

Demographic Effects in Latent Fingerprints Matching and their Relation to Image Quality

EMANUELA MARASCO, Center for Secure Information Systems, School of Computing and Engineering, George Mason University, USA

MENGLING HE, Department of Statistics and National Center for Forensic Science, University of Central Florida, USA

LARRY TANG, Department of Statistics and National Center for Forensic Science, University of Central Florida, USA

YUANTING TAO, Center for Secure Information Systems, School of Computing and Engineering, George Mason University, USA

Automatic processing of latent fingerprints can alleviate the subjectivity inherent in feature markups done by examiners that may be affected by factors such as visual perception, expertise and workload. Despite several benefits, the use of algorithmic decision systems is also associated with different risks for individuals, such as discrimination and unfair practices. The goal of this paper is to analyze the fairness of prediction and decision-making in forensics through discovering and mitigation of biases in automated algorithms operating on latent fingerprint images. Previous work analyzes biases only in match scores without considering the impact of image quality that is crucial to the reliability of the matching algorithm. Furthermore, due to their learning-based nature, quality predictors may be biased as well. In our previous work, we carried out an ROC regression analysis related to the demographic effects on latent fingerprint matching without considering the impact of image quality. In this paper, quality measures extracted from latent prints are considered in the predictive model as an additional covariate to the demographics. Experiments were carried out on the FBI WVU BioCop 2008 database that contains 469 right-thumb and 219 right-index latent fingerprint images with associated demographics. Quality is estimated using the latent fingerprint image quality (LFIQ) algorithm. The findings show that the proposed covariate-adjusted ROC curve conditioned on image quality and demographics is a more informative assessment scheme than an evaluation without quality.

CCS Concepts: • Security and privacy → Biometrics; • Computing methodologies → Machine learning; Computer vision.

Additional Key Words and Phrases: Biometrics, Machine learning, Computer vision

1 INTRODUCTION

Uncovering and clarifying biases in biometric technologies has recently gained increasing interest. A lot of attention has focused on how predictive models may be biased and how statistical and machine learning models' performance differs along social axes such as gender, age, and ethnicity [10, 19]. Algorithmic unfairness and its impact on the society is a critical concern for designers of AI systems including biometric technologies. Discrimination may result from different types of biases arising from the training data, technical constraints or societal biases. When given new input data, a machine learning model generates values based on a model trained using a certain dataset that may not account for variances in race, sexual orientation or identity, or age; thus, the outcomes may very negatively affect people's lives. In this regard, facial recognition systems

Authors' addresses: Emanuela Marasco, Center for Secure Information Systems, School of Computing and Engineering, George Mason University, Virginia, USA, emarasco@gmu.edu; Mengling He, Department of Statistics and National Center for Forensic Science, University of Central Florida, Florida, USA, menglinghe@knights.ucf.edu; Larry Tang, Department of Statistics and National Center for Forensic Science, University of Central Florida, Florida, USA, liansheng.tang@ucf.edu; Yuanting Tao, Center for Secure Information Systems, School of Computing and Engineering, George Mason University, Virginia, USA, ytao4@gmu.edu.

have been examined with respect to the demographic effects and their dependence on image acquisition and findings confirm that they are strongly related [5]. Fingerprint features have also been related to age group (individuals born at a similar time), gender (physical characteristics that distinguish males from females) and ethnicity (common culture and origin) [8, 13, 18]. Fingerprint texture offers one possible explanation for these differences.

Latent fingerprints are fingerprint impressions unintentionally left on surfaces at a crime scene. They are crucial in crime scene investigations for making identifications or exclusions of suspects. This data often represents an incomplete or distorted impression of a finger obtained when body's natural oils and sweat on the skin are left on a surface. Thus, processing latent fingerprints usually requires the involvement of experts in determining the value of the print as forensic evidence and only in the past few years scientists have released fully automated approaches. An interesting study by Yoon and Jain found that fingerprint matching scores vary with covariates in terms of demographic differentials [28]. Later, Marasco *et al.* utilized the idea of ROC regression techniques to incorporate demographic differentials in the ROC curve. The resulting covariate-specific ROC curves were able to successfully provide the interpretation of demographic bias by taking them into account when assessing the system's performance [17].

Recent research has only focused on gaining crucial insights about how demographic differentials influence biometric matching. Although image quality has been extensively used as a predictor of matching performance, automated quality estimators have not been investigated from a fairness perspective. The quality of a biometric signal expresses its utility to an automated system intended as its suitability to further process by the biometric system. Poor quality fingerprint images do not have a clear pattern of ridge and valleys and may result in spurious features and loss of identifiable information. Factors such as adverse skin conditions (e.g., dry, creased) can lower image quality and degrade performance. The importance of quality has been extensively highlighted for sensor-based fingerprints, its impact is even more significant in the presence of latent fingerprint evidence. For latent prints, quality is also an indicator of which type of evidence we are dealing with, which can help to use a more accurate model for that specific case, where standard ROC would not capture the variations in the sample. In the literature, it is well-known that matching performance is highly impacted by image quality [3]. Incorporating quality measures for latent fingerprints enables a more objective assessment and mitigation of algorithmic unfairness.

Previous work showed that demographic factors could influence the performance of automatic latent matching algorithms. However, the relationship between image quality and demographic effects has not been investigated. In the proposed study, we address the research question: "How do demographics affect matching of latent fingerprints of same image quality?". This work is inspired by recent studies in which ROC curves are used to evaluate the incremental effect of an additional marker in predicting a binary event. ROC curves are a standard way to evaluate the ability of a continuous marker to predict a binary outcome. The proposed approach uses ROC curves derived from regression models, where demographics and quality are both considered in the predictive model [2, 16, 23, 25]. The fitted values from the regression model are used to construct the ROC curve and compare it with the ROC curve derived from the regression excluding quality.

The contribution of this paper is three-fold: *i)* Investigate how demographic differentials affect the latent fingerprint image quality (LFIQ) algorithm, *ii)* Study how demographics impact automatic matching of latent fingerprints of the same quality, and *iii)* Discuss the use of covariate-specific ROC regression incorporating not only demographics but also LFIQ measures for a more fair assessment compared to the standard ROC. The

rest of the paper is organized as follows: Section 2 reviews research conducted on latent fingerprint image quality assessment, Section 3 describes the proposed quality-based mitigation approach, Section 4 presents the experimental results, Section 5 draws our conclusions and discusses future work.

2 LITERATURE REVIEW

Methods for demographics estimation from fingerprints have searched for gender clues in the ridge density structure that can be encoded by the local texture [18, 22]. An important study on an operational fingerprint database also revealed that fingerprint decision scores vary with subjects' demographic covariates [28]. The data used in Yoon's study were collected from records of the Michigan State Police with TenPrint cards as acquisition type. They experimentally confirmed that, for a given individual, genuine match scores decrease over time, impostor scores do not significantly vary, and that the accuracy remains stable. Image quality was considered as being the best covariate to explain the changes in the genuine match scores. Since Yoon's model only considers genuine and impostors separately and focuses on regression, our proposed framework could provide a more general approach applicable to adopt demographic information, as well as the quality of the latent prints.

