

000
001
002

Quality- and Demographic-Specific ROC Curves derived from Regression Models of Match Scores involving Latent Prints

003
004
005
006
007

008 Anonymous WACV submission

009
010
011
012
013014 Paper ID ****
015016
017
018
019
020

Abstract

021
022
023
024

Recent research has focused on gaining crucial insights about demographic differentials and their sources by analyzing how they influence matching. Previous work analyzes biases in matching algorithms without considering the impact of quality estimators that might also be biased due to their learning-based nature. In the proposed study, the research question we address is: "How do demographics affect matching of latent fingerprints of same image quality?".

025
026
027
028
029

We extend our previous approach that uses ROC curves derived from regression models by integrating quality measures extracted from latent prints. The contribution focuses on considering quality in the predictive model and not only demographics.

030
031
032
033
034
035
036
037
038
039

The evaluation is carried out on the FBI WVU BioCop 2008 database that contains 469 right-thumb and 219 right-index latent fingerprint images with associated demographics. Experimental results provide comprehensive empirical support related to the quality-based demographic differentials' impact on the performance of automatic latent fingerprint matching.

040
041
042
043
044
045
046
047
048
049
050
051
052
053

1. Introduction

Uncovering and clarifying biases in biometric technologies has recently gained increasing interest [6, 14]. 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. In this paper, we analyze fairness of prediction decision-making in forensics through discovering and mitigation of biases in automated algorithms operating on latent fingerprint images. 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

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

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 [25]. 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 [12].

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 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 [19]. Incorporating quality measures for latent fingerprints enables a more objective assessment and mitigation of algorithmic unfairness. 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. Most fingerprint matchers are based on minutiae points (e.g., ridge ending and ridge bifurcation). Image quality for query latent prints is defined through the extraction of Latent Fingerprint Image Quality (LFIQ) scores. Factors such as adverse skin conditions (e.g., dry, creased) can lower image quality and degrade performance.

ROC curves are a standard way to evaluate the ability of a continuous marker to predict a binary outcome. 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. The proposed approach uses

ROC curves derived from regression models, where demographics and quality are both considered in the predictive model [18, 22, 11, 2]. The fitted values from the regression model are used to construct the ROC curve and to 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 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

Fingerprint features have 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, 4, 13]. Methods for demographics estimation from fingerprints have searched for gender clues in the ridge density structure that can be encoded by the local texture [13, 17].

An important study on an operational fingerprint database also revealed that fingerprint decision scores vary with subjects' demographic covariates [25]. 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[3]. However, some unexplained variability by minutiae 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 de-

terminations. As a result, 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[24]. 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 global level, *ii*) reliability of minutiae, and *iii*) finger position estimation into the latent quality measurement (LFIQ) [23]. 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 image.

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. While the pooled ROC curve doesn't take any covariates into consideration, the covariate-specific ROC curve models the covariate effects on the ROC curves and are commonly used in medical diagnostics [26]. It provides error rates specific to the demographics of source subjects. Accounting for a set of covariates $X = (X_1, \dots, X_p)^T$ that may represent subject demographics, information on quality of the measurement on the process, etc., the ROC curve varies conditional on observed covariates $X = \mathbf{x}$.

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[26].

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

216 can be expressed as:
 217

$$218 \text{ROC}_{\mathbf{x}}(u) = 1 - G_1[G_0^{-1}(1-u)\frac{\sigma_0(\mathbf{x})}{\sigma_1(\mathbf{x})} - \frac{\mu_1(\mathbf{x}) - \mu_2(\mathbf{x})}{\sigma_1(\mathbf{x})}], \\ 219 \quad (3)$$

220 Here $G_i(z) = P(z \leq \epsilon_j)$ is the distribution function of the
 221 the regression error term which is independent of covariates
 222 for $i \in \{0, 1\}$. In this paper, we assume that $G(z)$ is the
 223 normal CDF, and the ROC curve in (3) is **so-called** binormal
 224 ROC curve. The derivation for this expression can be found
 225 in [16] and [21].

