

Provided for non-commercial research and education use.
Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/authorsrights>



Contents lists available at ScienceDirect

Pattern Recognition

journal homepage: www.elsevier.com/locate/pr

Likelihood ratio estimation in forensic identification using similarity and rarity



Yi Tang, Sargur N. Srihari*

University at Buffalo, The State University of New York, 338, Davis Hall, Buffalo, NY 14260, USA

ARTICLE INFO

Available online 7 August 2013

Keywords:

Likelihood ratio
Distance measure
Rarity
Forensic science
Impression evidence
Approximate inference

ABSTRACT

Forensic identification is the task of determining whether or not observed evidence arose from a known source. It is useful to associate probabilities with identification/exclusion opinions, either for presentation in court or to evaluate the discriminative power of a given set of attributes. At present, in most forensic domains outside of DNA evidence, it is not possible to make such a statement since the necessary probability distributions cannot be computed with reasonable accuracy, although the probabilistic approach itself is well-understood. In principle, it involves determining a likelihood ratio (LR) – the ratio of the joint probability of the evidence and source under the *identification* hypothesis (that the evidence came from the source) and under the *exclusion* hypothesis (that the evidence did not arise from the source). Evaluating the joint probability is computationally intractable when the number of variables is even moderately large. It is also statistically infeasible since the number of parameters to be determined from the data is exponential with the number of variables. An approximate method is to replace the joint probability by another probability: that of distance (or similarity) between evidence and object under the two hypotheses. While this reduces to linear complexity with the number of variables, it is an oversimplification leading to errors. We consider a third method which decomposes the LR into a product of two factors, one based on *distance* and the other on *rarity*. This result, which is exact for the univariate Gaussian case, has an intuitive appeal – forensic examiners assign higher importance to rare feature values in the evidence and low importance to common feature values. We generalize this approach to more complex data such as vectors and graphs, which makes LR estimation computationally tractable. Empirical evaluations of the three methods, done with several data types (continuous features, binary features, multinomial and graph) and several modalities (handwriting with binary features, handwriting with multinomial features and footwear impressions with continuous features), show that the distance and rarity method is significantly better than the distance only method.

© 2013 Elsevier Ltd. All rights reserved.

1. Introduction

Forensic identification concerns whether observed evidence arose from a known source. The probabilistic approach is to determine the likelihood ratio positive ($LR+$) [17,2,38,24,29] whose numerator is the joint probability of the evidence and source under the null hypothesis that the evidence arises from the source and the denominator is the joint probability under the alternate hypothesis that the evidence does not arise from the object. The evidence is deemed to have arisen from the source if $LR+ > 1$ and not from the source otherwise.

Determining the joint probability has high data requirements. Assume that the evidence and object are both characterized by n binary features, the joint distribution requires 2^{2n} probabilities or parameters. Even for small n this requires extremely large data sets for estimating parameters. Furthermore in forensic applications data sets are usually small making the approach infeasible.

A common solution is to use the probability of distance, or similarity, between the evidence and known instead of the joint probability [24,33]. The distance based method, which has a *constant* number of parameters, is simple to compute but there is a severe loss of information in going from a high-dimensional joint probability space to a one-dimensional distance space. This paper considers a third method based on a result of Lindley [21] for univariate Gaussian samples which combines the probability of distance with the probability of the mean of evidence and object, called *rarity*. Computing rarity exactly has the complexity of 2^n still, for which probabilistic graphical models (e.g., Bayesian networks) and mixture models are used to further simplify the computation.

1.1. Forensic evidence

Several forensic modalities can benefit from a characterisation of uncertainty, particularly *impression* evidence where “objects or materials retain characteristics of other objects or materials impressed against them”, e.g., latent fingerprints, footwear and

* Corresponding author. Tel.: +1 716 984 7258.

E-mail address: srihari@cedar.buffalo.edu (S.N. Srihari).

tire impressions [5], tool marks and bite marks. Handwriting can also be regarded as impression evidence.

The characteristics need an appropriate representation for both human and computer processing. Examples of such representations are feature vectors or graphs providing a structural description. Features in forensics are usually further divided into *class* and *individual* characteristics. Examples of class characteristics in the case of footwear evidence are those that result from the mass machine manufacturing and can associate the evidence only to one class, whereas *individual characteristics* can potentially associate the evidence to a specific instance of that class. Another example from latent print analysis where class characteristics are type 1 detail (whorls, arches, etc.) and individualizing characteristics are type 2 detail (minutiae).

1.2. Organization of paper

The rest of the paper is structured as follows. Section 2 defines the likelihood ratio (LR) under the identification scenario, Lindley's result and its generalization and the computational complexities of the three methods. Sections 3–5 discuss the application of the three methods to data with feature distributions that are continuous, discrete and graphs, respectively. A theoretical analysis using K-L divergence is given in Section 6. Concluding remarks are made in Section 7.

2. Likelihood ratio

Let $S = \{s_i\}$ be a set of sources. They correspond to, say, footwear outsoles. We denote the *evidence* by E , and the *object* by O . Let \mathbf{o} represent an object drawn from a source s_i , e.g., impression of a known footwear. Let \mathbf{e} represent the evidence drawn from a source s_j , e.g., a crime scene impression. The task is to determine the odds of whether \mathbf{o} and \mathbf{e} came from the same or different source. We can state two opposing hypotheses:

h^0 : \mathbf{o} and \mathbf{e} are from the *same* source ($i=j$); and
 h^1 : \mathbf{o} and \mathbf{e} are from *different* sources ($i \neq j$), which are the *identification* and *exclusion* hypotheses of forensics; some forensic statistics literature also refers to them as *prosecution* and *defense* hypotheses [2].

We can define two joint probability distributions $P(\mathbf{o}, \mathbf{e}|h^0)$ and $P(\mathbf{o}, \mathbf{e}|h^1)$ which specify as to how often each instance of the object and evidence occur together when they belong to the same source or to different sources. The relative strengths of evidence supporting the two hypotheses is quantified by the likelihood ratio

$$LR_j = LR(\mathbf{o}, \mathbf{e}) = \frac{P(\mathbf{o}, \mathbf{e}|h^0)}{P(\mathbf{o}, \mathbf{e}|h^1)} \quad (1)$$

The corresponding log-likelihood ratio, $LLR(\mathbf{o}, \mathbf{e}) = \ln P(\mathbf{o}, \mathbf{e}|h^0) - \ln P(\mathbf{o}, \mathbf{e}|h^1)$, has representational advantages: its sign is indicative of same or different source, it has a smaller range than LR, and additivity of contributions of independent features.¹

It is useful to convert LRs into probabilities of identification and exclusion using a Bayesian formulation.² Let the *prior* probabilities of the hypotheses be $P(h^0)$ and $P(h^1)$ with $P(h^0) + P(h^1) = 1$.

¹ In performing LLR additions, since LR values in the interval $(1, \infty)$ convert to positive LLRs and LR values in the interval $(0, 1)$ convert to negative LLRs, the precisions of LRs < 1 must be high, otherwise the ranges of positive and negative LLRs will not be symmetric.

² There has been recent discussion in the legal community regarding whether the terms identification and individualization, commonly used in expressing forensic opinion, are appropriate in view of the fact that there is uncertainty in identification [26]. In response the forensic community has proposed that probabilities can be associated with identification, and that when the probability approaches 1 it is referred to as individualization [39].

Defining the prior odds as $O_{prior} = P(h^0)/P(h^1)$, we can express the prior probability of the same source as

$$P(h^0) = O_{prior} / (1 + O_{prior}). \quad (2)$$

The prior odds can be converted into posterior odds as

$$O_{posterior} = \frac{P(h^0|\mathbf{o}, \mathbf{e})}{P(h^1|\mathbf{o}, \mathbf{e})} = O_{prior} \times LR(\mathbf{o}, \mathbf{e}). \quad (3)$$

Thus we can write the posterior probability of the same source as $P(h^0|\mathbf{o}, \mathbf{e}) = O_{posterior} / (1 + O_{posterior})$. The particular case of equal priors is of interest in forensics, as opinion without prior bias, i.e., no information other than the evidence is used. In this case we get a simple form for the probability of identification as

$$P(h^0|\mathbf{o}, \mathbf{e}) = \frac{LR(\mathbf{o}, \mathbf{e})}{1 + LR(\mathbf{o}, \mathbf{e})} = \frac{\exp(LLR(\mathbf{o}, \mathbf{e}))}{1 + \exp(LLR(\mathbf{o}, \mathbf{e}))}. \quad (4)$$

The probability of exclusion is $P(h^1|\mathbf{o}, \mathbf{e}) = 1 - P(h^0|\mathbf{o}, \mathbf{e}) = 1 / [1 + LR(\mathbf{o}, \mathbf{e})] = 1 / [1 + e^{LLR(\mathbf{o}, \mathbf{e})}]$.

Thus the key to determining the probability of identification is to determine LR defined by Eq. (1), which in turn requires the distributions $P(\mathbf{o}, \mathbf{e}|h^i)$ ($i = 0, 1$), defined over all possible values of objects and their evidential forms. If \mathbf{o} and \mathbf{e} are n -dimensional multinomial vectors with each feature taking K possible values, then $2K^{2n}$ parameters are needed to specify the joint distribution. Determining these distributions is computationally and statistically infeasible. Computationally, kernel density estimation [3] and finite mixture models [23] have been proposed, but they have limitations as well.³ More important is the statistical limitation of having a sufficient number of samples for so many parameters. Today, objects and evidence can be represented by ever finer features due to higher camera resolution and automatic feature extraction methods and their possible evidential forms are infinite.

One method of simplification is to use a (dis)similarity function between object and evidence. The approach is to define $d(\mathbf{o}, \mathbf{e})$ as a scalar *distance* between object and evidence and define another likelihood ratio as follows:

$$LR_D = \frac{P(d(\mathbf{o}, \mathbf{e})|h^0)}{P(d(\mathbf{o}, \mathbf{e})|h^1)} \quad (5)$$

The number of parameters needed to evaluate LR_D is constant, or $O(1)$, and is independent of the number of features n . Due to its simplicity, this method has been proposed for fingerprint identification [24], handwriting analysis [32], pharmaceutical tablet comparison [6], etc.

For certain feature spaces and distance measures, e.g., continuous features with Euclidean distance, this approach is equivalent to a *kernel* method [28]. The scalar distance d is just the magnitude of the vector difference \mathbf{d} . However, because it maps two distributions of $2n$ variables each into two scalar distributions there is severe loss of information (many pairs of \mathbf{o} and \mathbf{e} can have the same distance). A natural extension is to use vector difference \mathbf{d} , which quantifies the distribution of both the magnitude and the orientation of the difference between \mathbf{o} and \mathbf{e} , giving a much fine-grained characterization of the difference between \mathbf{o} and \mathbf{e} . While this likelihood ratio LR_{VD} provides the simplification of mapping two distributions of $2n$ variables each into two distributions of n variables each, there is still a loss of information in the many to one mappings.

³ Kernel density estimation is expensive in memory as it needs to store the entire training data set, and the cost of evaluating the density grows linearly with size of data set. For mixture models, an important issue is to select the number of components, also, the training algorithm such as EM may converge to a local optimum or the boundary of the parameter space [18].

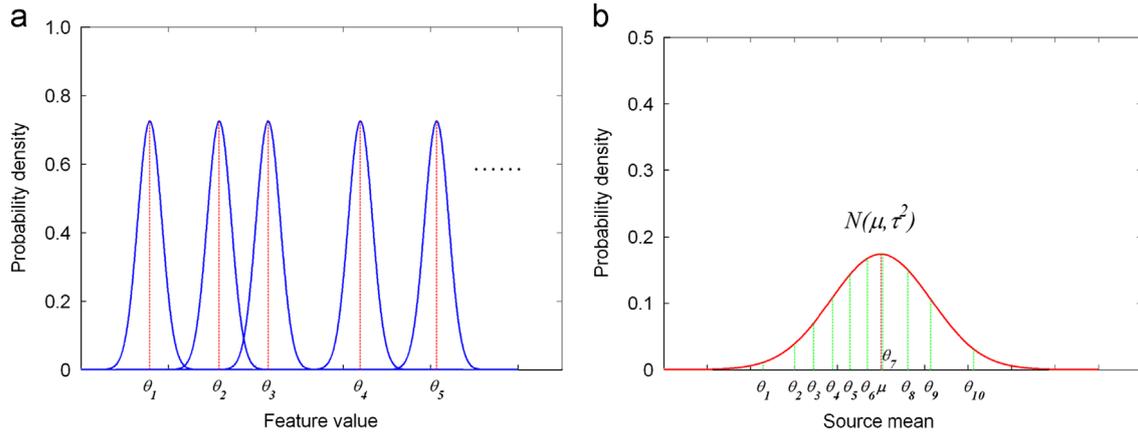


Fig. 1. Distributions of sources for object and evidence: (a) normally distributed sources, where each source s_i ($i = 1, 2, \dots$) is $\mathcal{N}(\theta_i, \sigma^2)$. Samples may come from a common source s_i or different sources s_i and s_j ($s_i \neq s_j$). (b) Source means $\{\theta_1, \theta_2, \dots\}$ are assumed to have distribution $\mathcal{N}(\mu, \tau^2)$, with $\tau \gg \sigma$. Samples coming from sources θ_1, θ_{10} are rarer (less frequent) than samples from θ_6 and θ_7 , suggesting that information about the distribution of source means is useful to assess strength of evidence.

