

21 Generative Models and Probability Evaluation for Forensic Evidence

Sargur N. Srihari¹ and Chang Su²

Abstract: Generative approaches to pattern recognition and machine learning involve two parts: first describing the underlying probability distributions and then using such models to compute probabilities or make classificatory decisions. We consider generative models for forensic evidence where the goal is to describe the distributions using graphical models and to use such models to compute probabilistic metrics for measuring the degree of individuality of a forensic modality or of a piece of evidence. The metrics are defined as variations of the probability of random correspondence (PRC) when evidence consists of a set of measurements and correspondence is within a tolerance. Three metrics are defined, the first two of which concern the modality and the third concerns evidence: 1) PRC of two samples, 2) PRC among a random set of n samples (n PRC), and 3) PRC between a specific sample among n others (specific n PRC). Computation of these probabilities are described using graphical models which makes all the variables explicit. The metrics are evaluated for several cases—some of which are illustrative (birthdays and heights) and others concern fingerprints. For birthdays, which are discrete-valued and exact, assuming uniformly distributed birthdays, while n PRC rapidly approaches unity with higher n (which is the well-known birthday paradox), specific n PRC grows much less rapidly. For human heights, which are continuous-valued scalars, a quantization is needed and results are gender-specific: assuming Gaussian distributed heights, the PRC for males is higher than for females due to lower variance. Two forms of fingerprint representation are considered: ridge flow and minutiae. Gaussian mixtures are used to model location and orientation of minutiae. With parameters estimated from standard databases (using expectation maximization), fingerprint PRCs are determined for given numbers of available and matching minutiae (which correspond to quantization), for the case of 36 available minutiae where 24 of them match, the PRC is 4.0×10^{-34} . The methodology put forward should be value to establish the relative value of different types of forensic evidence in courtroom scenarios.

Keywords:

¹The State University of New York at Buffalo, Buffalo, New York 14260. E-mail: srihari@cedar.buffalo.edu.

²The State University of New York at Buffalo, Buffalo, New York 14260. E-mail: srihari@cedar.buffalo.edu.

21.1 Introduction

In forensics there are several modalities for identifying an object or individual from evidence. Examples of modality classes (and specific instances of modalities) are impression evidence (latent prints, handwriting, shoe-prints), trace evidence (paint, hair, fiber), biological evidence (DNA, blood type), etc. It is useful to establish the strength of a given modality or of a given piece of evidence within a given modality. Not least of the reasons for such analysis being court rulings that require a scientific basis for evidence presented [1].

The terms class characterization and individualization are commonly used in forensics. In addition terminology from the biometric domain, such as verification and identification are also present. Thus it is useful to first define these terms. Class-characterization is the narrowing down of the evidence into a sub-class within the forensic modality, e.g., ethnicity. Individualization is sometimes defined as the exclusion of all other sources for the given evidence. Verification is the determination of whether a given evidence is from a given source and is a binary decision. Identification is the determination of the best match of the evidence given a finite set of sources for that evidence. Finally individuality of a forensic modality or of a piece of forensic evidence is the degree of distinctiveness of that type of evidence in a population.

In many modalities, particularly those based on visual patterns, the set of measurements made from the evidence have an inherent variability even for the same object or individual. Since evidence in many forensic modalities can be characterized by quantitative measurements [2] a natural metric for individuality is a probability or a probability distribution. Such a metric can be evaluated, for a particular modality, measurement (or feature vector) or or a given sample of data, by using either one of the classical approaches of machine learning: discriminative and generative [3, 4]. In both approaches, representative samples of evidence are used to construct models during a training phase.

In the discriminative approach to measuring individuality, samples are directly used to construct either a two-class or a multi-class classifier [5]. One such approach is that based on determining a similarity (or kernel) function s whose value is high when the input and a template have the same origin and low when they are not. Such a method is used, for instance, in automated fingerprint identification systems that determine the degree of *match* between two fingerprints [6]. By thresholding the value of s into binary classes: *same* and *different*, an accuracy measure can be estimated from known training samples, e.g., average probability of error or risk. The accuracy estimate itself provides a measure of individuality of measurement \mathbf{x} . Since the nature of the testing set is crucial, data from cohort groups such as twins are often used to determine error rates [7, 8].

Generative models are statistical models that represent the distribution of \mathbf{x} . They are referred to as being generative in that given the distribution, samples can be generated from them. In these models, a distribution of \mathbf{x}

is learnt through a training data set. Samples can be generated from this distribution to determine the probability of random correspondence. The training set used is immaterial as long as it is representative of the entire population.

Generative models contrast with discriminative models, in that a generative model is a full probability model of all variables, whereas a discriminative model provides a model only of the target variable(s) conditioned on the observed variables. A generative model can be used, for example, to simulate (i.e., generate) values of any variable in the model, whereas a discriminative model allows only sampling of the target variables conditioned on the observed quantities. If the observed data are truly sampled from the generative model, then fitting the parameters of the generative model to maximize the data likelihood is a common method.

Both approaches have their advantage and limitations. In the generative approach, the problems with realistic modeling of all the parameters of the measurement may become unsurmountable. But a good model would be a method of gaining insights into the fundamental accuracy bounds of the measurement; more typically, these models may be relatively less effective in predicting performance. In the discriminative approach, definition of decision thresholds based on minimizing risk is an issue [9]. Discriminative methods lead to higher performing automated systems, but do not readily lead to fundamental understanding of underlying issues [10].

This paper only focuses on generative models and defines probabilistic metrics that are useful for a given modality or given evidence within a modality (Section 21.2). These metrics are applied to several modalities: birthdays (Section 21.3), human heights of males and females (Section 21.4), and fingerprints using ridge flow and minutiae (Section 21.5). The modalities discussed provide a continuum in terms of problem complexity: birthdays are represented as discrete numbers with uniform distribution, heights are continuous-valued Gaussian distributed scalars and fingerprints involve multivariate continuous features and scores provided by fingerprint matching algorithms.

21.2 Generative Models of Individuality

The goal is to specify a model for randomly generating observed data from which the relevant metrics can be computed. Given the model values for the probability of two randomly chosen samples having the same value within some tolerance can be computed. The process involves the following steps:

- Step 1** Consider a generative model (proposal) for measurement \mathbf{x} ;
- Step 2** Formulate a method for estimating parameters of the model;
- Step 3** Evaluate the parameters from a data set;
- Step 4** Use the model to evaluate relevant individuality metrics.

The first four steps are to determine the distribution of the data. The final

step involves computing the probability of match or correspondence within some tolerance between two samples. Tolerance specification depends on the modality and type of measurement, e.g., in the case of a continuous scalar it can be specified as $\pm\varepsilon$.

As measures of individuality, three probabilities can be defined: PRC (Probability of Random Correspondence), n PRC (Probability of Random Correspondence given n samples) and Specific n PRC (Probability of Random Correspondence of a specific \mathbf{x} among n samples). These definitions are further described below.

- 1) PRC: probability that two randomly chosen samples have the same measured value \mathbf{x} within specified tolerance ε .
- 2) n PRC: the probability that among a set of n samples, some pair have the same value \mathbf{x} , within specified tolerance, where $n \geq 2$. Since there are $\binom{n}{2}$ pairs involved this probability is higher than PRC. Note that when $n = 2$, $\text{PRC} = n\text{PRC}$.
- 3) Specific n PRC: the probability that in a set of n samples, a specific one with value \mathbf{x} coincides with another sample, within specified tolerance. Since we are trying to match a specific value \mathbf{x} , this probability depends on the probability of \mathbf{x} and is generally smaller than PRC. The exact relationship with respect to PRC depends on the distribution of \mathbf{x} .

We note that the first two measures, PRC and n PRC characterize the forensic modality as described by a set of measurements and furthermore, the second is a function of the first. The third measure *specific* n PRC characterizes specific evidence, e.g., a specimen. It is dependent on the specific value as well as the distribution of the measurement.

21.2.1 Graphical Models

Since probability distributions and their computation play a central role in this discussion, it is useful to represent them using graphical models [4]. Directed graphical models for PRC, n PRC and specific n PRC are shown in Fig. 21.1. Each node represents a random variable (or a group of random variables), and the links express conditional probabilistic relationships between them. The variables are: \mathbf{x} is a random variable corresponding to the feature vector, \mathbf{y} is another random variable with the same distribution, $\{\mathbf{y}_i\}$ where $i = [1, \dots, n]$ represents a set of n random variables using the plate representation and \mathbf{y}_s where $s \in [1, \dots, n]$ is a random variable among $\{\mathbf{y}_i\}$.