226 We refer to the unknown source (i.e., the latent finger-
 227 print evidence) as Query Q while the one pertaining to the
 228 known source as Reference R . The demographic covariates
 229 for the source subjects are denoted similarly. For instance,
 230 the age variable pertaining to the unknown source is referred
 231 to as Age_Q while the one pertaining to the known source as
 232 Age_R . The same scheme is used for the gender and ethnic-
 233 ity covariates. We used two regression models to account
 234 for the effect of LFIQ on matching scores. The unique fea-
 235 ture of the regression models for the ROC curve is the inclu-
 236 sion of the label term to indicate whether a matching score
 237 is genuine or imposter, and also the inclusion of the interac-
 238 tion between covariates and the label[26]. The interaction
 239 terms ensure that genuine and imposter groups have differ-
 240 ent population means.
 241

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

$$244 Score \sim \beta_0 + \beta_L * label + \beta_1 * LFIQ \\ 245 + \beta_{L1} * LFIQ * label, \\ 246 \quad (Model A)$$

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

$$251 Score \sim \beta_0 + \beta_L * label + \beta_1 * LFIQ + \beta_2 * Age_Q \\ 252 + \beta_3 * Age_R + \beta_4 * Gender_Q + \beta_5 * Gender_R \\ 253 + \beta_6 * Ethnicity_Q + \beta_7 * Ethnicity_R \\ 254 + \beta_{L1} * LFIQ * label + \beta_{L2} * Age_Q * label \\ 255 + \beta_{L3} * Gender_Q * label \\ 256 + \beta_{L4} * Ethnicity_Q * label. \\ 257 \quad (Model B)$$

258 In the model, Age_Q and Age_R indicate the age category
 259 of the two images being compared, while $Gender_Q$ and
 260 $Gender_R$ represent their gender category (male or female),
 261 and $Ethnicity_Q$ and $Ethnicity_R$ represent their Ethnicity
 262 category (caucasian or non-caucasian). Specifically, when
 263 the subject is male then gender equals to '1', when the

264 subject is caucasian then ethnicity equals to '1'. We only
 265 consider the interaction term of label and the demograph-
 266 ics of the unknown source (e.g. Age_Q) because if label=1,
 267 the Score is from comparing the fingerprints of the same
 268 person, and hence the demographics of the unknown source
 269 and the reference are the same (e.g. $Age_Q = Age_R$). When
 270 label=0, then the interaction term is also equal to 0. There-
 271 fore, it is not meaningful to include both interaction terms in
 272 the model. After the model above is estimated, the regres-
 273 sion results can be used to compose the covariate-specific
 274 ROC curve in Eqn. (3), then the accuracy of computer al-
 275 gorithms can be estimated by the curves and can also be
 276 summarized using the area under the ROC curve (AUC).

3.1. Image Quality from Latent Prints

277 A deep network-based minutiae extractor, referred to as
 278 MinutiaeNet, is used to obtain the minutiae points from the
 279 input latent fingerprint image. LFIQ is then applied to the
 280 minutiae map extracted in the previous step to estimate im-
 281 age quality. The algorithms can only be applied to latent
 282 prints with more than 5 minutiae points. LFIQ is computed
 283 by three components, i.e., ridge quality expressed as local
 284 ridge continuity, minutiae reliability and finger position. A
 285 diagram illustrating the information flow of the proposed
 286 approach can be seen in Fig. 1.

3.1.1 MinutiaeNet

287 Fingerprint comparison is primarily based on minutiae
 288 points comparison. Several hand-crafted approaches have
 289 been used to augment the minutiae with their attributes to
 290 improve the recognition accuracy [9]. However, a robust
 291 automatic fingerprint minutiae extraction that's suitable for
 292 noisy fingerprint images, continues to be a bottleneck in fin-
 293 gerprint recognition technique. With rapid developments
 294 in computer technology, deep learning approach has been
 295 used by other researchers for minutiae extraction. Typically,
 296 minutiae extraction and matching involves pre-processing
 297 stages such as ridge extraction and ridge thinning men-
 298 tioned above, followed by minutiae extraction , and some
 299 heuristics to define minutiae attributes [9]. While such an
 300 approach performs well for high quality images, its per-
 301 formance does degrade for poor quality rolled/plain prints in
 302 particular for latent prints.

303 Extracting minutiae automatically from latent prints is
 304 a challenging task. Tang *et al.* utilized the idea of object
 305 detection (non-maximum suppression) to detect candidate
 306 minutiae location and orientation, but it suffers from two
 307 major weaknesses such as hard threshold to delete the can-
 308 didate patches, and the use of plain stacked CNN that suf-
 309 fers from vanishing gradient [20].