2.1. Lindley's result

Evaluating the likelihood ratio (positive) was considered by Lindley [21] when the object and evidence arise from univariate Gaussian distributions. This is illustrated in Fig. 1. The measurement considered was the refractive index of glass, with the object being fragment samples on a window and evidence being samples on a suspect's clothing.

Theorem 1 (Lindley). *If the object and evidence are continuous scalar random variables drawn from the same or different sources; samples are normally distributed about its source mean with a known constant variance σ^2 ; the source mean is normally distributed with mean μ and variance τ^2 with $\tau \gg \sigma$, and there are p object samples with mean o , and q evidence samples with mean e with $p=q$, then the likelihood ratio can be approximated by*

$$LR(o, e) = \frac{\tau}{\sigma\sqrt{2/p}} \exp\left\{-\frac{(o-e)^2}{4\sigma^2/p}\right\} \exp\left\{-\frac{(m-\mu)^2}{2\tau^2}\right\}, \quad (6)$$

where $m = (o + e)/2$ is the mean of o and e .

Proof. Let $o \sim \mathcal{N}(\theta_1, \sigma^2/p)$ and $e \sim \mathcal{N}(\theta_2, \sigma^2/q)$. To evaluate Eq. (1), consider first its denominator. Since o and e do not come from the same source, they are independent. Thus we have

$$P(o, e|h^1) = \int P(o|\theta_1)P(\theta_1) d\theta_1 \cdot \int P(e|\theta_2)P(\theta_2) d\theta_2,$$

which is a product of the unconditional distributions of o and e . These are, respectively, $\mathcal{N}(\mu, \tau^2 + \sigma^2/p)$ and $\mathcal{N}(\mu, \tau^2 + \sigma^2/q)$, and are independent.

Writing $\sigma_1^2 = \tau^2 + \sigma^2/p$, $\sigma_2^2 = \tau^2 + \sigma^2/q$, we alternately have that $P(o-e|h^1)$ is $\mathcal{N}(0, \sigma_1^2 + \sigma_2^2)$. Writing $z = (\sigma_2^2 e + \sigma_1^2 o)/(\sigma_1^2 + \sigma_2^2)$, we have $P(z|h^1)$ is $\mathcal{N}(\mu, \sigma_1^2 \sigma_2^2 / (\sigma_1^2 + \sigma_2^2))$. These are also independent.

As the Jacobian J of the transformation⁴ is 1, the denominator can be equivalently written by $P(o-e, z)|J|$, then expanded as

$$P(o, e|h^1) = \frac{1}{2\pi\sigma_1\sigma_2} \exp\left\{-\frac{(o-e)^2}{2(\sigma_1^2 + \sigma_2^2)}\right\} \exp\left\{-\frac{(z-\mu)^2(\sigma_1^2 + \sigma_2^2)}{2\sigma_1^2\sigma_2^2}\right\}. \quad (7)$$

⁴ Let X and Y be jointly distributed with p.d.f. $p(x, y)$, U and V be constructed from (X, Y) via a one-to-one transform $\begin{cases} U = U(X, Y) \\ V = V(X, Y) \end{cases}$. The Jacobian J of the inverse transform $\begin{cases} X = X(U, V) \\ Y = Y(U, V) \end{cases}$ is given by $J = \partial(x, y)/\partial(u, v) = \begin{vmatrix} \partial x/\partial u & \partial x/\partial v \\ \partial y/\partial u & \partial y/\partial v \end{vmatrix}$. Then the joint p.d.f. of (U, V) can be computed by $p(u, v) = p(x, y)|J|$ [40], where $|J|$ is the determinant of J .

The numerator of Eq. (1) is evaluated in a similar way. Let θ be the common value of θ_1 and θ_2 . Since the evidence and object are drawn independently from the same distribution

$$P(o, e|h^0) = \int P(o|\theta)P(e|\theta)P(\theta) d\theta. \quad (8)$$

The required joint unconditional distribution of o and e is bivariate normal with means μ , variances σ_1^2 and σ_2^2 and covariance τ_2 . Hence $P(o-e|h^0)$ is $\mathcal{N}(0, \sigma^2(1/p + 1/q))$. Writing $m = (po + qe)/(p + q)$, $P(m|h^0)$ is $\mathcal{N}(\mu, \tau^2 + \sigma^2/(p + q))$, and they are independent. Again with unit Jacobian, the numerator is

$$P(o, e|h^0) = \frac{1}{2\pi a\sigma\sigma_3} \exp\left\{-\frac{(o-e)^2}{2a^2\sigma^2} - \frac{(m-\mu)^2}{2\sigma_3^2}\right\}, \quad (9)$$

where $\sigma_3^2 = \tau^2 + \sigma^2/(p + q)$ and $a^2 = 1/p + 1/q$. Then the LR is the ratio of (9) to (7), given by

$$LR = \frac{\sigma_1\sigma_2}{a\sigma\sigma_3} \exp\left\{-\frac{(o-e)^2\tau^2}{a^2\sigma^2(\sigma_1^2 + \sigma_2^2)} - \frac{(m-\mu)^2}{2\sigma_3^2} + \frac{(z-\mu)^2(\sigma_1^2 + \sigma_2^2)}{2\sigma_1^2\sigma_2^2}\right\}.$$

Applying assumption of $p=q$, we have $\sigma_1^2 = \sigma_2^2$, and $m = z = (o + e)/2$. Thus,

$$LR = \frac{\sigma_1\sigma_2}{a\sigma\sigma_3} \exp\left\{-\frac{(o-e)^2\tau^2}{4\sigma^2\sigma_1^2/m} + \frac{(m-\mu)^2\tau^2}{2\sigma_1^2\sigma_3^2}\right\}. \quad (10)$$

From assumption of $\tau \gg \sigma$, we have $\sigma_1^2 \approx \sigma_2^2 \approx \sigma_3^2 \approx \tau^2$. Thus (10) can be finally rewritten as Eq. (6) □

When $p=q=1$, Eq. (6) takes the simple form

$$LR = \frac{\tau}{\sqrt{2}\sigma} \exp\left\{-\frac{(o-e)^2}{4\sigma^2}\right\} \exp\left\{-\frac{(m-\mu)^2}{2\tau^2}\right\}. \quad (11)$$

It is this form that is discussed (see Fig. 2), though similar remarks apply to the more general form of Eq. (6). The first exponential in Eq. (11), the *distance* factor, is a significance test of the null hypothesis of identity, based on comparing $|o-e|$ with the standard deviation $\sqrt{2}\sigma$. If $|o-e| = \lambda\sqrt{2}\sigma$, its value is $e^{-(1/2)\lambda^2}$. A small value of λ favors the hypothesis that the two sets came from the same source.

The second exponential, the *rarity* factor, measures how typical the two sets are. If $m = \frac{1}{2}(o + e)$ is μ , the most likely value, it is one in which case the first term alone (similarity) determines the strength of evidence. As $|m-\mu|$ increases relative to its standard deviation τ , i.e., increasing rarity, the second term increases. For a fixed similarity, characterized by difference $o-e$, rarity of the

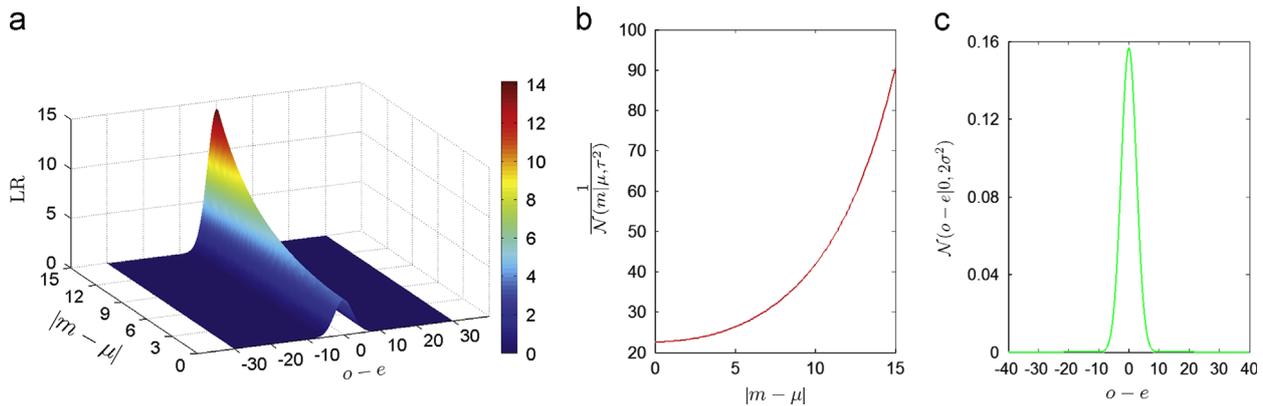


Fig. 2. Dependence of likelihood ratio on rarity and similarity in the univariate Gaussian case: (a) plot of LR against distance between known and evidence $o-e$ and rarity of evidence and known $|m-\mu|$, (b) when the difference $o-e$ is fixed, the rarity factor $1/\mathcal{N}(m|\mu, \tau^2)$ (and LR+) grows exponentially with rarity $|m-\mu|$. (c) When rarity $|m-\mu|$ is fixed, the difference factor $\mathcal{N}(o-e|0, 2\sigma^2)$ (and LR+) increases as $|o-e|$ the difference decreases. Parameters used are $o, e \in [1, 31]$, $\mu=16$, $\tau=9$, $\sigma=1.8$.

evidence is extremely important to establish the strength of evidence (Fig. 2(b)). When the rarity of the evidence is known, the strength of evidence decreases as dissimilarity increases (Fig. 2(c)).

Lindley's theorem justifies a forensic examiner's intuition that similarity between the known and the evidence carries more weight when the feature is rare. Also, since evaluating the LR using the joint distribution is computationally and statistically infeasible, Lindley's factorization provides a way out.

2.2. Generalization of Lindley's result

Here we generalize Lindley's theorem to the general cases of vectors, discrete features and graphs. We begin by rewriting Eq. (11) for the univariate Gaussian case as the product of two factors, the first as the density of the difference and the second as the reciprocal of the density of mean of object and evidence, as follows:

$$LR = \frac{\frac{1}{\sqrt{2\sigma}} \exp\left\{-\frac{(o-e)^2}{4\sigma^2}\right\}}{\frac{1}{\tau} \exp\left\{-\frac{(m-\mu)^2}{2\tau^2}\right\}} = \frac{\mathcal{N}(o-e|0, (\sqrt{2}\sigma)^2)}{\mathcal{N}(m|\mu, \tau^2)} = P(o-e|h^0) * \frac{1}{P(m)}. \tag{12}$$

In the generalization of Eq. (12), LR is approximated as the product of two factors, one based on difference and the other on rarity:

$$LR_{DR} = P(\mathbf{d}(\mathbf{o}, \mathbf{e})|h^0) * \frac{1}{P(\mathbf{m}(\mathbf{o}, \mathbf{e}))}, \tag{13}$$

where $\mathbf{d}(\mathbf{o}, \mathbf{e})$ is the difference between \mathbf{o} and \mathbf{e} , and $\mathbf{m}(\mathbf{o}, \mathbf{e})$ is the mean of \mathbf{o} and \mathbf{e} . The difference term in Eq. (13) can be defined and expanded as

$$T_d = P(\mathbf{d}(\mathbf{o}, \mathbf{e})|h^0) = \sum_{\mathbf{o}, \mathbf{e}} P(\mathbf{d}(\mathbf{o}, \mathbf{e})|\mathbf{o}, \mathbf{e}, h^0) P(\mathbf{o}, \mathbf{e}|h^0). \tag{14}$$

According to Eq. (14) the difference distribution is estimated from pairs of (object, evidence) samples from the same source. The dependence on the distribution of sources can be explicitly written as

$$P(\mathbf{o}, \mathbf{e}|h^0) = \sum_s P(\mathbf{o}, \mathbf{e}|h^0, s_i) P(s_i).$$

The rarity term in Eq. (13) can be expanded as

$$T_r = P(\mathbf{m}(\mathbf{o}, \mathbf{e})) = P(h^0)P(\mathbf{m}(\mathbf{o}, \mathbf{e})|h^0) + P(h^1)P(\mathbf{m}(\mathbf{o}, \mathbf{e})|h^1). \tag{15}$$

According to Eq. (15) rarity is estimated from representative pairs of object and evidence irrespective of source. To keep the

distribution unbiased, the number of sample pairs sampled under h^0 should almost equal the number of sample pairs that are drawn under h^1 .

The definitions of difference and mean with forensic data is not always clear, which we consider next.

2.2.1. Definition of difference

For continuous vector values of \mathbf{o} and \mathbf{e} the obvious choice is vector difference $\mathbf{d}(\mathbf{o}, \mathbf{e}) = \mathbf{o} - \mathbf{e}$. The difference between categorical features can be regarded as another categorical variable, which can be mapped to numbers for easy implementation, although those numbers do not indicate any ordering. Difference can be more generally regarded as a distance measure. Distance measures are closely related to similarity measures and usually considered to be reciprocals of each other. Specific distance measures for binary and multinomial features are given in Section 4.

