

Combining Evidence using Likelihood Ratios in Writer Verification

Sargur Srihari, Dmitry Kovalenko, Yi Tang and Gregory Ball
 CEDAR, Department of Computer Science and Engineering
 University at Buffalo, Buffalo NY 14260-2500, USA
 srihari@cedar.buffalo.edu

ABSTRACT

Forensic identification is the task of determining whether or not observed evidence arose from a known source. 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). In LR- based decision methods, particularly handwriting comparison, a variable number of input evidences is used. A decision based on many pieces of evidence can result in nearly the same LR as one based on few pieces of evidence. We consider methods for distinguishing between such situations. One of these is to provide confidence intervals together with the decisions and another is to combine the inputs using weights. We propose a new method that generalizes the Bayesian approach and uses an explicitly defined discount function. Empirical evaluation with several data sets including synthetically generated ones and handwriting comparison shows greater flexibility of the proposed method.

1. INTRODUCTION

Forensic identification concerns whether observed evidence arose from a known source. The probabilistic approach is to determine the likelihood ratio positive ($LR+$)¹⁻⁶ whose numerator is the joint probability of the evidence and source under the null, or *identification*, hypothesis that the evidence arises from the source and the denominator is the joint probability under the alternate, or *exclusion*, 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.

In this paper we consider the problem of combining LRs provided by several different pieces of evidence. The problem is illustrated with examples from handwriting comparison and a scoring approach is proposed.

2. LIKELIHOOD RATIO

Let $W = \{w_i\}$ be a set of individual writers. Let K be a writing sample from a known writer w_i . Let Q be a random sample from an unknown writer w_j , i.e., questioned writing. We can state two opposing hypotheses:

h^0 : K and Q are from the *same* writer ($i = j$), the *identification* hypothesis; and

h^1 : K and Q are from *different* writers ($i \neq j$), the *exclusion* hypothesis. The task is to determine the probability of identification (and its complement, exclusion).

Assuming that K and Q are represented by sets of measurements \mathbf{k} and \mathbf{q} respectively, we can define two joint probability distributions $P(\mathbf{k}, \mathbf{q}|h^0)$ and $P(\mathbf{k}, \mathbf{q}|h^1)$ which specify as to how often each instance of the known and questioned occur together when they belong to the same writer or to different writers. The relative strengths of evidence supporting the two hypotheses is quantified by the LR

$$LR_J = LR(\mathbf{k}, \mathbf{q}) = \frac{P(\mathbf{k}, \mathbf{q}|h^0)}{P(\mathbf{k}, \mathbf{q}|h^1)}. \quad (1)$$

The corresponding log-likelihood ratio, $LLR(\mathbf{k}, \mathbf{q}) = \ln P(\mathbf{k}, \mathbf{q}|h^0) - \ln P(\mathbf{k}, \mathbf{q}|h^1)$, 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*.

*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.

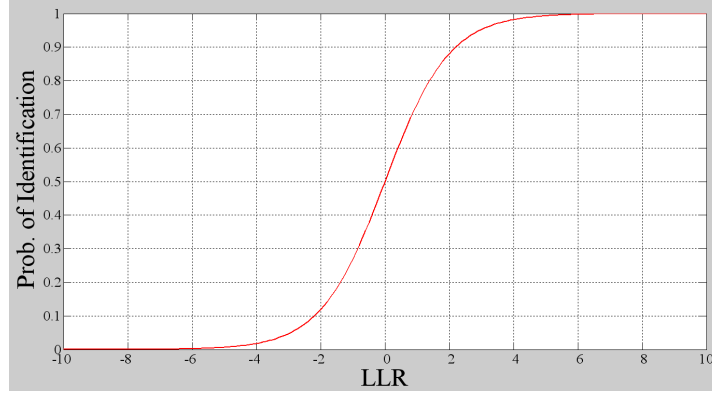


Figure 1. Probability of identification is a sigmoid function of the log-likelihood ratio.

Using a Bayesian formulation, we can convert LRs into probabilities of identification and exclusion. Let the *prior* probabilities of the hypotheses be $P(h^0)$ and $P(h^1)$ with

$$P(h^0) + P(h^1) = 1. \quad (2)$$

Defining the prior odds as $O_{prior} = \frac{P(h^0)}{P(h^1)}$, we can express the prior probability of the same source as $P(h^0) = O_{prior}/(1 + O_{prior})$. The prior odds can be converted into posterior odds as $O_{posterior} = \frac{P(h^0|\mathbf{k}, \mathbf{q})}{P(h^1|\mathbf{k}, \mathbf{q})} = O_{prior} \times LR(\mathbf{k}, \mathbf{q})$. Thus we can write the posterior probability of the same source as $P(h^0|\mathbf{k}, \mathbf{q}) = O_{posterior}/(1 + O_{posterior})$, or equivalently

$$P(h^0|\mathbf{k}, \mathbf{q}) = \frac{P(h^0) \cdot LR(\mathbf{k}, \mathbf{q})}{P(h^1) + P(h^0) \cdot LR(\mathbf{k}, \mathbf{q})} \quad (3)$$

