GAR_U(p_c) \quad (32)$$

where $GAR_c = 1 - FRR_c$, and reject it otherwise.

V. EXPERIMENTAL RESULTS

We evaluate the methodology developed in the previous sections for biometric authentication systems based on fingerprints. For evaluation purposes, it is necessary that the fingerprint databases consist of multiple impressions of a finger as well as impressions from several different fingers for each subject. Many publicly available databases do not meet these requirements and as a result, we focused on two databases that were appropriate for our purpose, namely, a database consisting of fingerprint impressions collected in our laboratory, and a different database obtained from West Virginia University.

The Michigan State University (MSU) database [8] consists of fingerprint impressions from 4 different fingers (the right index, right middle, left index and left middle fingers) of 160 users. A total of 4 impressions per finger were obtained; 2 impressions were obtained on the first day and the remaining two after a period of a week. The fingerprint images were acquired using a solid state sensor manufactured by Veridicom, Inc, with image sizes 300×300 and resolution 500 dpi. Figure 3 show all 4 impressions of 3 fingers in this database. The first two fingers (first two rows) are from the same subject whereas the images in the last row are from a different subject. A fingerprint similarity score was generated using an asymmetric matcher, described in [7]. All raw scores ranged between 0 and 1000, and thus, the transformation (9) with $a = 0$ and $b = 1000$ was used to convert the scores onto the real line. All subsequent analysis was performed on the transformed similarity scores. Thus, we have the following values for N and K (with $n = 160$, $c = d = 4$): $N = 160$ and dimensionality $K = 4 \times 4 \times 3 = 48$ for the set of genuine scores, $N = 160$ and $K = 4 \times 3 \times 4^2 = 192$ for the set of intra-subject impostor scores, and $N = 160 \times 159 = 25,440$ and $K = 4^2 \times 4^2 = 256$ for the set of inter-subject impostor scores. The number of parameters in the correlation matrices that need to be estimated for the genuine, intra-subject impostor and inter-subject impostor scores are, respectively, $(48 \times 47)/2 = 1128$, $(192 \times 191)/2 = 18,336$ and $(256 \times 255)/2 = 32,640$. The number of parameters far exceeds the total number of observations in each of the three sets of scores. In order to avoid overfitting, we reduce the value of K in each case. Instead of selecting all 4 fingers, we choose only $c = 2$, namely, the right index and right middle fingers, and use the $d = 2$ impressions per finger obtained on the first day. In this case, the number of parameters that need to be estimated are 6, 28 and 120 for the genuine, intra-subject and inter-subject impostor sets of scores, respectively.

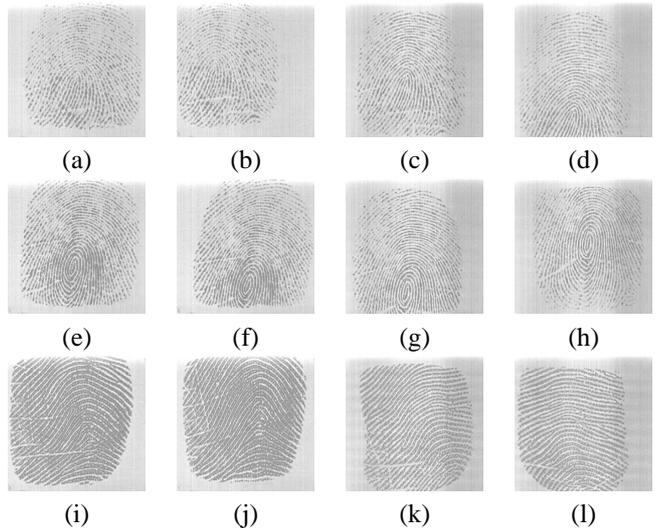


Fig. 3. Examples of fingerprint impressions from [8]: Each row gives the 4 impressions per finger collected. The first two rows are different fingers from the same subject, whereas the last row contains fingerprint impressions from a different subject.

Databases	n	c	d
MSU	160	2	2
WVU	263	1	2

TABLE II

VALUES OF n , c AND d FOR THE MSU AND WVU DATABASES USED IN THE EXPERIMENTS.

The West Virginia University (WVU) fingerprint database consists of fingerprint impressions from 263 different users. We used the first 2 impressions of the right index finger to obtain similarity scores with the same matcher as above; thus, $c = 1$ and $d = 2$ for the WVU database. Consequently, there is only one kind of impostor score, namely, the inter-subject impostor score for this database. Table II gives the number of subjects (n), as well as the values of c (number of different fingers per subject) and d (number of impressions per finger) for the MSU and WVU databases.

A. Estimating the joint distribution of similarity scores

In order to estimate the joint distribution, F , of similarity scores corresponding to the genuine, intra-subject and inter-subject impostor sets, we first need to estimate each marginal F_k , $k = 1, 2, \dots, K$ and correlation matrix R from observed data. The estimation of F_k and R are described in detail in the Appendix and in [6]. We show the results of the non-parametric estimation procedure for the first 2 marginal distributions corresponding to each of the genuine, intra-subject impostor and inter-subject impostor scores for the MSU database (see Figure 4). Note the very good agreement between the observed density histogram and the fitted density curve for each figure, especially at the tails of

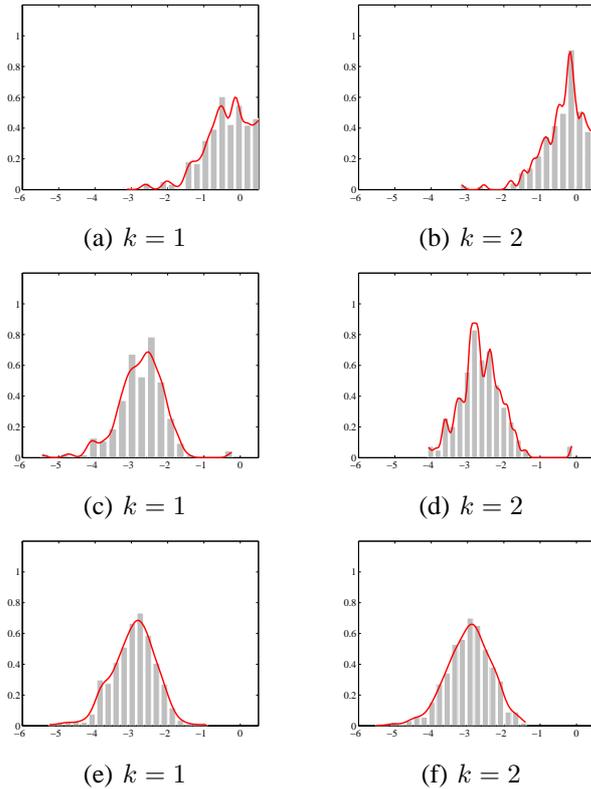


Fig. 4. Fitted density functions (solid line) for the genuine (a,b), intra-subject (c,d) and inter-subject (e,f) marginal distributions.

the distributions. A good fit at the tails is essential for the construction of a valid ROC curve that accurately reflects the authentication performance based on the observed data of similarity scores.

The estimate of the genuine correlation matrix (of dimension 4×4) is given by

$$\hat{R}_0 = \begin{pmatrix} 1.00 & 0.99 & 0.15 & 0.16 \\ 0.99 & 1.00 & 0.15 & 0.16 \\ 0.15 & 0.15 & 1.00 & 0.99 \\ 0.16 & 0.16 & 0.99 & 1.00 \end{pmatrix}. \quad (33)$$

The ordered row (and column) dimensions 1, 2, 3 and 4 respectively represents the scores $s(B_{1,1}, B_{1,2})$, $s(B_{1,2}, B_{1,1})$, $s(B_{2,1}, B_{2,2})$ and $s(B_{2,2}, B_{2,1})$; recall that $c = 2$ and $d = 2$. Consequently, the off-diagonal entries of (33) give the correlation between the corresponding row and column dimensions. For example, the entry 0.15 in the 2-nd row and 3-rd column of matrix \hat{R}_0 is the correlation between between $s(B_{1,1}, B_{1,2})$ and $s(B_{2,1}, B_{2,2})$. The off-diagonal entries of \hat{R}_0 indicate that there is a significant amount of correlation in the set of genuine similarity scores. We also obtained estimates of the intra-subject (of dimension 8×8) and inter-subject (of dimension 16×16) correlation matrices in a similar fashion (see the Appendix). We also developed an assessment of fit of the copula functions to the observed data and found that the estimated Gaussian copula functions are a good

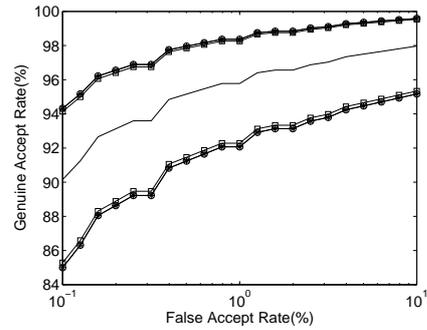


Fig. 5. Upper and lower ROC envelopes obtained using the three different methods: The non-parametric, semi-parametric bootstrap, and asymptotic envelopes are represented by the symbols \circ , \square , and $*$, respectively. The middle solid line is the non-parametric ROC curve.

fit to each of the genuine, intra-subject and inter-subject impostor sets of similarity scores. The methodology and related plots are presented in the Appendix.