2.2.2. Definition of mean

For continuous vector values the vector average is the obvious choice. For binary features, categorical features and graphs, a simple average does not have a meaningful interpretation, or even does not make any sense, e.g., average of the attributes of a line and the attributes of a circle. In this case, the weighted "mean" of the evidence and the object can be interpreted as the set of common features or matched characteristics, as they are the nearest candidate to represent the true mean of the object and the evidence. Thus we have for categorical discrete features or graphs, $\mathbf{m}(\mathbf{o}, \mathbf{e})$ is the common features between \mathbf{e} and \mathbf{o} . So both $\mathbf{d}(\mathbf{o}, \mathbf{e})$ and $\mathbf{m}(\mathbf{o}, \mathbf{e})$ have the same dimension.

2.2.3. Definition of rarity

Rarity is defined as the reciprocal of probability and therefore follows the same definitions as probability. In the discrete case, given a probability space (Ω, \mathcal{F}, P) , the rarity of a random event $\xi \in \mathcal{F}$ is defined by $R(\xi) = 1/P(\xi)$, where Ω is the sample space, $\mathcal{F} \subseteq 2^\Omega$ is the set of events, P is the probability measure, and $P(\xi) (\neq 0)$ is the probability of the event ξ . For the continuous case, let $\mathbf{x} = (x_1, \dots, x_n)^T$ be a continuous n -dimensional random vector with the p.d.f. $p(\mathbf{x})$ defined on a domain S . Suppose for every assignment of $\mathbf{x} \in S$, there is a confidence interval $(\mathbf{x} - \epsilon/2, \mathbf{x} + \epsilon/2)$ associated with \mathbf{x} at a given confidence level $1 - \alpha$, where $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T$, and α is a small positive value less than 1. This interval $(\mathbf{x} - \epsilon/2, \mathbf{x} + \epsilon/2)$ is a n -dimensional region whose volume is $\prod_{i=1}^n \epsilon_i$. Then the rarity of the event that \mathbf{x} takes value \mathbf{x}_0 is given by $R(\mathbf{x} = \mathbf{x}_0) = 1 / \int_{\mathbf{x}_0 - \epsilon/2}^{\mathbf{x}_0 + \epsilon/2} p(\mathbf{x}) d\mathbf{x}$. Since the magnitude of ϵ is usually small, the density in the region $(\mathbf{x}_0 - \epsilon/2, \mathbf{x}_0 + \epsilon/2)$ can be

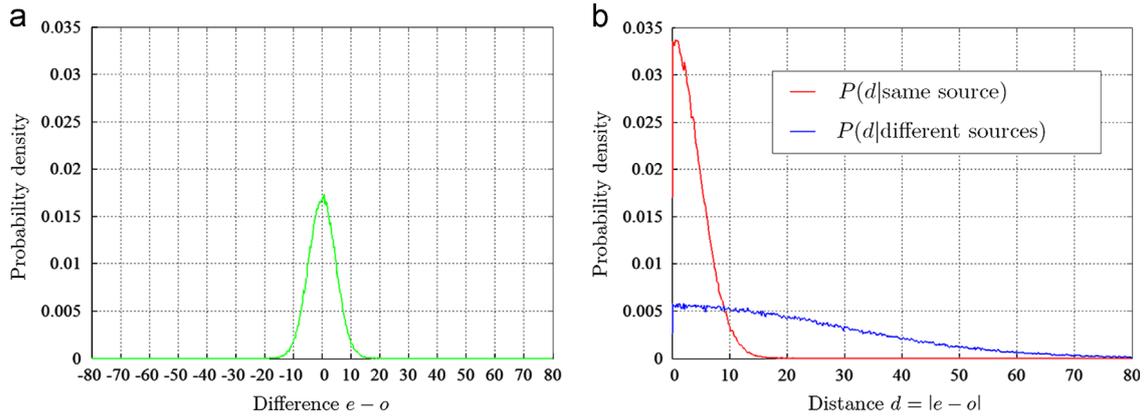


Fig. 3. Distribution of difference when observation and evidence have a univariate Gaussian distributions: (a) within-source distribution of the signed difference: $P(e-o|h^0) = \mathcal{N}(0, 2\sigma^2)$; (b) within/between-source distributions of the absolute difference: $d = |e-o|$, where the red curve shows $P(d|h^0) = 2\mathcal{N}(0, 2\sigma^2)$, $d \in [0, +\infty)$, and the blue curve shows $P(d|h^1) = 2\mathcal{N}(0, 2\tau^2 + 2\sigma^2)$, $d \in [0, +\infty)$. Parameters used are $\mu=200$, $\tau=20$, $\sigma=3.3$. Each distribution is built from 20,000 samples independently drawn from Gaussian distributions. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this article.)

considered as constant, therefore rarity can be approximated by $R(\mathbf{x} = \mathbf{x}_0) = 1/p(\mathbf{x}_0) \prod_{i=1}^n \epsilon_i$, $\forall \mathbf{x}_0 \in S$.

In the case of categorical variables and graphs, rarity is the reciprocal of the probability (or p.d.f.) of observing these common features in the general population.

2.2.4. Computational considerations

The distributions of difference and mean each involve n variables. The complexity of computing the two distributions exactly for non-Gaussian distributed data is still exponential with the number of features, with the number of parameters required being $2K^n$. However, we are dealing with two separate n -dimensional distributions rather than a single $2n$ -dimensional distribution as in the case of the joint distribution method. For determining the n -dimensional distributions, approximation techniques based on mixture models or probabilistic graphical models, e.g., Bayesian networks [30], can be used. Finally, the explicit factorization into difference and rarity matches human intuition thereby having explanatory power which is essential in forensics.

2.3. Evaluating LR methods

We now have three methods for computing LR: exact LR_j defined by Eq. (1), the commonly used LR_D defined by Eq. (5), and the approximate LR_{DR} defined by Eq. (13). It is noted that all these values, including exact LR, correspond to the particular feature representation chosen.

In order to compare the performance different methods of determining LR we need an appropriate metric. A method can be evaluated assuming a binary decision of *identification/exclusion* using the decision rule: if $LR > 1$, then the known and evidence have the same source (identification), i.e., $s_i = s_j$, and if $LR < 1$ we have exclusion, i.e., $s_i \neq s_j$. Denoting the probability of misclassifying as $s_i \neq s_j$ when $s_i = s_j$ as ϵ_1 , and the probability of misclassifying as $s_i = s_j$ when $s_i \neq s_j$ as ϵ_2 , the average error rate is $\frac{1}{2}(\epsilon_1 + \epsilon_2)$.

We note a disadvantage of the average error rate measure: two very different LRs, e.g., 2 and 20,000,000 mean the same. If under h^0 we have $LR=2$ then it is more likely to be inclusive and $LR=20,000,000$ means a high probability of identification for a wide range of priors. But if under h^1 , $LR=2$ is an acceptable “error” but $LR=20,000,000$ implies an unreliable technique. So the “cost” of different assigned LRs with different magnitudes cannot be the same. Distributions of LRs under h^0 and h^1 plotted as cumulative distributions, known as Tippet plots, show probabilities of mis-leading evidences $LR > 1$ in h^0 , $LR < 1$ in h^1 and the magnitude of

errors; Tippet plots are useful for visual comparison of different techniques/experiments [13,32]. Scoring rules, such as *ClIr* in [8], give a single metric (as with average error rate) but penalizing errors and the magnitude of the errors. Since LR values can always be thresholded, so that small magnitudes of LR are deemed inconclusive and not included in the evaluation, the average error rate is a reasonable baseline measure.

In the next two sections we specialize each of the three LR methods to different types of underlying feature values and provide experimental evaluations, first considering fixed length vector representations that are continuous, then the discrete case (binary and multinomial) and finally graph representations.

3. Continuous case

3.1. Univariate Gaussian

Here LR_j is given by Eq. (6). To derive an expression for LR_D , consider $P(e-o|h^0) \sim \mathcal{N}(0, (a\sigma)^2)$, where $a^2 = (1/p + 1/q)$ is defined by the number of evidence and observation samples; shown for $p=q=1$ in Fig. 3(a). The absolute difference $d = |e-o|$ follows a *folded normal distribution*⁵ expressed by

$$P(d|h^0) = \begin{cases} \mathcal{N}(0, a^2\sigma^2) & \text{if } d = 0, \\ 2\mathcal{N}(0, a^2\sigma^2) & \text{if } d > 0. \end{cases}$$

Similarly, $P(e-o|h^1) \sim \mathcal{N}(0, 2\tau^2 + a^2\sigma^2)$, which leads to

$$P(d|h^1) = \begin{cases} \mathcal{N}(0, 2\tau^2 + a^2\sigma^2) & \text{if } d = 0 \\ 2\mathcal{N}(0, 2\tau^2 + a^2\sigma^2) & \text{if } d > 0. \end{cases}$$

The two distance distributions are illustrated in Fig. 3(b). Then a good approximation is $LR_D = (\sqrt{2}\tau/a\sigma) \exp\{-d^2/2a^2\sigma^2\}$, as typically $\tau \gg \sigma$. When $p=q=1$, we further have $LR_D = (\tau/\sigma) \exp\{-d^2/4\sigma^2\}$, which can be substituted into (11) to get

$$LR_j = LR_D * \frac{1}{\sqrt{2} \exp\left\{-\frac{(m-\mu)^2}{2\tau^2}\right\}}, \quad (16)$$

which indicates that LR_D differs from LR_j by a factor of un-normalized rarity. The expression for LR_{DR} is obtained by

⁵ Given $X \sim \mathcal{N}(\mu, \sigma^2)$, the random variable $Y = |X|$ follows a folded normal distribution, which is also called *half-normal* distribution if $\mu = 0$.

rewriting LR_J in (11) as

$$LR_{DR} = \frac{1}{\sqrt{2\pi}2^{1/2}\sigma} \exp\left\{-\frac{(e-o)^2}{2(2^{1/2}\sigma)^2}\right\} * \frac{1}{\sqrt{2\pi}\tau} \exp\left\{-\frac{(m-\mu)^2}{2\tau^2}\right\}, \quad (17)$$

which decomposes LR_J as the product of two factors, i.e., *difference* and *rarity*, where difference factor is quantified by the probability density of the distribution of the difference $e-o$ given h^0 , i.e., $e-o|h^0 \sim \mathcal{N}(0, (\sqrt{2}\sigma)^2)$, and rarity is quantified by the reciprocal of the p.d.f. of $m = (e+o)/2$, which assesses how unusual m is, i.e., m strictly follows $\mathcal{N}(\mu, \tau^2 + \sigma^2/2)$ that can be reduced to $\mathcal{N}(\mu, \tau^2)$ as $\tau \gg \sigma$. If $m = \mu$, the most likely value, the rarity is 1, and it increases as $|m-\mu|$ increases. In other words, the likelihood of hypothesis h^0 being true is greater the more unusual are the evidence and the object. For the multiple-sample case, this decomposition still holds by defining difference and rarity by $\mathcal{N}(0, (\sqrt{2/m}\sigma)^2)$ and $\mathcal{N}(\mu, \tau^2)$, respectively.

Experimental evaluation of the two cases, $LR_J = LR_{DR}$ and LR_D , was performed with two univariate Gaussian data sets: synthetic data and paper thickness data. In the synthetic cases, shown in Fig. 4, source means θ_i ($i=1,2,\dots$) were drawn from Gaussian $\mathcal{N}(\mu, \tau^2)$, then for each source mean θ_i , samples were drawn from distribution $\mathcal{N}(\theta_i, \sigma^2)$. Average error rate plotted against ratio τ/σ in Fig. 4(a) shows that: (i) the error rates of both methods decrease as τ/σ increases, because large value of τ/σ indicates samples from different sources are well separated; and (ii) the distance and rarity method consistently outperforms the distance method.

Paper thickness data. We applied the distance and rarity method to data consisting of measurements of the thickness of two pieces of paper (Table 1), where the task is to decide if they come from the same source. Assuming a Gaussian distribution, the mean and variance of the measurements of each page are estimated by a maximum likelihood estimator. The mean of the thickness of a piece of paper is estimated to be 0.1 mm [1], or 0.0039 in., while its standard deviation is approximated as 0.00048 in. (by averaging the variances shown in the original Fig. 3 in [16]). Computing T_d and T_r for the univariate case described in Section 2, it is concluded that the two pieces of paper come from *different* sources. As a cross-check of the result, we also performed a *likelihood ratio test* (LRT, a standard method of hypothesis testing [12]), where the likelihood functions in both the numerator and denominator are computed using maximum likelihood based on a Gaussian assumption. The result of LRT is consistent with the proposed method.

3.2. Multivariate Gaussian

Aitken and Lucy [3] extended Lindley's work to multivariate data in *vector* form. Bolck et al. [6] compared the distance and joint distribution methods for LR computation and showed that the latter performs better. Their methods were tested with only low-dimensional data and normally distributed features. Further, both failed to reveal clearly two key terms in the expression of the LR.