A Latent Fingerprint Quality assessment was tasked by the FBI to assess the quality of friction ridge images for use by latent print examiners. Latent print examiners were tasked with providing an assessment of overall image usefulness, pattern classification, and anticipated difficulty for over 1,000 exemplar fingerprint images[6]. However, some unexplained variability by minutiae could be due to the lack of reproducibility of determinations among examiners. Thus, Yoon *et al.* proposed a method of defining quality measure for latent fingerprints, namely Latent Fingerprint Image Quality (LFIQ), which can be used as a predictor of latent matching performance. Based on the minutiae from latent image estimation, the earlier LFIQ method defines a latent quality measure by combining a qualitative quality feature (*i.e.*, the average ridge clarity) and a quantitative quality feature (*i.e.*, the number of minutiae) to estimate the objective target quality[29]. However, the earlier method that features minutiae count is not a good measurement for latent quality estimation in the presence of unreliable minutiae. Later on, Yoon further incorporated *i*) the connectivity of good ridge quality regions at the global level, *ii*) the reliability of minutiae, and *iii*) finger position estimation into the latent quality measurement (LFIQ) [27]. The experimental results show that the modified LFIQ has a high correlation with latent matching accuracy, and the model can be effectively used to reveal the quality measurement of latent fingerprint images.

LQMetric, an objective and automated algorithm to measure the quality of latent prints, was designed to predict AFIS performance as well as to augment or replace the informal subjective assessments of quality used by latent print examiners [6, 15]. LQMetric has been widely used since its release in 2014, it's a multipurpose tool that could provide assistance to examiners in determining which latents are appropriate for image-only searching (as opposed to requiring human-marked minutiae); replacing the informal "good", "bad", and "ugly" categories; and in research evaluating the efficacy of new latent print processing or development methods. Ezeobiejesi proposed an automatic region-of-interest based latent fingerprint quality assessment technique using deep learning. The experimental result has shown that in terms of predicting latent AFIS performance, the quality prediction by their deep learning model performs better than then state-of-art latent fingerprint value determination model [7].

A study by Dr. Ulery and Hicklin *et al.* proposed an introduction on assessment of fingerprint value and minutia count, and how it related to image clarity and feature content. The relationships were modeled between features and determined values by 21 certified latent print examiners in 1850 latent impressions [4]. The value of the print being suitable for individualization was assessed using EFS categories of value (VID), limited (VEO), and no

value (NV). A “no value” impression states that a comparison cannot be made solely based on individualized determination, regardless of quality. An association was found between minutiae count and value assessments; there was a high count of VEO and NV determinations and a low count of VID [4]. However, unexplained variability by minutia could be due to the lack of reproducibility of determinations among examiners. It would be beneficial to develop quality and quantity metrics in order to better analyze relationships and comparison determinations.

Later on, Hicklin *et al.* discussed how to correctly discern friction ridges in finger and hand prints, particularly its clarity. The researchers used an analysis/comparison approach to understand the difference between the value determinations made by latent print examiners and examiners’ annotation of minutiae, features, and image quality. The results yield a strong association between minutiae count and value determinations. This defined method can be used in metrics, data interchanging, or as an aid in fingerprint matching [11].

3 COVARIATE-SPECIFIC ROC CURVE

The Receiver Operating Characteristic (ROC) curve is a popular way to evaluate and compare the accuracy of classification markers when the outputs are continuous. Let Y be the continuous variable representing these outputs, and we wish to model $Pr[Y < j|X]$, the probability of a response in value j or lower for subjects with covariates X . While the pooled ROC curve does not take any covariates into consideration, the covariate-specific ROC curve models the covariate effects on the ROC curves[30].

Accounting for a set of covariates $X = (X_1, \dots, X_p)^T$ that may represent subject demographics or image quality may provide more specific error rates. We define the covariate-specific ROC curve as

$$\text{ROC}_{\mathbf{x}}(u) = 1 - F_{1,\mathbf{x}}(F_{0,\mathbf{x}}^{-1}(1-u)), \quad u \in (0, 1) \quad (1)$$

where $F_{1,\mathbf{x}}(t) = P(T \leq t | D = 1, X = \mathbf{x})$, is the distribution of scores in the genuine group conditional on the covariates, and $F_{0,\mathbf{x}}^{-1}(u) := \inf\{t \in \mathbb{R} : F_{0,\mathbf{x}}(t) \geq u\}$ is the quantile function of the imposter group conditional on the covariates[30].

The conditional distributions $F_{j,\mathbf{x}}(t)$ and $F_{j,\mathbf{x}}^{-1}(u)$, for $j \in \{0, 1\}$, are estimated from linear regression models on genuine and imposter match scores as follows:

$$T_j = \mu_j(\mathbf{x}) + \sigma_j(\mathbf{x})\epsilon_j, \quad j \in \{0, 1\}, \quad (2)$$

where the conditional mean and the conditional variance of T are $\mu_j(\mathbf{x}) = E(T|D = j, X = \mathbf{x})$ and $\sigma_j^2(\mathbf{x}) = \text{var}(T|D = j, X = \mathbf{x})$ given observed covariates $X = \mathbf{x}$, respectively. And the error term ϵ_j is independent of \mathbf{x} . Then, for a given covariate \mathbf{x} , the covariate-specific ROC curve can be expressed as:

$$\text{ROC}_{\mathbf{x}}(u) = 1 - G_1[G_0^{-1}(1-u)\frac{\sigma_0(\mathbf{x})}{\sigma_1(\mathbf{x})} - \frac{\mu_1(\mathbf{x}) - \mu_0(\mathbf{x})}{\sigma_1(\mathbf{x})}], \quad (3)$$

Here $G_i(z) = P(z \leq \epsilon_i)$ is the distribution function of the the regression error term which is independent of covariates for $i \in \{0, 1\}$. In this paper, we assume that $G(z)$ is the normal CDF, and the ROC curve in (3) is the so-called binormal ROC curve. The derivation for this expression can be found in [21] and [26].

We refer to the unknown source (i.e., the latent fingerprint evidence) as Query Q while the one pertaining to the known source as Reference R . The demographic covariates for the source subjects are denoted similarly. For instance, the age variable pertaining to the unknown source is referred to as Age_Q while the one pertaining to the known source as Age_R . The same scheme is used for the gender and ethnicity covariates. We used two

regression models to account for the effect of LFIQ on matching scores. The unique feature of the regression models for the ROC curve is the inclusion of the label term to indicate whether a matching score is genuine or imposter, and also the inclusion of the interaction between covariates and the label[30]. The interaction terms ensure that genuine and imposter groups have different population means.