The particular case of equal priors is of interest in forensics, as opinion without prior bias. In this case we get a simple form for the probability of identification as

$$P(h^0|\mathbf{k}, \mathbf{q}) = \frac{LR(\mathbf{k}, \mathbf{q})}{1 + LR(\mathbf{k}, \mathbf{q})} = \frac{\exp(LLR(\mathbf{k}, \mathbf{q}))}{1 + \exp(LLR(\mathbf{k}, \mathbf{q}))} = \sigma[LLR(\mathbf{k}, \mathbf{q})]. \quad (4)$$

where σ is the *sigmoid* function $\sigma(a) = \frac{1}{1 + e^{-a}}$. The probability of exclusion is $P(h^1|\mathbf{k}, \mathbf{q}) = 1 - P(h^0|\mathbf{k}, \mathbf{q}) = 1 - \sigma[LLR(\mathbf{k}, \mathbf{q})]$,

The probability of identification with respect to the LLR follows a sigmoid function as shown in Figure 1. Since this function reaches either 1 or zero very quickly, the LLR has to be computed quite precisely for the probability of identification to provide meaningful discrimination. Determining the LR or LLR as defined by Eq. 1 requires the distributions $P(\mathbf{k}, \mathbf{q}|h^i)$ ($i = 0, 1$), defined over all possible values of objects and their evidential forms. If \mathbf{k} and \mathbf{q} consist of n measurements (characteristics of writing) each, with each measurement taking r possible values, then $2r^{2n}$ parameters are needed to specify LR_j defined by Eq. 1. Since this is exponential in n , which can be quite high when writing samples are compared, determining these distributions is computationally and statistically infeasible.

3. DISTANCE APPROXIMATION

A common solution to the difficulty of determining the joint distribution in forensics is to use the probability of distance, or similarity, between the evidence and known instead of the joint probability.^{4,7} The distance-based approach assumes that $\mathbf{k} = [k_1, \dots, k_n]$, and $\mathbf{q} = [q_1, \dots, q_n]$ are points in a fixed-dimensional vector-space and that there exists a scalar *distance* $D(\mathbf{k}, \mathbf{q})$ between them. A likelihood ratio is defined as follows

$$LR_D = \frac{P(D(\mathbf{k}, \mathbf{q})|h^0)}{P(D(\mathbf{k}, \mathbf{q})|h^1)}. \quad (5)$$

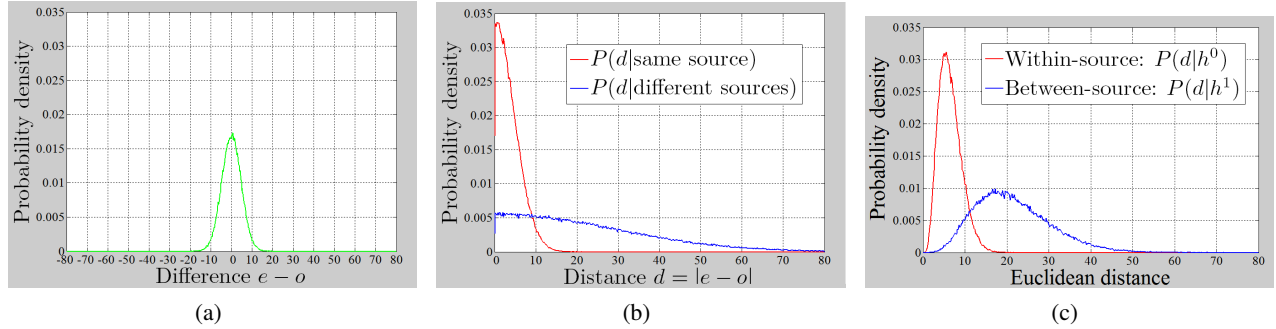


Figure 2. Distribution of distance when features are Gaussian: (a) univariate samples - distribution of the signed difference: $P(e - o|h^0) = \mathcal{N}(0, 2\sigma^2)$; (b) univariate samples- absolute difference, $d = |e - o|$, follows a folded-normal distribution, and (c) multivariate 5-dimensional Gaussian samples- distances follow Gamma distributions.

The number of parameters needed to characterize the probability distributions in Eq. 5 is independent of the number of measurements. While LR_D , which has a *constant* number of parameters, is simple to compute, there is a severe loss of information in going from probabilities of points in $2n$ -dimensional space to probabilities of a scalar value. However, due to its simplicity, the distance method is been widely used in fingerprint identification,⁴ speaker verification, pharmaceutical tablet comparison,⁸ etc.

The performance of the distance based method, when thresholded to give a yes-no decision, is worse than both the joint distribution method and a method that takes into account a rarity factor.⁶ A compromise solution⁹ assumes that $D(\mathbf{k}, \mathbf{q})$ consists of a set of features $F = \{f_i\}_{i=1}^n$ where $f_i = D_i(k_i, q_i)$ are scalar distances along a set of n measurements. If we assume that the features (distances) are independent we can write

$$LR_F(\mathbf{k}, \mathbf{q}) = \prod_i LR(f_i) \quad (6)$$

where $LR(f_i) = \frac{P(f_i|h^0)}{P(f_i|h^1)}$ whose numerator and denominator are probabilities of distances along individual measurements, conditioned on the respective hypotheses. Since this method has more degrees of freedom than LR_D , it performs better than a scalar distance⁶ †. It also has the advantage over scalar distance that the measurements need not form a vector space, be commensurate, etc.