We take measurements on samples drawn from the object and evidence, where each measurement is a *vector* of k real values. Measurements on multiple samples from the same source are normally distributed about its true value with a constant covariance matrix \mathbf{C} . Let \mathbf{e} be the *mean* of p measurements on the evidence, $\mathbf{e} \sim \mathcal{N}(\theta_1, p^{-1}\mathbf{C})$. Similarly, let \mathbf{o} be the *mean* of q measurements on the object, $\mathbf{o} \sim \mathcal{N}(\theta_2, q^{-1}\mathbf{C})$. True values θ_i ($i=1,2,\dots$) are also normally distributed about mean μ with covariance Σ . Since between-source variation is usually much greater than within-source variation [21,3], we have $|\Sigma| \gg |\mathbf{C}|$. The likelihood ratio then is $LR_J = P(\mathbf{e}, \mathbf{o}|h^0)/P(\mathbf{e}, \mathbf{o}|h^1)$, whose denominator and numerator are expanded as follows.

Denominator: It is a product of two marginal distributions of \mathbf{e} and \mathbf{o} given by $\mathcal{N}(\mu, \Sigma + p^{-1}\mathbf{C})$ and $\mathcal{N}(\mu, \Sigma + q^{-1}\mathbf{C})$, respectively. Denote $\Sigma + p^{-1}\mathbf{C}$ by Σ_1 , and $\Sigma + q^{-1}\mathbf{C}$ by Σ_2 , we have $\mathbf{e}-\mathbf{o}|h^1 \sim \mathcal{N}(\mathbf{0}, \Sigma_1 + \Sigma_2)$, and $(\Sigma_1 + \Sigma_2)^{-1}(\Sigma_2\mathbf{e} + \Sigma_1\mathbf{o})$ (denoted by \mathbf{z}) $\sim \mathcal{N}(\mu, \Sigma_1\Sigma_2(\Sigma_1 + \Sigma_2)^{-1})$. So the denominator is

$$\frac{1}{(2\pi)^k |\Sigma_1\Sigma_2|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{e}-\mathbf{o})^T(\Sigma_1 + \Sigma_2)^{-1}(\mathbf{e}-\mathbf{o})\right\} \exp\left\{-\frac{1}{2}(\mathbf{z}-\mu)^T(\Sigma_1 + \Sigma_2)^{-1}(\mathbf{z}-\mu)\right\}. \quad (18)$$

Numerator: Given h^0 , $\mathbf{e}-\mathbf{o} \sim \mathcal{N}(\mathbf{0}, a^2\mathbf{C})$, where $a^2 = 1/p + 1/q$, and $\mathbf{w} = (p\mathbf{e} + q\mathbf{o})/(p+q) \sim \mathcal{N}(\mu, \Sigma + (p+q)^{-1}\mathbf{C})$. Denote $\Sigma + (p+q)^{-1}\mathbf{C}$ by Σ_3 , the numerator is written by

$$\frac{1}{(2\pi)^k |a|\mathbf{C}\Sigma_3|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{e}-\mathbf{o})^T(a^2\mathbf{C})^{-1}(\mathbf{e}-\mathbf{o})\right\} \exp\left\{-\frac{1}{2}(\mathbf{w}-\mu)^T\Sigma_3^{-1}(\mathbf{w}-\mu)\right\}. \quad (19)$$

When $p=q$, $\mathbf{w} = \mathbf{z} = \frac{1}{2}(\mathbf{e} + \mathbf{o})$. As $|\Sigma| \gg |\mathbf{C}|$, the likelihood ratio is approximated by

$$LR_J = a^{-1}|\mathbf{C}\Sigma^{-1}|^{1/2} * \exp\left\{-\frac{1}{2}(\mathbf{e}-\mathbf{o})^T(a^2\mathbf{C})^{-1}(\mathbf{e}-\mathbf{o})\right\} * \exp\left\{\frac{1}{2}(\mathbf{w}-\mu)^T\Sigma^{-1}(\mathbf{w}-\mu)\right\}. \quad (20)$$

Table 1
Paper thickness measurements in inches.

Page 1	0.0042	0.0042	0.0042	0.0043	0.0042	0.0042	0.0041	0.0041
Page 2	0.0043	0.0043	0.0043	0.0043	0.0044	0.0044	0.0044	0.0043

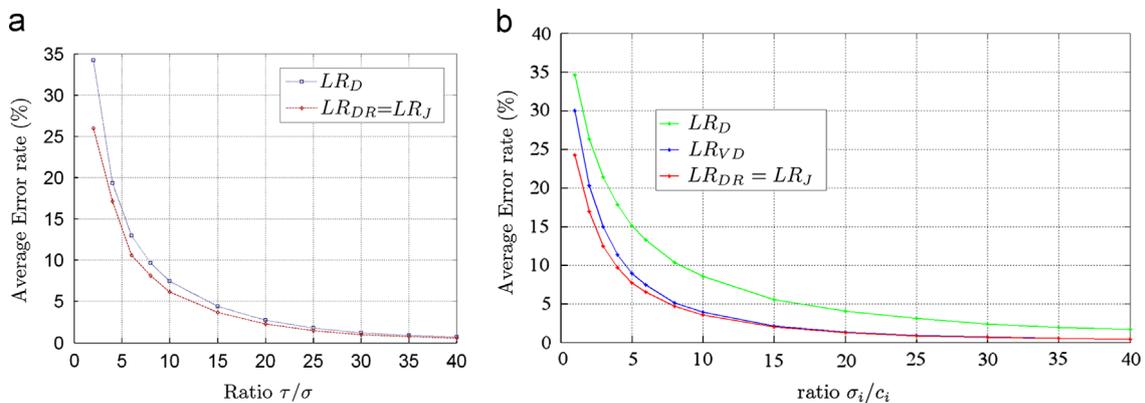


Fig. 4. Average error rates on synthetic Gaussian data: (a) univariate, and (b) multivariate (five dimensions). In both plots error rate is lowest with LR_{DR} . (a) Univariate and (b) five dimensions.

The first exponential term represents the difference between evidence and object, while the second term accounts for rarity.

Thus when features are either of low dimension or follow simple parametric distributions, $P(\mathbf{e}, \mathbf{o}|h^i)$ and LR_j have a closed form solution. If the feature distribution is either multimodal or of high dimensions, LR_j is difficult to evaluate, thus making LR_D a possible solution.

For LR_D , the obvious dissimilarity measure is the Euclidean distance. The distributions $P(d|h^i)$ ($i=0,1$) can be represented by Gamma distributions (see Fig. 5), where their shape and scale parameters are estimated from samples [32].

For LR_{DR} , the difference factor T_d is defined by the p.d.f. of the difference $\mathbf{e}-\mathbf{o}$ given h^0 , and rarity T_r is computed by the reciprocal of the p.d.f. of \mathbf{w} , where \mathbf{w} is the weighted mean of \mathbf{e} and \mathbf{o} .

In particular, when features are multivariate Gaussian, we can rewrite (20) as

$$LR_{DR} = \mathcal{N}(\mathbf{e}-\mathbf{o}|\mathbf{0}, a^2\mathbf{C}) * \frac{1}{\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})}, \quad (21)$$

where $a^2 = 1/p + 1/q$. It can be found that difference factor T_d is given by Gaussian $\mathcal{N}(\mathbf{e}-\mathbf{o}|\mathbf{0}, a^2\mathbf{C})$, and rarity T_r is given by the reciprocal of Gaussian $\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

We used synthetic data to evaluate the LR methods. Because $LR_j=LR_{DR}$, we list only the results of LR_D and LR_{DR} . As shown in Fig. 5, source means θ_i ($i=1, \dots, n$) were drawn from a 5-dimensional Gaussian $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma} = \text{diag}(\sigma_i^2)$. Then for each source mean θ_i , samples were drawn from Gaussian $\mathcal{N}(\theta_i, \mathbf{C})$, where $\mathbf{C} = \text{diag}(c_i^2)$. For each pair of samples (\mathbf{e}, \mathbf{o}), LR_D and LR_{DR} were computed to decide whether to classify them as having come from the same/different source(s) in each method. The results with different ratios σ_i/c_i are shown in Fig. 4(b). We conclude that (i) LR_{DR} performs consistently better than LR_D ; and (ii) as the ratio σ_i/c_i goes up, the error rates of both methods go down. A close view of Fig. 4(a) and (b) reveals that, as the dimension gets higher, the difference in performance increases.

4. Discrete features

Evaluation and performance of the three likelihood ratio methods with binary and multinomial features is discussed next.

4.1. Binary features

Here $x = [x_1, \dots, x_n]$ where $x_i \in \{0, 1\}$. We first consider independent x_i and then the dependent case. In order to define LR_D and LR_{DR} we need an appropriate difference or distance measure. The correlation distance between two binary vectors \mathbf{x} and \mathbf{y} is

$$d = \frac{1}{2} \frac{S_{11}S_{00} - S_{10}S_{01}}{2((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}}, \quad (22)$$

where S_{ij} ($i, j \in \{0, 1\}$) is the number of occurrences of matches with i in the first vector and j in the second vector at corresponding positions.

Binary feature vectors from the domain of forensic handwriting identification were used in the evaluation [41]. Binary features, known as GSC [32], represent a handwritten word image; it consists of 1024 bits corresponding to gradient (384 bits), structural (384 bits), and concavity (256 bits) features. Two data sets were used, the first consisting of 3125 extracted images of *th* written by 499 writers, where each writer wrote *th* 3 to 10 times. The second data set consisted of 10,111 images of *and* extracted from handwritten documents written by 1555 authors.

Histograms of $P(d|h^i)$ ($i=0,1$) using the *th* data set with GSC features (Fig. 6(a and b)) are unimodal. The nearly symmetric distribution $P(d|h^0)$ can be approximated by a Gaussian $\mathcal{N}(\mu_s, \sigma_s^2)$,

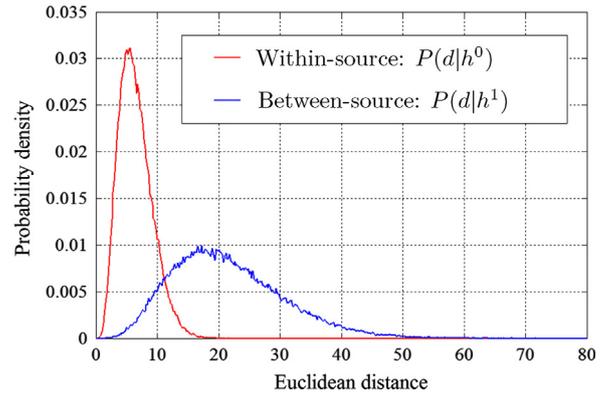


Fig. 5. Distributions of the Euclidean distance between samples that follow 5-dimensional Gaussian distribution, which are modeled by two Gamma distributions. Parameters used are $\boldsymbol{\mu} = [276 \ 221 \ 223 \ 268 \ 234]^T$, $\boldsymbol{\Sigma} = \text{diag}(\sigma_i) = \text{diag}([105 \ 94.5 \ 31.5 \ 10.5 \ 3.5])$, ratio $r = \sigma_i/c_i = 10$, $\mathbf{C} = \boldsymbol{\Sigma}/r$.

with maximum likelihood estimated parameters $\mu_s = 0.325$ and $\sigma_s = 0.0748$. For distribution $P(d|h^1)$, similarity, namely $1-d$ (Fig. 6 (c)), can be modeled by a gamma distribution

$$P(d|h^1) = \frac{1}{\Gamma(\alpha)\beta^\alpha} (1-d)^{\alpha-1} \exp\left\{-\frac{1-d}{\beta}\right\}, \quad (23)$$

where $\Gamma(\cdot)$ is the gamma function, α and β are estimated from all the training samples: $\alpha = 125.032$, $\beta = 0.00483$. The distance method LR_D can be improved: if for each writer, $P(d|h^1)$ is estimated from the distances between the samples of that writer and all other writers (denoted by $LR_D(\text{spec.})$) [31].

4.1.1. Independent binary features

For determining LR_j , both the numerator and denominator have a simple form, $P(\mathbf{x}, \mathbf{y}|h^i) = \prod_{j=1}^n P(x_j, y_j|h^i)$. For each feature $j \in \{1, \dots, n\}$, $i \in \{0, 1\}$, $P(x_j, y_j|h^i)$ can be estimated separately. For LR_{DR} the distance factor is computed by multiplying the probability of each bit difference and rarity is modeled by a mixture of Bernoulli distributions [37]; an example of a mixture component of *th* images represented using GSC features is shown in Fig. 7.

Performance of five methods LR_j , LR_{VD} , LR_{DR} and two distance methods LR_D are given in Fig. 8(a), which shows both average error rate and computation time; in the evaluation two-third of the samples were used for training and the rest for testing.

4.1.2. Dependent binary features

Modeling full dependence between all variables is intractable, e.g., for n binary variables, LR_j needs $2(2^n - 1)$ free parameters, and LR_{DR} needs $2(2^n - 1)$ parameters. A compromise is provided by Bayesian networks (BNs), i.e., directed probabilistic graphical models, which allow including salient dependencies [19]. BN structures can be learnt using chi-squared independence tests and log-likelihood scoring [36]. For LR_j the *and* binary data is represented by 2048 binary variables (combination of object and evidence). For LR_{DR} there were two BNs of 1024 variables each for the numerator and denominator of the LLR. Performances of the three methods are in Fig. 8(b).