B. Construction of the ROC confidence bands

The 95% ROC confidence bands are constructed based on the semi-parametric bootstrap, asymptotic and the non-parametric bootstrap approaches for the MSU and WVU databases. The resulting upper and lower bounds of all the three approaches closely match with each other for the two databases; due to space restrictions, we only show the bands for the MSU database in Figure 5. Figure 5 shows that the semi-parametric bootstrap and the asymptotic approaches give good approximations to the true upper and lower confidence bands even for moderate sample sizes.

C. Effects of correlation on the ROC confidence bands

Our next set of experiments consist of studying the effect of correlation among the multiple impressions of a user on the width of the ROC confidence band. Since this requires varying the correlation, this experiment is not possible using real data since real data would give only one estimate of correlation for each of the sets of genuine, intra-subject and inter-subject impostor similarity scores. Instead, our experiment is based on simulated sets of genuine, inter-subject impostor and intra-subject impostor similarity scores from the multivariate Gaussian K -copula models with Toeplitz forms for the correlation matrix. Let

$$R_*(\rho) = \begin{pmatrix} 1 & \rho & \rho & \cdots & \rho \\ \rho & 1 & \rho & \cdots & \rho \\ \rho & \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \rho & \cdots & 1 \end{pmatrix} \quad (34)$$

denote the correlation matrix with all off-diagonal entries equal to ρ . The dimension of $R_*(\rho)$ will be different according to whether the sets of scores are genuine, intra-subject or inter-subject impostor scores.

Sets/Estimates	$\hat{\rho}_1$	$\hat{\rho}_2$	$\dim R_*(\rho_1)$	$\dim R_*(\rho_2)$
Genuine	0.15	0.99	c	$d(d-1)$
Intra-Subject Impostor	0.80	0.27	$c(c-1)$	d^2
Inter-Subject Impostor	0.26	0.55	c^2	d^2

TABLE III

DIFFERENT VALUES OF $\hat{\rho}_1$ AND $\hat{\rho}_2$ FOR THE GENUINE, INTRA-SUBJECT IMPOSTOR AND INTER-SUBJECT IMPOSTOR SIMILARITY SCORES, AS WELL AS THE DIFFERENT DIMENSIONS OF $R(\rho_1)$ AND $R(\rho_2)$ FOR AN ASYMMETRIC MATCHER.

For a genuine set, the parameterization of the correlation matrix as $R \equiv R_*(\rho_1) \otimes R_*(\rho_2)$ implies that the correlation between any two components of $\underline{s}(f)$ corresponding to finger f is ρ_2 , and the correlation between a component of $\underline{s}(f)$ and a component of $\underline{s}(f')$ for two different fingers, $f \neq f'$, is $\rho_1 \cdot \rho_2$. For an intra-subject impostor set, the parameterization of the correlation matrix implies that the correlation between any two components of $\underline{s}(f, f')$ for each pair (f, f') is ρ_2 , and the correlation between a component of $\underline{s}(f, f')$ and a component of $\underline{s}(g, g')$ for two different pairs, $(f, f') \neq (g, g')$, is $\rho_1 \cdot \rho_2$. For an inter-subject impostor set, the parameterization implies that the correlation between any two pairs of components in $\underline{s}(f, f')$ is ρ_2 , and the correlation between a component of $\underline{s}(f, f')$ and a component of $\underline{s}(g, g')$ for two different pairs, $(f, f') \neq (g, g')$, is $\rho_1 \cdot \rho_2$.

One advantage of selecting correlation matrices to be of the form $R \equiv R_*(\rho_1) \otimes R_*(\rho_2)$ is that the matrices can be determined from specifying only two real numbers, ρ_1 and ρ_2 , and is therefore, easy to use for illustrative purposes. For a given estimated correlation matrix \hat{R} , we find the values of ρ_1 and ρ_2 that minimize the sum of Euclidean distances between the entries of \hat{R} and $R_*(\rho_1) \otimes R_*(\rho_2)$,

$$\|\hat{R} - R_*(\rho_1) \otimes R_*(\rho_2)\|^2, \quad (35)$$

where $R_*(\rho_1)$ and $R_*(\rho_2)$ are as in (34) with ρ_1 and ρ_2 plugged in for ρ , respectively, and \otimes is the Kronecker delta product. The minimizers of ρ_1 and ρ_2 , $\hat{\rho}_1$ and $\hat{\rho}_2$, for each of the genuine, intra-subject impostor and inter-subject impostor sets of scores, as well as the dimensions of each of $R_*(\rho_1)$ and $R_*(\rho_2)$ are given in Table III for the MSU database. For the WVU database, the estimated values of ρ_2 was found to be 0.99 and 0.39, respectively, for the genuine and impostor sets of similarity scores.

In order to show the effects of increasing correlation on the ROC confidence bands, four combinations of (ρ_1, ρ_2) were selected. The first three combinations are (i) $(\rho_1 = 0, \rho_2 = 0)$, (ii) $(\rho_1 = 0, \rho_2 = \hat{\rho}_2)$, and (iii) $(\rho_1 = \hat{\rho}_1, \rho_2 = \hat{\rho}_2)$, where $\hat{\rho}_1$ and $\hat{\rho}_2$ are selected according to the entries of Table III for each set of genuine, intra-subject impostor and inter-subject impostor similarity scores. The fourth combination (iv) is obtained by setting the genuine ρ_1 to 0.6 and the remaining ρ_1 s and ρ_2 s selected according to the entries in

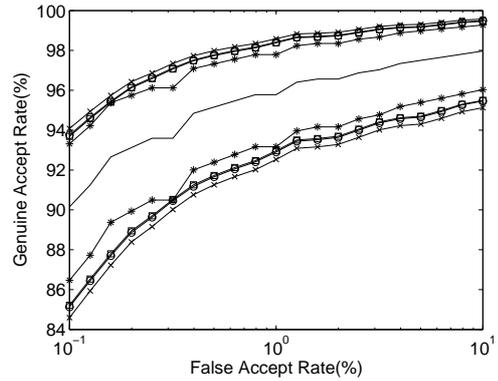


Fig. 6. Effects of correlation on the ROC confidence bands. The lines with $*$, \square , \circ and \times , respectively, denote the four different combinations of intra-finger and inter-finger correlations (i), (ii), (iii) and (iv).

Table III. The 95% ($\alpha = 0.05$) level confidence bands for the ROC curve were constructed based on $B^* = 1,000$ bootstrap resamples. Figure 6 gives the ROC confidence bands based on the semi-parametric bootstrap. Note that the width of the confidence bands generally increases as we move from combination (i) to (iv). The median widths of the confidence bands for the four combinations are 4.62, 5.41, 5.51, 6.06, respectively. The effects of correlation on the confidence bands using the asymptotic approach and for the WVU database were similar to the bootstrap approach, and therefore, are not presented here.

D. Validation of the ROC confidence bands

We conducted several tests to validate the ROC confidence bands at a specified confidence level. Recall that the $100(1 - \alpha)\%$ ROC confidence bands, by definition, cover the true ROC curve with a probability of *at least* $100(1 - \alpha)\%$ on repeated sampling from the underlying population of similarity scores. Treating the entire MSU database with $n = 160$ subjects as the underlying population, we selected a subset of 120 subjects from this population for constructing the semi-parametric bootstrap ROC confidence bands; a subset of 120 subjects (as opposed to smaller subsets of the data) is selected so that estimation of the non-parametric marginal distributions can be performed reliably. We then determined if the population ROC curve (the empirical ROC curve for the 160 subjects) was within the constructed confidence bands. This procedure was repeated 200 times (with different subsets of 120 subjects from the population of 160), and each time, we determined if the population ROC curve was within the constructed ROC confidence bands. The percentage of coverage based on this validation procedure should be at least $100(1 - \alpha)\%$. In our experiments we selected $\alpha = 0.05$ for the 95% ROC confidence bands, and obtained a coverage proportion of 99.5%. For the WVU database, validation of the ROC confidence bands was carried out with sub-samples of 198 users. The procedure of constructing the ROC

confidence bands was repeated 500 times. The empirical ROC curve (ROC curve based on the 263 users) was found to be inside the 95% confidence bands in 497 (out of the 500) trials, resulting in a coverage probability of 99.4%.