310 For this paper, a fully automatic minutiae extractor called
 311 MinutiaeNet is used for minutiae extraction; specifically,

312 313 314 315 316 317 318 319 320 321 322 323

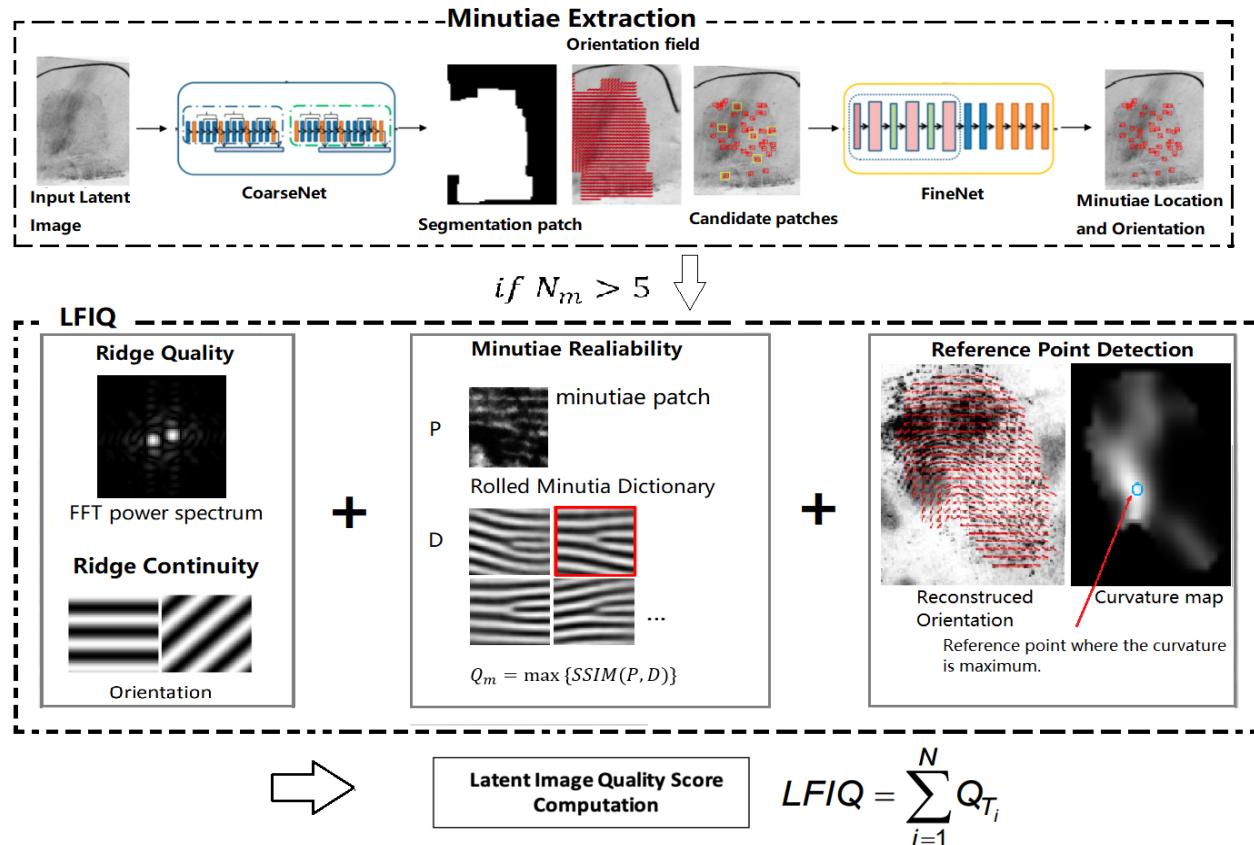


Figure 1. The LFIQ extraction process. Portions adapted from slides of [10].

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 [5]. 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 [15]. 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 finger print as a primary input and the corresponding enhanced image, segmentation map, and orientation field as 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 [15].

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 number of candidates minutiae for each latent fingerprint. This adaptive threshold will obtain all candidate points over 45% threshold, plus additional best ranking candidates until at least 20 minutiae points are provided.