A binary-valued random variable z in Fig. 21.1 (a) indicates that if \mathbf{x} is the same as \mathbf{y} within a tolerance ε , in which case $z = 1$ and $z = 0$ otherwise. This leads to the following definition

$$p(z|\mathbf{x}, \mathbf{y}) = \begin{cases} p(z = 1|\mathbf{x}, \mathbf{y}) = 1 \text{ and } p(z = 0|\mathbf{x}, \mathbf{y}) = 0, & \text{if } \mathbf{x} = \mathbf{y}, \\ p(z = 1|\mathbf{x}, \mathbf{y}) = 0 \text{ and } p(z = 0|\mathbf{x}, \mathbf{y}) = 1, & \text{if } \mathbf{x} \neq \mathbf{y}. \end{cases} \quad (21.1)$$

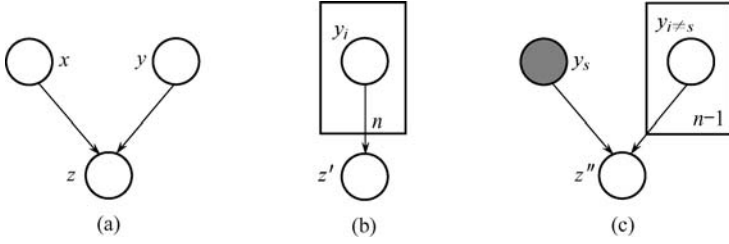


Fig. 21.1 Graphical models for computing coincidence probabilities. The models correspond to: (a) PRC, the probability of two samples having the same value, (b) n PRC, the probability of at least two samples among n having the same value, and (c) specific n PRC, the probability that a specific sample is among n . Here \mathbf{x} , \mathbf{y} , \mathbf{y}_s , and \mathbf{y}_i are feature vectors with identical distribution. We are interested in the distributions of z and z' which express probabilities of match/non-match. Note that \mathbf{y}_s is shaded indicating that its value is observed.

By marginalizing the joint distribution as

$$p(z) = \sum_{\mathbf{x}} \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}, z), \tag{21.2}$$

applying Bayes rule to write

$$p(z) = \sum_{\mathbf{x}} \sum_{\mathbf{y}} p(z|\mathbf{x}, \mathbf{y})p(\mathbf{x}, \mathbf{y}), \tag{21.3}$$

and noting the independence of \mathbf{x} and \mathbf{y}

$$p(z) = \sum_{\mathbf{x}} \sum_{\mathbf{y}} p(z|\mathbf{x}, \mathbf{y})p(\mathbf{x})p(\mathbf{y}), \tag{21.4}$$

we obtain the PRC as $\rho \equiv p(z = 1)$.

Next we consider the case of n PRC where there are n identically distributed random variables $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n]$. This is shown using the plate representation in Fig. 21.1 (b), where z' is the state indicating if in a set of n random variables at least one value \mathbf{y}_i is the same as another value \mathbf{y}_j in which case $z' = 1$ and $z' = 0$ otherwise. This leads to the following definition

$$p(z'|\mathbf{Y}) = \begin{cases} p(z' = 1|\mathbf{Y}) = 1 \text{ and } p(z' = 0|\mathbf{Y}) = 0, & \text{if } \exists \mathbf{y}_i, \mathbf{y}_j \in \mathbf{Y} \ni \mathbf{y}_i = \mathbf{y}_j, \\ p(z' = 1|\mathbf{Y}) = 0 \text{ and } p(z' = 0|\mathbf{Y}) = 1, & \text{if } \forall \mathbf{y}_i, \mathbf{y}_j \in \mathbf{Y} \ni \mathbf{y}_i \neq \mathbf{y}_j. \end{cases} \tag{21.5}$$

The marginal distribution of z' is obtained as

$$p(z') = \sum_{\mathbf{Y}} p(z'|\mathbf{Y})p(\mathbf{Y}). \tag{21.6}$$

The specific instance of $p(z' = 1)$, which is the n PRC, can be written as

$$p(z' = 1) = 1 - p(z' = 0) = 1 - \left[1 - \sum_{\mathbf{x}} \sum_{\mathbf{y}} p(z = 1 | \mathbf{x}, \mathbf{y}) p(\mathbf{x}) p(\mathbf{y}) \right]^{\frac{n(n-1)}{2}}. \tag{21.7}$$

Denoting the n PRC value as $\rho[n] \equiv p(z' = 1)$ and using Eq. (21.4), we have a relationship between PRC and n PRC as

$$\rho[n] = 1 - (1 - \rho)^{\frac{n(n-1)}{2}}. \tag{21.8}$$

Finally we consider the case where the value of $\mathbf{y}_s \in \mathbf{Y}$ is known as represented by the shaded node in Fig. 21.1 (c). Given a specific value \mathbf{y}_s , we define z'' as

$$p(z'' | \mathbf{y}_s, \mathbf{Y}) = \begin{cases} p(z'' = 1 | \mathbf{y}_s, \mathbf{Y}) = 1 \text{ and } p(z'' = 0 | \mathbf{y}_s, \mathbf{Y}) = 0, \\ \qquad \qquad \qquad \text{if } \exists \mathbf{y}_i \in \mathbf{Y} \ni \mathbf{y}_i = \mathbf{y}_s, \\ p(z'' = 1 | \mathbf{y}_s, \mathbf{Y}) = 0 \text{ and } p(z'' = 0 | \mathbf{y}_s, \mathbf{Y}) = 1, \\ \qquad \qquad \qquad \text{if } \forall \mathbf{y}_i \in \mathbf{Y} \ni \mathbf{y}_i \neq \mathbf{y}_s, \end{cases} \tag{21.9}$$

the specific n PRC is then given by the marginal probability

$$p(z'' = 1 | \mathbf{y}_s) = \sum_{\mathbf{Y}'} p(z'' = 1, \mathbf{Y}' | \mathbf{y}_s) = \sum_{\mathbf{Y}'} p(z'' = 1 | \mathbf{y}_s, \mathbf{Y}') p(\mathbf{Y}'), \tag{21.10}$$

where $\mathbf{Y}' = [\mathbf{y}_1, \dots, \mathbf{y}_{s-1}, \mathbf{y}_{s+1}, \dots, \mathbf{y}_n]$ and $p(\mathbf{Y}')$ is the joint probability of the $n - 1$ individuals. Since n PRC is dependent on the specific value $\mathbf{y} = \mathbf{y}_s$, it is useful to compute the expected value of n PRC, where the expectation is computed with respect to the distribution of \mathbf{y} .

The rest of this paper concerns the application of the developed methods to three cases: birthdays, human heights and fingerprints represented by minutiae and ridges. In the case of birthdays, which are discrete-valued, the tolerance is zero, i.e., requiring exact match. In the case of heights, which are continuous-valued scalars, a tolerance is specified in terms of the height differences that are considered to be the same. In the case of fingerprints, the measurement is a variable-length vector with each element represented as a triple, and tolerance is specified within the fingerprint matching algorithm.

21.3 Application to Birthdays

Birthdays are a standard case for the probability of coincidences [11]. As is typical, we disregard leap years and assume that the 365 possible birthdays are equally likely. We use a uniform density to model the birthday distribution and the probability that a person has any specific birthday is $1/365$. As this

case is for illustrative purpose only we will forego the goodness of fit test against a database of birthdays. Also we only consider the case of exact matches rather than almost matches, e.g., birthdays coincide within one day or k days.

The PRC, n PRC and specific n PRC are determined as follows.

21.3.1 PRC

PRC is the probability that any two persons have the same birthday. This is given by Eq. (21.4) and the definition in Eq. (21.1)

$$\begin{aligned} p_\varepsilon = p(z = 1) &= \sum_{\mathbf{x}=1}^{365} \sum_{\mathbf{y}=1}^{365} p(z = 1 | \mathbf{x}, \mathbf{y}) \left(\frac{1}{365} \right) \left(\frac{2\varepsilon + 1}{365} \right) \\ &= \sum_1^{365} \left(\frac{2\varepsilon + 1}{365^2} \right) = \frac{2\varepsilon + 1}{365}, \end{aligned} \quad (21.11)$$

where ε is the tolerance for birthdays. In the special case when $\varepsilon = 0$, the case of exactly matching birthdays, we have $p_\varepsilon = 1/365 = 0.0027$. Note that when $\varepsilon = 1$ the two birthdays can differ in three ways.

21.3.2 n PRC

The value of n PRC for birthdays is the probability that in a group of n persons, some pair of them have the same birthday. To compute it, it is easier to first calculate the probability $\bar{p}(n)$ that all n persons have different birthdays. The n PRC is complementary to $\bar{p}(n)$ and is given by

$$p(n) = 1 - \bar{p}(n) = 1 - (1 - p_\varepsilon)^{\frac{n(n-1)}{2}}, \quad (21.12)$$

where p_ε is the PRC value. Thus n PRC increases very rapidly with n , in fact exponentially with an exponent which is $O(n^2)$.

We can see how n PRC increases monotonically with n in Fig. 21.2. As an example consider the case of 40 randomly chosen individuals, or $n = 40$. We see the probability that two individuals have the same birthday in this group is 0.9. This is the well-known birthday paradox [11] which has been cited as a reason against claiming individualization in forensic evidence [12]. We shall see that there is a counter-paradox when we consider specific n PRCs.

21.3.3 Specific n PRC

To evaluate the probability of at least one individual sharing a birthday with a given individual among n individuals, the specific n PRCs are calculated.

The event of at least one of the n persons having the same birthday with the given individual is complementary to all the n people having the different birthday with the specific people. Assuming that $p(b)$ is the probability of an individual has the birthday b , the specific n PRC is $p(b, n)$ which is given by

$$p(b, n) = 1 - \bar{p}(b, n) = 1 - (1 - p(b))^{n-1}, \tag{21.13}$$

where $\bar{p}(b, n)$ is the probability that no other person has birthday within ϵ days of b .

The expected value of specific n PRC is given by

$$\begin{aligned} E[p(b, n)] &= \sum_{b=1}^{365} p(b)p(b, n) \\ &= \sum_{b=1}^{365} p(b)[1 - (1 - p(b))^{n-1}]. \end{aligned} \tag{21.14}$$

Assuming a uniform distribution of birthdays $p(b) = 1/365$ and allowing for birthdays to match when they are ϵ of each other, the expected value of specific n PRC is $1 - (1 - p_\epsilon)^{n-1}$. This increases exponentially as in Eq. (21) but with an exponent which is $O(n)$.