The first model accounting for the effect of only LFIQ score is given by

$$Score \sim \beta_0 + \beta_L * label + \beta_1 * LFIQ + \beta_{L1} * LFIQ * label, \quad (\text{Model A})$$

where the variable *label* equals '1' if the two images are from the same subjects, while it equals to '0' otherwise. To account for the demographics, the second model includes age, gender and ethnicity as follows:

$$\begin{aligned} Score \sim & \beta_0 + \beta_L * label + \beta_1 * LFIQ + \beta_2 * Age_Q + \beta_3 * Age_R + \beta_4 * Gender_Q + \beta_5 * Gender_R \\ & + \beta_6 * Ethnicity_Q + \beta_7 * Ethnicity_R + \beta_{L1} * LFIQ * label + \beta_{L2} * Age_Q * label \\ & + \beta_{L3} * Gender_Q * label + \beta_{L4} * Ethnicity_Q * label. \end{aligned} \quad (\text{Model B})$$

In the model, Age_Q and Age_R indicate the age of the two images being compared, while $Gender_Q$ and $Gender_R$ represent their gender category (male or female), and $Ethnicity_Q$ and $Ethnicity_R$ represent their Ethnicity category (caucasian or non-caucasian). Specifically, when the subject is male, then gender equals '1', when the subject is caucasian, then ethnicity equals '1'. We only consider the interaction term of label and the demographics of the unknown source (e.g., Age_Q) because if $label=1$, the Score is from comparing the fingerprints of the same person, and hence the demographics of the unknown source and the reference are the same (e.g., $Age_Q = Age_R$). When $label=0$, then the interaction term is also equal to 0. Therefore, it is not meaningful to include both interaction terms in the model. After the model above is estimated, the regression results can be used to compose the covariate-specific ROC curve in Eqn. (3), then the accuracy of computer algorithms can be estimated by the curves and can also be summarized using the area under the ROC curve (AUC).

3.1 Image Quality from Latent Prints

A deep network-based minutiae extractor, referred to as MinutiaeNet, is used to obtain the minutiae points from the input latent fingerprint image. LFIQ is then applied to the minutiae map extracted in the previous step to estimate image quality. The algorithms can only be applied to latent prints with more than five minutiae points. LFIQ is computed by three components, i.e., ridge quality expressed as local ridge continuity, minutiae reliability and finger position. A diagram illustrating the information flow of the proposed approach can be seen in Fig. 1.

3.1.1 MinutiaeNet. Fingerprint comparison is primarily based on minutiae points comparison. Several hand-crafted approaches have been used to augment the minutiae with their attributes to improve the recognition accuracy [14]. However, a robust automatic fingerprint minutiae extraction that's suitable for noisy fingerprint images, continues to be a bottleneck in fingerprint recognition technique. With rapid developments in computer technology, deep learning approach has been used by other researchers for minutiae extraction. Typically, minutiae extraction and matching involves pre-processing stages such as ridge extraction and ridge thinning mentioned above, followed by minutiae extraction, and some heuristics to define minutiae attributes [14]. While such an approach performs well for high-quality images, its performance does degrade for poor quality rolled/plain prints, in particular for latent prints.

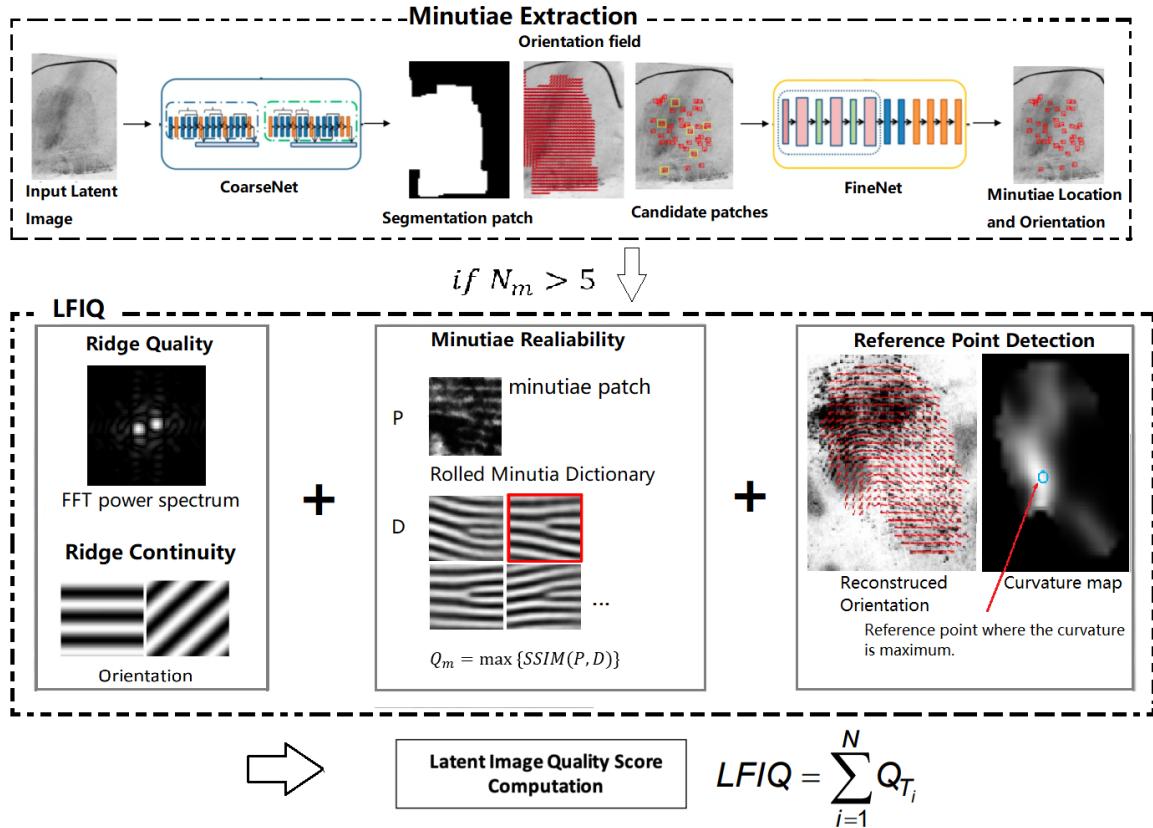


Fig. 1. The LFIQ extraction process.

Extracting minutiae automatically from latent prints is a challenging task. Tang *et al.* utilized the idea of object detection (non-maximum suppression) to detect candidate minutiae location and orientation, but it suffers from two major weaknesses such as hard threshold to delete the candidate patches, and the use of plain stacked CNN that suffers from vanishing gradient [24].

For this paper, a fully automatic minutiae extractor called MinutiaeNet is used for minutiae extraction; specifically, we use the MinutiaeNet introduced by Cao and Nguyen *et al.* [1]. This tool consists of a robust patch-based minutiae classifier that significantly boosts the precision and recall of candidate patches. This approach could provide reliable minutiae location and orientation without using hard threshold or fine-tuning. This method uses residual net instead of just plain stacked convolutional layers to make the classifier more precise [9]. The experimental results show that the MinutiaeNet is robust and has superior performance in terms of precision, recall and F1 values over published state-of-the-art on both benchmark datasets, namely FVC 2004 and NIST SD27 [20]. The architecture is based on two deep neural networks called CoarseNet and FineNet. CoarseNet uses a residual learning-based convolutional neural network with fingerprint domain knowledge to predict the minutiae score map and minutiae orientation. This minutiae score map is generated using latent fingerprint as a primary input and the corresponding enhanced image, segmentation map, and orientation field as a

secondary input. FineNet, on the other hand, is a minutiae classifier based on an inception residual network that processes each candidate patch to improve the minutiae score map and approximate minutiae orientation using regression [20].