Assuming independence of *and* data (GSC) produced results similar to that of *th* shown in Fig. 8(a), i.e., error rates have the ordering $LR_{DR} < LR_{VD} < LR_D < LR_j$. Our interest was to see if we could improve the accuracy by removing the independence assumption, particularly using BNs to model the dependence between individual binary features, where each sample in the *and* data set contains more identifying characteristics than a sample *th* and serve as a good test set for BN learning.

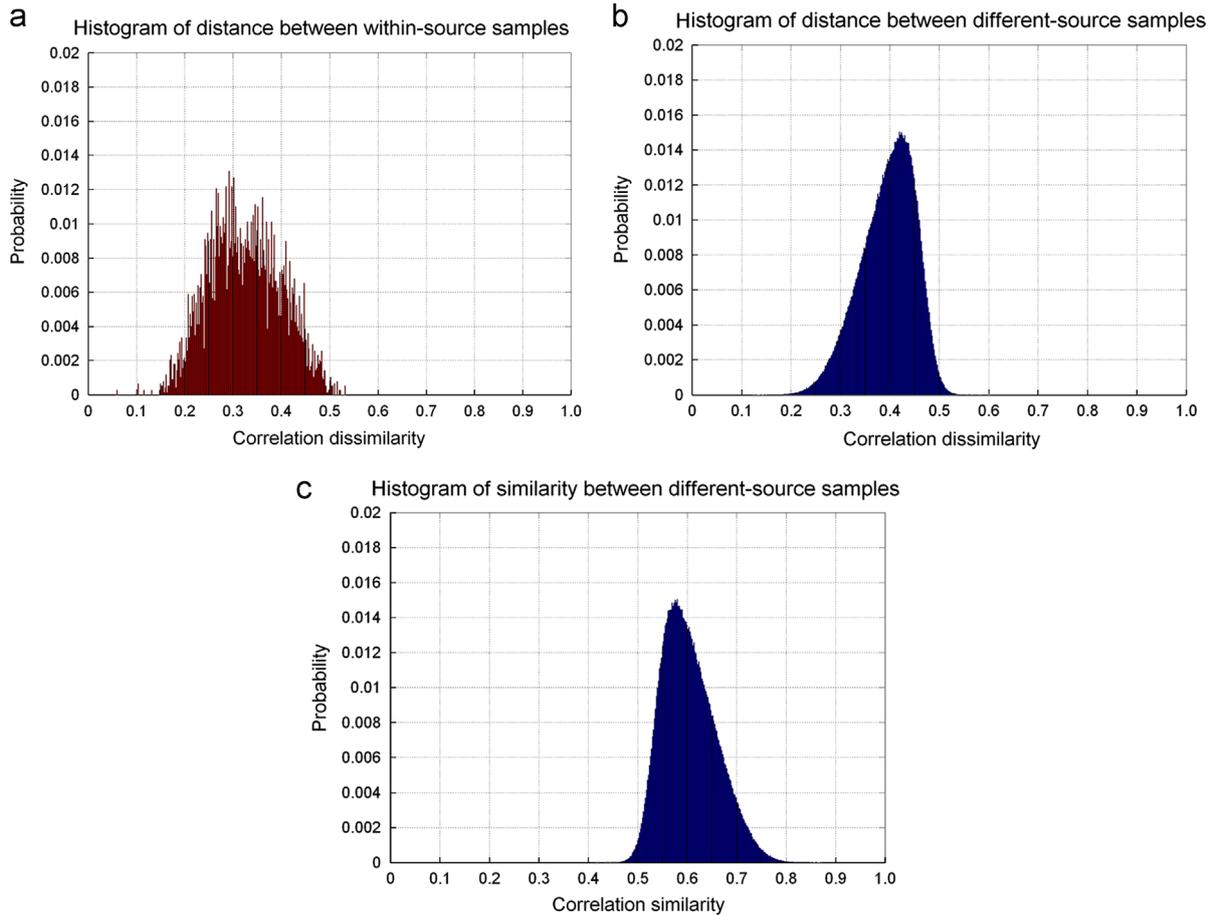


Fig. 6. Distributions of distance d between th represented by GSC features: (a) when written by the same writer (hypotheses h^0), (b) when written by different writers (hypothesis h^1), and (c) distribution of the similarity $(1-d)$ given h^1 .

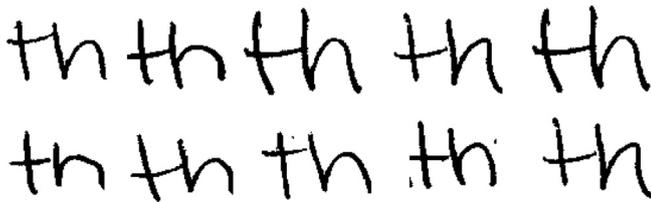


Fig. 7. Sample cluster found using a mixture of Bernoulli.

4.2. Multinomial features

Here $\mathbf{x} = [x_1, \dots, x_n]$ where x_i takes one of the n_i values $\{x_i^0, x_i^1, \dots, x_i^{n_i}\}$. As an example, characteristics of th defined by forensic document examiners [22] are given in Fig. 9, where $n=6$, $n_1=4$, $n_2=4$, $n_3=3$, $n_4=4$, $n_5=4$, $n_6=5$. As in the binary case, computing LR_j is intractable. If we assume independent x_i , we get the simple form $LR_j = \prod_{i=1}^n \prod_{k=1}^{n_i} P(x_j^k, y_j^k | h^0) / P(x_j^k, y_j^k | h^1)$. For LR_D we can use one of the several data-driven similarity measures: we consider *Lin* [20], *Occurrence Frequency* (OF) and *Goodall3* [7], abbreviated as L , O and G , respectively.

Handwriting images of th as described earlier, but with these human-determined multinomial features, x_i ($i=1, \dots, 6$), were used in this experiment. Since distance, which is continuously distributed in $[0, 1]$, is multimodal (Fig. 10), a Gaussian mixture model (GMM) [4] was used. For each distribution, the E-M algorithm [15]

together with the *Bayesian Information Criterion* (BIC) [4] was used to estimate parameters. Then the likelihood ratios with three distance measures, $LR_{D(O)}$, $LR_{D(G)}$, and $LR_{D(L)}$, were computed by the ratio of two mixture densities under alternative hypotheses h^1 ; corresponding performances are shown in the three middle columns of Fig. 8(c).

Distance and rarity method. The difference between categorical variables can be interpreted as another *categorical variable*. Suppose there is a 6-dimensional vector $(n_1, \dots, n_6)^T$, where n_k ($k=1, \dots, 6$) is the number of values taken by the k th variable x_k . Categorical values can be easily mapped to numbers so we might as well assume that x_k can take any integer between 0 and n_k-1 . Thus the difference between categorical variables x_k and y_k , $d_k(x_k, y_k)$, can be defined as follows:

$$d_k(x_k, y_k) = \begin{cases} x_k & \text{if } x_k = y_k, \\ n_k + \min(x_k, y_k) & \text{otherwise.} \end{cases} \quad (24)$$

To illustrate, suppose x_k or y_k can take values in the set $\{0, 1, 2, 3\}$, it is easy to verify that the definition in (24) can map (x_k, y_k) to integers in the range $[0, 6]$, as shown in Table 2. Now the variables d_k ($k=1, \dots, 6$) constitute the *vertices* of the Bayesian network (BN) for modeling the joint distribution of the K differences given h^0 , i.e., $P(\mathbf{d} | h^0)$, where $\mathbf{d} = \{d_1, \dots, d_6\}$. For modeling rarity, a BN consisting of a set of vertices $\mathbf{R} = \{R_k | k=1, \dots, 6\}$ is constructed, where each vertex R_k represents a tuple (x_k, y_k) , and each edge from R_k to $R_{k'}$ indicates that the variable (x_k, y_k) has a direct effect on $(x_{k'}, y_{k'})$. A BN structure learning algorithm [36] obtained the

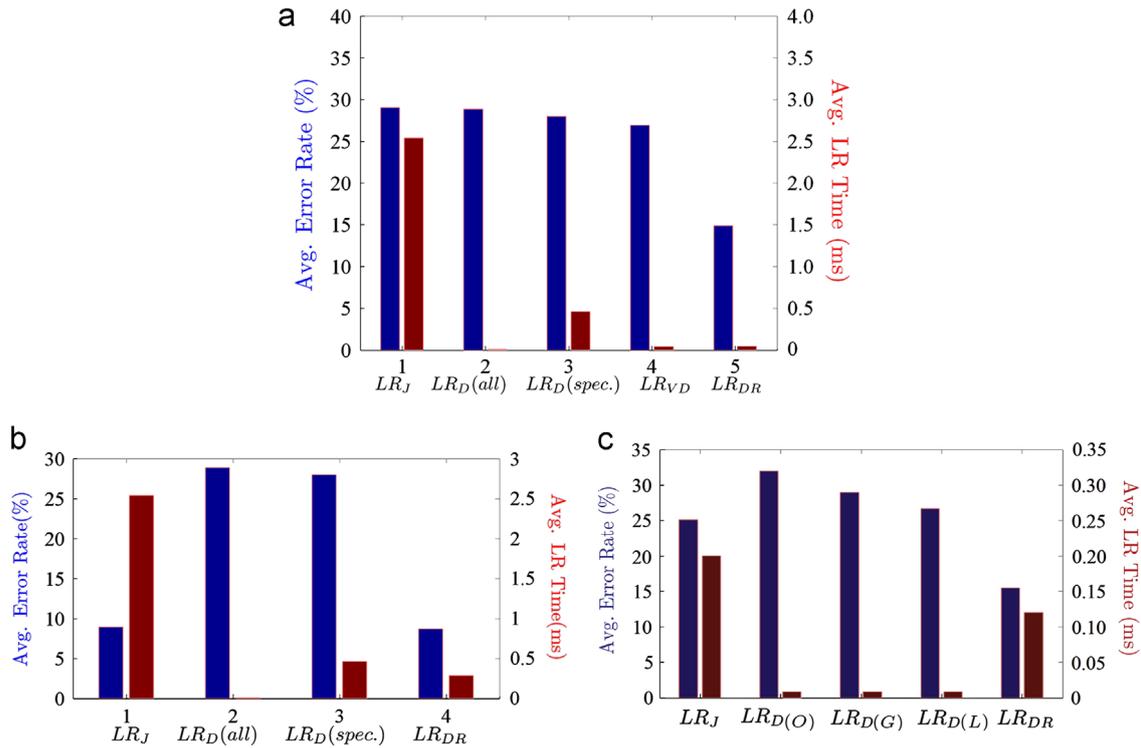


Fig. 8. Performance of likelihood ratio methods LR_J , LR_D , and LR_{DR} with discrete data (binary and multinomial). For each method left bar corresponds to average error rate and right bar to average computation time per sample: (a) binary features ($n=1024$ GSC bits) of th data, assuming statistical independence where two versions of LR_D (all and specific) and LR_{VD} are included, (b) binary features (GSC) of and data where dependence between features is modeled using Bayesian networks, and (c) multinomial features of th data, where three similarity measures O , G and L were used for LR_D and Bayesian networks for LR_{DR} .

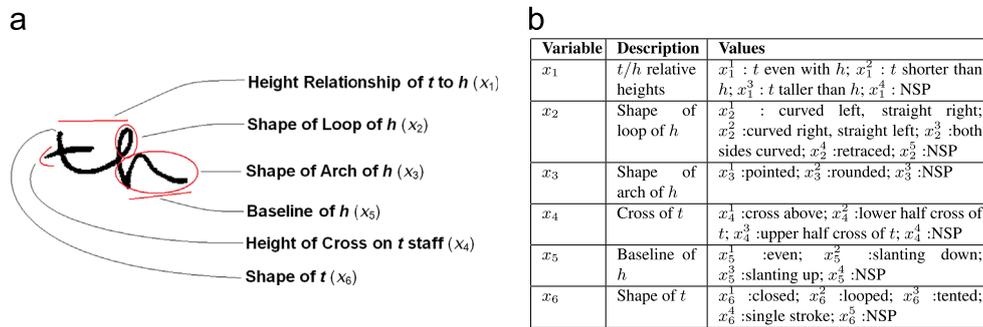


Fig. 9. Multinomial features: characteristics of th given by FDEs (human): (a) definitions of six variables, and (b) discrete values for each variable (where NSP stands for no set pattern).

BNs shown in Fig. 11(a and b). LR_{DR} is given by the ratio of the two joint probabilities given by the two networks. The resulting performance is shown in the rightmost columns of Fig. 8(c).

5. Application to graphs

5.1. Graph representation

Graphs, as a non-linear data structure, are widely used for representing objects in pattern recognition and machine learning due to several advantages over vector representations. They can model the relationships (binary) between different objects, or different parts of the same object [10].