Fig. 21.2 shows the expected value of specific n PRC for some values of n . We see that given a specific birthday, the probability of a second individual having that same birthday among 40 individuals is about 0.1, which is significantly smaller than the n PRC value which is 0.9. Even with $n = 1000$, which

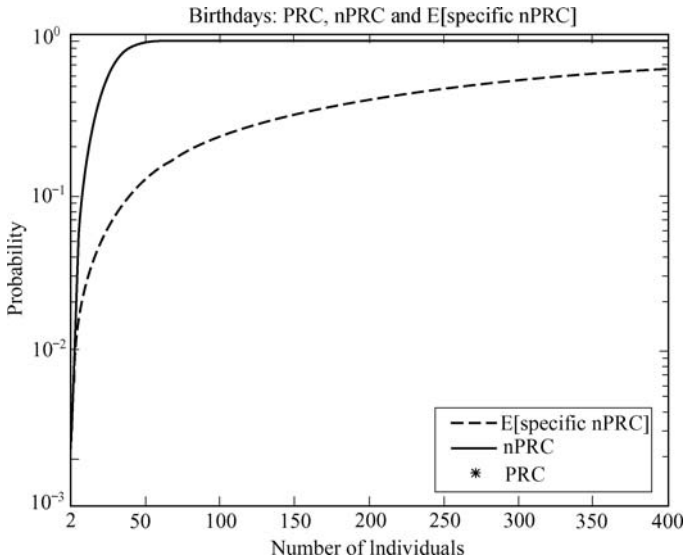


Fig. 21.2 Individuality of Birthdays: PRC, n PRC and Expectation of specific n PRC.

in nearly three times the number of days in a year, there is more than 6% probability that there is no other person with a specific birthday, say June 1st. In view of the birthday paradox, this is an interesting counter-paradox!

The expectation of specific n PRC is also plotted in Fig. 21.2 where it is seen that specific n PRC grows much slower than n PRC. In the case of birthdays the following inequality holds:

$$\text{PRC} \leq E[\text{Specific } n\text{PRC}] \leq n\text{PRC}$$

with equality holding when $n = 2$.

21.4 Application to Human Heights

We next consider a simple human measurement, height, which is continuous-valued and for which data is abundantly available. We use distributions of heights for males and females in the United States as produced by the CDC [13]. The Gaussian density is a reasonable model to fit the distribution. The Gaussian probability density functions for heights (in inches) for males and females aged 20 years and over is given in Fig. 21.3. The parameters for these distributions are as follows: mean $\mu_f = 63.8$, $\mu_m = 69.3$ and standard deviation $\sigma_f = 4.2$, $\sigma_m = 3.3$. Note that the standard deviation for females is much larger than for males. Again we forego the goodness of fit for the case of heights as it is also an illustrative case.

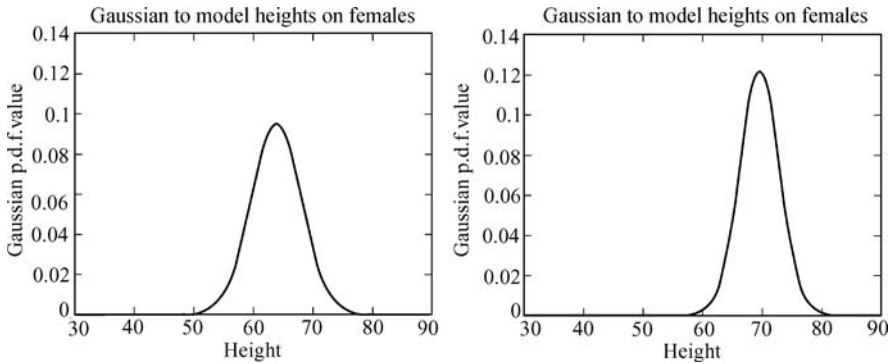


Fig. 21.3 Gaussian density models of human height: (a) female and (b) male, using US-CDC data. Females have higher variance.

We derive a value for the probability of random correspondence (PRC) for heights in terms of the known Gaussian distributions as follows. The probability of two individuals having the same height with some tolerance $\pm \epsilon$, can be written as

$$p_\epsilon = \int_{-\infty}^{\infty} \int_{y-\epsilon}^{y+\epsilon} P(x|\mu, \sigma)P(y|\mu, \sigma)dx dy, \tag{21.15}$$

where

$$P(h|\mu, \sigma) \sim \mathcal{N}(\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(h-\mu)^2}{2\sigma^2}}.$$

Eq. (21.15) can be numerically evaluated for values of μ and σ corresponding to the male and female height distributions. Using a tolerance of $\epsilon = 0.1$ inches we can calculate the PRCs. The PRC for female height with a mean height of 63.8 inches and standard deviation of 4.2 inches is $p_\epsilon = 0.0134$ and the PRC for male height with a mean height of 69.3 inches and standard deviation of 3.3 inches is $p_\epsilon = 0.0173$ (tolerance of 0.1 inch). The PRC increases with ϵ as shown in Fig. 20.4.

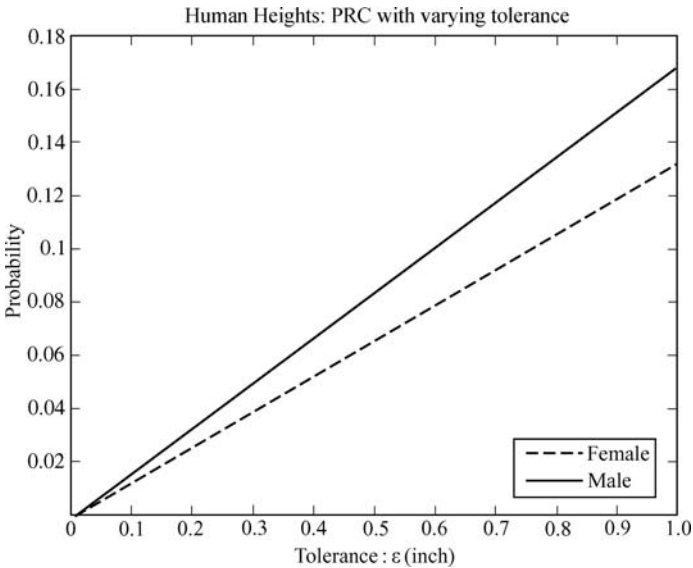


Fig. 21.4 Human Heights: PRC with varying tolerance ϵ .

The value of n PRC can be computed from the PRC value. This is the probability that in a randomly chosen group some pair will have the same height. The continuous variable can be handled by the implicit discreteness due to the tolerance ϵ . It is easier to first calculate the probability $\bar{p}(n)$ that all n heights are different from each other. The n PRC value is complementary to $\bar{p}(n)$ and given by

$$p(n) = 1 - \bar{p}(n) = 1 - (1 - p_\epsilon)^{\frac{n(n-1)}{2}}, \tag{21.16}$$

where p_ϵ is the PRC value. Values of n PRC with different n are shown in the second column of Tables 21.1 and 21.2.

Since the variance of heights of females is higher than for males, n PRC grows slower for females than for males. If $n = 40$ it is almost certain that there will be two persons with the same height within a tolerance of 0.1 inch.

As we shall see these probabilities will be much smaller in the case of specific heights.

Table 21.1 Female Height Probabilities: n PRC, Specific n PRC (57 inches) and Expectation of Specific n PRC

n	PRC $p(n)$	Specific n PRC ($h = 57$ inches) $p(h, n)$	Expectation of Specific n PRC $E[p(h, n)]$
2	0.0134	0.0051	0.0134
5	0.1262	0.0203	0.0525
10	0.4551	0.0452	0.1137
15	0.7574	0.0694	0.1703
20	0.9229	0.0930	0.2227
40	0.99995	0.1815	0.3965
80	1	0.3335	0.6214
160	1	0.5581	0.8233
320	1	0.8057	0.9316
800	1	0.9835	0.9784
1600	1	0.9997	0.9904

Table 21.2 Male Height Probabilities: n PRC, specific n PRC (68 inches) and expectation of specific n PRC

n	PRC $p(n)$	Specific n PRC ($h = 68$ inches) $p(h, n)$	Expectation of Specific n PRC $E[p(h, n)]$
2	0.0173	0.0224	0.0173
5	0.1601	0.0865	0.0664
10	0.5440	0.1842	0.1423
15	0.8400	0.2715	0.2111
20	0.9637	0.3494	0.2734
40	0.9999997	0.5862	0.4704
80	1	0.8326	0.6991
160	1	0.9726	0.8723
320	1	0.9993	0.9316
800	1	0.9999992	0.9838
1600	1	0.999999999995	0.9927

To evaluate the probability that at least one among n individuals has the same height as a given individual, the specific n PRCs are estimated. To compute it, it is easier to first calculate the probability $\bar{p}(h, n)$ that all n heights are different from the given height h . The event of at least one of the n persons having the same height with the given individual is complementary to all n heights being different. Assuming that $p(h)$ is the probability of an individual at a height of h , the specific PRC $p(h, n)$ is given by

$$p(h, n) = 1 - \bar{p}(h, n) = 1 - (1 - p(h))^{n-1}, \tag{21.17}$$

where $p(h) = \int_{h-\varepsilon}^{h+\varepsilon} P(l|\mu, \sigma)dl$, $P(l|\mu, \sigma)$ is the generative model for height.