E. Sample size requirements

As correlated multiple biometric observations affect the width of the ROC confidence bands, we now proceed to determine the number of users, n^* , required by a system to report a $100(1 - \alpha)\%$ ROC confidence band with a width of at most w . We take $w = 1\%$. Our results are based on simulation with correlations selected according to combinations (i-iv) in Section V-C. Thus, the results reported here can be generalized to real data which exhibit different degrees of intra-finger and inter-finger correlations. The values of n^* are given for different combinations of c and d , and therefore, varying dimensionality of the genuine, intra-subject and inter-subject sets of similarity scores. Consequently, we assume a common marginal for each of the three sets given by the mixture over component scores. We selected $C_0 = 0.1\%$, $C_1 = 10\%$ and $M = 21$ here, and $p_m = 10^{(-1+0.1(m-1))}$, $m = 1, 2, \dots, M$. For each $m = 1, 2, \dots, M$, the width of the ROC confidence band at each $FAR = p_m$ (see equation (29)) is given by

$$\begin{aligned} w(p_m) &= e_U(p_m) - e_L(p_m) \\ &= \frac{4z_{1-\alpha, M} \sqrt{W(p_m)(1 - W(p_m))}}{\sqrt{n}} \end{aligned} \quad (36)$$

for large $n (= N_0)$, where $z_{1-\alpha, M}$ is the $100(1 - \alpha)\%$ percentile of the distribution of z_M defined in (28); the second equality is from applying the delta method [15] to $e_U(p_m) - e_L(p_m)$ in (29). In order to determine $z_{1-\alpha, M}$, we first estimate the covariance matrices Γ_0 and Γ_1 (see equation (59) in the Appendix) as accurately as possible. This estimation is performed based on 1000 simulated samples from each of the correlation combinations (i-iv) for $n = 1000$ subjects. To achieve a width of w for the confidence band implies that $w(p_m) \leq w$ for all p_m , $m = 1, 2, \dots, M$. Thus, the minimum number of users required is given by the formula $n^* = n_0 + 1$ where n_0 is the greatest integer less than or equal to

$$\max_{1 \leq m \leq M} \left(\frac{4z_{1-\alpha, M} \sqrt{W(p_m)(1 - W(p_m))}}{w(p_m)} \right)^2. \quad (37)$$

We also compare the minimum sample size requirements given by our method to that of the subset bootstrap approach [3]. One important point is that [3] obtains confidence rectangles, and not confidence bands, at each threshold value on the ROC curve. In order to perform a valid band to band comparison of the two methods, we applied the subset bootstrap procedure to the alternative parametrization of the ROC curve given in (20). As mentioned earlier, the subset bootstrap is not able to

give an overall confidence level of $100(1 - \alpha)\%$ using M individual $100(1 - \alpha)\%$ confidence intervals. To guarantee a $100(1 - \alpha)\%$ confidence level, the level of each individual confidence interval would have to be $100(1 - \alpha/M)\%$ using Bonferroni's inequality. For $m = 1, 2, \dots, M$, the minimum sample size requirement, $n_{sb}(m)$, for the m -th confidence interval can be obtained using similar asymptotic arguments as in Section IV-B with $C_0 = C_1 = p_m$. It follows that the minimum sample size required to achieve the pre-specified width for all M confidence intervals is given by

$$n_{sb}^* = \max_{1 \leq m \leq M} n_{sb}(m). \quad (38)$$

Table IV reports the average n^* and n_{sb}^* over 10 simulation runs with the numbers below n^* (respectively, n_{sb}^*) representing the average total number of observations n^*cd (n_{sb}^*cd). The numbers in the parenthesis are the corresponding standard deviations over the 10 runs. If a biometric authentication system was tested based on n users, where n is chosen according to the n^* entries in Table IV, we will be 95% certain that the true ROC curve will lie in the interval $[\hat{W} - 0.5, \hat{W} + 0.5]$. Table IV indicates that as the correlation among the multiple impressions of a finger increases for each fixed c and d , the total number of observations needed to achieve the width w for the confidence band increases. The same holds true when c and d values are increased for each correlation combination. Thus, in the presence of non-zero correlation, we are better off selecting a larger number of users rather than increasing the number of acquisitions per user. Note that the sample sizes required by our method, n^* , is smaller compared to n_{sb}^* for achieving the same overall confidence level.

We also obtained the minimum sample sizes determined by the "Rule of 3" [20] and the "Rule of 30" [14] (see Appendix for their derivation). For the fingerprint database [8], n_3 was approximately 150 for all pairs of correlation combination, c and d , while n_{30} was approximately 770. Comparing the values of n_3 and n_{30} with n^*cd , we see that both n_3 and n_{30} grossly underestimate the total number of biometric acquisitions required to achieve a desired width. The underestimation becomes more prominent when significant correlation is present between multiple acquisitions of the biometric templates from a subject.

To illustrate the effects of correlation on the sample size requirement for the WVU database, we choose three combinations of the genuine and impostor within finger correlations, namely, $(\rho_2^{gen}, \rho_2^{imp}) = (0, 0), (0.49, 0.19)$ and $(0.99, 0.39)$ to reflect the no correlation (or, independence), intermediate and high correlation states. Table V reports the average n^* and n_{sb}^* over 10 simulation runs for the width $w = 1\%$, with the average total number of observations, n^*d and n_{sb}^*d given by the entries directly below the n^* s. The numbers in the parenthesis are the corresponding standard deviations over the 10 runs. Note

	Values of c and d					
	$c = 1, d = 2$		$c = 2, d = 2$		$c = 2, d = 3$	
Correlations (ρ_1, ρ_2)	n^*	n_{sb}^*	n^*	n_{sb}^*	n^*	n_{sb}^*
	mean (sd)	mean (sd)	mean (sd)	mean (sd)	mean (sd)	mean (sd)
(0,0)	11,443 (246)	48,674 (600)	5,809 (148)	24,201 (373)	1,967 (31)	8,143 (136)
	22,885 (492)	97,350 (1,200)	23,235 (590)	96,810 (1,493)	11,801 (190)	48,860 (814)
(0, $\hat{\rho}_2$)	20,439 (790)	90,725 (315)	10,476 (279)	46,209 (837)	9,505 (263)	43,500 (455)
	40,877 (1,581)	181,450 (630)	41,905 (1,115)	184,840 (3,346)	57,028 (1,580)	261,000 (2,729)
($\hat{\rho}_1, \hat{\rho}_2$)	21,403 (1,004)	90,477 (407)	11,056 (346)	47,855 (430)	9,749 (163)	46,269 (968)
	42,806 (2,008)	180,950 (813)	44,223 (1,382)	191,420 (1,720)	58,492 (977)	277,620 (5,811)
(0.6, $\hat{\rho}_2$)	19,015 (503)	89,993 (429)	13,321 (506)	61,394 (884)	11,558 (423)	56,723 (826)
	38,029 (1,006)	179,990 (858)	53,285 (2,026)	245,570 (3,536)	69,346 (2,540)	340,340 (4,956)

TABLE IV

MEAN n^* AND n_{sb}^* VALUES FOR ACHIEVING A WIDTH OF 1% FOR THE 95% CONFIDENCE BAND. THE TOTAL NUMBER OF OBSERVATIONS, n^*cd AND n_{sb}^*cd , ARE GIVEN BELOW THE n^* AND n_{sb}^* ENTRIES, RESPECTIVELY. ENTRIES ARE CALCULATED AS THE MEANS OF 10 SIMULATION RUNS. THE CORRESPONDING STANDARD DEVIATIONS ARE GIVEN IN PARENTHESIS.

	Values of c and d					
	$c = 1, d = 2$		$c = 1, d = 3$		$c = 1, d = 4$	
Correlations ($\rho_2^{gen}, \rho_2^{imp}$)	n^*	n_{sb}^*	n^*	n_{sb}^*	n^*	n_{sb}^*
	mean (sd)	mean (sd)	mean (sd)	mean (sd)	mean (sd)	mean (sd)
(0,0)	12,875 (283)	47,526 (655)	4,251 (77)	16,170 (280)	2,103 (37)	8,144 (169)
	25,749 (477)	95,050 (1,310)	12,754 (231)	48,510 (841)	8,412 (148)	32,580 (676)
(0.49, 0.19)	15,215 (513)	61,195 (1,074)	7,719 (215)	35,053 (697)	6,200 (299)	29,149 (940)
	30,430 (1,025)	122,390 (2,148)	23,158 (645)	105,160 (2,091)	24,799 (1,197)	116,600 (3,761)
(0.99, 0.39)	23,802 (886)	90,334 (170)	20,898 (414)	86,357 (400)	18,748 (698)	84,478 (766)
	47,604 (1,772)	180,670 (304)	62,693 (1,244)	259,070 (1,200)	74,991 (2,793)	337,910 (3,064)

TABLE V

MEAN n^* AND n_{sb}^* VALUES FOR ACHIEVING A WIDTH OF 1% FOR THE 95% CONFIDENCE BAND BASED ON THE WEST VIRGINIA UNIVERSITY DATABASE. THE TOTAL NUMBER OF OBSERVATIONS, n^*cd AND n_{sb}^*cd , ARE GIVEN BELOW THE n^* AND n_{sb}^* ENTRIES, RESPECTIVELY. ENTRIES ARE CALCULATED AS THE MEANS OF 10 SIMULATION RUNS. THE CORRESPONDING STANDARD DEVIATIONS ARE GIVEN IN PARENTHESIS.

here, again, that n^* is smaller compared to n_{sb}^* for achieving the same overall confidence level.