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

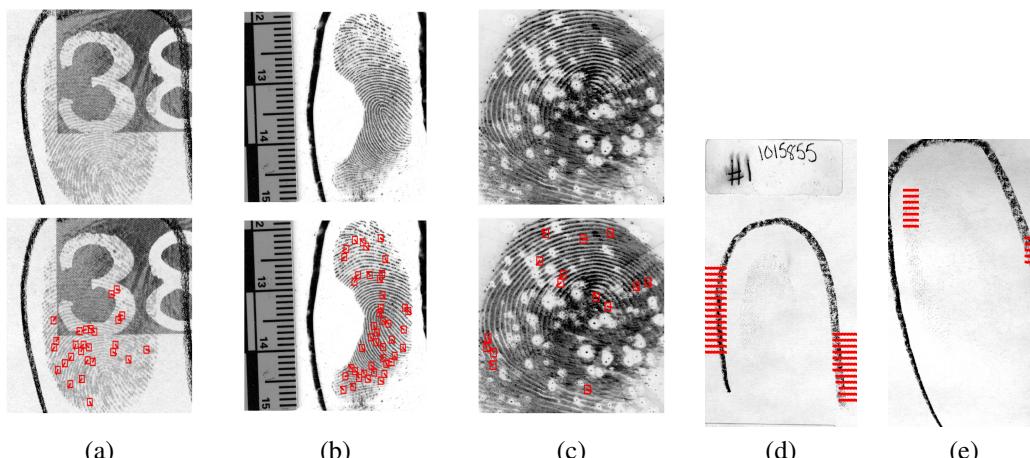


Figure 2. Sample images result from MinutiaeNet on right thumb finger from the FBI Biometric Collection (BioCoP) Next Generation Identification Phase 1 (2008 - 2009) dataset[7]: (a) occluded (b) partial and (c) distorted latent fingerprint impressions; (d) and (e) are wrong orientation field (red area) obtained from the low contrast dry-latent print (unsurprisingly LFIQ = 0).

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 the texture of the ridges becomes indistinguishable even by a professional examiner's eye.

3.1.2 Latent Fingerprint Image Quality (LFIQ)

In which, 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 the model can be effectively used to *i*) automatize quality measurement of latent fingerprint image and *ii*) assist latent examiners in their value determination [23].

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 archetype fingerprints) and assigning high weights to minutiae in the central regions of the finger in the LFIQ computation [23]

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

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

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) [7]. 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 pertain 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

540 18-29 was highest accounting for 74% percent of people,
 541 8% between 30-39 years old, 7% between 40-49 years old
 542 and 11% above 50 years old. Among the ethnicity, Cau-
 543 casians accounted for 79% of the people, only 6.2% Asian,
 544 3.8% Asian Indian, 3.7% African American, 2.4% African,
 545 2.1% Hispanic. For the experiments of this paper, we use
 546 right index and right thumb.
 547

548 The latent fingerprints collection was carried out by
 549 gloving the subject's hands with nitrile gloves that induce
 550 sweating required for the development of the first latent fin-
 551 gerprints. Three quality sets were needed, good, bad and
 552 ugly, so that three whole or partial impressions for each fin-
 553 ger were made on each of the substrates. Three different
 554 substrates were used: paper, plastic, and glass/porcelain.
 555 The items were separated based on substrate type and pro-
 556 cessed in one of three ways: *i*) chemical (ninhydrin) pro-
 557 cessing, *ii*) cyanoacrylate processing, *iii*) lift cards (pro-
 558 cessed with black fingerprint powder at the collection site).
 559 All fingerprints processed with cyanoacrylate were digitally
 560 photographed, while all ninhydrin and black powder finger-
 561 prints were scanned.
 562

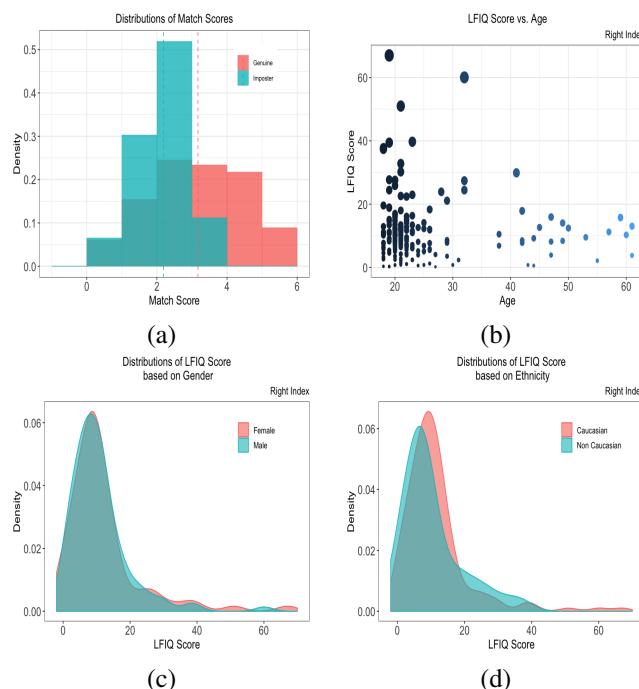
4.2. Latent-to-reference print comparison

563 The match scores in this paper were obtained using the
 564 end-to-end latent fingerprint search system recently pub-
 565 lished by Cao *et al.*. The algorithm does include automated
 566 ridge structure cropping, latent image pre-processing, fea-
 567 ture extraction, feature comparison, and outputs a candidate
 568 list. The model is robust to poor quality latents by gen-
 569 erated set of virtual minutiae to construct a texture template.
 570 This fully automated latent search system was evaluated on
 571 NIST SD27 (258 latents); MSP (1,200 latents), WVU (449
 572 latents) and N2N (10,000 latents) against a background set
 573 of 100K rolled prints, which includes the true rolled mates
 574 of the latents with rank-1 retrieval rates of 65.7%, 69.4%,
 575 65.5%, and 7.6% respectively[1].
 576