The expected value of specific n PRC is given by

$$E[p(h, n)] = \int_0^\infty p(h, n)p(h)dh \tag{21.18}$$

The third and fourth columns of Tables 21.1 and 21.2 shows specific n PRCs and expectation of specific n PRC for some values of n . For example, among 40 randomly chosen females, the probability that one of them has height $h = 57$ inches (or 4ft 9in) is 0.18 and the expectation that one of them has same height as a given one is 0.39. Similarly, among 40 randomly chosen males, the probability that one of them has height $h = 68$ inches (or 6ft 8in) is 0.58 and the expectation that one of them has same height as a given one is 0.47. These probabilities are much smaller than the corresponding values for n PRC, which are close to certainty.

The values will depend specific heights h . PRC, n PRC, expected value of specific n PRC and specific n PRC of 5'9" and 5'2" on male height over varying number of individuals n are shown in Fig. 21.5. The value of n PRC increases rapidly with n . The rate of increase of specific n PRC depends on the specific height — with a more likely height (5'9") having a higher value than a less likely value (5'2"), i.e., an unusual height will have a low specific n PRC. The expected value has an intermediate value.

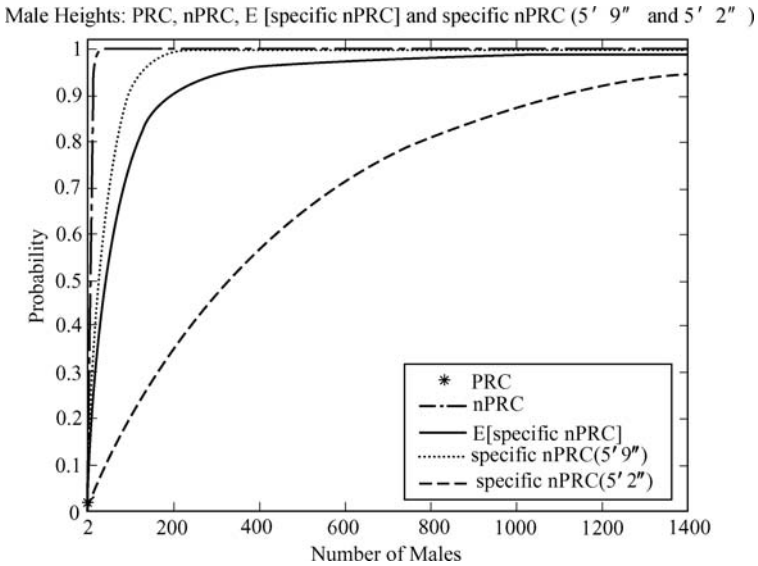


Fig. 21.5 Probabilities of male heights coinciding within 0.1 inch: PRC, n PRC, Expectation of specific n PRC and Specific n PRC of 5'9" and 5'2".

In the case of human heights the following inequality holds:

$$\text{PRC} \leq \text{E}[\text{Specific } n\text{PRC}] \leq n\text{PRC}$$
 with equality holding when $n = 2$.

21.5 Application to Fingerprints

We consider next the evaluation of the relevant probabilities in a commonly used forensic modality — that of fingerprints. Their use in human identification has been based on two premises, that, (i) they do not change with time and (ii) they are unique for each individual. While in the past, identification based on latent fingerprints had been accepted by courts, more recently their use has been questioned under the basis that the premises stated above have not been objectively tested and error rates not been scientifically established [14]. Though the first premise has been accepted, the second one on individuality is being challenged.

Fingerprint individuality studies date back to the late 1800s. More than twenty models have been proposed to establish the improbability of two random people (or fingers) have the same fingerprint [5]. These models can be classified into five different categories: grid-based [16, 17], ridge-based [17], fixed probability [18], relative measurement [19, 20] and generative [21–24]. All models try to quantify the uniqueness property, e.g., the probability of false correspondence. A match here does not necessarily mean an exact match but a match within given tolerance levels.

Features for representing fingerprints are classified into three types [25]. Level 1 features provide class-characterization of fingerprints based on ridge flow. They are divided into five primary classes: whorl, left loop, right loop, arch and tent (Fig. 21.6). Some of the primary classes have secondary classes resulting in a total of eight class types.

Level 2 features, which are more useful for identification, are also known as minutiae. Fingerprints such as those shown in Fig. 21.6 are first aligned. This is done manually where core points are identified and then the image is centered. The minutiae correspond to ridge endings and ridge bifurcations. Automatic fingerprint matching algorithms use minutiae as the salient features, e.g., [6], since they are stable and are reliably extracted. A minutia is represented by its location and direction; direction is determined by the ridge ending at the location (Fig. 21.7). The type of minutiae (either bifurcation or ending) is not distinguished since this information is not as reliable as the information on location and direction. Level 3 features, such as pores and scars are ancillary features.

Our goal is to model the distribution of fingerprints based on features. The distributions considered are those based on ridge flow types (Section 21.5.1) and minutiae (Section 21.5.2). Methods for estimating distribution parameters are discussed in Section 21.5.3. PRCs calculation algorithms are proposed in Section 21.5.4. Section 21.5.5 shows the experiments and results.

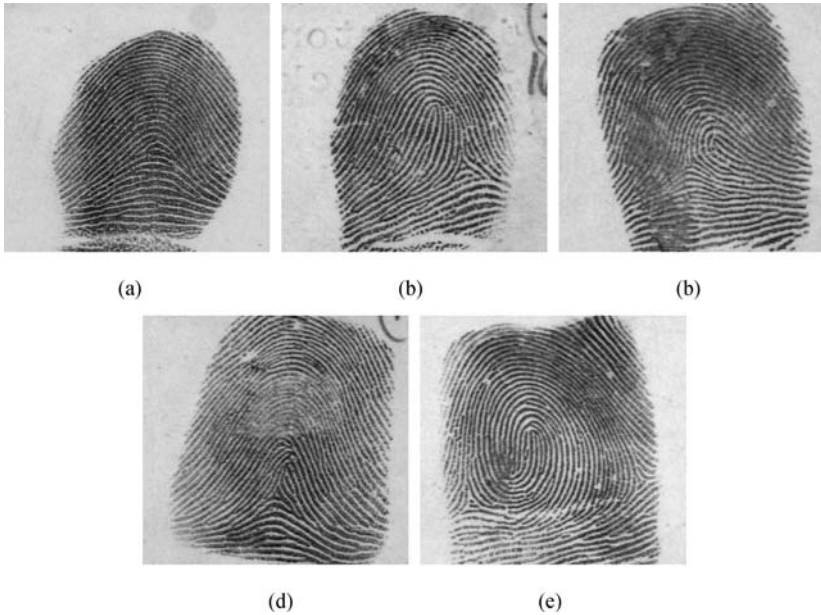


Fig. 21.6 Examples of five main types of ridge flow in fingerprints, referred to as Level 1 features: (a) arch, (b) left loop, (c) right loop. (d) tented arch, and (e) whorl. From NIST Special Database 4.

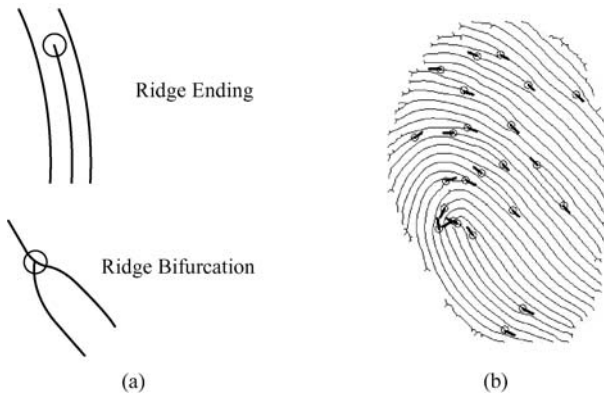


Fig. 21.7 Representation of fingerprints using minutiae: (a) locations of ridge endings and ridge bifurcations are indicated by circles, and (b) minutiae directions are indicated by line segments in a skeletonized fingerprint image.

21.5.1 Distribution of Ridge Flow Type

A simple distribution of the Level 1 ridge flow types is obtained by counting the relative frequency of each of the primary and secondary types in a

fingerprint database. In one such evaluation [8] loops account for 64% of the fingers, with the secondary types being: 30% left loops, 27% right loops and 7% double loops. Arches account for 18% of the primary types, with the secondary types being: plain arches (13%) and tented arches (5%). Whorls account for the remainder of the Level 1 types (19%).

Level 1 features are clearly broad class characteristics which are useful for exclusion of individual fingers but not by themselves useful for the tasks of verification, identification and individualization.

Assuming that fingerprints are distinguished by 6 secondary types, the PRC for ridge flow types is calculated by Eq. 21.3. Given above type frequencies, we have PRC value $p_\epsilon = 0.2233$. The n PRC and specific n PRC can be calculated by Eqs. 21.8 and 21.9. Tables 21.3 and 21.4 show the n PRC and Specific n PRC with different n .