VI. CONCLUSION

With the growing deployment of biometric systems in several government and commercial applications, it has become even more important to validate the performance levels of a system claimed by a vendor. We present a flexible semi-parametric approach for estimating both the genuine and impostor distributions of similarity scores using multivariate Gaussian copula functions with non-parametric marginals. Confidence bands for the ROC curve are constructed using (i) semi-parametric bootstrap re-samples, and (ii) asymptotic approximations derived from the estimated models. We also determine the minimum required number of subjects needed to achieve a desired width for the confidence band of the ROC curve. Currently, the implementation of the ROC validation procedure and the estimation of required number of samples are based on fingerprint databases with a small number of subjects. We plan to test our methodology on larger databases as they become available. We will also focus on extending the current framework to validate reported performances of multimodal systems.

ACKNOWLEDGMENT

The authors wish to thank Karthik Nandakumar, Arun Ross, Umut Uludag and Yi Chen for their help when we were conducting our experiments. This research is partially supported by the NSF ITR grant 0312646.

REFERENCES

- [1] J. R. Beveridge, K. She, and B. A. Draper, "A nonparametric statistical comparison of principal component and linear discriminant subspaces for face recognition," In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR 2001)*, Hawaii, December, 2001.
- [2] R. Bolle, J. Connell, S. Pankanti, N. Ratha, and A. Senior, *Guide to Biometrics*, Springer, 2004.
- [3] R. Bolle, N. Ratha, and S. Pankanti, "Error analysis of pattern recognition systems: The subsets bootstrap," *Computer Vision and Image Understanding*, vol. 93, no. 1, pp. 1–33, January, 2004.
- [4] R. M. Bolle, S. Pankanti, and N. K. Ratha, "Evaluation techniques for biometrics-based authentication systems (FRR)," In *Proceedings of the 14th International Conference on Pattern Recognition, ICPR*, pp. 2831–2837, August, 2000.
- [5] U. Cherubini, E. Luciano, and W. Vecchiato, *Copula Methods in Finance*, Wiley, 2004.
- [6] S. C. Dass, Y. Zhu, and A. K. Jain, "Validating a biometric authentication system: Sample size requirements," *Technical Report MSU-CSE-05-23*, Computer Science and Engineering, Michigan State University, East Lansing, Michigan, August 2005.
- [7] A. K. Jain, L. Hong, and R. Bolle, "On-line fingerprint verification," *IEEE Transactions on Pattern Recognition and Machine Intelligence*, vol. 19, no. 4, pp. 302–314, 1997.
- [8] A. K. Jain, S. Prabhakar, and A. Ross, "Fingerprint matching: Data acquisition and performance evaluation," *Technical Report TR99-14*, Michigan State University, 1999.
- [9] R. A. Johnson and D. W. Wichern, *Applied Multivariate Statistical Analysis*, Prentice Hall, Englewood Cliffs, NJ, 1988.
- [10] R. J. Micheals and T. E. Boulton, "Efficient evaluation of classification and recognition systems," In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR 2001)*, Hawaii, December, 2001.
- [11] R. G. Miller, *Simultaneous Statistical Inference*, Springer-Verlag, NY, 1981.
- [12] D. F. Morrison, *Multivariate Statistical Methods*, McGraw-Hill, NY, 1990.
- [13] R. B. Nelsen, *An Introduction to Copulas*, Springer, 1999.
- [14] J. Porter, "On the '30 error criterion'," In *National Biometric Center Collected Works*, eds. J. Wayman, pp. 51–56, 2000.
- [15] C. R. Rao, *Linear Statistical Inference And Its Applications*, Wiley, 1991.
- [16] M. E. Schuckers, "Using the beta-binomial distribution to assess performance of a biometric identification device," *International Journal of Image and Graphics (Special Issue on Biometrics)*, vol. 3, no. 3, pp. 523–529, July, 2003.
- [17] U. K. Biometrics Working Group, "Best practices in testing and reporting performance of biometric devices", 2000. Online: www.cesg.gov.uk/technology/biometrics.
- [18] J. Wayman, "Technical testing and evaluation of biometric identification devices," In *Biometrics: Personal Identification in Networked Society*, eds. A. K. Jain, R. Bolle, and S. Pankanti, Kluwer Academic Publishers, 1999.
- [19] J. Wayman, "Confidence interval and test size estimation for biometric data," In *National Biometric Center Collected Works*, eds. J. Wayman, pp. 91–95, 2000.
- [20] J. Wayman, "Technical testing and evaluation of biometric identification devices," In *National Biometric Center Collected Works*, eds. J. Wayman, pp. 67–89, 2000.

APPENDIX I SIMULATION FROM F

We first describe how to simulate samples from F assuming that F is of the form (14). This simulation procedure will be needed for the estimation of the marginals F_k and generating bootstrap samples from F to construct the ROC confidence bands. The following steps outline how to generate N samples from F : (1) Generate a vector $Z = (Z_1, Z_2, \dots, Z_K)^T$ from Φ_R^K , the K -dimensional multivariate normal with mean 0, variance 1, and correlation matrix R , (2) Obtain the vector $U = (U_1, U_2, \dots, U_K)^T$ by letting $U_k = \Phi(Z_k)$ for $k = 1, 2, \dots, K$, and (3) Obtain the vector $S^* = (s_1^*, s_2^*, \dots, s_K^*)^T$ using $s_k^* = F_k^{-1}(U_k)$ for $k = 1, 2, \dots, K$, where F_k^{-1} is the inverse of F_k . It follows that S^* is a sample from F . In order to obtain a sample of size N , steps (1-3) are repeated N times resulting in the simulated samples $\{s^*(i, k), k = 1, 2, \dots, K\}$ for $i = 1, 2, \dots, N$. In practice, one difficulty is that the marginal distributions and the correlation matrices for the genuine and impostor similarity scores will generally be unknown, and will have to be estimated from the observed scores (this is discussed in the subsequent section). Once they have been estimated, we can follow steps (1-3) to obtain samples from the fitted copula models.

A. Estimation of F_k and R

The marginal distribution functions, F_k , and the correlation matrix R are generally unknown and have to be estimated from the observed vector of similarity scores, $\{S_i, i = 1, 2, \dots, N\}$. The empirical distribution function for the k -th marginal is given by

$$E_k(s) = \frac{1}{N} \sum_{i=1}^N I\{s(i, k) \leq s\}, \quad (39)$$

where $I(A)$ is the indicator function of the set A ; $I(A) = 1$ if A is true, and 0 otherwise. Note that $E_k(s) = 0$ for all $s < s_{min}$ and $E_k(s) = 1$ for all $s \geq s_{max}$, where s_{min} and s_{max} , respectively, are the minimum and maximum of the observations $\{s(i, k) : i = 1, 2, \dots, N\}$. Next, we define $\mathcal{H}(s) \equiv -\log(1 - E_k(s))$, and note that discontinuity points of $E_k(s)$ will also be points of discontinuity of $\mathcal{H}(s)$. In order to obtain a continuous estimate of $\mathcal{H}(s)$, the following procedure is adopted: For an M -partition $s_{min} \equiv s_0 < s_1 < \dots < s_M \equiv s_{max}$ of $[s_{min}, s_{max}]$, the value of $\mathcal{H}(s)$ at a point $s \in [s_{min}, s_{max}]$ is redefined via the linear interpolation formula

$$\hat{\mathcal{H}}(s) = \mathcal{H}(s_m) + (\mathcal{H}(s_{m+1}) - \mathcal{H}(s_m)) \cdot \frac{s - s_m}{s_{m+1} - s_m} \quad (40)$$

whenever $s_m \leq s \leq s_{m+1}$ and subsequently, the estimated distribution function, $\hat{F}_k(s)$, of $F_k(s)$ is obtained as

$$\hat{F}_k(s) = 1 - \exp\{-\hat{\mathcal{H}}(s)\}. \quad (41)$$

It follows that each $\hat{F}_k(s)$ is a continuous distribution function. Next we generate B^* samples from \hat{F}_k : (1) Generate a uniform random variable U in $[0, 1]$, (2) Define $V = -\log(1 - U)$, and (3) Find the value V^* such that $\hat{\mathcal{H}}(V^*) = V$. It follows that V^* is a random variable with distribution function \hat{F}_k . To generate B^* independent realizations from \hat{F}_k , we repeat the steps (1-3) B^* times. Finally, a non-parametric density estimate of F_k is obtained based on the simulated samples using a Gaussian kernel.