In cases of small amount of candidate minutiae given, an adaptive threshold was applied in minutiae classifier - FineNet for determining final minutiae:

$$\text{Threshold} = \begin{cases} 0.45, & \text{if } N_m > 20 \\ 0.45 - n, & \text{if } N_m \leq 20 \\ n \in \{0.05, 0.1, 0.15 \dots 0.40\} \end{cases}$$

Where N_m is the number of candidates minutiae for each latent fingerprint. This adaptive threshold will obtain all candidate points over 45% threshold, plus additional top-ranking candidates until either minimum threshold of 0.05 is reached or at least 20 minutiae points are provided.

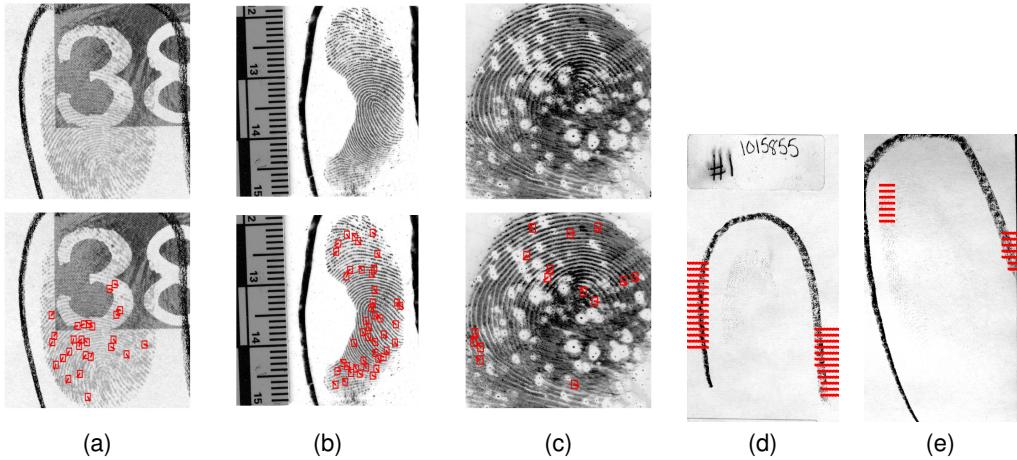


Fig. 2. Sample images result from MinutiaeNet on right thumb finger from the FBI Biometric Collection (BioCoP) Next Generation Identification Phase 1 (2008 - 2009) dataset[12]: (a) occluded (b) partial and (c) distorted latent fingerprint impressions; (d) and (e) are wrong orientation field (red area) obtained from the low contrast dry-latent print (unsurprisingly LFIQ = 0).

The result is shown in Fig. 2, minutiae points extracted from right index and right thumb with different background noise are illustrated. The algorithm works well in difficult situations such as Fig.2(a) distorted, Fig.2(b) occluded, and Fig.2(c) partial latent fingerprints which demonstrated its robustness to noisy background. Orientation field of two latent impressions marked in Fig.2(d) and 2(e) show the CoarseNet mapping suffers poor performance from dry latent that are extremely low contrast in grey scale (compared to background), due to the discontinuity of broken ridges. The algorithm failed to detect the correct center of a latent at the beginning, this is caused by the weak continuity of true ridges has been further coarsened into background noise while each level of residual net is fused to get the final minutiae score map. In fact, these phenomena commonly happen to latent fingerprints that belong to the no value (NV) category due to the texture of the ridges becomes indistinguishable even by a professional examiner's eye.

3.1.2 Latent Fingerprint Image Quality (LFIQ). LFIQ score can be an independent latent fingerprint's quality indicator for evaluating the matching performance of any existing Automated Fingerprint Identification System (AFIS). This proposed model can be effectively used to *i*) automatize quality measurement of latent fingerprint image and *ii*) assist latent examiners in their value determination [27].

The objective quality metric for latent fingerprints is the Latent Fingerprint Image Quality (LFIQ). It can be utilized to successfully distinguish high-quality latent fingerprints that don't require human intervention, as well as compensate for the subjective aspect of value determination by latent examiners.

The LFIQ is determined by three parameters that affect the quality of latent fingerprints: (i) ridge quality, (ii) minutiae reliability, and (iii) finger position. Local ridge clarity and the friction ridge regions with high ridge clarity determine ridge quality. The reliability of minutia is determined by its likelihood of being a genuine minutia. The position of a finger is determined by detecting the reference point (e.g., core point(s) or maximum curvature point for arch-type fingerprints) and assigning high weights to minutiae in the central regions of the finger in the LFIQ computation [27]

3.1.3 Latent Quality Score Computation. For each triangle T_i

$$Q_{T_i} = Q_{r_i} \sum_{j=1}^3 Q_{m_{ij}} W_{m_{ij}},$$

Where Q_{r_i} is the average ridge quality in T_i , $Q_{m_{ij}}$ is the Reliability of j-th minutia of T_i , $W_{m_{ij}}$ is the weight based on the finger position

The quality score of a latent is computed as follows:

$$\text{LFIQ} = \sum_{i=1}^N Q_{T_i}$$

N is the number of triangles in the latent print.

LFIQ requires the least amount of 5 candidate minutiae to achieve any meaningful result ($\text{LFIQ} > 0$).

4 EXPERIMENTAL RESULTS

4.1 Dataset

The dataset used in this study is a subset of the FBI Biometric Collection of People (BioCoP) Next Generation Identification Phase 1 (2008 - 2009) [12]. The data collection involved the acquisition of latent-deposited fingerprints on common materials as well as standard ink and paper methods. The ink and paper data was used as an exemplar set for both electronic capture performed using BioCOP and the latent substrate capture. Each scanned image is saved as a grayscale type image with a resolution of 1000 ppi. These fingerprints pertained to a total of 1504 subjects and were collected at West Virginia University. There was a nearly equal amount of male to female participants with 52% to 48% ratio. Also, among the participants, the age group between 18-29 was highest accounting for 74% percent of people, 8% between 30-39 years old, 7% between 40-49 years old and 11% above 50 years old. Among the ethnicity, Caucasians accounted for 79% of the people, only 6.2% Asian, 3.8% Asian Indian, 3.7% African American, 2.4% African, 2.1% Hispanic. For the experiments of this paper, we use the right index and right thumb.

The latent fingerprints collection was carried out by gloving the subject's hands with nitrile gloves that induce sweating required for the development of the first latent fingerprints. Three quality sets were needed, good,

bad and ugly, so that three whole or partial impressions for each finger were made on each of the substrates. Three different substrates were used: paper, plastic, and glass/porcelain. The items were separated based on substrate type and processed in one of three ways: *i*) chemical (ninhydrin) processing, *ii*) cyanoacrylate processing, *iii*) lift cards (processed with black fingerprint powder at the collection site). All fingerprints processed with cyanoacrylate were digitally photographed, while all ninhydrin and black powder fingerprints were scanned.