Table 21.3 Ridge Flow Types: n PRC

n	2	3	4	5	6	7	8
n PRC	0.2233	0.5314	0.7805	0.9201	0.9774	1	1

Table 21.4 Ridge Flow Types: Specific n PRC belong to class LL = left loop, RL = right loop, DL = double loop, PA = plain arch, TA = tented arch and W = whorl.

n	Specific n PRC					
	LL	RL	DL	PA	TA	W
2	0.3000	0.2700	0.0700	0.1300	0.0500	0.1900
4	0.6570	0.6110	0.1956	0.3415	0.1426	0.4686
6	0.8319	0.7927	0.3043	0.5016	0.2262	0.6513
8	0.9176	0.8895	0.3983	0.6227	0.3017	0.7712
10	0.9596	0.9411	0.4796	0.7145	0.3698	0.8499
12	0.9802	0.9686	0.5499	0.7839	0.4312	0.9015
14	0.9903	0.9833	0.6107	0.8364	0.4867	0.9354
16	0.9953	0.9911	0.6633	0.8762	0.5367	0.9576
18	0.9977	0.9953	0.7088	0.9063	0.5819	0.9722
20	0.9989	0.9975	0.7481	0.9291	0.6226	0.9818

21.5.2 Distribution of Minutiae

Each minutia is represented as $\mathbf{x} = (s, \theta)$ where $s = (x_1, x_2)$ is its location and θ its direction. The distribution of minutiae location conditioned on ridge flow is shown in Fig. 21.8 where there were 400 fingerprints of each type. In the model we develop the combined distribution over all types is used (Figure 8 (f)).

Since minutia location has a multimodal distribution, a mixture of K Gaussians is a natural approach. For the data set considered a value of $K = 3$ provided a good fit, as validated by a goodness of fit test. Values of $K = 4, 5$ do not fit the data as well. A Gaussian mixture with $k = 3$ is shown in Fig.

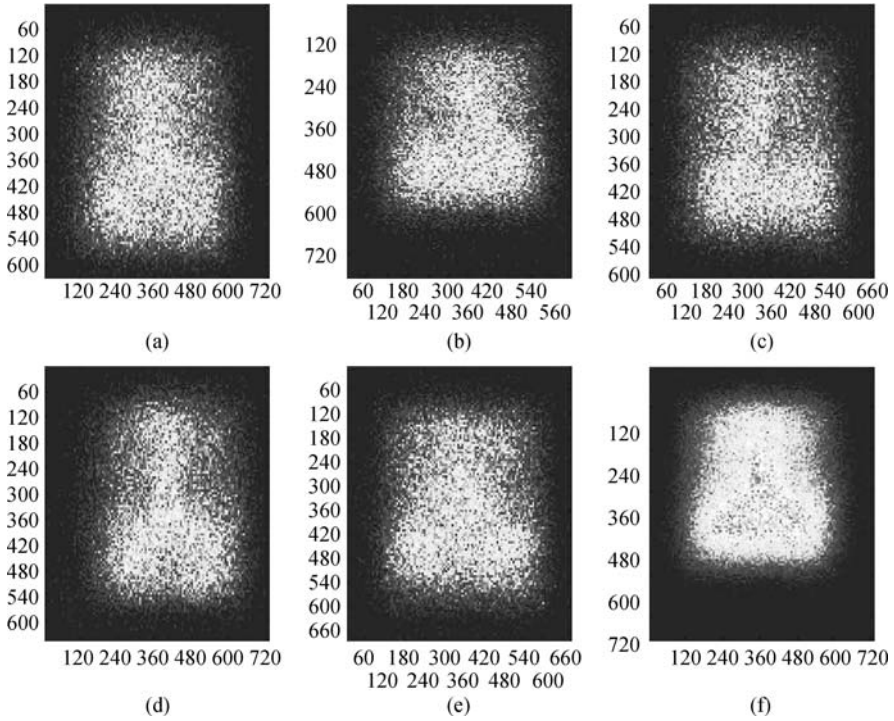


Fig. 21.8 Distribution of minutia location for different types of ridge flow: (a) arch, (b) left loop, (c) right loop, (d) tent, (e) whorl, and (e) all types combined.

21.9.

Since minutiae orientation is a periodic variable, it is modeled by a *circular normal* or von Mises distribution which itself is derived from the Gaussian [4, 26]. Such a model is better than mixtures of hyper-geometric and binomial distributions [21, 22].

Such a model for minutiae distributions involves a random variable z that represents the particular mixture component from which the minutia is drawn. In this model both minutiae location and orientation depend on the component they belong to. Minutiae location and orientation are conditionally independent given the component. This is represented by

$$p(\mathbf{x}|z) = p(s, \theta|z) = p(s|z)p(\theta|z), \tag{21.19}$$

whose graphical model is shown in Fig. 21.10 from which we have the joint distribution

$$p(\mathbf{x}, z) = p(z)p(\mathbf{x}|z). \tag{21.20}$$

Marginalizing over the components, we have the distribution of minutiae as

$$p(\mathbf{x}) = \sum_z p(z)p(\mathbf{x}|z). \tag{21.21}$$

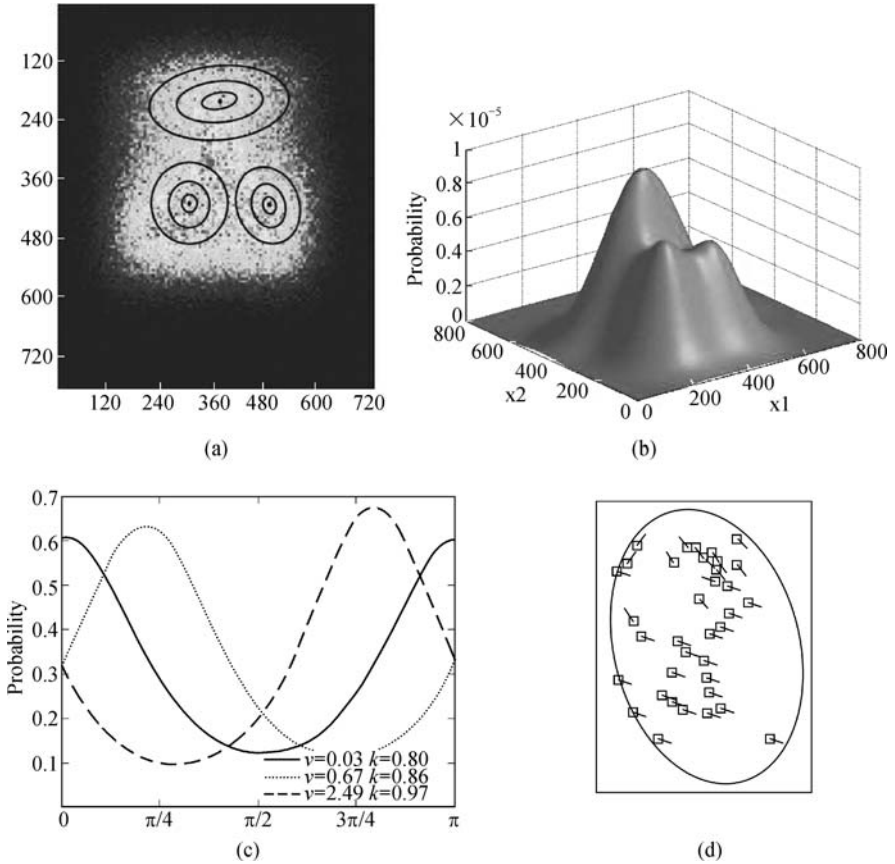
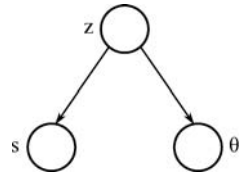


Fig. 21.9 Model for minutiae distribution using Gaussian mixture for location and von Mises for direction: (a) Gaussian mixture model for minutia location with three components, (b) three-dimensional plot of mixture model, (c) von Mises distributions of minutiae orientation for each of the three components, where the green curve corresponds to the upper cluster, blue the lower left cluster and red the lower right cluster, and (d) sample generated from model.

Fig. 21.10 Mixture model for distribution of minutiae $\mathbf{x} = (s, \theta)$. The joint distribution of minutia location and orientation is given by $p((s, \theta), z) = p(z)p(s|z)p(\theta|z)$.



Substituting Eq. (21.19) in Eq. (21.21) we have

$$p(\mathbf{x}) = \sum_z p(z)p(s|z)p(\theta|z). \tag{21.22}$$

Since minutiae location within each component is Gaussian and minutiae orientation within each component is von Mises we can write

$$p(\mathbf{x}|\Theta) = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(s|\mu_k, \Sigma_k) \cdot \mathcal{V}(\theta|\nu_k, \kappa_k), \quad (21.23)$$

where K is the number of mixture components, π_k are non-negative component weights that sum to one, $\mathcal{N}(s|\mu_k, \Sigma_k)$ is the bivariate Gaussian probability density function of minutia with mean μ_k and covariance matrix Σ_k , $\mathcal{V}(\theta|\nu_k, \kappa_k)$ is the von Mises probability density function of minutiae orientation with mean angle ν_k and precision (inverse variance) κ_k , and $\Theta = \{\pi_k, \mu_k, \Sigma_k, \nu_k, \kappa_k, \rho_k\}$ where $k = 1, 2, \dots, K$ is the set of all parameters of the k Gaussian and von Mises distributions.