The estimate of R based on the observed similarity score vectors $\{\mathcal{S}_i : i = 1, 2, \dots, N\}$ is obtained in the following way: Define a new vector $\mathcal{Z}_i = (Z(i, 1), Z(i, 2), \dots, Z(i, K))^T$ where

$$Z(i, k) = \Phi^{-1}(E_k(s(i, k))), \quad (42)$$

for $k = 1, 2, \dots, K$. The mean vector $\bar{\mathcal{Z}}$ is then obtained by averaging over the vectors \mathcal{Z}_i , that is,

$$\bar{\mathcal{Z}} = \frac{1}{N} \sum_{i=1}^N \mathcal{Z}_i \quad (43)$$

and the covariance matrix is defined as

$$J = \frac{1}{N} \sum_{i=1}^N (\mathcal{Z}_i - \bar{\mathcal{Z}}) \cdot (\mathcal{Z}_i - \bar{\mathcal{Z}})^T. \quad (44)$$

The estimate of $\rho_{kk'}$ is given by

$$\hat{\rho}_{kk'} = \frac{\sigma_{kk'}}{\sqrt{\sigma_{kk}\sigma_{k'k'}}}, \quad (45)$$

where $\sigma_{kk'}$ is the (k, k') -th entry of J in (44), and the estimated correlation matrix is given by $\hat{R} = ((\hat{\rho}_{kk'}))$. The total number of correlation parameters that need to be estimated is $K(K-1)/2$; thus, it is necessary to have $K(K-1)/2$ much smaller than N to avoid over-fitting.

B. Assessing the Goodness of Fit

We present a method here for graphically assessing the goodness of fit of the estimated multivariate Gaussian K -copula model to the observed data. We first give the general methodology, and then apply it to the observed genuine and impostor similarity scores. Lower dimensional marginals of a K -copula function $C(u_1, u_2, \dots, u_K)$ can be obtained by fixing the irrelevant u_k s to be equal to one: For example, if we require the 2-dimensional copula function in the dimensions of k and k' , where $k \neq k'$, $k, k' = 1, 2, \dots, K$, this can be obtained by setting the other u_j s ($j \neq k, j \neq k'$) to 1, that is,

$$C_{k,k'}(u_k, u_{k'}) \equiv C(1, 1, \dots, u_k, 1, \dots, 1, u_{k'}, 1, \dots, 1). \quad (46)$$

It follows that all lower k -dimensional ($k < K$) marginals of the multivariate Gaussian K -copula are Gaussian k -copulas. In particular, for $k = 2$, we obtain $\binom{K}{2}$ bivariate Gaussian copulas from a single Gaussian K -copula as in (13). Each bivariate Gaussian copula is characterized by a single correlation parameter; for dimensions k and k' , this parameter is $\rho_{kk'}$ of matrix R .

The bivariate empirical copula based on N independent bivariate observations (X_i, Y_i) , $i = 1, 2, \dots, N$ is defined as follows: For each $0 \leq x \leq 1$ and $0 \leq y \leq 1$,

$$C_{emp}(x, y) = \frac{1}{N} \sum_{i=1}^N I\{X_i \leq X_{([Nx])}, Y_i \leq Y_{([Ny])}\}, \quad (47)$$

where $X_{([Nx])}$ (respectively, $Y_{([Ny])}$) is the $[Nx]$ -th ($[Ny]$ -th) element in the ordered list of X (Y) samples, and the notation $[u]$ represents the greatest integer less than or equal to u . The empirical copula function gives the best approximation to the true but unknown copula function that generated the observed data (X_i, Y_i) , $i = 1, 2, \dots, N$.

Our graphical test for checking goodness of fit consists of the following steps: (i) Obtain the $\binom{K}{2}$ 2-dimensional marginal copulas based on \hat{R} . For the dimension pair (k, k') , we obtain the contour plot of $C_{k,k'}(u_k, u_{k'})$ given by

$$C_{k,k'}(u_k, u_{k'}) = \Phi_{\hat{\rho}_{kk'}}^2(\Phi^{-1}(u_k), \Phi^{-1}(u_{k'})). \quad (48)$$

(ii) Obtain the empirical copula based on the score vectors $(s(i, k), s(i, k'))^T$ for $i = 1, 2, \dots, N$ using equation (47); here $s(i, k)$ are the X samples and $s(i, k')$ are the Y samples.

C. Results for the fingerprint database [8]

The estimates of the intra-subject impostor correlation matrix (of dimension (8×8)) is given by $\hat{R}_{11} =$

$$\begin{pmatrix} 1.00 & 0.58 & 0.52 & 0.42 & 0.90 & 0.53 & 0.54 & 0.41 \\ 0.58 & 1.00 & 0.44 & 0.47 & 0.58 & 0.46 & 0.88 & 0.46 \\ 0.52 & 0.44 & 1.00 & 0.45 & 0.50 & 0.86 & 0.37 & 0.42 \\ 0.42 & 0.47 & 0.45 & 1.00 & 0.41 & 0.41 & 0.43 & 0.87 \\ 0.90 & 0.58 & 0.50 & 0.41 & 1.00 & 0.53 & 0.55 & 0.41 \\ 0.53 & 0.46 & 0.86 & 0.41 & 0.53 & 1.00 & 0.40 & 0.42 \\ 0.54 & 0.88 & 0.37 & 0.43 & 0.55 & 0.40 & 1.00 & 0.44 \\ 0.41 & 0.46 & 0.42 & 0.87 & 0.41 & 0.42 & 0.44 & 1.00 \end{pmatrix}. \quad (49)$$

We also obtained the estimate of the inter-subject impostor correlation matrix, \hat{R}_{12} , which is of dimension 16×16 . Due to the large dimensionality associated with this matrix, we do not present it here.

For assessing the goodness of fit, the total number of pairs of components for the sets of genuine, intra-subject and inter-subject scores are, respectively, $\binom{4}{2} = 6$, $\binom{8}{2} = 28$, and $\binom{16}{2} = 120$. Figures 7, 8 and 9 respectively give the plots of 6 component pairs for the genuine, intra-subject impostor and inter-subject impostor sets in this case. Note that the figures indicate that there is a good agreement between the empirical and the proposed Gaussian copula functions. We checked all of the pairwise copula plots and found that there were no major discrepancies between the empirical contours and the fitted Gaussian copula contours. Thus, we conclude that the proposed Gaussian copula functions are good models for representing the correlation structures in all of the genuine, intra-subject and inter-subject sets of scores. There is always a problem of quantitatively assessing the quality of a model fit to the observed data when the sample size is very large (as in the case of the genuine and impostor sets of similarity scores here). A small discrepancy between the observed data and model fit will magnify due to the large sample size and cause a quantitative goodness of fit test to be statistically significant. The point to note here is that the test can potentially be statistically significant even if the models are a good fit to the observed data set.

D. Rules of 3 and 30

Recall that the Rule of 3 and the Rule of 30 are rules of thumb to select the sample size, n , for the reliable estimation of an error probability, p , based on n independent binary observations, x_1, x_2, \dots, x_n , with $P(x_i = 1) = 1 - P(x_i = 0) = p$ (see [20] and [14] for details). Since both the rules were derived for setting up confidence intervals for specific values of FAR and GAR (and not confidence bands for a range of FAR and GAR values), we were required to modify them slightly to suit the present case. For the Rule of 3, we computed the quantity $FRR_m = 1 - GAR(p_m)$ for $m = 1, 2, \dots, M$ and derived the minimum sample size as

$$n_3 = \max_{1 \leq m \leq M} \frac{3}{FRR_m}. \quad (50)$$

The smallest sample size based on the Rule of 30 was obtained using the formula

$$n_{30} = \max_{1 \leq m \leq M} \frac{(2 * 1.96)^2}{FRR_m}. \quad (51)$$

E. Asymptotic Theory

We derive several results below to validate the asymptotic representation of z in equation (28). In proving these results, we assume that the biometric entities considered are the different subjects, and the matcher S is asymmetric. Recall that the total number of subjects was denoted by n , and d impressions of c fingers for each subject were acquired for validating a vendor's claim. In this case, $N_0 = n$, $K_0 = cd(d-1)$, $N_{11} = n$, $K_{11} = c(c-1)d^2$, $N_{12} = n(n-1)$ and $K_{12} = c^2d^2$. The asymptotic results presented here will be for $n \rightarrow \infty$ with c and d fixed.