4.3. Results

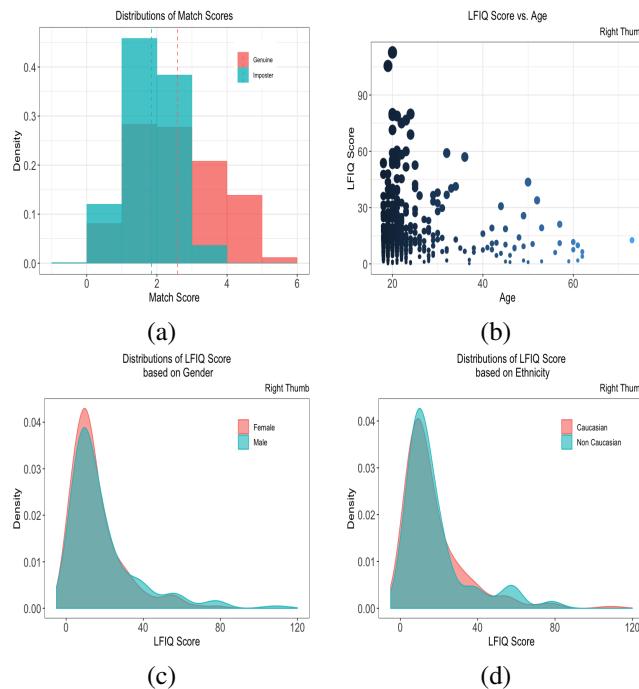
577 Fig. 3 (a) and Fig. 4 (a) show the distributions of match
 578 scores generated using the MSU identification system for
 579 both right index and right thumb.
 580

581 Fig. 3 also illustrates how demographics impact LFIQ
 582 measures for right index. Fig. 3 (b) points out instead that
 583 LFIQ differs across different age groups with younger pop-
 584 ulation achieving the highest scores ($LFIQ > 60$) while the
 585 subjects greater than 41 years old exhibit an upper bound of
 586 only $LFIQ = 20$. Fig. 3 (c) suggests that there is no trend
 587 with respect to gender which may be due to lack of tex-
 588 tural information capture by the LFIQ algorithm. From the
 589 literature, gender estimation from fingerprints does exploit
 590 textural differences between males and females. Similarly,
 591 Fig. 3 (d) highlights that higher LFIQ scores are more likely
 592 to be achieved by the Non-Caucasian population. Fig. 4 il-
 593



594
 595
 596
 597
 598
 599
 600
 601
 602
 603
 604
 605
 606
 607
 608
 609
 610
 611
 612
 613
 614
 615
 616
 617
 618
 619
 620
 621
 622
 623
 624
 625
 626
 627
 628
 629
 630
 631
 632
 633
 634
 635
 636
 637
 638
 639
 640
 641
 642
 643
 644
 645
 646

Figure 3. Distributions of the Right Index Finger: (a) Distributions of the Genuine and Impostor Match Scores, (b) Scatter Plot of the LFIQ Scores vs. Age, (c) Density Plot of the LFIQ Scores by Gender and (d) Density Plot of the LFIQ Scores by Ethnicity

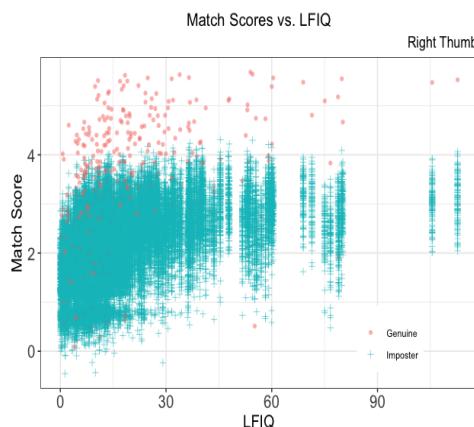


642
 643
 644
 645
 646

Figure 4. Distributions of the Right Thumb Finger: (a) Distributions of the Genuine and Impostor Match Scores, (b) Scatter Plot of the LFIQ Scores vs. Age, (c) Density Plot of the LFIQ Scores by Gender and (d) Density Plot of the LFIQ Scores by Ethnicity

648 illustrates how demographics impact LFIQ measures for right
 649 thumb which shows same behavior as right index.
 650