We can condition the hypotheses explicitly on the features and rewrite Eq. 2 as: $P(h^0|F) + P(h^1|F) = 1$, and Eq. 3 as:

$$P(h^0|F) = \frac{P(h^0) \cdot \prod_i LR(f_i)}{P(h^1) + P(h^0) \cdot \prod_i LR(f_i)}. \quad (7)$$

When Q and K are characterized by a single measurement that is Gaussian, $P(q - k|h^0) \sim \mathcal{N}(0, (a\sigma)^2)$, the signed distance is Gaussian as shown in Fig. 2(a). The absolute difference $d = |q - k|$ follows a *folded normal distribution* as shown in Fig. 2(b). In the multivariate case, we take measurements on samples drawn from the object and evidence, where each measurement is a *vector* of k real values. When the Euclidean distance is used to measure the distance between two points, the distributions $P(d|h^i)$ ($i = 0, 1$) can be represented by Gamma distributions⁹ as as shown in Fig. 2(c).

4. THE SCORING PROBLEM

Consider two simple problems for identification, where the priors are equal, i.e., $P(h^0) = P(h^1)$. In the first problem a single feature is to be used for identification/exclusion, whereas in the second problem nine features available, with the following associated LR values:

†The performance of the three methods LR_J , LR_D and LR_F depends on the relative magnitude of between-source variance and within-source variance. In the case of multivariate Gaussian measurements with five features, when the ratio r of between-source variance and within-source variance is 10, the error rates of the three methods were 0.035, 0.04 and 0.085 respectively. As the ratio becomes larger, the difference between LR_J and LR_F becomes increasingly small. The method LR_{VD} for vector-distance in⁶ applies to both independent features and dependent features. If assuming independent features (distances), then LR_{VD} exactly equals LR_F . Yet, when features are correlated, results of LR_F and LR_{VD} were similar.

1. A single feature $F = f_1$ with $LR(f_1) = 96$.
2. Nine features $F = \{f_i\}_{i=1}^9$ with $\{LR(f_i), i = 1, \dots, 9\} = \{3, 4, 2, \frac{1}{4}, 2, 2, \frac{1}{3}, 6, 2\}$.

In both cases, by applying Eq.7, we get: $P(h^0|F) = \frac{96}{97}$. In the first case the result is based on one feature (a discriminating characteristic) and could be happenstance, whereas in the second case the support of h^0 is more trustworthy since it is based on nine features.

As another example, consider the use of the distance approach to two different handwriting comparison problems shown in Figures 3 and 4. In the first case only a single word is being compared. In the second case, a whole paragraph (address block) is being compared. The total LLR score in both cases is quite similar. We would like the second case to produce a stronger result than the first since there is much more evidence.

In the rest of this paper we consider methods for distinguishing between such cases. This problem is especially important when probabilities of features are estimated approximately.

5. METHODS FOR RELATIVE SCORING OF HYPOTHESES

One approach is to use not a single numerical value for posterior probability of h^0 but a credible interval.¹⁰ When there are no parameters of h^0 then we can assume that there is a distribution of likelihood ratios and estimate a credible interval for it. A crucial part of this approach is to define a proper prior probability for likelihood ratios.

Another way of dealing with the problem at hand is to use statistical hypothesis testing:¹¹ to define likelihood ratio confidence interval with a given level of confidence. If we assume a log-normal distribution with unknown mean and variance as a distribution of likelihood ratios than we can calculate confidence interval for log-likelihood ratio mean and use it in equation 7.

A completely different approach to the problem is to use a score instead of posterior probability of h^0 .¹² This approach allows us not to make any assumptions about likelihood ratio distribution. We have designed such scoring method based on Bayesian approach with weighted likelihoods and a discount function that diminishes the score when features in a comparison have a small sum of weights.

The problem is to get a score for hypothesis h^0 that reflects how likely h^0 is true in comparison with its alternative and that explicitly includes how trustworthy a hypotheses comparison is.

Here we assume that all features F are independent from each other. A common way to solve the problem is to use Bayesian approach (see Eq. 7). To adapt its ideas in order to explicitly include a hypotheses comparison trustworthiness we discuss several existing methods and suggest a new one.