We will first derive the asymptotic theory for $\sqrt{N_0}(\hat{W}(p) - W(p))$, and then extend it to the quantity $\sqrt{N_0}(\sin^{-1}\sqrt{\hat{W}(p)} - \sin^{-1}\sqrt{W(p)})$. We denote the densities of G_0 and G_1 , assuming they exist, by g_0 and g_1 , respectively. The quantity $\sqrt{N_0}(\sin^{-1}\sqrt{\hat{W}(p)} - \sin^{-1}\sqrt{W(p)})$ is a continuous function of $p \in [C_0, C_1]$ since the component marginals and their estimates for the genuine, intra-subject impostor and inter-subject impostor joint distributions are continuous. In order to find the asymptotic distribution of $\sqrt{N_0} \max_{C_0 \leq p \leq C_1} |\sin^{-1}\sqrt{\hat{W}(p)} - \sin^{-1}\sqrt{W(p)}|$, we first define a partition of $[C_0, C_1]$: $C_0 \equiv p_1 < p_2 < \dots < p_M \equiv C_1$. Defining $z(p) = \sin^{-1}\sqrt{\hat{W}(p)} - \sin^{-1}\sqrt{W(p)}$, we have

$$\sqrt{N_0} \max_{C_0 \leq p \leq C_1} |z(p)| \approx \sqrt{N_0} \max_{1 \leq m \leq M} |z(p_m)| \quad (52)$$

for large M . Thus, we first derive the joint asymptotic distribution of the M -dimensional vector $\sqrt{N_0} z(p_m)$, $m = 1, 2, \dots, M$, and then obtain the distribution of the maximum of the absolute values of these m components. Note that by Taylor's expansion, we have

$$\sqrt{N_0} z(p) \approx D(p) \sqrt{N_0} (\hat{W}(p) - W(p)) \quad (53)$$

for large N_0 , where $D(p) = \frac{1}{\sqrt{4W(p)(1-W(p))}}$. In other words, we require to find the distribution of $D_M \cdot \hat{W}_M$ where

$$\hat{W}_M \equiv \begin{pmatrix} \sqrt{N_0} (\hat{W}(p_1) - W(p_1)), \hat{W}(p_2) - W(p_2), \\ \dots, \hat{W}(p_M) - W(p_M) \end{pmatrix}^T \quad (54)$$

is an M -dimensional vector and D_M is the diagonal matrix with the (m, m) -th entry given by $D(p_m)$. We introduce some notation before stating the main results. For $m = 1, 2, \dots, M$, define ξ_m and $\hat{\xi}_m$ to be the p_m -th upper quantiles of G_1 and \hat{G}_1 , respectively, that is

$$\xi_m \equiv G_1^{-1}(p_m) \quad \text{and} \quad \hat{\xi}_m \equiv \hat{G}_1^{-1}(p_m). \quad (55)$$

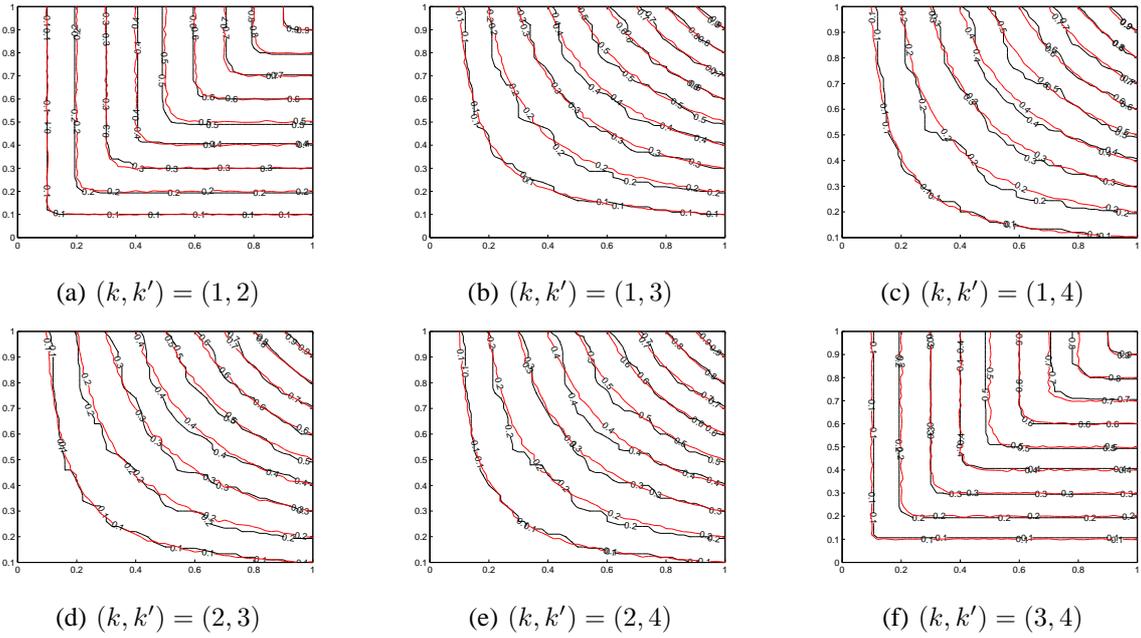


Fig. 7. Nine level curves (at levels 0.1, 0.2, ..., 0.9) indicating a good match between the empirical copula (black lines) and the estimated bivariate Gaussian copula (red lines) along dimensions k and k' for the genuine scores.

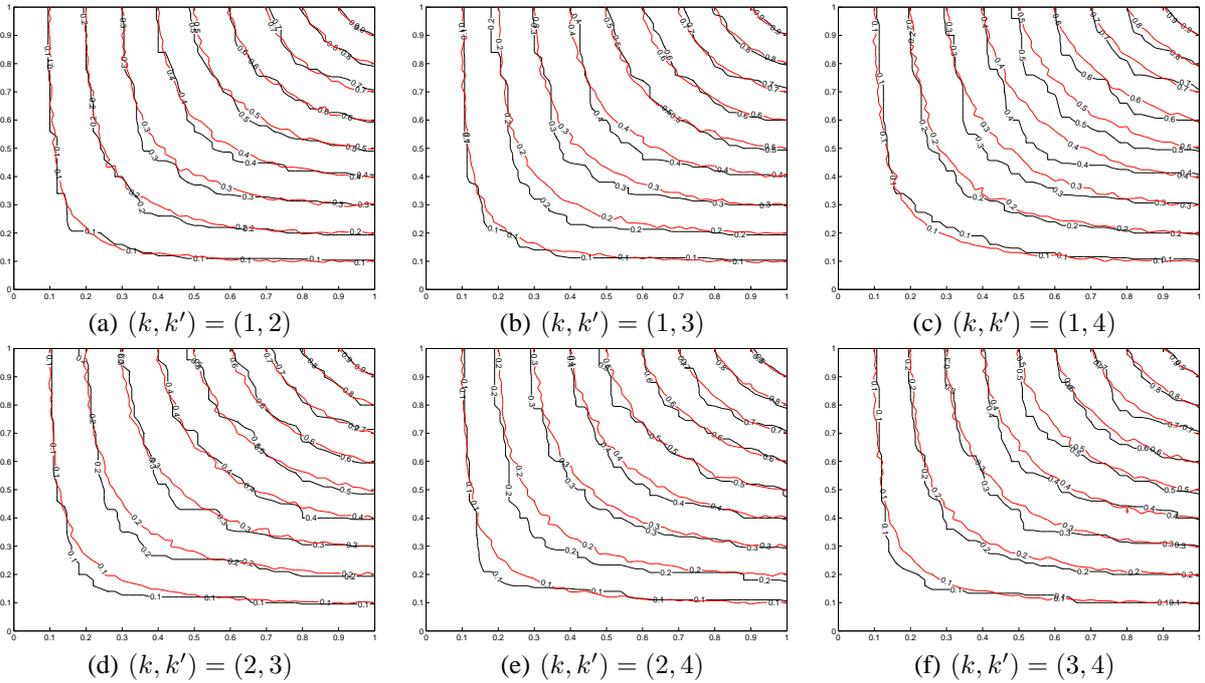


Fig. 8. Nine level curves (at levels 0.1, 0.2, ..., 0.9) indicating a good match between the empirical copula (black lines) and the estimated bivariate Gaussian copula (red lines) along dimensions k and k' for the intra-subject impostor scores.