651 In Fig.5, the scatter plot of match scores versus LFIQ
 652 measures pertaining to right thumb shows a clear improvement
 653 of the separability between genuine and impostors for
 654 LFIQ>60. A similar trend has been observed for right index.
 655



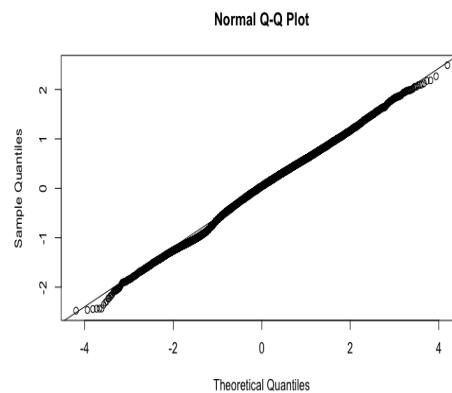
666 Figure 5. Scatter plot of Match Scores vs. LFIQ scores for Right
 667 Thumb.

668 We examine how the ROC curve varies conditioning on
 669 observed covariates by reporting the regression results for
 670 the two covariate-specific ROC models described in Section
 671 3. Various classifiers are capable of predicting the
 672 likelihood of a sample belonging to a class. A probabilistic
 673 classifier is implemented by setting a threshold which
 674 divides the entire data into different classes. The results
 675 shown in Table 1 was used to compute the sensitivities of
 676 the right index and the right thumb fingers. The QQ plots
 677 of the model were used to check the model diagnostics,
 678 Fig.6 shows the plot of model B for right index which
 679 indicates that the model fits the data very well. The similar
 680 results have been observed in right thumb data and also in
 681 Model A. The significant interaction terms indicate mean
 682 difference between genuine and imposter scores changes
 683 with covariates. For example, Model B for the right index
 684 fingers has significant interaction terms, LFIQ*label,
 685 Gender_Q*label and Ethnicity_Q*label. This indicates that
 686 the difference in mean scores between genuine and imposter
 687 groups change when LFIQ scores or age change. The differ-
 688 ence in mean matching scores between genuine and im-
 689 poster groups also change with ethnicity levels. For the
 690 right thumb, significant interaction terms indicate the differ-
 691 ence in mean matching scores between genuine and im-
 692 poster groups change as LFIQ scores and for different gen-
 693 der change. The mean difference is a main component in
 694 the ROC curve expression. Significant interactions imply
 695 when the covariates change, the ROC curve adjusting for
 696 the covariates tends to have significant changes.

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*
<i>Age_Q</i>	-	0.0015*	-	-0.0025*
<i>Age_R</i>	-	0.0021*	-	0.0020*
<i>Gender_Q</i>	-	0.0797*	-	0.0167*
<i>Gender_R</i>	-	0.1057*	-	0.0930*
<i>Ethnicity_Q</i>	-	0.0624*	-	-0.0827*
<i>Ethnicity_R</i>	-	0.0251*	-	0.0117
LFIQ*label	0.0257*	0.0266*	0.0216*	0.0210*
<i>Age_Q</i> *label	-	0.0091	-	0.0034
<i>Gender_Q</i> *label	-	0.4077*	-	0.2535*
<i>Ethnicity_Q</i> *label	-	0.3453*	-	-0.0742

* p-value<0.05

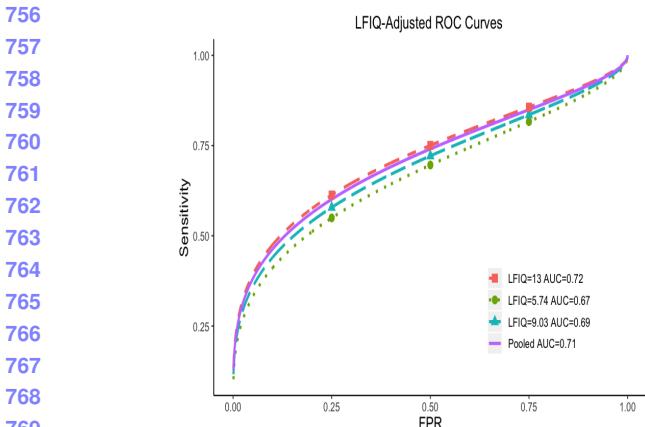
710 Table 1. Result of the designed covariate-specific ROC Models
 711