4.2 Latent-to-reference print comparison

The match scores in this paper were obtained using the end-to-end latent fingerprint search system recently published by Cao *et al.*. The algorithm does include automated ridge structure cropping, latent image pre-processing, feature extraction, feature comparison, and outputs a candidate list. The model is robust to poor quality latents by generated set of virtual minutiae to construct a texture template. This fully automated latent search system was evaluated on NIST SD27 (258 latents); MSP (1,200 latents), WVU (449 latents) and N2N (10,000 latents) against a background set of 100K rolled prints, which includes the true rolled mates of the latents with rank-1 retrieval rates of 65.7%, 69.4%, 65.5%, and 7.6% respectively[1].

4.3 Results

Fig. 3 shows the distributions of the genuine and impostors match scores generated using the MSU identification system for right index and right thumb.

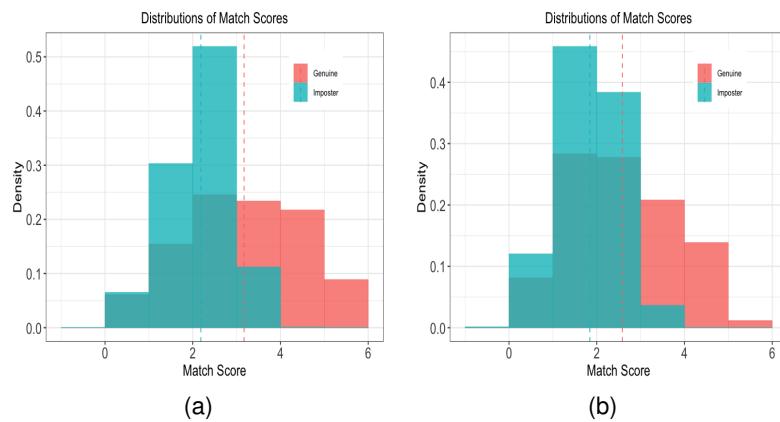


Fig. 3. Distributions of the Genuine and Impostor Match Scores: (a) Right Index and (b) Right Thumb.

Fig. 4 illustrates how demographics impact LFIQ measures. Fig. 4 (a) points out instead that LFIQ differs across different age groups with the younger population achieving the highest scores ($LFIQ > 60$) while the subjects greater than 41 years old exhibit an upper bound of only $LFIQ = 20$. Fig. 4 (b) suggests that there is no trend with respect to gender, which may be due to a lack of textural information capture by the LFIQ algorithm. From the literature, gender estimation from fingerprints does exploit textural differences between males and females. Similarly, Fig. 4 (c) highlights that higher LFIQ scores are more likely to be achieved by the Non-Caucasian population. The same behavior can be found in Fig. 4 for the right thumb finger as we can see in (d), (e) and (f).

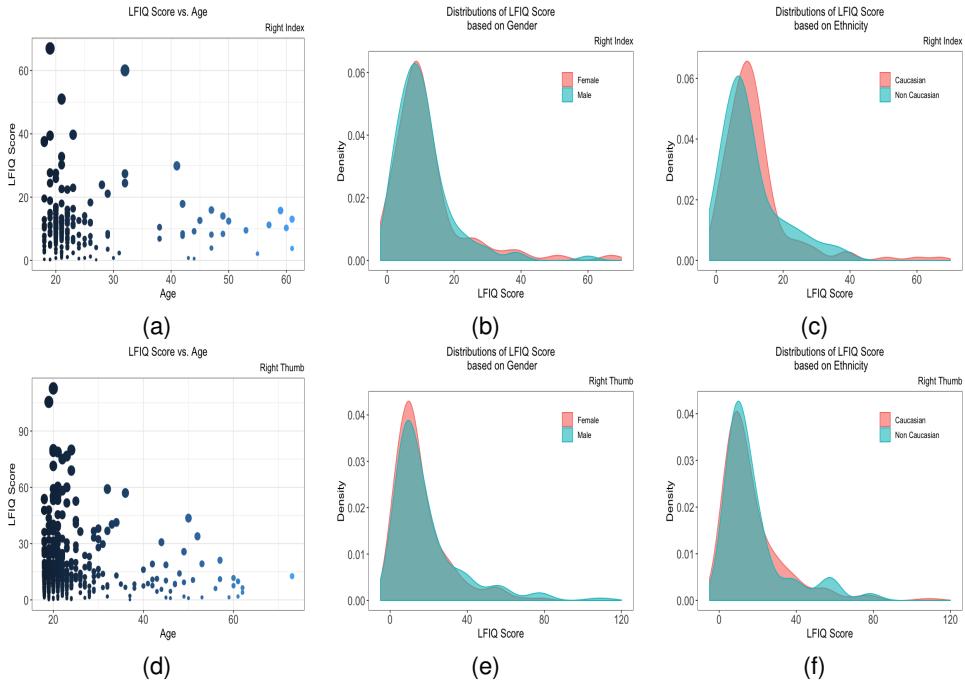


Fig. 4. Distributions of LFIQ Score vs. Demographics: (a), (b) and (c) are from right index finger; (d), (e) and (f) are from right thumb finger.

In Fig. 5, the scatter plot of match scores versus LFIQ measures pertaining to right thumb shows a clear improvement of the separability between genuine and impostors for $LFIQ > 60$. A similar trend has been observed for the right index.

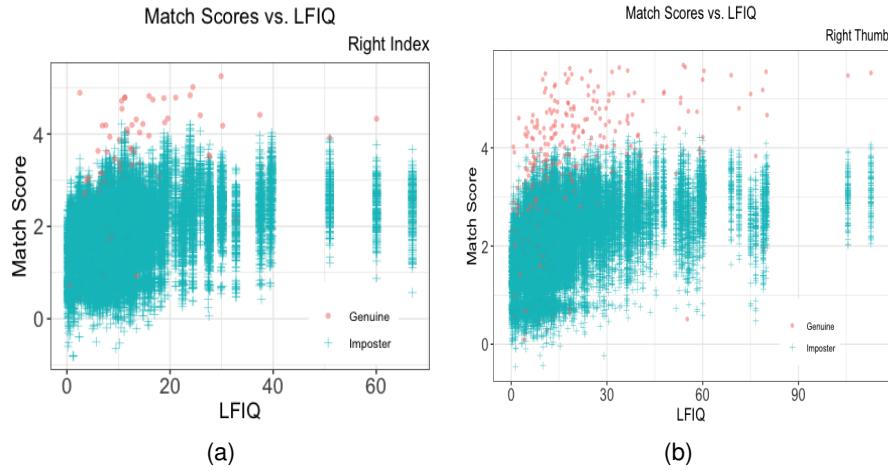


Fig. 5. Scatter plot of Match Scores vs. LFIQ score.