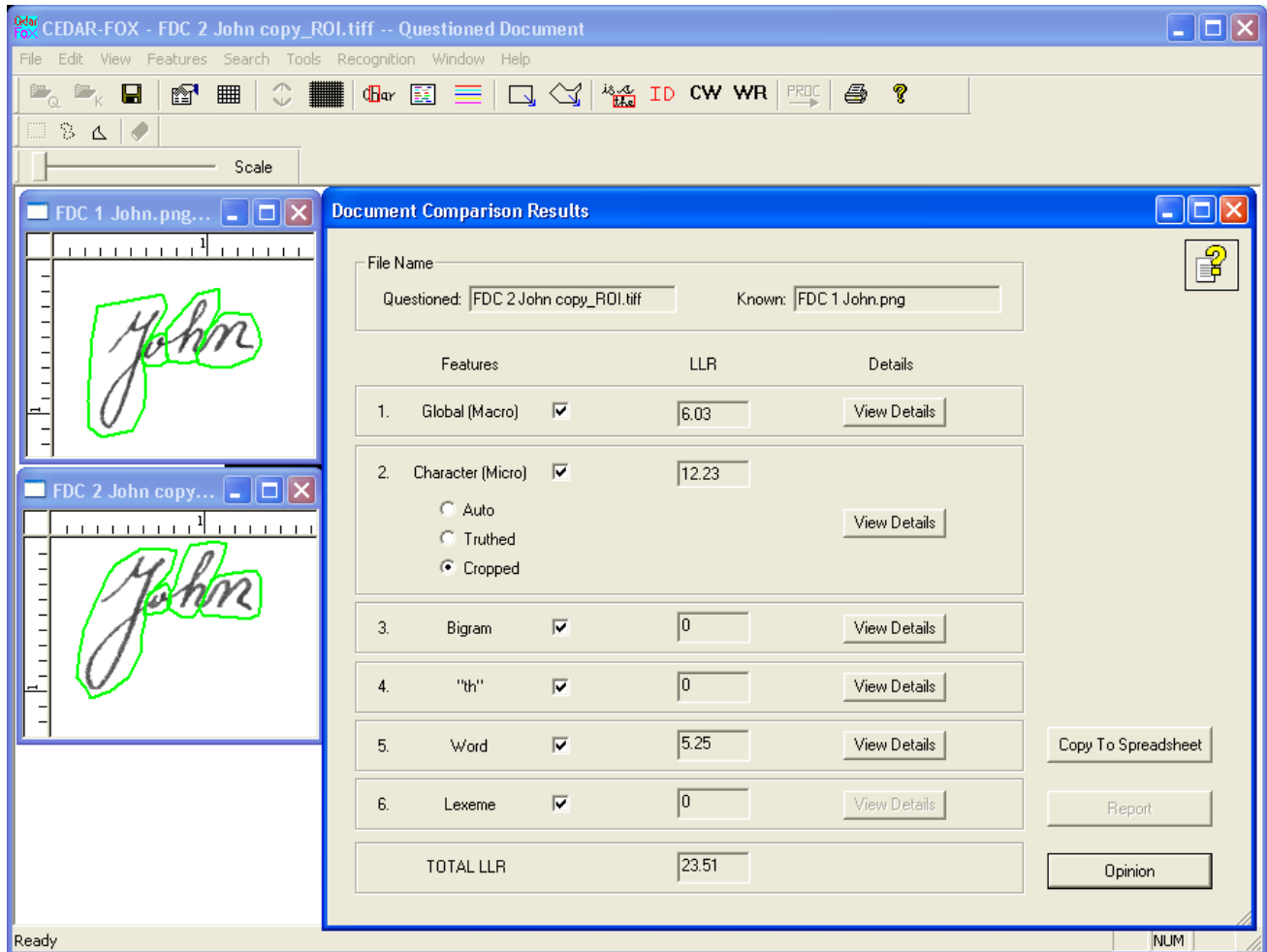
Assume that LR's of features are themselves i.i.d. samples of some distribution. For example, we can assume that it is a log-normal distribution. Then it is possible to identify a credible interval:¹³ an $\alpha \cdot 100\%$ credible interval is a set C such that $P(C) = \alpha$. There are different ways to evaluate credible intervals approximately or exactly when it is possible.¹³ Values of Eq. 7 determined on credible interval boundaries together with level α define a result for the stated problem.

An analogy of a credible interval in statistical hypothesis testing is a confidence interval.¹¹ It is an interval $(L_\alpha(F), R_\alpha(F))$ evaluated from a given set of features F , that frequently includes the parameter of interest, if the experiment is repeated. The frequency here is determined by confidence level α . For example, the parameter of interest could be mean of LLRs if they comply with the log-normal distribution. Values of Eq. 7 determined on confidence interval boundaries along with confidence level α define a solution for the stated problem.

A useful property of using credible or confidence intervals is that in most cases with increase in the number of features the distance between interval boundaries is decreasing. A problem of using credible or confidence intervals is a necessity of combining interval boundaries and level α into one score if only one single value is required by the application (as a score for hypothesis h^0).

To drop any assumptions on likelihood ratio distributions, we have designed a new method to solve the problem. To include credibility of the hypotheses comparison into the score we suggest the following expression instead of the posterior probability ratio:

$$S \left(\frac{h^0|F}{h^1|F} \right) = \frac{P(h^0)}{P(h^1)} \cdot \left(\prod_{i=1}^n \left(\frac{P(f_i|h^0)}{P(f_i|h^1)} \right)^{w_i} \right)^{d(\sum_i w_i)} \quad (8)$$



(a)

Macro Features

Macro Features	Enable	Q-Feat	K-Feat	Difference	Log Same	Log Diff	LLR
Entropy	<input type="checkbox"/>	0	0	0	0	0	0
Threshold	<input type="checkbox"/>	0	0	0	0	0	0
Black Pixels	<input type="checkbox"/>	0	0	0	0	0	0
Ext. Cont.	<input checked="" type="checkbox"/>	1	1	0	-1.47	-2.5	1.03
Int. Cont.	<input checked="" type="checkbox"/>	0	3	3	-2.13	-2.5	0.37
H. Slope	<input checked="" type="checkbox"/>	0.16	0.2	0.03	2.41	2.29	0.12
P. Slope	<input checked="" type="checkbox"/>	0.37	0.36	0.01	3.42	1.31	2.11
V. Slope	<input checked="" type="checkbox"/>	0.36	0.33	0.02	2.94	1.93	1.01
N. Slope	<input checked="" type="checkbox"/>	0.09	0.09	0	4.58	2.04	2.54
Stroke Width	<input type="checkbox"/>	0	0	0	0	0	0
Avg. Slant	<input checked="" type="checkbox"/>	23.45	24.3	0.84	-1.04	-3.49	2.45
Avg. Height	<input checked="" type="checkbox"/>	0	0	0	-2.02	-3.01	0.99
Avg. Word Gap	<input checked="" type="checkbox"/>	0	0	0	-3.51	1.08	-4.59
Sum:					3.18	-2.85	6.03

(b)

Micro (character) Features

Type of Character: Manually Cropped

Enter Character ID: [Display Character](#)