Let \mathbf{E} be the set of features for the evidence graph, i.e., $\mathbf{E} = \{\mathbf{f}_1, \dots, \mathbf{f}_{n_1}\}$, where each feature \mathbf{f}_i ($i=1, \dots, n_1$) is a n -dimensional vector of attributes, i.e., $\mathbf{f}_i = (f_{i1}, \dots, f_{in})^T$. Similarly, let \mathbf{O} be the set of features for the object graph, i.e., $\mathbf{O} = \{\mathbf{g}_1, \dots, \mathbf{g}_{n_2}\}$, where each feature \mathbf{g}_j is represented by $\mathbf{g}_j = (g_{j1}, \dots, g_{jn})^T$, $j=1, \dots, n_2$. We can

represent \mathbf{E} and \mathbf{O} as two attributed relational graphs (ARGs) [27], where each ARG is formally defined as a three-tuple (V, R, n) where $V = \{V_i | 1 \leq i \leq n\}$ is the set of nodes and $R = \{R_{ij} | 1 \leq i, j \leq n\}$ is the set of relations between nodes. Each node has a weight and an attribute vector, $V_i = (w_i, \mathbf{v}_i)$, and each relation R_{ij} has an attribute vector \mathbf{r}_{ij} .

Graph matching is the process of finding a correspondence between the nodes and edges of two graphs that satisfies some constraints ensuring that similar substructures in one graph are mapped to similar substructures in the other [14]. Associated with inexact graph matching is a distance measure or a cost function between two graphs. Examples of graph distance measures are graph edit distance [27,9], graph distance metric based on the maximal common subgraph (MCS) [11].

For the rarity term, we need to determine the mean between object and evidence when each is represented as a graph. We first determine corresponding nodes between them, and induce two sub-graphs such that the two sub-graphs will have an equal number of nodes and edges, then construct a feature vector from

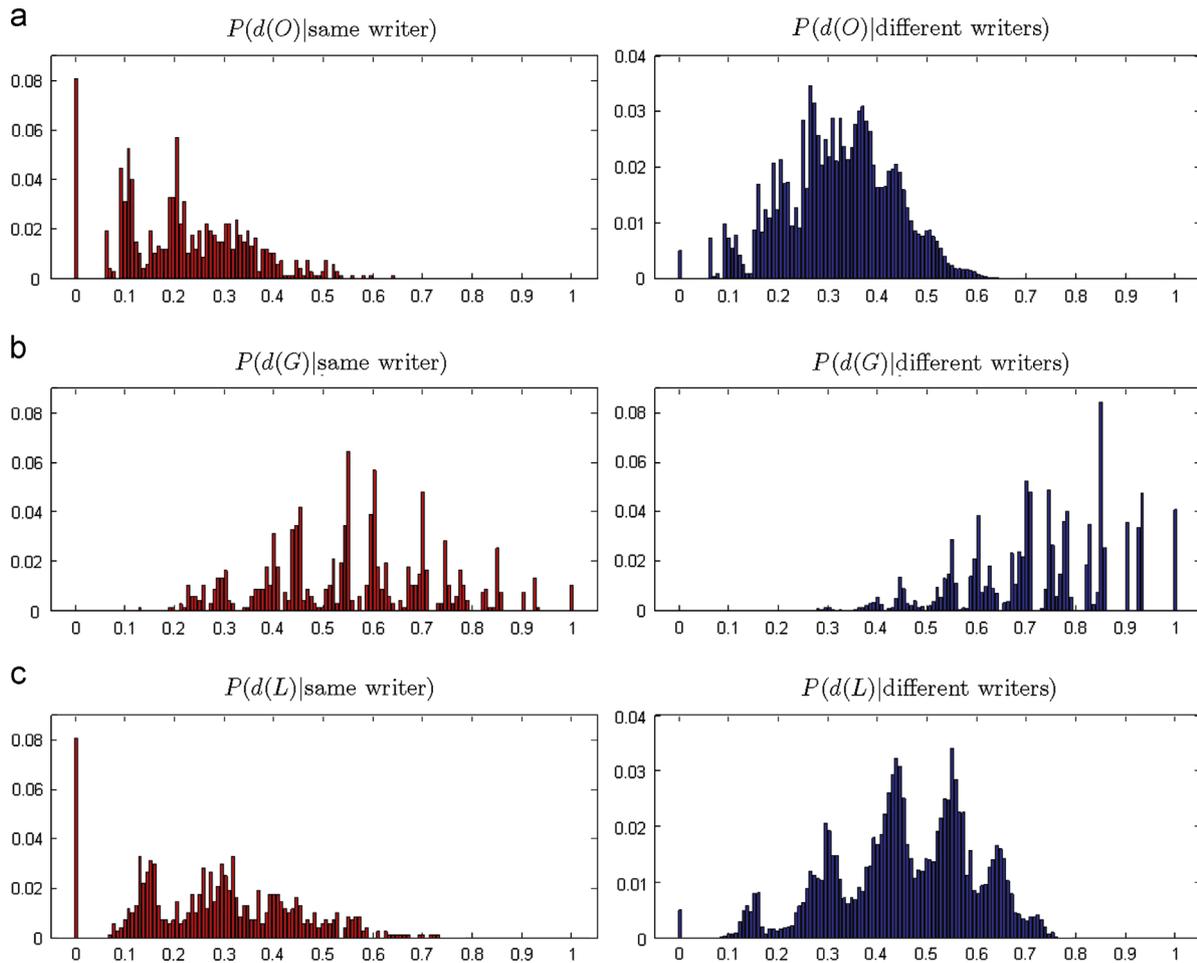


Fig. 10. Distributions of distance measures O , G and L for th , based on the six multinomial features, given hypotheses h^i ($i=0,1$).

Table 2
An example of difference between categorical variables.

x_k	y_k	$d_k(x_k, y_k)$
0	0	0
1	1	1
2	2	2
3	3	3
0	1	4
0	2	4
0	3	4
1	2	5
1	3	5
2	3	6

each of the two sub-graphs. Since the number of matching nodes is variable, evaluating the distribution of the mean requires an exponential number of parameters. It may also not be normalized. Instead we use the *magnitude* of the mean vector, which is scalar.

The scalar distance (if Euclidean distance is used) can be regarded as the magnitude of the vector difference. With analogy, we could use the distribution of the *magnitude* of the mean vector as a substitute for the distribution of the mean vector itself, i.e.,

$$LR_{DR} = \frac{P(d(\mathbf{E}, \mathbf{O})|h^0)}{P(\mathbf{m}(\mathbf{E}, \mathbf{O}))} \approx \frac{P(d(\mathbf{E}, \mathbf{O})|h^0)}{P(|\mathbf{m}(\mathbf{E}, \mathbf{O})|)} \quad (25)$$

In this way, the distributions in both numerator and denominator of the LR_{DR} have been reduced from the original high dimensional

feature space to a scalar space. This approximation has an intuitive appealing: two graphs that have more matched features will have a greater value of $|\mathbf{m}(\mathbf{E}, \mathbf{O})|$ than another two graphs that have less matched features. By mapping the distribution of mean vector $\mathbf{m}(\mathbf{E}, \mathbf{O})$ of varied length to the distribution of its magnitude, we have overcome the difficulty of defining the distribution of the difference $\mathbf{E} - \mathbf{O}$, avoided the normalization issue mentioned before, and made both numerator and denominator having the same dimension.

5.2. Application to footwear impression evidence

Footwear impressions are commonly found at crime scenes and are valuable for establishing links between suspect(s) and the crime scene. Identification involves extracting and comparing features from the impression (evidence) and a known impression (object). Each impression is represented by an *attributed relational graph* (ARG) [27,34,35] in which attributed nodes correspond to basic shapes, and attributed edges describe relationships between nodes. The Earth Mover's Distance (EMD) [25] can efficiently handle high-dimensional feature spaces and is close to perceptual similarity [25]. Derived from EMD, *footwear print distance* (FPD) [34,35] was used to compute the distance between two footwear prints represented by ARGs. Fig. 12(a–h) shows an example.

We used 1060 degraded footwear prints with ground truth to build distributions $P(d|h^i)$ ($i=0,1$), as shown in Fig. 13(a), where $P(d|h^0)$ was approximated by a Gaussian $\mathcal{N}(\mu_0, \sigma_0^2)$, and $P(d|h^1)$ was approximated by a mixture of Gaussians, where $LR_D = P(d|h^0)/P(d|h^1)$. In the

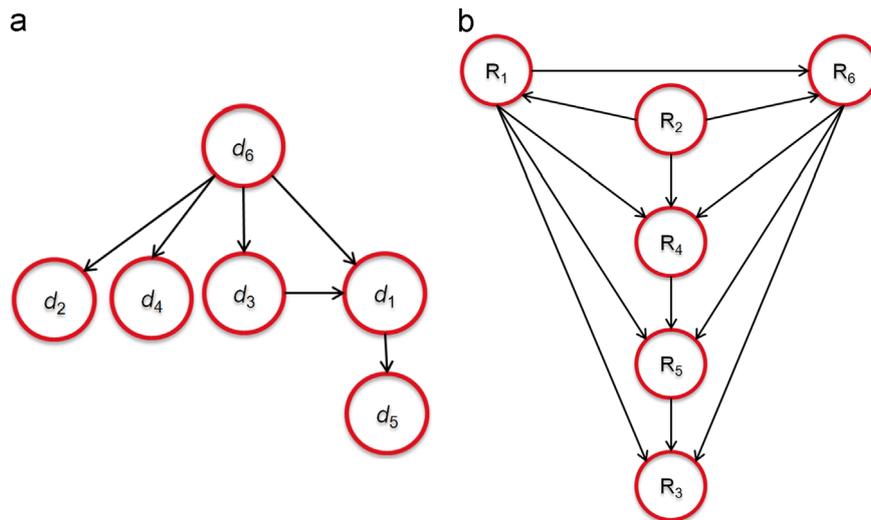


Fig. 11. Bayesian networks for LR_{DR} learnt from multinomial data: (a) difference distribution $P(d_k|h^0)$, used to evaluate the numerator, and (b) rarity distribution $P(R_k|h^0)$, used to evaluate the denominator.

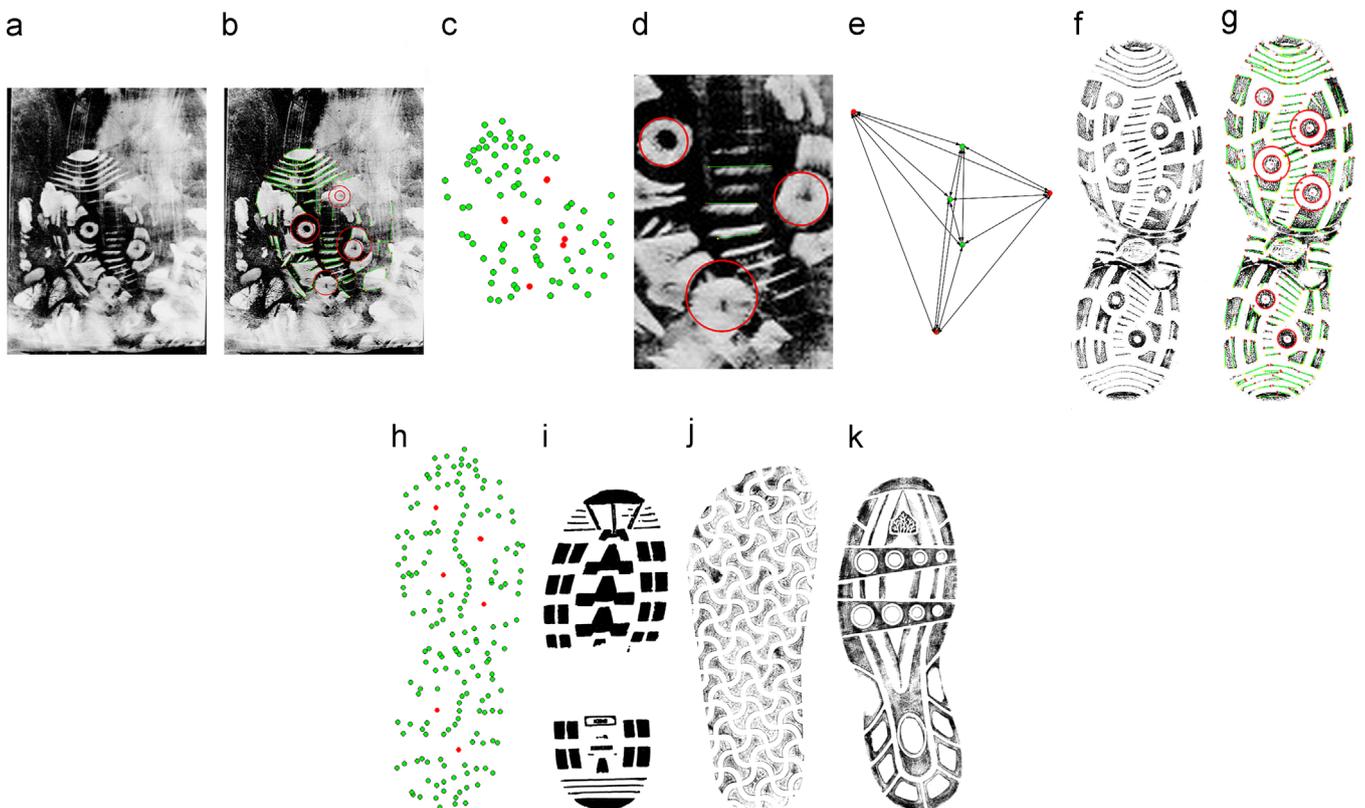


Fig. 12. Footwear impression evidence analysis using graph representation: (a) crime scene impression, (b) extracted features which are either circles (dark) or lines (light); for clarity, edges are omitted, (c) fully connected attribute relational graph (ARG) where nodes represent detected circles (dark) or lines (light); for clarity, edges are omitted, (d) enlarged sub-image with features, (e) sub-graph corresponding to sub-image, (f) a known print in database (serial number #16360), (g) features of known print, (h) ARG of (g): the distance between ARGs in (c) and (h) is 0.0835, and (i,j,k) are three other database prints illustrating common primitive shapes. Best viewed in color. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this article.)

experiments, we computed the FPD between all pairs of footwear prints in the training data set (1060 prints), retrieved matching nodes for each pair, calculated the magnitude of the mean vector of the matching features, and built the distribution $P(|\mathbf{m}(\mathbf{E}, \mathbf{O})|)$, as shown in Fig. 13(b). Two test data sets were used in evaluation: 1200 synthetic scene impressions (set **A**, generated by degrading known prints) and 200 crime scene impressions (set **B**, provided by a forensic agency along with the ground truth), where the reference data set (set **R**)

consists of 5034 known prints. Average error rate and average computation time were used to measure the accuracy and efficiency of the two methods, and the results are shown in Fig. 13(c).