We examine how the ROC curve varies conditioning on observed covariates by reporting the regression results for the two covariate-specific ROC models: **Model A** and **Model B**. Various classifiers are capable of predicting the likelihood of a sample belonging to a class. A probabilistic classifier is implemented by setting a threshold which divides the entire data into different classes. The results shown in Table 1 were used to compute the sensitivities of the right index and the right thumb fingers. The residual plots of the model in Fig.6 were used to check the model diagnostics, the residual plots of **Model B** for the right index indicate that the linear model fits the data very well and also validates the model's normality assumption. The similar results have been observed in right thumb data and also in **Model A**. The significant interaction terms, as we can see in Table 1, indicate the mean difference between genuine and imposter scores changes with covariates. For example, **Model B** for the right index fingers has significant interaction terms, $LFIQ^*\text{label}$, $Gender_Q^*\text{label}$ and $Ethnicity_Q^*\text{label}$. This indicates that the difference in mean scores between genuine and imposter groups change when LFIQ scores or age change. The difference in mean matching scores between genuine and imposter groups also change with ethnicity levels. For the right thumb, significant interaction terms indicate the difference in mean matching scores between genuine and imposter groups change as LFIQ scores and for different gender change. The mean difference is a main component in the ROC curve expression. Significant interactions imply when the covariates change, the ROC curve adjusting for the covariates tends to have significant changes.

Table 1. Result of the designed covariate-specific ROC Models

Model	Right Index		Right Thumb	
	(A)	(B)	(A)	(B)
Intercept	1.5756*	1.3017*	1.8731*	1.8833*
Label	0.4419*	-0.3022	0.6196*	0.4667*
$LFIQ$	0.0227*	0.0229*	0.0177*	0.0175*
Age_Q	-	0.0015*	-	-0.0025*
Age_R	-	0.0021*	-	0.0020*
$Gender_Q$	-	0.0797*	-	0.0167*
$Gender_R$	-	0.1057*	-	0.0930*
$Ethnicity_Q$	-	0.0624*	-	-0.0827*
$Ethnicity_R$	-	0.0251*	-	0.0117
$LFIQ^*\text{label}$	0.0257*	0.0266*	0.0216*	0.0210*
$Age_Q^*\text{label}$	-	0.0091	-	0.0034
$Gender_Q^*\text{label}$	-	0.4077 *	-	0.2535*
$Ethnicity_Q^*\text{label}$	-	0.3453*	-	-0.0742

* p-value<0.05

Fig.7 shows the ROC curves based on **Model A** only considering LFIQ score, and the corresponding AUC values are also reported in the figure. Here, the 25th, 50th, and 75th percentiles of LFIQ score were chosen to compute the curves. The purple curve in the figure is the pooled ROC curve used for comparison. We can observe that as LFIQ score increases, the model's identifying ability increases.

Fig.8 and Fig.9 show the ROC curves based on **Model B** considering demographics and also LFIQ score for index and thumb finger, respectively. We can see that when the demographics are constant, the increase in LFIQ will increase the model's identifying ability, which is consistent with the univariate **Model A**.

In Fig.8, we can also observe that for the same LFIQ score and demographics, as age increases the model has a much better identifying ability for right index, while the increase in identifying ability for right thumb in Fig.9 is much smaller. Fig.8 also shows that when LFIQ scores and other demographics are adjusted, male subjects perform better in identification than females; we can find similar behavior in Fig.9. In Fig.8, it is observed that

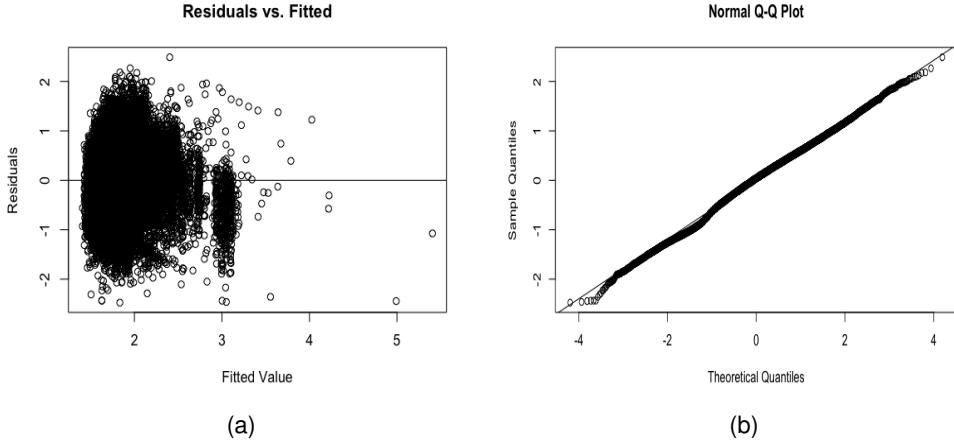


Fig. 6. Model Diagnostics of Model B for Right Index.

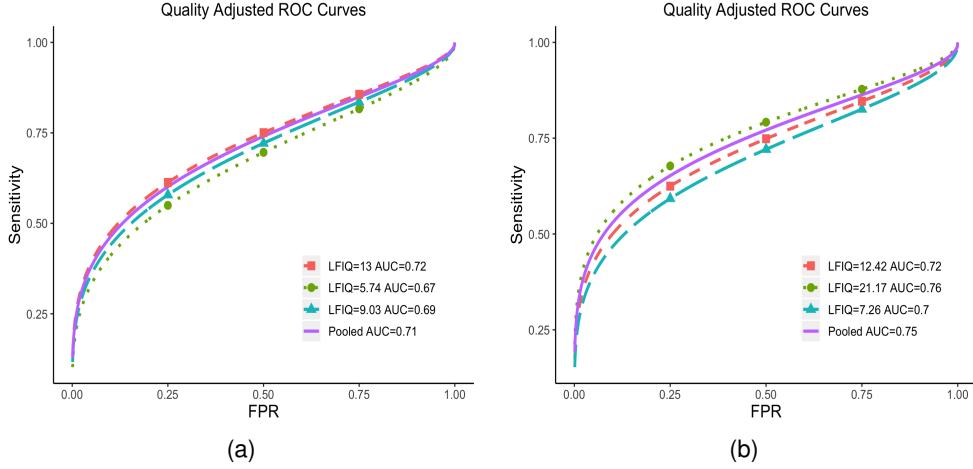


Fig. 7. Covariate-specific ROC curves conditioned on LFIQ scores: (a) Right Index, (b) Right Thumb.

when the LFIQ score and other demographics are identical, the caucasian group has a better identifying ability compared with the non-caucasian group for the right index. While for the right thumb, as we can see in Fig.9, the non-caucasian group performs better in identification. And also, based on the regression results in Table 1, we can know that Gender has the most significant impact on the algorithm's performance, while Age has the slightest influence among the demographics.

In Fig.8 (c) and (d), we can see that although LFIQ is not able to provide insights about gender differentials, when incorporated as an additional covariate into the proposed demographic-adjusted ROC regression, it can contribute to improving the performance that was not obtained in previous research based only on demographic covariates. Specifically, from Fig.8 (d) we can see that for males we have a better performance using the quality-based adjusted-ROC with higher LFIQ values.