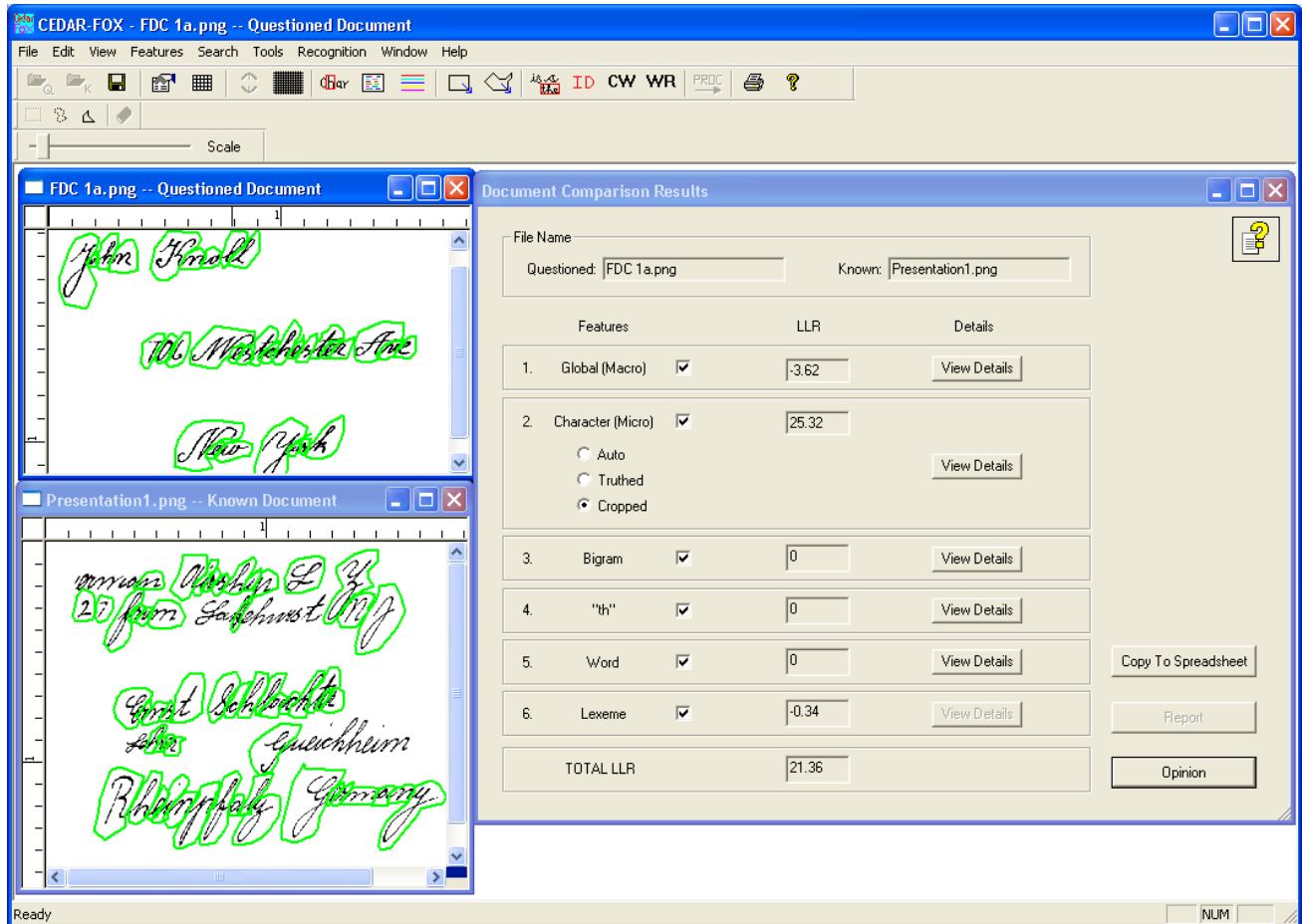
Enabled	No. in Q	No. in K	No. Dist	SCORE	Enabled	No. in Q	No. in K	No. Dist	SCORE	Enabled	No. in Q	No. in K	No. Dist	SCORE			
<input type="checkbox"/>	0	0	0	0.00	<input type="checkbox"/>	l	0	0	0	0.00	<input type="checkbox"/>	F	0	0	0	0.00	
<input type="checkbox"/>	1	0	0	0.00	<input type="checkbox"/>	m	0	0	0	0.00	<input type="checkbox"/>	G	0	0	0	0.00	
<input type="checkbox"/>	2	0	0	0.00	<input checked="" type="checkbox"/>	n	1	1	2	4.10	<input type="checkbox"/>	H	0	0	0	0.00	
<input type="checkbox"/>	3	0	0	0.00	<input checked="" type="checkbox"/>	o	1	1	3	1.39	<input type="checkbox"/>	I	0	0	0	0.00	
<input type="checkbox"/>	4	0	0	0.00	<input type="checkbox"/>	p	0	0	0	0.00	<input checked="" type="checkbox"/>	J	1	1	4	3.73	
<input type="checkbox"/>	5	0	0	0.00	<input type="checkbox"/>	q	0	0	0	0.00	<input type="checkbox"/>	K	0	0	0	0.00	
<input type="checkbox"/>	6	0	0	0.00	<input type="checkbox"/>	r	0	0	0	0.00	<input type="checkbox"/>	L	0	0	0	0.00	
<input type="checkbox"/>	7	0	0	0.00	<input type="checkbox"/>	s	0	0	0	0.00	<input type="checkbox"/>	M	0	0	0	0.00	
<input type="checkbox"/>	8	0	0	0.00	<input type="checkbox"/>	t	0	0	0	0.00	<input type="checkbox"/>	N	0	0	0	0.00	
<input type="checkbox"/>	9	0	0	0.00	<input type="checkbox"/>	u	0	0	0	0.00	<input type="checkbox"/>	O	0	0	0	0.00	
<input type="checkbox"/>	a	0	0	0.00	<input type="checkbox"/>	v	0	0	0	0.00	<input type="checkbox"/>	P	0	0	0	0.00	
<input type="checkbox"/>	b	0	0	0.00	<input type="checkbox"/>	w	0	0	0	0.00	<input type="checkbox"/>	Q	0	0	0	0.00	
<input type="checkbox"/>	c	0	0	0.00	<input type="checkbox"/>	x	0	0	0	0.00	<input type="checkbox"/>	R	0	0	0	0.00	
<input type="checkbox"/>	d	0	0	0.00	<input type="checkbox"/>	y	0	0	0	0.00	<input type="checkbox"/>	S	0	0	0	0.00	
<input type="checkbox"/>	e	0	0	0.00	<input type="checkbox"/>	z	0	0	0	0.00	<input type="checkbox"/>	T	0	0	0	0.00	
<input type="checkbox"/>	f	0	0	0.00	<input type="checkbox"/>	A	0	0	0	0.00	<input type="checkbox"/>	U	0	0	0	0.00	
<input type="checkbox"/>	g	0	0	0.00	<input type="checkbox"/>	B	0	0	0	0.00	<input type="checkbox"/>	V	0	0	0	0.00	
<input checked="" type="checkbox"/>	h	1	1	1	3.01	<input type="checkbox"/>	C	0	0	0	0.00	<input type="checkbox"/>	W	0	0	0	0.00
<input type="checkbox"/>	i	0	0	0.00	<input type="checkbox"/>	D	0	0	0	0.00	<input type="checkbox"/>	X	0	0	0	0.00	
<input type="checkbox"/>	j	0	0	0.00	<input type="checkbox"/>	E	0	0	0	0.00	<input type="checkbox"/>	Y	0	0	0	0.00	
<input type="checkbox"/>	k	0	0	0.00							<input type="checkbox"/>	Z	0	0	0	0.00	