720 Figure 6. QQ plot of Model B for Right Index.
 721

722 Fig.7 shows the ROC curves based on **Model A** only con-
 723 sidering LFIQ score for right index, and the correspond-
 724 ing AUC values are also reported in the figure. Here, the
 725 25th, 50th, and 75th percentiles of LFIQ score were
 726 chosen to compute the curves. The purple curve is the pooled
 727 ROC curve which is used for comparison. We can observe
 728 that as LFIQ score increases, the model's identifying abil-
 729 ity increases. A similar behaviour has been found also for
 730 right thumb with a slightly lower performance compared to
 731 right index. Fig.8 and Fig.9 show the ROC curves based
 732 on **Model B** considering demographics and also LFIQ score
 733 for right index and right thumb, respectively. We can see
 734 that when the demographics are constant, the increase in
 735 LFIQ will increase the model's identifying ability, which is
 736 consistent with the univariate **Model A**.
 737

738 In Fig.8 we can also observe that for the same LFIQ
 739 score and demographics, as age increases the model has a
 740 much better identifying ability for right index, while we do
 741 not see much increase in identifying ability for right thumb
 742 in Fig.9 . Fig.8 also shows that when LFIQ scores and other
 743 demographics are adjusted, male subjects **performs** better in
 744



770 Figure 7. Covariate-specific ROC curves conditioned on LFIQ
771 scores for Right Index.

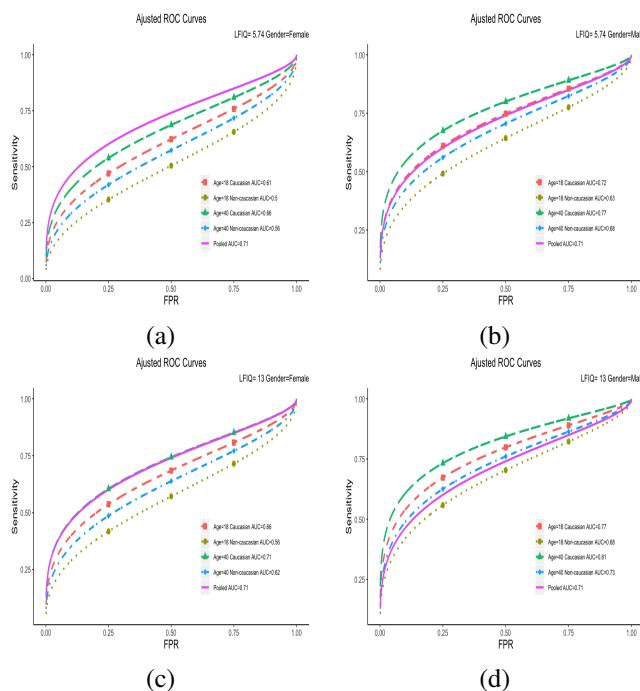


Figure 8. Covariate-specific ROC curves conditioned on LFIQ scores and demographics for right index.

identification than female; we can find the similar behaviour in Fig.9. In Fig.8 it is observed that when LFIQ score and other demographics are same, the caucasian group has a better identifying ability compared with the non-caucasian group for right index. While for right thumb, as we can see in Fig.9, the non-caucasian group performs better in identification.

In Fig.8 (c) and (d), we can see that although LFIQ is not able 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 previ-

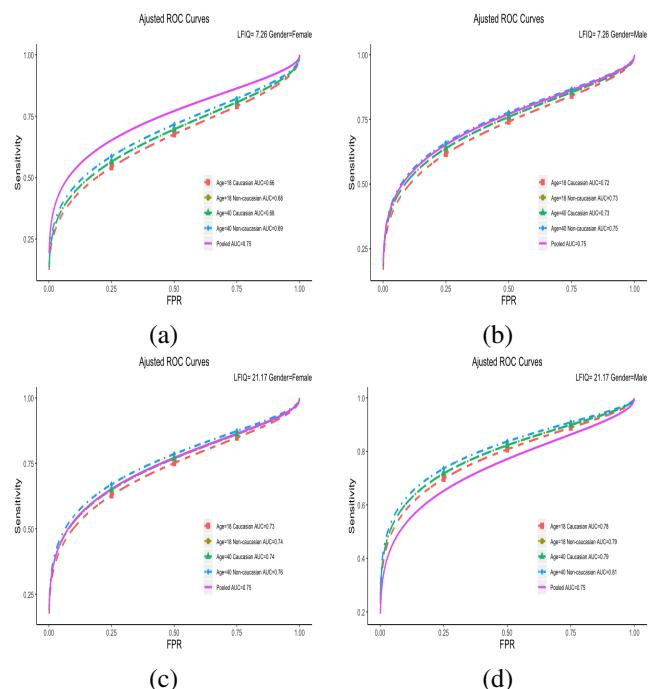


Figure 9. Covariate-specific ROC curves conditioned on LFIQ scores and demographics for right thumb.

ous 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.