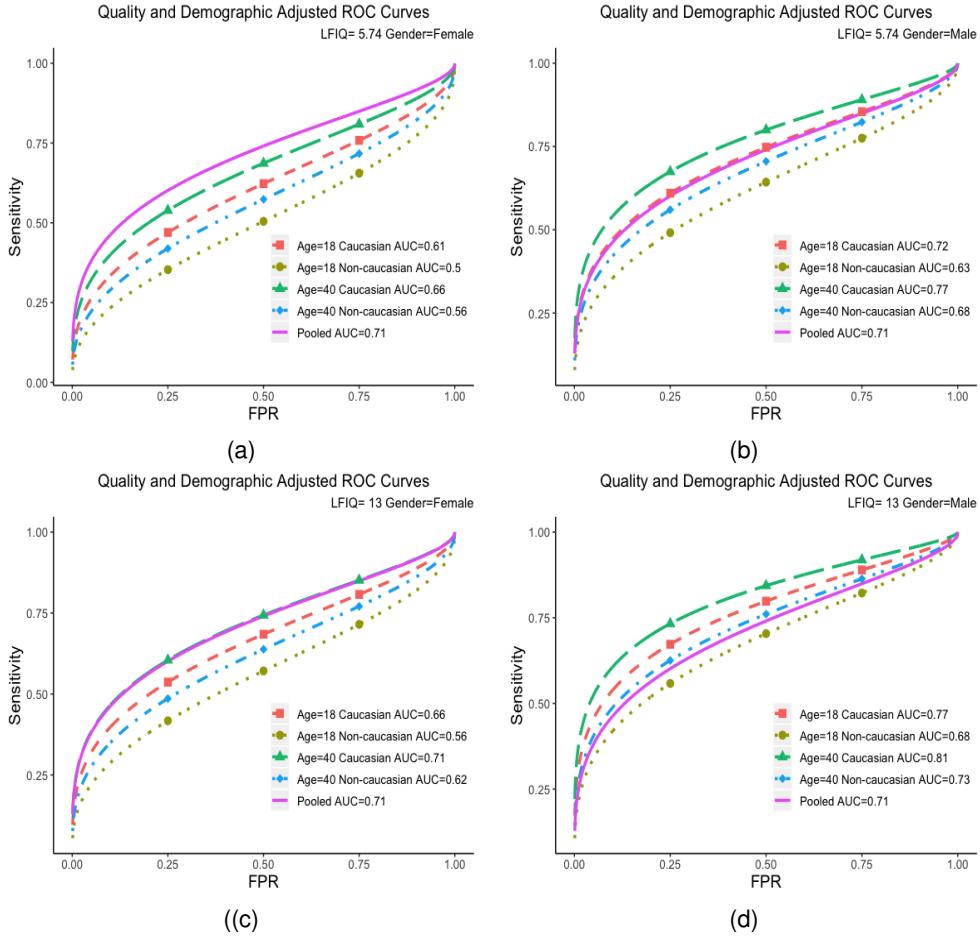


Fig. 8. Covariate-specific ROC curves conditioned on LFIQ scores and demographics for right index.

5 CONCLUSIONS

In most cases, latent fingerprint images are obtained under non-ideal acquisition conditions, resulting in partial or distorted impressions, contaminated by background noise. Reliable latent quality assessment can prevent any evidence of value from being discarded; contrarily, by rejecting poor quality fingerprints, an AFIS can reduce incidents of false accepts or false rejects. The identification ability of an automatic matching algorithm is expected to increase in the presence of a high-quality image, and it could be further improved if the demographics of each subject can be obtained.

In this paper, we evaluate a regression model derived from match scores in which both image quality and demographic covariates are incorporated. The regression model includes the interaction between covariates and group status (genuine or impostor) so that all scores are analyzed in the same model. The covariate-adjusted ROC curve is then obtained from the regression coefficients. Our findings show that the proposed covariate-adjusted assessment scheme conditioned on image quality and demographics is more informative than the

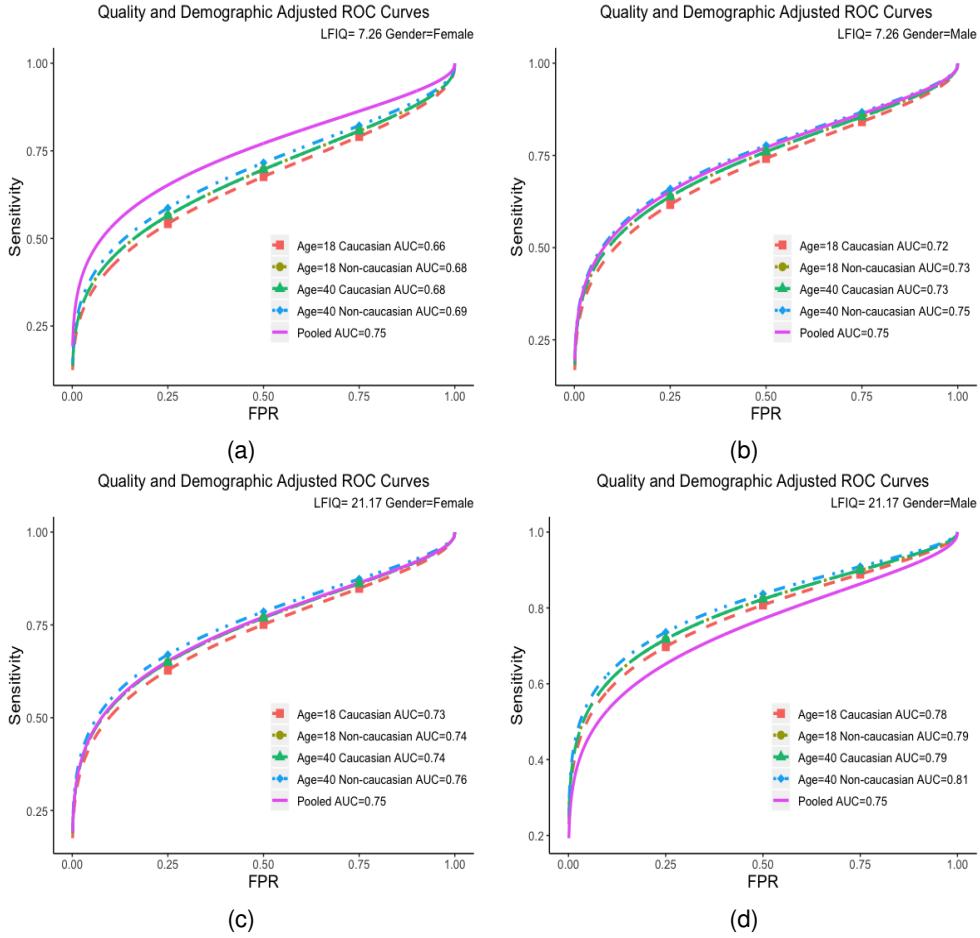


Fig. 9. Covariate-specific ROC curves conditioned on LFIQ scores and demographics for right thumb.

traditional ROC curve. Specifically, between high quality and low quality the impact of demographics on the ROC curves changes. The accuracy increases as age increases for right index and right thumb. Also, the matching algorithm tends to have higher accuracy in male subjects than females for both types of fingers. The accuracy is higher for the Caucasian group than the non-Caucasian group for the right index, but opposite for the thumb.