Total Number of Distances: 10
 Total LLR Score: 12.23

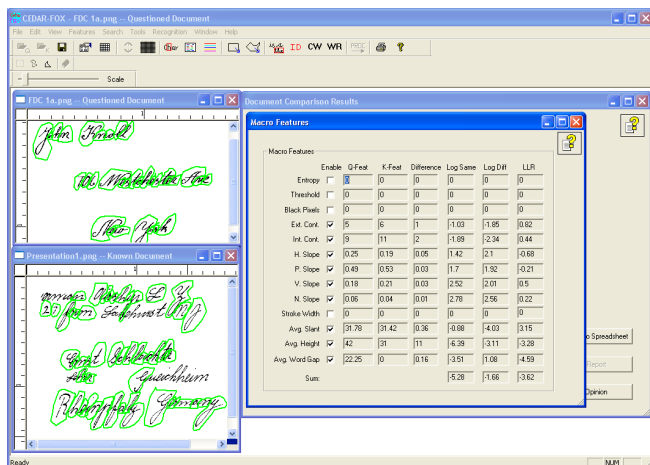
[OK](#)

(c)

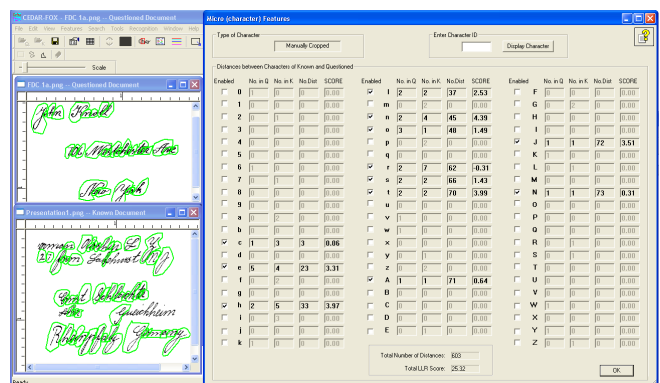
Figure 3. LLRs of a single word "John": (a) total LLR is 23.51 which includes macro LLR of 6.03, micro LLR of 12.23 and word LLR of 5.25, (b) breakdown over 13 macro features that sum to 6.03, (c) micro feature LLRs (62) that sum to 12.23.



(a)



(b)



(c)

Figure 4. LLRs of an address block: (a) showing total LLR score of 21.36. (b) macro feature LLRs that sum to -3.62, and (c) micro feature LLRs that sum to 25.32.

where: $S \left(\frac{h^0|F}{h^1|F} \right)$ is the relative score of hypothesis h^0 against h^1 . The term

$$\prod_{i=1}^n \left(\frac{P(f_i|h^0)}{P(f_i|h^1)} \right)^{w_i} \quad (9)$$

is a weighted likelihood ratio,¹⁴ where weights $w_i \in \mathbb{R}^+$. A trivial set of weights is:

$$w_i = 1, \quad i = \overline{1, n} \quad (10)$$

which defines that all of the used features are equal in their contribution to the overall value. A more discriminating approach is to use higher weights for more accurate and influential features.