Rather than using the standard form of the von Mises distribution for the range $[0, 2\pi]$, since minutiae orientations are represented as being in the range $[0, \pi)$, we use the alternate form [26] as follows:

$$\begin{aligned} \mathcal{V}(\theta|\nu_k, \kappa_k, \rho_k) &= \rho_k v(\theta) \cdot I\{0 \leq \theta < \pi\} \\ &\quad + (1 - \rho_k) v(\theta - \pi) \cdot I\{\pi \leq \theta < 2\pi\} \end{aligned} \quad (21.24)$$

where $I\{A\}$ is the indicator function of the condition A ,

$$v(\theta) \equiv v(\theta|\nu_k, \kappa_k) = \frac{2}{I_0(\kappa_k)} \exp[\kappa_k \cos 2(\theta - \nu_k)], \quad (21.25)$$

minutiae arising from the k th component have directions that are either θ or $\theta + \pi$ and the probabilities associated with these two occurrences are ρ_k and $1 - \rho_k$ respectively.

Since fingerprint ridges flow smoothly with very slow direction changes, direction of neighboring minutiae are strongly correlated, i.e., minutiae that are spatially close tend to have similar directions with each other. However, minutiae in different regions of a fingerprint tend to be associated with different region-specific minutiae directions thereby demonstrating independence [27, 28]. The model allows ridge orientations to be different at different regions (different regions can be denoted by different components) while it makes sure that nearby minutiae have similar orientations (as nearby minutiae will belong to the same component).

Since several minutiae are observed in a finger, a joint distribution model is needed. The minutiae are assumed to be independent of each other, with each minutiae, consisting of an (s, θ) pair, being distributed according to a mixture of component densities. This is discussed further when we consider parameter estimation in Section 21.5.3.

21.5.3 Parameter Estimation

We now develop an equivalent formulation of the mixture distribution given in Eq. (21.23) by involving an explicit latent variable. This will allow us to formulate the problem of parameter estimation in terms of the expectation maximization (EM) algorithm.

We define the joint distribution $p(\mathbf{x}, \mathbf{z})$ in terms of a marginal distribution $p(\mathbf{z})$ and a conditional distribution $p(\mathbf{x}|\mathbf{z})$, corresponding to the graphical model in Fig. 21.11 (a).

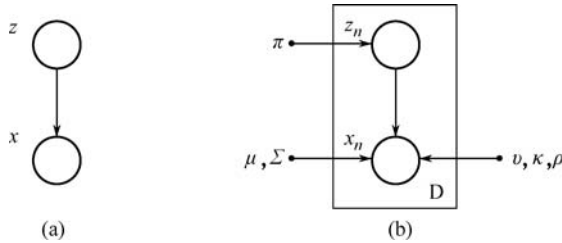


Fig. 21.11 Graphical models representing mixture for: (a) single minutia whose distribution is expressed as $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$, (b) set of D identically distributed minutiae with corresponding latent points (z_n) , where $(n = 1, \dots, D)$.

Given that the total number of minutiae observed in a finger is (D) , a joint distribution model is needed. The (D) minutiae are assumed to be independent of each other, with each minutiae, consisting of an $\mathbf{x}(s, \theta)$ pair, being distributed according to a mixture of component densities. This is shown in Fig. 21.11 (b).

The K -dimensional random variable \mathbf{z} has a 1-of- K representation in which a particular element z_k is equal to 1 and all other elements are equal to 0, we can write

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \tag{21.26}$$

Similarly the conditional distribution of \mathbf{x} given a particular value for \mathbf{z} is given by

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k) \cdot \mathcal{V}(\theta|\nu_k, \kappa_k, \rho_k) \tag{21.27}$$

which can also be written in the form

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)^{z_k} \cdot \mathcal{V}(\theta|\nu_k, \kappa_k, \rho_k)^{z_k}. \tag{21.28}$$

The joint distribution is given by $p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$, and the marginal distribution of \mathbf{x} is obtained by summing the joint distribution over all possible states

of \mathbf{z} to give

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(s|\mu_k, \Sigma_k) \cdot \mathcal{V}(\theta|\nu_k, \kappa_k, \rho_k), \quad (21.29)$$

where we have made use of Eqs. (21.27) and (21.29). Thus the marginal distribution of \mathbf{x} is a mixture of the form (24). If we have several observed minutiae $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_D$ then, because we have represented the marginal distribution in the form $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$, it follows that for every observed minutia \mathbf{x}_n , there is a corresponding latent variable \mathbf{z}_n .

We are now able to work with the joint distribution $p(\mathbf{x}, \mathbf{z})$ instead of the marginal distribution $p(\mathbf{x})$. To estimate the unknown parameters using the maximum likelihood approach, we use the EM algorithm. The number of components K for the mixture model was found after validation using k -means clustering.

E-Step Using γ_{dk} is to denote the responsibility of component k for minutiae \mathbf{x}_d , its value can be found using Bayes's theorem

$$\begin{aligned} \gamma_{dk} \equiv p(z_k = 1|\mathbf{x}_d) &= \frac{p(z_k = 1)p(\mathbf{x}_d|z_k = 1)}{\sum_{k=1}^K p(z_k = 1)p(\mathbf{x}_d|z_k = 1)} \\ &= \frac{\pi_k \mathcal{N}(s_d|\mu_k, \sigma_k) \mathcal{V}(\theta_d|\nu_k, \kappa_k, \rho_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(s_d|\mu_k, \sigma_k) \mathcal{V}(\theta_d|\nu_k, \kappa_k, \rho_k)}. \end{aligned} \quad (21.30)$$

M-Step The estimates of the Gaussian distribution parameters π_k, μ_{mk} , and Σ_{mk} at the $(n + 1)$ th iteration, are given by

$$\pi_k^{(n+1)} = \frac{1}{D} \sum_{d=1}^D \gamma_{dk}^{(n)}, \quad (21.31)$$

$$\mu_{mk}^{(n+1)} = \frac{\sum_{d=1}^D \gamma_{dk}^{(n)} s_m}{\sum_{d=1}^D \gamma_{dk}^{(n)}}, \quad (21.32)$$

$$\Sigma_{mk}^{(n+1)} = \frac{\sum_{d=1}^D \gamma_{dk}^{(n)} (s_m - \mu_{mk}^{(n+1)})(s_m - \mu_{mk}^{(n+1)})^T}{\sum_{d=1}^D \gamma_{dk}^{(n)}}. \quad (21.33)$$

The parameters for orientation distributions are obtained using expectation maximization for the von Mises distribution [29]. The estimates of ν_{mk}

and κ_{mk} at the $(n + 1)$ th iteration are given by

$$\nu_{mk}^{(n+1)} = \frac{1}{2} \tan^{-1} \left(\frac{\sum_{d=1}^D \gamma_{dk}^{(n)} \sin 2\psi_d}{\sum_{d=1}^D \gamma_{dk}^{(n)} \cos 2\psi_d} \right), \tag{21.34}$$

$$\frac{I'_0(\kappa_{mk}^{(n+1)})}{I_0(\kappa_{mk}^{(n+1)})} = \frac{\sum_{d=1}^D r_{dk}^{(n)} \cos 2(\psi_d - \nu_k^{(n+1)})}{\sum_{d=1}^D r_{dk}^{(n)}}. \tag{21.35}$$

The solution for Eq. (21.35), which involves Bessel functions, obtained using an iterative method gives the estimate for κ_{mk} . The estimate of ρ_{mk} is then obtained as

$$\rho_{mk}^{(n+1)} = \frac{\sum_{d=1}^D I\{c_d^{(n+1)} = k, \theta_d \in [0, \pi)\}}{\sum_{d=1}^D I\{c_d^{(n+1)} = k\}}, \tag{21.36}$$

where $c_d^{(n+1)} = \operatorname{argmax}_k \gamma_{dk}^{(n+1)}$ is the component label for the observation d at the $(n + 1)$ th iteration, ψ_j is the orientation of the minutiae m_j .

21.5.4 Evaluation of PRCs

To compute the PRCs, we first define correspondence, or match, between two minutiae. Let $\mathbf{x}_a = (s_a, \theta_a)$ and $\mathbf{x}_b = (s_b, \theta_b)$ be a pair of minutiae. The minutiae are said to correspond if for tolerance $\varepsilon = [\varepsilon_s, \varepsilon_\theta]$,

$$|s_a - s_b| \leq \varepsilon_s \text{ and } |\theta_a - \theta_b| \leq \varepsilon_\theta \tag{21.37}$$

where $|s_a - s_b|$, the Euclidean distance between the minutiae location $s_a = (x_{a1}, x_{a2})$ and $s_b = (x_{b1}, x_{b2})$, is given by

$$|s_a - s_b| = \sqrt{(x_{a1} - x_{b1})^2 + (x_{a2} - x_{b2})^2}. \tag{21.38}$$

Then, the probability that a random minutia \mathbf{x}_a would match a random minutia \mathbf{x}_b is given by

$$p_\varepsilon(\mathbf{x}) = p(|\mathbf{x}_a - \mathbf{x}_b| \leq \varepsilon | \Theta) = \int \int_{\mathbf{x}_a | |\mathbf{x}_a - \mathbf{x}_b| \leq \varepsilon} p(\mathbf{x}_a | \Theta) p(\mathbf{x}_b | \Theta) d\mathbf{x}_a d\mathbf{x}_b, \tag{21.39}$$

where Θ is the set of parameters describing the distribution of the minutiae location and direction.