In this paper, the match scores were extracted using Cao's 2018 algorithm based on comparing minutiae and texture templates. A more "tailored" analysis of the quality covariate can be obtained by accessing the minutiae template using by the matching algorithm. Furthermore, the MinutiaeNet CNN-based minutiae extractor used in this work was trained on the FVC 2002 dataset that does not contain examples of dry fingers. Thus, the minutiae extracted from dry fingers in the WVU database are often not enough for being processed through LFIQ.

In future work, we will: *i*) Enhance the robustness of the CNN-based minutiae extractor to dry fingers by fine-tuning it on the WVU database; *ii*) Access to the minutiae templates extracted by the MSU identification

system and apply LFIQ to those; and *iii*) Enhance the LFIQ algorithm by fusing multi-layer of minutiae maps from additional extractors for increased robustness in the presence of poor image quality.

REFERENCES

- [1] Kai Cao, Dinh-Luan Nguyen, Cori Tymoszek, and Anil K. Jain. 2018. End-to-End Latent Fingerprint Search. *arXiv preprint arXiv:1812.10213* (2018).
- [2] Nancy R Cook. 2008. Statistical evaluation of prognostic versus diagnostic models: beyond the ROC curve. *Clin Chem*. Jan;54(1) (2008), 17–23.
- [3] Tabassi Elham and Grotherand Patrick. 2009. *Fingerprint Image Quality*. Springer US, Boston, MA, 482–490. https://doi.org/10.1007/978-0-387-73003-5_52
- [4] Bradford T. Ulery et al. 2013. Understanding the sufficiency of information for latent fingerprint value determinations. *Forensic Science International* 226 (2013), 99–106.
- [5] Cynthia M. Cook et al. 2019. Demographic Effects in Facial Recognition and Their Dependence on Image Acquisition: An Evaluation of Eleven Commercial Systems. *IEEE Transactions on Biometrics, Behavior, and Identity Science* vol.1, no.1 (2019), 32–41.
- [6] Hicklin et al. 2011. Latent Fingerprint Quality: A Survey of Examiners. *Journal of Forensic Identification* 61(4) (2011), 385–419.
- [7] Jude Ezeobiejesi and Bir Bhanu. 2018. Latent Fingerprint Image Quality Assessment Using Deep Learning. *Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* (2018), 621–629.
- [8] P Gnanasivam and Dr S Muttan. 2012. Estimation of Age Through Fingerprints Using Wavelet Transform and Singular Value Decomposition. *International Journal of Biometrics and Bioinformatics (IJBB)* 6, 2 (2012), 58–67.
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. *IEEE CVPR* (2016), 770–778.
- [10] Deborah Hellman. 2020. Measuring Algorithmic Fairness. *Va. L. Rev.* 106 (2020), 811.
- [11] R Austin Hicklin, JoAnn Buscaglia, and Maria Antonia Roberts. 2013. Assessing the clarity of friction ridge impressions. *Forensic science international* 226, 1-3 (2013), 106–117.
- [12] L. Hornak, W. LaRue, B. Cukic, A. Ross, K. Morris, J. Dawson, S. Crihalmeanu, N. Kalka, and N. Kayal. 2009. FBI Biometric Collection of People (BioCoP): Next Generation Identification Phase 1 (2008 - 2009). *2008 Biometric Collection Project 08-06-2008 to 12-31-2009 FINAL REPORT* (2009).
- [13] Crystal Huynh, Erica Brunelle, Lenka Halamkova, Juliana Agudelo, and Jan Halamek. 2015. Forensic identification of gender from fingerprints. *Analytical chemistry* 87, 22 (2015), 11531–11536.
- [14] Anil Jain, Lin Hong, and Ruud Bolle. 1997. On-line fingerprint verification. *IEEE Trans PAMI*,19(4) (1997), 302–314.
- [15] Nathan D Kalka, Michael Beachler, and R. Austin Hicklin. 2020. LQMetric: A Latent Fingerprint Quality Metric for Predicting AFIS Performance and Assessing the Value of Latent Fingerprints. *Journal of Forensic Identification* 70(4) (2020), 443–463.
- [16] Michael W. Kattan. 2003. Judging new markers by their ability to improve predictive accuracy. *PubMed*. May 7;95(9) (2003), 634–5.
- [17] Emanuela Marasco, Mengling He, Larry Tang, and Sumanth Sriram. 2020. Accounting for Demographic Differentials in Forensic Error Rate Assessment of Latent Prints via Covariate-Specific ROC Regression. *CVIP 2020 CCIS* 1376 (2020), 338–350.
- [18] Emanuela Marasco, Luca Lugini, and Bojan Cukic. 2014. Exploiting Quality and Texture Features to Estimate Age and Gender from Fingerprints. *SPIE Defense and Security* (2014).
- [19] Shira Mitchell, Eric Potash, Solon Barocas, Alexander D'Amour, and Kristian Lum. 2020. Algorithmic Fairness: Choices, Assumptions, and Definitions. *Annual Review of Statistics and Its Application* 8 (2020), 141–63.
- [20] Dinh-Luan Nguyen, Kai Cao, and Anil K. Jain. 2017. Robust Minutiae Extractor: Integrating Deep Networks and Fingerprint Domain Knowledge. *CoRR* abs/1712.09401 (2017). arXiv:1712.09401 <http://arxiv.org/abs/1712.09401>
- [21] Margaret Sullivan Pepe. 2000. An Interpretation for the ROC Curve and Inference using GLM Procedures. *Biometrics* 56, 2 (2000), 352–359.
- [22] Ajita Rattani, Cunjian Chen, and Arun Ross. 2014. Evaluation of texture descriptors for automated gender estimation from fingerprints. (2014), 764–777.
- [23] Venkatraman E. Seshan, Mithat Gönen, and Colin B. Begg. 2013. Comparing ROC Curves Derived From Regression Models. *Stat Med*. April 30; 32(9): (2013), 1483–1493.
- [24] Yao Tang, Fei Gao, and Jufu Feng. 2017. Latent fingerprint minutia extraction using fully convolutional network. *IEEE IJCB* (2017).
- [25] AN Tosteson, Milton C Weinstein, Jack Wittenberg, and Colin B Begg. 1994. ROC curve regression analysis: the use of ordinal regression models for diagnostic test assessment. *PubMed*. 102 Suppl 8(Suppl 8) (1994), 73–80.
- [26] Anna N Angelos Tosteson and Colin B Begg. 1988. A General Regression Methodology for ROC Curve Estimation. *Medical Decision Making* 8, 3 (1988), 204–215.
- [27] Soweon Yoon, Kai Cao, Eryun Liu, and Anil K. Jain. 2013. LFIQ: Latent fingerprint image quality. In *2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems (BTAS)*. 1–8.

- [28] Soweon Yoon and Anil K Jain. 2015. Longitudinal study of fingerprint recognition. *Proceedings of the National Academy of Sciences* 112, 28 (2015), 8555–8560.
- [29] Soweon Yoon, Eryun Liu, and Anil K Jain. 2012. On latent fingerprint image quality. *Computational Forensics* (2012), 67–82.
- [30] Xiao-Hua Zhou, Donna K McClish, and Nancy A Obuchowski. 2009. *Statistical Methods in Diagnostic Medicine*. John Wiley & Sons.