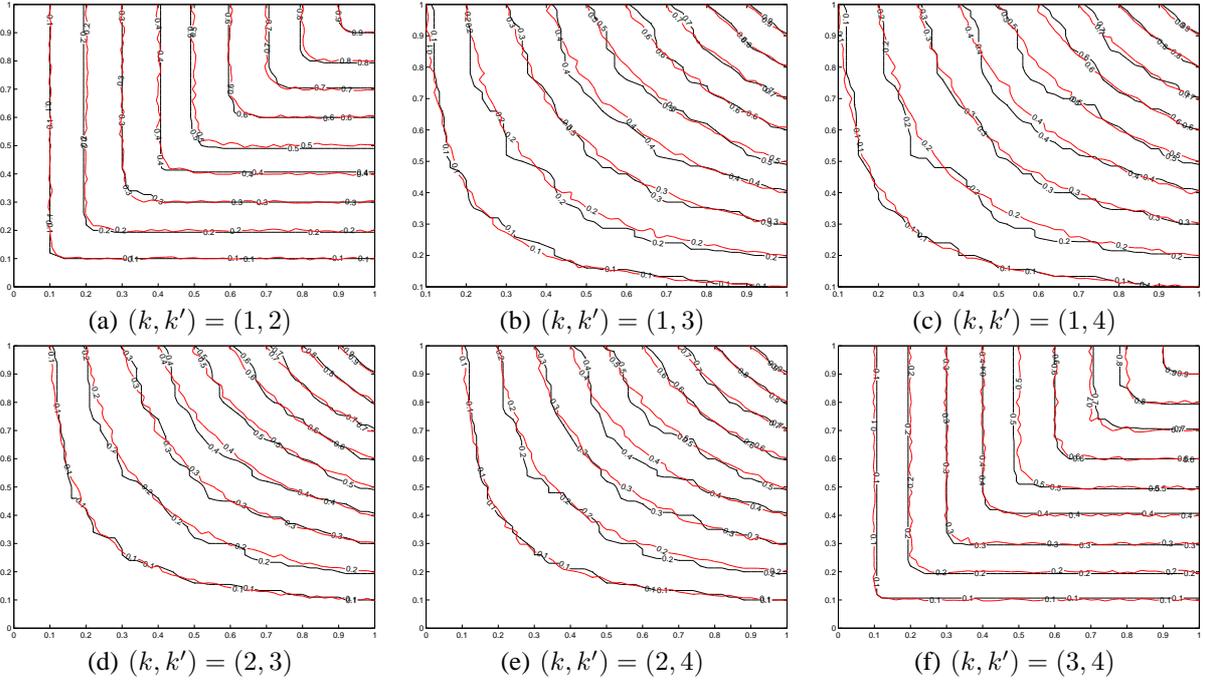


Fig. 9. Nine level curves (at levels 0.1, 0.2, ..., 0.9) indicating a good match between the empirical copula (black lines) and the estimated bivariate Gaussian copula (red lines) along dimensions k and k' for the inter-subject impostor scores.

Since $\hat{G}_1 - G_1$ converges almost surely to 0, we have $\hat{\xi}_m - \xi_m \rightarrow 0$ as $N_0 \rightarrow \infty$. Also, denoting $\hat{G}_{0,M} \equiv \sqrt{N_0} (\hat{G}_0(\hat{\xi}_1) - G_0(\hat{\xi}_1), \hat{G}_0(\hat{\xi}_2) - G_0(\hat{\xi}_2), \dots, \hat{G}_0(\hat{\xi}_M) - G_0(\hat{\xi}_M))^T$ and $\hat{G}_{1,M} \equiv \sqrt{N_0} (G_0(\hat{\xi}_1) - G_0(\xi_1), G_0(\hat{\xi}_2) - G_0(\xi_2), \dots, G_0(\hat{\xi}_M) - G_0(\xi_M))^T$, we have

$$\hat{W}_M = \hat{G}_{0,M} + \hat{G}_{1,M}. \quad (56)$$

Lemmas 1 - 4 in Appendix II can be used to show that $\hat{G}_{0,M}$ and $\hat{G}_{1,M}$ are asymptotically independent, and the limiting distributions of $\hat{G}_{0,M}$ and $\hat{G}_{1,M}$ are multivariate normals with means 0 and covariance matrices given by Θ_0 and $\frac{N_0}{N_1}\Theta_1$, respectively; see Lemmas 2 and 3 for the forms of Θ_0 and Θ_1 , respectively. Thus, it follows that for the M -partition $C_0 \equiv p_1 < p_2 < \dots < p_M \equiv C_1$, the distribution of $\sqrt{N_0}(z(p_m), m = 1, 2, \dots, M)$ is given by

$$D_M \cdot \hat{W}_M = D_M \cdot \hat{G}_{0,M} + D_M \cdot \hat{G}_{1,M}. \quad (57)$$

Since $\hat{G}_{0,M}$ and $\hat{G}_{1,M}$ are asymptotically independent, it follows that $D_M \cdot \hat{G}_{0,M}$ and $D_M \cdot \hat{G}_{1,M}$ are also asymptotically independent, and the limiting distributions of $D_M \cdot \hat{G}_{0,M}$ and $D_M \cdot \hat{G}_{1,M}$ are multivariate normals with means 0 and covariance matrices given by

$$\Gamma_0 = D_M \Theta_0 D_M^T \quad \text{and} \quad \Gamma_1 = \frac{N_0}{N_1} D_M \Theta_1 D_M^T, \quad (58)$$

respectively. Since the covariance matrices above depend on unknown parameters, they will, in practice, be determined by plugging in parameter estimates in place of the unknown parameters; for example, the (m, m) -th

entry of D_M , $D_M(p_m) = \frac{1}{\sqrt{4W(p_m)(1-W(p_m))}}$, will be estimated by plugging in $\hat{W}(p_m)$ in place of $W(p_m)$.

APPENDIX II LEMMAS

We now state and prove the required lemmas. Define $G_{11}(\lambda) = \frac{1}{K_{11}} \sum_{k=1}^{K_{11}} P\{s_{11}(1, k) > \lambda\}$ and $G_{12}(\lambda) = \frac{1}{K_{12}} \sum_{k=1}^{K_{12}} P\{s_{12}(1, k) > \lambda\}$. It follows then, that $G_1(\lambda) = \frac{N_{11}K_{11}}{N_1} G_{11}(\lambda) + \frac{N_{12}K_{12}}{N_1} G_{12}(\lambda)$. For $m = 1, 2, \dots, M$, define $\xi_{12,m} = G_{12}^{-1}(p_m)$. We introduce a few notations for the subsequent discussion: Let $\beta_H(k, m) = P\{s_H(1, k) > \xi_{H,m}\}$ and $\beta_H(k, k', m, m') = P\{s_H(1, k) > \xi_{H,m}, s_H(1, k') > \xi_{H,m'}\}$ for the sets $H = \{0, 11, 12\}$, respectively, denoting the genuine, intra-subject impostor and inter-subject impostor cases.

We state

Lemma 1: The M -dimensional vector

$$\sqrt{N_{12}} \left(\frac{g_1(\xi_1)(\hat{\xi}_1 - \xi_1)}{\sqrt{p_1(1-p_1)}}, \frac{g_1(\xi_2)(\hat{\xi}_2 - \xi_2)}{\sqrt{p_2(1-p_2)}}, \dots, \frac{g_1(\xi_M)(\hat{\xi}_M - \xi_M)}{\sqrt{p_M(1-p_M)}} \right)^T \rightarrow Z_M \quad (59)$$

where Z_M is a multivariate normal random variable with zero means, unit variances and correlation matrix given by

$$\Theta_{12}(m, m') = \frac{1}{K_{12}^2} \sum_{k=1}^{K_{12}} \sum_{k'=1}^{K_{12}} \theta_{12}(k, k', m, m') \quad (60)$$

where

$$\theta_{12}(k, k', m, m') = \frac{\beta_{12}(k, k', m, m') - \beta_{12}(k, m)\beta_{12}(k', m')}{\sqrt{p_m(1-p_m)} \cdot \sqrt{p_{m'}(1-p_{m'})}}. \quad (61)$$