Function $d(\sum_i w_i) : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is a discount function which reflects a credibility of the comparison. It ought to give a support to a comparison with big sum of features weights and discount the score for a comparison with the small sum. A trivial discount function is:

$$d_0 \left(\sum_i^n w_i \right) = \frac{n}{\sum_i^n w_i} \quad (11)$$

which does not make any discount but normalizes weighting scheme. An example of discount function which makes a discount for a small sum of weights is:

$$d_1 \left(\sum_i^n w_i \right) = \frac{n}{\max(m, \sum_i^n w_i)} \quad (12)$$

where: m is a constant threshold. When trivial weights are used and a number of features is more than a threshold then score S is equal to the ratio of posterior probabilities defined in equation 7.

Another example of discount function applies different discounts for different values of sum of weights:

$$d_2 \left(\sum_i^n w_i \right) = \beta \cdot \log \left(\sum_i^n w_i \right) \quad (13)$$

where: β – is a magnitude parameter.

Discounted weighted likelihood ratio (DWLR) is non-negative since original LRs are non-negative:

$$\forall f_i : \frac{P(f_i|h^0)}{P(f_i|h^1)} \geq 0 \Rightarrow \forall \{w_i : w_i \in \mathbb{R}^+\}, \forall \{d(\sum_i w_i) : \mathbb{R}^+ \rightarrow \mathbb{R}^+\} : \left(\prod_{i=1}^n \left(\frac{P(f_i|h^0)}{P(f_i|h^1)} \right)^{w_i} \right)^{d(\sum_i w_i)} \geq 0 \quad (14)$$

Substituting DWLR to equation 7 instead of LR we get a posterior score for hypothesis h^0 :

$$P^s(h^0|F) = \frac{P(h^0) \cdot \left(\prod_{i=1}^n \left(\frac{P(f_i|h^0)}{P(f_i|h^1)} \right)^{w_i} \right)^{d(\sum_i w_i)}}{P(h^1) + P(h^0) \cdot \left(\prod_{i=1}^n \left(\frac{P(f_i|h^0)}{P(f_i|h^1)} \right)^{w_i} \right)^{d(\sum_i w_i)}} \quad (15)$$

Range of $P^s(h|F)$ is $[0, 1]$ since DWLR is non-negative and $P(h)$ is prior probability of hypothesis h . Moreover, with trivial weights (equation 10) and trivial discount function (equation 11) equation 15 coincide with equation 7. Hence, posterior probability $P(h|F)$ is a special case of posterior score $P^s(h|F)$.

Accuracy of using score S instead of posterior probabilities ratio can be at least the same since the ratio is a special case of S (with trivial weights and the trivial discount function). Another advantage of using score S is absence of additional assumptions in comparison with using credible or confidence intervals. The disadvantage however is that in general it results in a score for the hypotheses of interest in comparison with its alternative, but not in the exact posterior probabilities ratio.

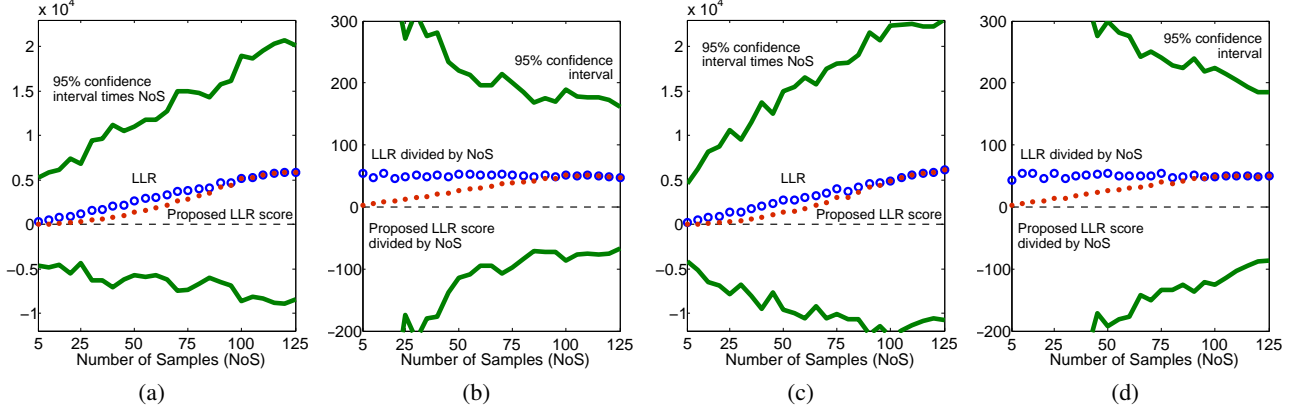


Figure 5. Scores on normal distribution $\mathcal{N}(50, 25)$ (a-b) and uniform distribution $\mathcal{U}[0, 100]$ (c-d): i) Green solid lines – 95% confidence intervals for the log-likelihood ratio mean; ii) Blue circles – log-likelihood ratio; iii) Red dots – discounted weighted log-likelihood ratio with trivial weights (eq. 10) and discount function d_1 (eq. 12) and threshold $m = 100$.

6. EXPERIMENTAL RESULTS

We evaluated the evidence combination method in three test scenarios: (i) synthetic data, (ii) simple problem of Section 4: one with a single feature and the second with nine, both with the same LLR value, and (iii) handwriting verification: comparison of single word and comparison of address blocks (Figs. 3 and 4).