Finally, the PRC, or the probability of matching at least \hat{m} pairs of minutiae within ϵ between two randomly chosen fingerprint f_1 and f_2 is calculated as

$$p_\epsilon(\hat{m}, m_1, m_2) = \binom{m_1}{\hat{m}} \binom{m_2}{\hat{m}} \hat{m}! \cdot p_\epsilon(\mathbf{x})^{\hat{m}} (1 - p_\epsilon(\mathbf{x}))^{(m_1 - \hat{m}) \cdot (m_2 - \hat{m})}, \quad (21.40)$$

where m_1 and m_2 are numbers of minutiae in fingerprints f_1 and f_2 , $p_\epsilon(\mathbf{x})^{\hat{m}}$ is the probability of matching \hat{m} specific pairs of minutiae between f_1 and f_2 , $(1 - p_\epsilon(\mathbf{x}))^{(m_1 - \hat{m}) \cdot (m_2 - \hat{m})}$ is the probability that none of minutiae pair would match between the rest of minutiae in f_1 and f_2 and $\binom{m_1}{\hat{m}} \binom{m_2}{\hat{m}} \hat{m}!$ is the number of different match sets that can be paired up.

Given n fingerprints and assuming that the number of minutiae in a fingerprint m can be modeled by the distribution $p(m)$, the general PRCs $p(n)$ is given by

$$p(n) = 1 - \bar{p}(n) = 1 - (1 - p_\epsilon)^{\frac{n(n-1)}{2}}, \quad (21.41)$$

where p_ϵ is the probability of matching two random fingerprint from n fingerprints. If we set the tolerance in terms of number of matching minutiae to \hat{m} , p_ϵ is calculated by

$$p_\epsilon = \sum_{m'_1 \in M_1} \sum_{m'_2 \in M_2} p(m'_1) p(m'_2) p_\epsilon(\hat{m}, m'_1, m'_2), \quad (21.42)$$

where M_1 and M_2 contain all possible numbers of minutiae in one fingerprint among n fingerprints, and $p_\epsilon(\hat{m}, m'_1, m'_2)$ can be calculated by Eq. 21.40.

Given a specific fingerprint f , the specific n PRCs can be computed by

$$p(f, n) = 1 - (1 - p(f))^{n-1}, \quad (21.43)$$

where $p(f)$ is the probability that \hat{m} pairs of minutiae are matched between the given fingerprint f and a randomly chosen fingerprint from n fingerprints.

$$p(f) = \sum_{m' \in M} p(m') \binom{m'}{\hat{m}} p(f_i), \quad (21.44)$$

where M contains all possible numbers of minutiae in one fingerprint among n fingerprints, $p(m')$ is the probability of a fingerprint having m' minutiae in n fingerprints, m_f is the number of minutiae in the given fingerprint f , minutiae set $f_i = (\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{i\hat{m}})$ is the subset of the minutiae set of given fingerprint and $p(f_i)$ is the joint probability of minutiae set f_i based on learned generative model.

21.5.5 Evaluation with Fingerprint Databases

Parameters of the fingerprint generative model introduced in Sections 21.5.2 was evaluated using the NIST fingerprint database. The NIST fingerprint database, NIST Special Database 4, contains 8-bit gray scale images of randomly selected fingerprints. Each print is 512×512 pixels with 32 rows of white space at the bottom of the print. The entire database contains fingerprints taken from 2000 different fingers with 2 impression of the same finger. Thus, there are a total of 2 000 fingerprints using which the model has been developed. The fingerprints are classified into one of five categories (left loop, whirl, right loop, tented arch, and arch) with an equal number of prints from each class (400). The number of components K for the mixture model was found after validation using k -means clustering.

Values of PRC p_ϵ are calculated using the formula introduced in Section 21.5.4. The tolerance is set at $\epsilon_s = 10$ pixels and $\epsilon_\theta = \pi/8$. For comparison, the empirical PRC $\hat{p}_\epsilon(\mathbf{x})$ was calculated with the same tolerance. To compute $\hat{p}_\epsilon(\mathbf{x})$, the empirical probabilities of matching a minutiae pair between imposter fingerprints are calculated first by

$$\hat{p}_\epsilon(\mathbf{x}) = \frac{1}{I} \sum_{i=1}^I \frac{\hat{m}_i}{m_i \times m'_i} \quad (21.45)$$

where I is the number of the imposter fingerprints pairs, \hat{m}_i is the number of matched minutiae pairs and m_i and m'_i are the numbers of minutiae pairs in each of the two fingerprints. Then, the empirical PRC \hat{p}_ϵ can be calculated by Eq. 21.41.

Both the theoretical and empirical PRCs are given in Table 21.5. The PRCs are calculated through varying number of minutiae in two randomly chosen fingerprint f_1 and f_2 and the number of matches between them. We can see that more minutiae the template and input fingerprint have, higher the PRC is. In experiments conducted on the NIST 4 dataset, the PRC values obtain here are smaller than the results in [22, 24]. The differences mainly result from use of different ways to evaluate PRC from generative models which is described in Section 21.5.4. Different matching tolerance, which $p_\epsilon(\mathbf{x})$ depends on may cause the differences also. Note that the theoretical PRC based on our model are close and have the same trend to empirical PRC. The consistency between the theoretical probabilities and empirical probabilities shows the validation of our generative model. The PRCs for the different m_1 and m_2 with 6, 26, 56, and 76 matches are shown in Fig. 21.12. It is obvious to note that, when \hat{m} decreases or m_1 and m_2 increase, the probability of matching two random fingerprints is more.

Based on the PRC value, n PRC can be computed. Table 21.6 shows the n PRCs through varying numbers of matched minutiae pairs \hat{m} in different numbers of fingerprints.

Table 21.5 PRC for different fingerprint matches with varying m_1 (number of minutiae in fingerprint f_1), m_2 (number of minutiae in fingerprint f_2) and \hat{m} (number of matched minutiae). p_ϵ is the theoretical PRC for the general population and \hat{p}_ϵ is the empirical PRC for NIST 4 dataset.

m_1	m_2	\hat{m}	p_ϵ	\hat{p}_ϵ
16	16	6	6.3×10^{-9}	1.0×10^{-8}
		12	3.4×10^{-23}	1.0×10^{-22}
		16	1.2×10^{-37}	5.4×10^{-37}
36	36	6	2.1×10^{-4}	3.4×10^{-4}
		12	1.1×10^{-11}	3.1×10^{-11}
		24	4.0×10^{-34}	3.6×10^{-33}
		36	3.6×10^{-72}	9.9×10^{-71}
56	56	6	1.8×10^{-2}	2.6×10^{-2}
		12	7.8×10^{-7}	2.1×10^{-6}
		24	2.6×10^{-21}	2.2×10^{-20}
		36	1.7×10^{-42}	4.5×10^{-41}
		48	5.0×10^{-72}	4.1×10^{-70}
76	76	56	1.1×10^{-101}	2.0×10^{-99}
		6	1.4×10^{-1}	1.8×10^{-1}
		12	5.0×10^{-4}	1.1×10^{-3}
		24	5.6×10^{-14}	4.2×10^{-13}
		36	4.4×10^{-29}	1.1×10^{-27}
		48	3.6×10^{-49}	2.8×10^{-47}
		60	3.6×10^{-75}	8.8×10^{-73}
		76	4.9×10^{-128}	5.4×10^{-125}

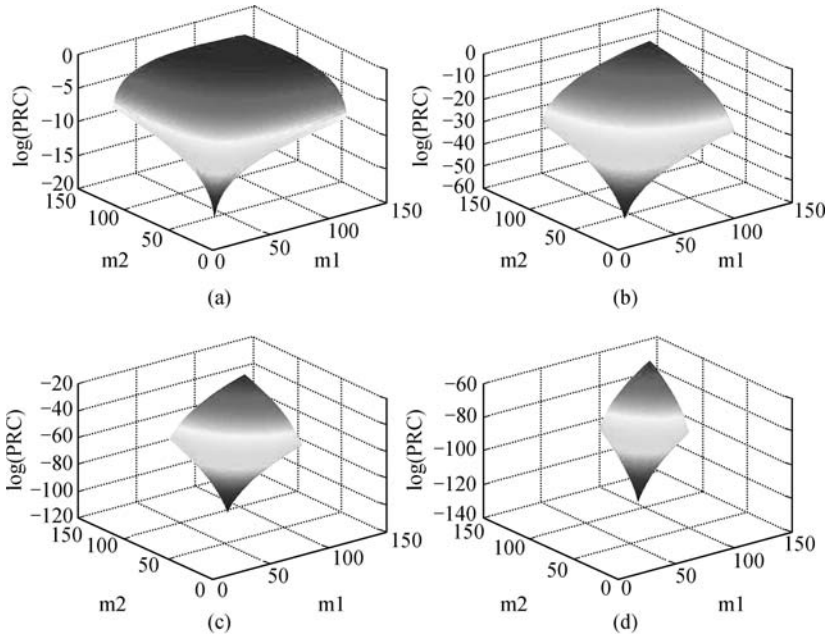


Fig. 21.12 PRCs with different number of the matched minutiae for (a) $\hat{m} = 6$, (b) $\hat{m} = 26$, (c) $\hat{m} = 56$, and (d) $\hat{m} = 76$.