Proof: Consider the expression

$$\begin{aligned} & P \left\{ \sqrt{N_{12}} \frac{g_1(\xi_m)(\hat{\xi}_m - \xi_m)}{\sqrt{p_m(1-p_m)}} \leq x_m, \right. \\ & \quad \left. 1 \leq m \leq M \right\} \\ &= P \left\{ \hat{\xi}_m \leq \xi_m + \frac{x_m}{g_1(\xi_m)} \sqrt{\frac{p_m(1-p_m)}{N_{12}}}, \right. \\ & \quad \left. 1 \leq m \leq M \right\} \\ &= P \left\{ \hat{G}_1 \left(\xi_m + \frac{x_m}{g_1(\xi_m)} \sqrt{\frac{p_m(1-p_m)}{N_{12}}} \right) > p_m, \right. \\ & \quad \left. 1 \leq m \leq M \right\} \\ &= P \{ K_{11}X_{11} + K_{12}X_{12} > N_{12}p_m, \\ & \quad 1 \leq m \leq M \}, \end{aligned}$$

where X_H is a Binomial random variable with parameters N_H for the total number of trials and $p_H^m \equiv G_H(\xi_m + \frac{x_m}{g_1(\xi_m)} \sqrt{\frac{p_m(1-p_m)}{N_{12}}})$ as the probability of success in each trial, for $H = \{11\}$ and $\{12\}$. It follows that the last expression above can be re-written as $P\{K_{12}Z_{12}^m > Q^m, m = 1, 2, \dots, M\}$ where

$$\begin{aligned} Q^m &= \frac{1}{\sqrt{N_{12}p_{12}^m(1-p_{12}^m)}} \left[N_{11}p_m - \right. \\ & \quad \left. N_{11}G_1 \left(\xi_m + \frac{x_m}{g_1(\xi_m)} \sqrt{\frac{p_m(1-p_m)}{N_{12}}} \right) - \right. \\ & \quad \left. K_{11}Z_{11}^m \sqrt{N_{11}p_{11}^m(1-p_{11}^m)} \right], \end{aligned}$$

$Z_{11}^m = (X_{11} - N_{11}p_{11}^m)/\sqrt{N_{11}p_{11}^m(1-p_{11}^m)}$, and $Z_{12}^m = (X_{12} - N_{12}p_{12}^m)/\sqrt{N_{12}p_{12}^m(1-p_{12}^m)}$. As $n \rightarrow \infty$, using the Taylor's expansion for $G_1 \left(\xi_m + \frac{x_m}{g_1(\xi_m)} \sqrt{\frac{p_m(1-p_m)}{N_{12}}} \right)$ and the facts that $N_{11}/N_{12} \rightarrow 0$, $N_1/N_{12} \rightarrow K_{12}$ and $p_{12}^m \rightarrow p_m$, we get $Q^m \rightarrow -K_{12}x_m$. The limiting distributions of each Z_H^m is normal with mean 0 and variance 1, for $u = \{11\}$ and $\{12\}$. Further, a computation of the covariance gives the expression (60) for the covariance between Z_{12}^m and $Z_{12}^{m'}$. QED.

For the next lemma, define $\theta_0(k, k', m, m')$ by

$$\theta_0(k, k', m, m') = \beta_0(k, k', m, m') - \beta_0(k, m)\beta_0(k', m'), \quad (62)$$

and let Θ_0 be the $M \times M$ matrix whose (m, m') -th entry is given by

$$\Theta_0(m, m') = \frac{1}{K_0^2} \sum_{k=1}^{K_0} \sum_{k'=1}^{K_0} \theta_0(k, k', m, m'). \quad (63)$$

We state

Lemma 2: Let $\underline{t} = (t_1, t_2, \dots, t_M)^T$. If $\hat{\varphi}_0(\underline{t})$ denotes the characteristic function of $\hat{G}_{0,M}$, and $\varphi_0(\underline{t}) \equiv \exp\{-\frac{1}{2}\underline{t}^T \Theta_0 \underline{t}\}$ is the characteristic function of an M -dimensional normal with mean 0 and covariance matrix Θ_0 , then

$$|\hat{\varphi}_0(\underline{t}) - \varphi_0(\underline{t})| \rightarrow 0 \quad (64)$$

as $n \rightarrow \infty$.

Proof: The proof of Lemma 2 will first involve conditioning on $\hat{\xi}_m$ for $m = 1, 2, \dots, M$. Using the multivariate Central Limit Theorem [15], it follows that $\sqrt{N_0}(\hat{G}_0(\hat{\xi}_m) - G_0(\hat{\xi}_m))$ converges to an M -variate normal distribution with zero means and covariance matrix given by Θ_0 , where $\hat{\Theta}_0$ is the matrix Θ_0 in (63) with $\hat{\xi}_m$ used in place of $\xi_{12,m}$. But, note that, $\hat{\xi}_m \rightarrow \xi_{12,m}$ so that $\hat{\Theta}_0 \rightarrow \Theta_0$. Lemma 2 follows. QED.

For the next lemma, let Θ_1 denote the $M \times M$ matrix whose (m, m') -th entry is given by

$$\sigma_{12}(m, m') = J(m) \cdot \Theta_{12}(m, m') \cdot J(m'), \quad (65)$$

where $\Theta_{12}(m, m')$ is as given in (60) and

$$J(m) \equiv \sqrt{p_m(1-p_m)} \cdot \frac{g_0(\xi_m)}{g_1(\xi_m)}.$$

We state

Lemma 3: Let $\underline{u} = (u_1, u_2, \dots, u_M)^T$. If $\hat{\varphi}_1(\underline{u})$ denotes the characteristic function of $\sqrt{\frac{N_{12}}{N_0}} \hat{G}_{1,M}$ and $\varphi_1(\underline{u}) \equiv \exp\{-\frac{1}{2}\underline{u}^T \Theta_1 \underline{u}\}$, then

$$|\hat{\varphi}_1(\underline{u}) - \varphi_1(\underline{u})| \rightarrow 0 \quad (66)$$

as $n \rightarrow \infty$.

Proof: The m -th component of $\hat{G}_{1,M}$, $\sqrt{N_0}(G_0(\hat{\xi}_m) - G_0(\xi_m))$, can be written as $\sqrt{N_0}g_0(\xi_m)(\hat{\xi}_m - \xi_m)$ using Taylor's expansion for large n since $\hat{\xi}_m - \xi_m \rightarrow 0$. We can re-write this as

$$\sqrt{\frac{N_0}{N_{12}}} \frac{g_0(\xi_m)}{g_1(\xi_m)} \cdot \sqrt{p_m(1-p_m)} \cdot \left(\sqrt{N_{12}} \frac{(\hat{\xi}_m - \xi_m)}{\sqrt{p_m(1-p_m)}} \right). \quad (67)$$

Lemma 3 follows from applying Lemma 1 to (67). QED.

The next lemma is

Lemma 4: Let $\varphi_{0,1}(\underline{t}, \underline{u}) \equiv E(e^{i\underline{t}^T \hat{G}_{0,M} + i\underline{u}^T \hat{G}_{1,M}})$ be the characteristic function of $(\hat{G}_{0,M}, \hat{G}_{1,M})$. Then,

$$|\varphi_{0,1}(\underline{t}, \underline{u}) - \varphi_0(\underline{t}) \cdot \varphi_1(\sqrt{\frac{N_0}{N_{12}}} \underline{u})| \rightarrow 0 \quad (68)$$

as $n \rightarrow \infty$, where $\varphi_0(\underline{t})$ and $\varphi_1(\underline{u})$ are as defined in Lemmas 2 and 3, respectively.

Proof: We first condition on all the impostor similarity scores. Thus, we have

$$\begin{aligned} & \varphi_{0,1}(\underline{t}, \underline{u}) \\ &= E(e^{i\underline{t}^T \hat{G}_{0,M} + i\underline{u}^T \hat{G}_{1,M}}) \\ &= E(e^{i\underline{u}^T \hat{G}_{1,M}} E(e^{i\underline{t}^T \hat{G}_{0,M}} | \mathcal{S}_{11} \cup \mathcal{S}_{12})) \\ &= E(e^{i\underline{u}^T \hat{G}_{1,M}} \varphi_0^*(\underline{t})), \end{aligned}$$

where $\varphi_0^*(\underline{t})$ is $\varphi_0(\underline{t})$ with Θ_0 replaced by $\hat{\Theta}_0$. Next, we have

$$\begin{aligned} & |\varphi_{0,1}(\underline{t}, \underline{u}) - \varphi_0(\underline{t})\varphi_1(\sqrt{\frac{N_0}{N_{12}}}\underline{u})| \\ = & |M_1 + M_2| \leq |M_1| + |M_2| \end{aligned}$$

where $M_1 = E(e^{i\underline{u}^T \hat{G}_{1,M}}(\varphi_0^*(\underline{t}) - \varphi_0(\underline{t})))$ and $M_2 = E(e^{i\underline{u}^T \hat{G}_{1,M}}\varphi_0(\underline{t}) - \varphi_0(\underline{t})\varphi_1(\sqrt{\frac{N_0}{N_{12}}}\underline{u}))$. Note that $|M_1| \leq E|\varphi_0^*(\underline{t}) - \varphi_0(\underline{t})| \rightarrow 0$ as $n \rightarrow \infty$ (since $\varphi_0^*(\underline{t})$ and $\varphi_0(\underline{t})$ are bounded functions by Lemma 2, and point-wise convergence implies convergence in expectation). Also $|M_2| \leq |\hat{\varphi}_1(\sqrt{\frac{N_0}{N_{12}}}\underline{u}) - \varphi_1(\sqrt{\frac{N_0}{N_{12}}}\underline{u})| \rightarrow 0$ using Lemma 3. Lemma 4 follows. QED.