5. Conclusions

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 some circumstances, such as when the subject is male. Our findings show that the proposed covariate-adjusted assessment scheme conditioned on image quality and demographics is more informative than the traditional ROC curve. It shows that demographics including age, gender and ethnicity tend to have impacts on the discriminating abilities of the matching algorithm.

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 to 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

864 CNN-based minutiae extractor to dry fingers by fine-tuning
865 it on the WVU database; *ii*) Access to the minutiae tem-
866 plates extracted by the MSU identification system and ap-
867 ply LFIQ to those; and *iii*) Enhance the LFIQ algorithm by
868 fusing multi-layer of minutiae maps from additional extrac-
869 tors for increased robustness in the presence of poor image
870 quality.

References

- 871
- 872
- 873 [1] K. Cao, D. Nguyen, C. Tymoszek, and A. Jain. End-to-End
874 Latent Fingerprint Search. *arXiv preprint arXiv:1812.10213*,
875 2018. 4, 6
- 876 [2] Nancy R Cook. Statistical evaluation of prognostic versus
877 diagnostic models: beyond the roc curve. *Clin Chem.*,
878 Jan;54(1):17–23, 2008. 2
- 879 [3] Hicklin et al. Latent fingerprint quality: A survey of examin-
880 ers. *Journal of Forensic Identification*, 61(4):385–419, 2011.
881 2
- 882 [4] P Gnanasivam and Dr S Muttan. Estimation of age through
883 fingerprints using wavelet transform and singular value de-
884 composition. *International Journal of Biometrics and Bioin-
885 formatics (IJBB)*, 6(2):58–67, 2012. 2
- 886 [5] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning
887 for image recognition. *IEEE CVPR*, pages 770–778, 2016. 4
- 888 [6] Deborah Hellman. Measuring Algorithmic Fairness. *Va. L.
889 Rev.*, 106:811, 2020. 1
- 890 [7] L. Hornak, W. LaRue, B. Cukic, A. Ross, K. Morris, J. Dawson,
891 S. Crihalmeanu, N. Kalka, and N. Kayal. FBI Biometric
892 Collection of People (BioCoP): Next Generation Identifi-
893 cation Phase 1 (2008 - 2009). *2008 Biometric Collection
894 Project 08-06-2008 to 12-31-2009 FINAL REPORT*, 2009. 5
- 895 [8] Crystal Huynh, Erica Brunelle, Lenka Halamkova, Juliana
896 Agudelo, and Jan Halamek. Forensic identification of gen-
897 der from fingerprints. *Analytical chemistry*, 87(22):11531–
898 11536, 2015. 2
- 899 [9] A. Jain, L. Hong, and R. Bolle. On-line fingerprint verifica-
900 tion. *IEEE Trans. PAMI*, 19(4):302–314, 1997. 3
- 901 [10] Anil Jain, Soweon Yoon, Kai Cao, and Eryun Liu. Latent
902 fingerprint quality(lfiq). Lecture PowerPoint - biomet-
903 rics.cse.msu.edu. 4
- 904 [11] Michael W. Kattan. Judging new markers by their ability to
905 improve predictive accuracy. *PubMed.*, May 7;95(9):634–5.,
906 2003. 2
- 907 [12] E. Marasco, Mengling He, L. Tang, and Sumanth Sriram.
908 Accounting for demographic differentials in forensic error
909 rate assessment of latent prints via covariate-specific roc
910 regression. *CIVP 2020*, CCIS 1376:338–350, 2020. 1
- 911 [13] E. Marasco, L. Lugini, and B. Cukic. Exploiting Quality and
912 Texture Features to Estimate Age and Gender from Finger-
913 prints. *SPIE Defense and Security*, 2014. 2
- 914 [14] Shira Mitchell, Eric Potash, Solon Barocas, Alexander
915 D’Amour, and Kristian Lum. Algorithmic fairness: Choices,
916 assumptions, and definitions. *Annual Review of Statistics
917 and Its Application*, 8:141–63, 2020. 1
- 918 [15] Dinh-Luan Nguyen, Kai Cao, and Anil K. Jain. Robust
919 minutiae extractor: Integrating deep networks and finger-
920 print domain knowledge. *CoRR*, abs/1712.09401, 2017. 4
- 921
- 922
- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942
- 943
- 944
- 945
- 946
- 947
- 948
- 949
- 950
- 951
- 952
- 953
- 954
- 955
- 956
- 957
- 958
- 959
- 960
- 961
- 962
- 963
- 964
- 965
- 966
- 967
- 968
- 969
- 970
- 971