6. Summary

The experimental performance and computational complexity of the three methods for estimating the likelihood ratio can now

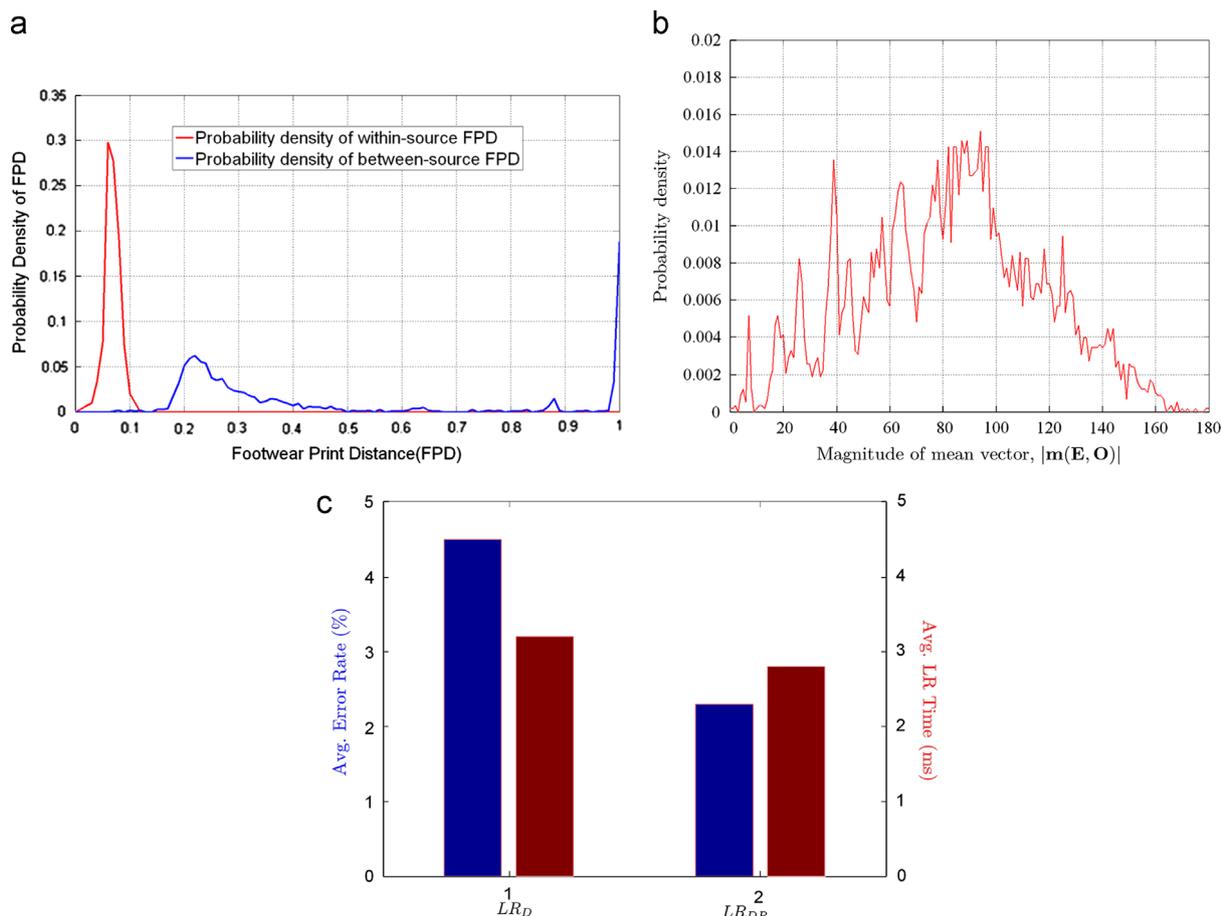


Fig. 13. Analysis of forensic data represented as a graph (footwear): (a) within-source and between-source distributions of FPD $P(d|h^i)$ ($i=0,1$) built using 1060 footwear prints, (b) distribution of the magnitude of the mean vector $\mathbf{m}(\mathbf{E}, \mathbf{O})$, i.e., $P(|\mathbf{m}(\mathbf{E}, \mathbf{O})|)$, where \mathbf{E} and \mathbf{O} represent evidence graph and object graph – distribution built with 5829 pairs of samples, and (c) average error rates of LR_D and LR_{DR} together with computation time.

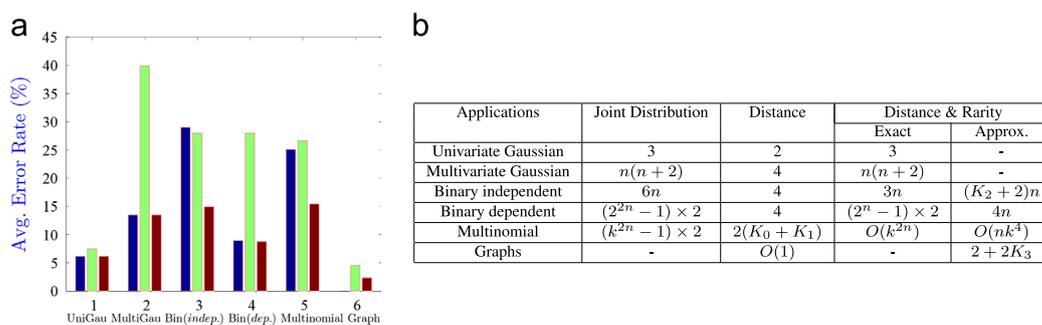


Fig. 14. Comparison of LR methods (LR_j , LR_D and LR_{DR}) over six feature types (univariate Gaussian, multivariate Gaussian, independent binary, dependent binary, multinomial and graph): (a) bars indicate error rates for different applications (for graph only LR_D and LR_{DR} are included) and (b) number of parameters needed in each method where n is the no. of variables, $\{K_i\}$, $i=0,1,2,3$ are no. of mixture components in distributions, and k is the maximum number of parents in the BN.

be summarized. We also provide an information-theoretic comparison of the methods.

6.1. Performance

Error rates for data sets of different feature types are summarized in Fig. 14(a). When features are normally distributed, LR_{DR} has the same error rate as LR_j . When feature distributions are non-Gaussian, LR_{DR} performs best. Independence assumptions make LR_j perform worse than other methods.

6.2. Model complexity of LR methods

The number of parameters needed by each method is given in Fig. 14(b). For univariate Gaussian, both LR_j and LR_{DR} need three parameters μ , τ and σ , while LR_D needs only 2 parameters τ and σ since the means are zero. For multivariate Gaussian, LR_j needs n , $n(n+1)/2$ and $n(n+1)/2$ free parameters for μ , Σ and \mathbf{C} in (20), respectively, totaling n^2+2n . Similarly for LR_{DR} in Eq. (21). LR_D needs 2 parameters for each gamma distribution $P(d|h^i)$ ($i=0,1$), totaling 4.

For independent n -dimensional binary features, LR_j needs 6 free parameters for each bit, totaling $6n$; LR_D needs only 4 parameters;

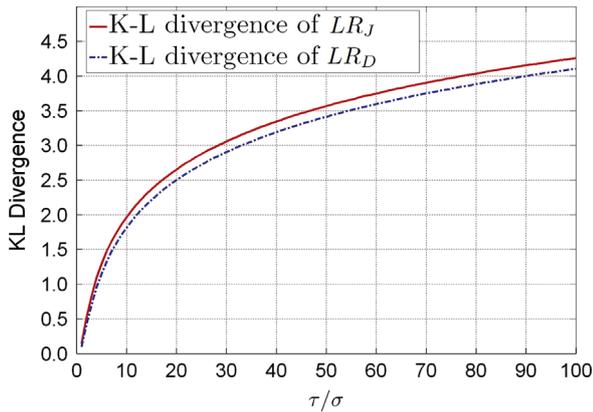


Fig. 15. Comparison of discriminative power of LR methods in univariate Gaussian case: K-L divergence of joint distribution method (red solid line) and distance method (blue dash-dotted) versus ratio τ/σ . The expectation of the difference is about 0.15, which suggests that, on average, for a single normally distributed feature, the joint distribution method is $e^{0.1532} - 1$ ($\approx 16.6\%$) more powerful in inferring the true state of hypothesis H . (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this article.)

LR_{DR} needs $2n$ and n parameters for exact computation of T_d and T_r , whereas approximating T_r using a mixture of Bernoulli needs Kn parameters, where K_2 is the number of clusters. For dependent binary features using exact approach, LR_J needs $(2^{2n} - 1) \times 2$ parameters, LR_{DR} needs $(2^n - 1) \times 2$ parameters; if using approximation with Bayesian networks, LR_{DR} needs $2n$ parameters each for T_d and T_r . For both exact and approximate approaches, LR_D needs 4 parameters.

For n -dimensional multinomial vectors with each feature taking k possible values, LR_J needs $(k^{2n} - 1) \times 2$; LR_D needs $2(K_0 + K_1)$, where K_0 and K_1 are the number of clusters for modeling $P(d|h^0)$ and $P(d|h^1)$, respectively; and the exact approach of LR_{DR} needs $O(k^n)$ for T_d , and $k^{2n} - 1$ for T_r , with the total of $O(k^n + k^{2n}) = O(k^{2n})$, while approximation using Bayesian networks reduces the number to $O(nk^4)$.

For graphs, directly computing LR_J is infeasible; LR_{DR} needs 2 parameters for modeling T_d and $2K_3$ parameters for modeling T_r , where K_3 is the number of Gaussian mixture components in $P(\mathbf{m}(\mathbf{E}, \mathbf{O}))$; LR_d needs a constant number of parameters.

6.3. Information theoretic analysis

The discriminative capability of each LR method can be analyzed using the *Kullback–Leibler (KL) divergence* between the alternative distributions. The KL divergence between two distributions $p(\mathbf{x})$ and $q(\mathbf{y})$ for the continuous case is $KL(p(\mathbf{x}) \| q(\mathbf{x})) = \int p(\mathbf{x}) \ln\{p(\mathbf{x})/q(\mathbf{x})\} d\mathbf{x}$. For discrete variables with distributions $P(\mathbf{x})$ and $Q(\mathbf{y})$ it is $KL(P \| Q) = \sum P(\mathbf{x}) \ln P(\mathbf{x})/Q(\mathbf{x})$. We can compare the KL distances between likelihood functions used by each method.

For LR_J , the KL divergence between two likelihood functions, $P(\mathbf{o}, \mathbf{e}|h^i)$ ($i=0,1$), can be interpreted as the expected *discrimination information* per sample for h^0 over h^1 . Below we present analytical results on univariate and multivariate Gaussian distributions.

6.3.1. Univariate Gaussian

From Section 2.1, we have

$$P(e, o|h^0) \sim \mathcal{N}\left(\begin{bmatrix} \mu \\ \mu \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \tau^2 \\ \tau^2 & \sigma_2^2 \end{bmatrix}\right) \quad (26)$$

and

$$P(e, o|h^1) \sim \mathcal{N}\left(\begin{bmatrix} \mu \\ \mu \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}\right), \quad (27)$$

where $\sigma_1^2 = \sigma_2^2 = \tau^2 + \sigma^2$. The joint distribution being bivariate, the KL divergence can be shown to be

$$KL(P(e, o|h^0) \| P(e, o|h^1)) = \frac{1}{2} \ln \frac{(\tau^2 + \sigma^2)^2}{\sigma^2(2\tau^2 + \sigma^2)}. \quad (28)$$

Similarly, the KL divergence between two distributions of distance $d = |o - e|$ can be computed by

$$KL(P(d|h^0) \| P(d|h^1)) = \frac{1}{2} \ln \left\{ \frac{\tau^2 + \sigma^2}{\sigma^2} \right\} - \frac{\tau^2}{2(\tau^2 + \sigma^2)}. \quad (29)$$

Denoting τ/σ by r , the difference between (28) and (29) is $\Delta KL = \frac{1}{2} \ln\{(r^2 + 1)/(2r^2 + 1)\} + r^2/(2(r^2 + 1))$, which is plotted against r in Fig. 15. Assuming a uniform distribution of r , i.e., $p(r) = U(a, b)$, the expectation of ΔKL is given by $\mathbb{E}[\Delta KL] = (1/2(b-a))F(r)|_a^b$, where $F(r) = r \ln(r^2 + 1)/(2r^2 + 1) + r + \arctan r - \sqrt{2} \arctan \sqrt{2}r$. For $a=1, b=1000, \mathbb{E}[\Delta KL] = 0.1532$, i.e., the discriminative power of LR_J is $e^{0.1532} - 1$ ($\approx 16.6\%$) greater than LR_D .