The specific n PRCs are also computed by Eq. (21.44) and given by Table 21.7. Here three fingerprints are chosen as query prints and they are shown in Fig. 21.13. The first one is a full print in good quality, the second one is a full print in low quality and the third one is a partial print. The specific n PRCs are calculated through varying number of minutiae in each template fingerprint (m) and the number of matches (\hat{m}), assuming that the number of fingerprints in template database (n) is 100 000. The numbers of minutiae in 3 given query fingerprint m_f are 71, 42, and 18. In 100 000 randomly chosen fingerprints there is only $5.137\ 1 \times 10^{-39}$ probability that one of them have 12 matched minutiae with the fingerprint F_1 .



Fig. 21.13 Three specific fingerprints (from the same finger) used to calculate probabilities: (a) good quality full print F_1 , (b) low quality full print F_2 , and (c) partial print F_3 .

Table 21.6 Fingerprint Probabilities: n PRCs with varying n and \hat{m}

Number of fingerprints n	Number of matches \hat{m}	$p(n)$
1000	16	1
	36	$4.173\ 1 \times 10^{-5}$
	56	$1.721\ 9 \times 10^{-18}$
	76	$4.854\ 6 \times 10^{-38}$
	96	$4.236\ 9 \times 10^{-64}$
100 000	16	1
	36	$3.414\ 6 \times 10^{-1}$
	56	$1.723\ 6 \times 10^{-14}$
	76	$4.859\ 4 \times 10^{-34}$
	96	$4.241\ 1 \times 10^{-60}$
10 000 000	16	1
	36	1
	56	$1.723\ 6 \times 10^{-10}$
	76	$4.859\ 4 \times 10^{-30}$
	96	$4.241\ 1 \times 10^{-56}$

Table 21.7 Fingerprint Probabilities: Specific n PRCs for fingerprints in Fig. 21.13 with $n = 100,000$ fingerprints

Fingerprint f	Number of minutiae m_f	Number of matches \hat{m}	$p(n, f)$
F_1	71	55	8.5815×10^{-276}
		45	1.8829×10^{-211}
		35	2.0781×10^{-153}
		25	3.0360×10^{-102}
		12	5.1371×10^{-39}
		6	6.5876×10^{-15}
F_2	42	40	5.1783×10^{-212}
		35	2.4791×10^{-183}
		25	8.4136×10^{-115}
		12	1.6059×10^{-43}
		6	4.9787×10^{-16}
F_3	18	18	4.0256×10^{-100}
		12	3.4509×10^{-58}
		6	1.4068×10^{-22}

21.6 Summary

While forensic evidence of many modalities have long been used in the judicial system, e.g., impression evidence (latent prints, handwriting, shoe-prints), trace evidence (paint flakes, pollen, fibers, glass, hair), etc., characterizing their accuracy in identification is still needed to provide a scientific basis. The degree of individuality of a forensic modality can be established quantitatively by using either discriminative or generative approaches. In the former an error rate is determined after a suitable classifier is constructed. In the latter a probability distribution is determined from which different types of probabilities of random correspondence (PRC) are evaluated.

Generative models of individuality attempt to model the distribution of features and then use the models to determine the probability of random correspondence. We have proposed such models of individuality for birthdays, heights and fingerprints. Individuality is evaluated in terms of three probability measures: probability of random correspondence (PRC) between two individuals, general probability of random correspondence (n PRC) between two individuals among a group of n individuals and specific probability of random correspondence (specific n PRC) which is the probability of matching a given individual among n individuals.

The generative model for birthdays is based on assuming a uniform distribution where each date is assumed to be equally likely. The PRC for birthdays is 0.002 7. For 40 individuals, the n PRC is 0.891 2 and specific n PRC (any date) is 0.101 5. The n PRC value is the well-known birthday paradox where with $n = 23$ the value reaches a high value of 0.5.

The heights of female and male individuals are modeled by Gaussian distribution, where the parameters are learned from the health statistic data. The PRCs for female and male heights are 0.013 4 and 0.017 3. In a group of 40 people, the n PRCs for females and males are 0.999 95 and 0.999 999 97 and specific n PRCs for females (57 inches) and male(68 inches) are 0.181 5 and 0.586 2.

Models for fingerprint individuality have been proposed for ridge flows and minutiae. A mixture distribution was proposed to model minutiae information. The new generative model is compared by implementation and experiments with the empirical results on the NIST 4 dataset. The PRC obtained for a fingerprint template and input with 36 minutiae each with 12 matching minutiae is 1.1×10^{-11} . This probability is very close to the empirical result which is 3.1×10^{-11} . n PRC and specific n PRC of fingerprints are computed also. Considering the case of 100 000 fingerprints, the n PRC with minutiae information where 56 minutiae pairs are matched is $1.723 6 \times 10^{-14}$. Given a specific fingerprint with 71 minutiae, the specific n PRC with minutiae information where 55 out of 71 minutiae are matched is $8.581 5 \times 10^{-276}$. The proposed generative model offers a reasonable and accurate fingerprint representation. The results provide a much stronger argument for the individuality of fingerprints in forensics than previous generative models.

Acknowledgements

This work was supported by a grant from the Department of Justice, Office of Justice Programs Grant NIJ 2005-DD-BX-K012. The opinions expressed are those of the authors and not of the DOJ.

References

- [1] Aitken C, Taroni F (2004) *Statistics and the Evaluation of Evidence for Forensic Scientists*. Wiley, New York
- [2] Ashbaugh D R (1999) *Quantitative-Qualitative Friction Ridge Analysis: An Introduction to Basic and Advanced Ridgeology*. CRC Press, Boca Raton, FL
- [3] Banerjee A, Dhillon I, Ghosh J et al (2003) Generative Model-based Clustering of Directional Data. In *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*
- [4] Bishop C (2006) *Pattern Recognition and Machine Learning*. Springer, New York
- [5] Bolle R M (2003) *Guide to Biometrics*. LLC, 2003. Springer, New York
- [6] Champod C, Margot P (1996) Computer Assisted Analysis of Minutiae Occurrences on Fingerprints. In *Almog J and Spinger E, Editors, Proc. International Symposium on Fingerprint Detection and Identification*, page 305. Israel National Police, Jerusalem

- [7] Chen J, Moon Y s (2007) A Minutiae-based Fingerprint Individuality Model. *Computer Vision and Pattern Recognition*, 2007. CVPR '07
- [8] Chen Y, Jain A K (2009) Beyond Minutiae: A Fingerprint Individuality Model with Pattern, Ridge and pore Features. In Tistarelli M and Nixon M S, editors, *ICB*, vol 5558 of *Lecture Notes in Computer Science*, pages 523–533. Springer
- [9] Dass S, Zhu Y, Jain A K (2005) Statistical Models for Assessing the Individuality of Fingerprints. *Fourth IEEE Workshop on Automatic Identification Advanced Technologies*, pages 3–9
- [10] Diaconia P, Mosteller F (1989) *Methods for Studying Coincidences*. *J. Am. Stat. Association*, 84: 853–857
- [11] Duda R O, Hart P E, Stork D G (2001) *Pattern Classification*, Second Edition. New York: John Wiley and Sons
- [12] Galton F (1892) *Finger Prints*. McMillan, London
- [13] Henry E (1900) *Classification and Uses of FingerPrints*. Routledge & Sons London
- [14] Hsu R V, Martin B (2008) *An Analysis of Minutiae Neighborhood Probabilities*. *Biometrics: Theory, Applications and Systems*, 2008.
- [15] Jebarra T (2004) *Machine Learning: Discriminative and Generative*. Kluwer Academic
- [16] Mardia K V, Jupp P E (2000) *Directional Statistics*. Wiley
- [17] McDowell M A, Fryar C D, R H et al (2005) *Anthropometric Reference Data for Children and Adults: United S Tates Population, 1999 to 2002*. *Advance data from Vital and Health Statistics*; Number 361
- [18] Pankanti S, Prabhakar S, Jain A K (2002) On the Individuality of Fingerprints. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(8)
- [19] Roxburgh T (1933) Galton's work on the Evidential Value of Fingerprints. *Indian Journal of Statistics*, 1: 62
- [20] Saks F, Koehler J J (2008) The Individualization Fallacy in Forensic Science Evidence. *Vanderbilt Law Review*, 61: 199–219
- [21] Srihari S N, Cha S, Arora H (2002) Individuality of Handwriting. *Journal of Forensic Sciences*, 47(4): 856–872
- [22] Srihari S N, Cha S, Arora H (2008) Discriminability of Fingerprints of twins. *Journal of Forensic Identification*, 58(1): 109–127
- [23] Srihari S N, Huang C, Srinivasan H (2008) On the Discriminability of the Handwriting of twins. *Journal of Forensic Sciences*, 53(2): 430–446
- [24] Stoney D A (2001) Measurement of Fingerprint Individuality. In H. Lee and R. Gaensslen, editors, *Advances in Fingerprint Technology*. CRC Press
- [25] Su C, Srihari S N (2008) Generative Models for Fingerprint Individuality Using Ridge Models. In *Proceedings of International Conference on Pattern Recognition*. IEEE Computer Society Press
- [26] Trauring M (1963) Automatic Comparison of Finger-ridge Patterns. *Nature*, page 197
- [27] United States Court of Appeals for the Third Circuit. *USA v. Byron Mitchell*, 2003. No. 02-2859.
- [28] United States Supreme Court. *Daubert et. al. v Merrell Dow Pharmaceuticals*, 1993. 509 U.S. 579.
- [29] Watson C, Garris M, Tabassi E (2004) *User's Guide to NIST Fingerprint Image Software 2 (NFIS2)*. NIST