6.3.2. Multivariate Gaussian

Based on earlier discussion, we have (assuming $p=q=1$)

$$P(\mathbf{e}, \mathbf{o}|h^0) \sim \mathcal{N}\left(\begin{bmatrix} \mu \\ \mu \end{bmatrix}, \begin{bmatrix} \Sigma + \mathbf{C} & \Sigma \\ \Sigma & \Sigma + \mathbf{C} \end{bmatrix}\right), \quad (30)$$

$$P(\mathbf{e}, \mathbf{o}|h^1) \sim \mathcal{N}\left(\begin{bmatrix} \mu \\ \mu \end{bmatrix}, \begin{bmatrix} \Sigma + \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \Sigma + \mathbf{C} \end{bmatrix}\right), \quad (31)$$

where covariance \mathbf{C} takes the form of $\text{diag}(c_i^2)$ as features are independent; similarly $\Sigma = \text{diag}(\sigma_i^2)$.

The KL divergence between (30) and (31) is

$$KL(P(\mathbf{e}, \mathbf{o}|h^0) \| P(\mathbf{e}, \mathbf{o}|h^1)) = \frac{1}{2} \ln \left\{ \frac{\prod_{i=1}^n (\sigma_i^2 + c_i^2)^2}{\prod_{i=1}^n (2\sigma_i^2 + c_i^2)c_i^2} \right\} + \frac{n}{2}.$$

For the distance method based on the Euclidean distance, it reduces to

$$KL(P(d|h^0) \| P(d|h^1)) = \ln \left\{ \frac{c}{\sqrt{\sigma^2 + c^2}} \right\}. \quad (32)$$

Then the difference between the K-L divergences of the two methods can be shown to be

$$\Delta KL = \frac{n}{2} \ln \left\{ \frac{(r^2 + 1)^2}{2r^2 + 1} \right\} + \frac{1}{2} \ln\{r^2 + 1\} + \frac{n}{2},$$

where $r = \sigma/c$, which indicates that when r is fixed, the difference between the K-L divergence of the two methods increases approximately *linearly* w.r.t. the number of features n .

Under the Gaussian assumption $LR_{DR} = LR_J$. Thus under the Gaussian assumption LR_{DR} discriminates better than LR_D . For all three methods, as the number of features n increases, the discriminative capability improves. However, as n increases, the difference (i.e., the gap) between their discriminative capabilities increases. The non-Gaussian case has been addressed experimentally in the earlier discussion.

7. Conclusion

We have considered the problem of estimating the likelihood ratio in forensic scenarios. For a given set of features exact evaluation involves computing a likelihood ratio under alternative hypotheses of identification and exclusion. The joint distribution method for non-Gaussian distributed features is computationally intractable (requiring $O(2^{2n})$ parameters) and impossible even when the number of variables n is moderately large. A commonly used simplification is the distance method, which reduces the number of parameters from $O(2^{2n})$ to $O(1)$ by approximating the

ratio of the original joint distributions by the ratio of the distributions of the distance, and inevitably leads to severe information loss and degraded accuracy. However, the joint distribution method for Gaussian distributed features has a closed-form solution which can be decomposed into the product of two factors, one based on the distance/difference and the other on rarity (distribution of the mean of object and evidence), where each factor has only n variables. This decomposition also has the intuitive value that human examiners look for rare features. We have explored generalizing this formulation with non-Gaussian distributed features arising in several forensic scenarios. Since the distance and rarity method still has complexity of $O(2^n)$, approximation using Bayesian networks has been proposed. Experiments with six different data sets (univariate Gaussian, multivariate Gaussian, independent binary, dependent binary, multinomial, and graphs) show that the distance and rarity method performs consistently better than the commonly used distance method in terms of accuracy, while being more efficient than the joint distribution method.

Conflict of interest

None declared.

References

- [1] <http://hypertextbook.com/facts/2001/JuliaSherlis.shtml>.
- [2] C. Aitken, F. Taroni, *Statistics and the Evaluation of Evidence for Forensic Scientists*, second edition, Wiley, 2004.
- [3] C.G.G. Aitken, D. Lucy, Evaluation of trace evidence in the form of multivariate data, *Journal of the Royal Statistical Society* 53 (1) (2004) 109–122.
- [4] C.M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [5] W.J. Bodziak, *Footwear Impression Evidence: Detection, Recovery and Examination*, second edition, CRC Press, 2000.
- [6] A. Bolck, C. Weyermann, L. Dujourdy, P. Esseiva, J. Berg, Different likelihood ratio approaches to evaluate the strength of evidence of MDMA tablet comparisons, *Forensic Science International* 191 (1) (2009) 42–51.
- [7] S. Boriah, V. Chandola, V. Kumar, Similarity measures for categorical data: a comparative evaluation, in: *SIAM Data Mining Conference*, 2008, pp. 243–254.
- [8] N. Brummer, J. du Preez, Application independent evaluation of speaker detection, *Computer Speech & Language* 20 (2–3) (2006) 230–275.
- [9] H. Bunke, On a relation between graph edit distance and maximum common subgraph, *Physical Review Letters* 18 (8) (1997) 689–694.
- [10] Horst Bunke, Christophe Irniger, Michel Neuhaus, Graph matching—challenges and potential solutions, in: *ICIAP'05*, 2005, pp. 1–10.
- [11] Horst Bunke, Kim Shearer, A graph distance metric based on the maximal common subgraph, *Pattern Recognition Letters* 19 (3–4) (1998) 255–259.
- [12] G. Casella, R. Berger, *Statistical Inference*, Duxbury, 2001.
- [13] C.F. Tippett, V.J. Emerson, F. Lawton, S.M. Lampert, The evidential value of the comparison of paint flakes from sources other than vehicles, *Journal of the Forensic Science Society* 8 (1968) 61–65.
- [14] Donatello Conte, Pasquale Foggia, Carlo Sansone, Mario Vento, Thirty years of graph matching in pattern recognition, *International Journal of Pattern Recognition and Artificial Intelligence* 18 (3) (2004) 265–298.
- [15] A.P. Dempster, N.M. Laird, D.B. Rubin, Maximum likelihood from incomplete data via the EM algorithm, *Journal of the Royal Statistical Society, Series B (Methodological)* 39 (1) (1977) 1–38.
- [16] C.T.J. Dodson, Y. Oba, W.W. Sampson, On the Distributions of Mass, Thickness and Density in Paper, 2000.
- [17] I.W. Evett, J.A. Lambert, J.S. Buckleton, A Bayesian approach to interpreting footwear marks in forensic casework, *Science and Justice* 38 (4) (1998) 241–247.
- [18] M.A.T. Figueiredo, A.K. Jain, Unsupervised learning of finite mixture models, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24 (2002) 381–396.
- [19] D. Koller, N. Friedman, *Probabilistic Graphical Models: Principles and Techniques*, MIT Press, 2009.
- [20] D. Lin, An information-theoretic definition of similarity, in: *Proceedings of ICML '98*, 1998, pp. 296–304.
- [21] D.V. Lindley, A problem in forensic science, *Biometrika* 64 (2) (1977) 207–213.
- [22] R.J. Muehlberger, K.W. Newman, J. Regent, J.G. Wichmann, A statistical examination of selected handwriting characteristics, *Journal of Forensic Sciences* 22 (1977) 206–210.
- [23] K. Nandakumar, Y. Chen, S.C. Dass, A.K. Jain, Likelihood ratio-based biometric score fusion, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 30 (2) (2008) 342–347.
- [24] C. Neumann, C. Champod, R. Puch-Solis, N. Egli, A. Anthonioz, A. Bromage-Griffiths, Computation of likelihood ratios in fingerprint identification for configurations of any number of minutiae, *Journal of Forensic Sciences* 52 (1) (2007) 54–64.
- [25] Y. Rubner, C. Tomasi, L.J. Guibas, The Earth Mover's Distance as a metric for image retrieval, *International Journal of Computer Vision* 40 (2000) 99–121.
- [26] J. Saks, R. Kohler, The individualization fallacy in forensic science evidence, *Vanderbilt Law Review* (2007). Available at: http://works.bepress.com/michael_saks/1.
- [27] A. Sanfeliu, K.S. Fu, A distance measure between attributed relational graphs for pattern recognition, *IEEE Transactions on Systems, Man, and Cybernetics* 13 (1983) 353–362.
- [28] B. Scholkopf, The kernel trick for distances, *Advances in Neural Information Processing Systems* 13 (2001).
- [29] S.N. Srihari, *Computational Methods for Handwritten Questioned Document Examination*, National Criminal Justice Research Report, 2010.
- [30] S.N. Srihari, Evaluating the rarity of handwriting formations, in: *Proceedings of the International Conference on Document Analysis and Recognition (ICDAR)*, 2011, pp. 618–622.
- [31] S.N. Srihari, G. Ball, Comparison of statistical models for writer verification, in: *Document Recognition and Retrieval*, 2009, pp. 1–10.
- [32] S.N. Srihari, C. Huang, H. Srinivasan, On the discriminability of the handwriting of twins, *Journal of Forensic Sciences* 53 (2) (2008) 430–446.
- [33] S.N. Srihari, H. Srinivasan, Comparison of ROC and likelihood decision methods in automatic fingerprint verification, *International Journal of Pattern Recognition and Artificial Intelligence (IJPRAI)* 22 (3) (2008) 535–553.
- [34] Y. Tang, S.N. Srihari, H. Kasiviswanathan, Similarity and clustering of footwear prints, in: *Proceedings of the International Conference on Granular Computing*, 2010, pp. 459–464.
- [35] Y. Tang, S.N. Srihari, H. Kasiviswanathan, J.J. Corso, Footwear print retrieval system for real crime scene marks, in: *IWCF*, 2010, pp. 88–100.
- [36] Yi Tang, Sargur N. Srihari, An efficient algorithm to learn Bayesian networks for multinomials, in: *Proceedings of International Conference on Pattern Recognition*, Tsukuba, Japan, 2012.
- [37] Yi Tang, Sargur N. Srihari, Handwriting individualization using distance and rarity, in: *Proceedings of Document Retrieval and Recognition*, San Francisco, CA, 2012.
- [38] F. Taroni, C.G.G. Aitken, P. Garbolino, A. Biedermann, *Bayesian Networks and Probabilistic Inference in Forensic Science*, Wiley, 2006.
- [39] M. Taylor, et al., *Latent Print Examination and Human Factors: Improving the Practice through a Systems Approach*, NIST, 2002.
- [40] L. Wasserman, *All of Statistics: A Concise Course in Statistical Inference*, Springer, 2003.
- [41] B. Zhang, S.N. Srihari, Binary vector dissimilarity measures for handwriting identification, in: *SPIE, Document Recognition and Retrieval X*, 2003.

Yi Tang obtained his Ph.D. in Computer Science from University at Buffalo, The State University of New York in 2012. He received his BE in Computer Science from Beijing University of Posts and Telecommunications, China, in 2006. His research interests include pattern recognition, machine learning, computer vision and computational forensics.

Srihari is a computer scientist who has made scientific contributions in the areas of pattern recognition, machine learning, and computational forensics. He is a SUNY Distinguished Professor in the Department of Computer Science and Engineering at the University at Buffalo, The State University of New York.

Srihari is the founding director of the Center of Excellence for Document Analysis and Recognition (CEDAR) at the University at Buffalo. Work at the center led to the first automated systems for reading handwritten postal addresses which were deployed in the United States, Australia and the United Kingdom.

Srihari served on the Board of Scientific Counselors of the National Library of Medicine for six years (2001–2007) and the National Academy of Sciences Committee on Identifying the Needs of the Forensic Science Community (2006–2009).

Srihari's honors include Fellow of the Institute of Electronics and Telecommunications Engineers (IETE, India) in 1992, Fellow of the Institute of Electrical and Electronics Engineers (IEEE) in 1995, Fellow of the International Association for Pattern Recognition in 1996 and distinguished alumnus of the Ohio State University College of Engineering in 1999.

Srihari received a B.Sc. in Physics and Mathematics from the Bangalore University (National College) in 1967, a B.E. in Electrical Communication Engineering from the Indian Institute of Science, Bangalore in 1970, and a Ph.D. in Computer and Information Science from the Ohio State University, Columbus in 1976.