6.1 Test Scenario 1: Synthetic Data

Synthetically created data sets were used since it allows comparing methods on different (initially specified) sample distributions and with a wide range of data sets sizes. Two types of synthetically generated data sets were used:

- I) Data set consists of values derived from normal distribution $\mathcal{N}(50, 25)$ with mean $\mu = 50$ and variance $\sigma^2 = 25$.
- II) Data set consists of values derived from uniform distribution $\mathcal{U}[0, 100]$ defined on interval $[0, 100]$.

Data sets of each type consisted of 5 through 125 samples with step size of 5. Every sample in each data set represents a log-likelihood ratio value for one feature.

Throughout testing we use a discrete uniform prior: $P(h^0) = P(h^1) = 0.5$. Hence, we focus on a score for log-likelihood ratio (LLR) rather than on the posterior score for the hypotheses h^0 . The following methods were used in the comparison:

- I) Original Bayesian approach where score is defined as posterior probabilities ratio of the hypotheses.
- II) An approach based on estimation of confidence interval for LLR distribution mean. We assume that likelihood ratios have a normal distribution $\mathcal{N}(\mu, \sigma^2)$ with unknown mean and variance.
- III) Discounted weighted likelihood ratio (DWLR) with trivial weights (see equation 10) and discount function d_1 (see equation 12) and threshold $m = 100$.

We did not include an approach based on credible intervals since if we assume that log-likelihood has a normal distribution with unknown mean and precision and use a normal-gamma distribution as a conjugate prior then (with certain parameters) credible and confidence intervals coincide.¹⁵

Results on data sets generated from normal distribution $\mathcal{N}(50, 25)$ are shown in figures 5(a) and 5(b). Results on data sets generated from uniform distribution $\mathcal{U}[0, 100]$ are shown in figures 5(c) and 5(d). Blue circles show log-likelihood ratio (method I); green solid lines and blue circles on figures 5(b) and 5(d) show 95% confidence intervals and log-likelihood ratio average accordingly (method II); red dots show discounted weighted log-likelihood ratio (method III).

In both sets of tests method III converges to log-likelihood ratio for large numbers of samples and gives a discounted value for small number of them. In particular, the score average is close to zero for less than 10 samples while LLR average is the same as for large number of samples (see figures 5(b) and 5(d)).

The obtained results show that there are the following disadvantages of using confidence intervals: i) a necessity to make an assumption about LLR distribution; ii) confidence interval boundaries are wide (even for 95% confidence level) and especially in case when real data distribution is different from the assumed one (see figure 5(c-d)); iii) there is no common way to combine confidence interval boundaries and confidence level into one score (see figures 5(a),5(c): if we use confidence interval boundaries times the number of samples then the result scores become even more distant).

The advantages of using DWLR are: i) it generalizes Bayesian approach and allows to explicitly set a discount for hypotheses comparison based on a few features; ii) no need to make additional assumptions on likelihood ratios distribution; iii) flexibility of choosing weights and discount function allows tuning the approach for every particular application. However, there are disadvantages also: i) the likelihood score results not in a posterior probability but in a relative score; ii) the flexibility of choosing weights and a discount function makes it non-trivial to find the best settings for a particular application.

6.2 Test Scenario 2: Simple Problems

Let us refer to the example from Section 4: comparison of two cases (in the first case one feature is given and in the second - nine features). Applying the Bayesian approach we get:

$$P_1(h^0|F) = P_2(h^0|F) = 0.99 \quad (16)$$

where: $P_1(h^0|F)$ – posterior probability of h^0 in the first case, which is the same as posterior probability $P_2(h^0|F)$ of h^0 in the second case.

Table 1. Posterior scores for hypothesis h^0 in example from Section 4 obtained using trivial weights and different discount functions.

Discount function:	$P_1^s(h^0 F)$	$P_2^s(h^0 F)$	$ P_2^s(h^0 F) - P_1^s(h^0 F) $
$d_0 (\sum_i^n w_i) = \frac{n}{\sum_i^n w_i}$	0.99	0.99	0
$d_1 (\sum_i^n w_i) = \frac{n}{\max(m, \sum_i^n w_i)}$	0.56	0.89	0.33
$d_2 (\sum_i^n w_i) = \beta \cdot \log (\sum_i^n w_i)$	0.50	0.88	0.38

Posterior probability scores for the first and the second cases obtained using DWLR with trivial weights and discount functions d_0, d_1 with threshold $m = 20$, d_2 with magnitude parameter $\beta = 0.2$ are shown in Table 1. As one can see from Table 1 posterior scores with trivial weights and trivial discount function d_0 coincide with posterior probabilities. The use of DWLR with discount function different from the trivial one enables distinguishing between the two cases (using one feature for comparison or using nine features) though total LLR is the same in these cases.

6.3 Test Scenario 3: Handwriting Comparison

Let us now apply the proposed DWLR to evaluate posterior score for examples in Figures 3 and 4. The total $LLR(\mathbf{k}, \mathbf{q})$ for comparison of two words (figure 3) is 23.51. Using Eq. 4 we get the following value of the posterior probability $P_{word}(h^0|\mathbf{k}, \mathbf{q})$ of hypothesis h^0 (K and Q are from the *same* writer):

$$P_{word}(h^0|\mathbf{k}, \mathbf{q}) = 1 - 6 \cdot 10^{-11} \quad (17)$$

Posterior probability of h^0 in the second example (figure 4) with comparison of address blocks is:

$$P_{block}(h^0|\mathbf{k}, \mathbf{q}) = 1 - 5 \cdot 10^{-10} \quad (18)$$

The difference between them seems to be quite small:

$$P_{word}(h^0|\mathbf{k}, \mathbf{q}) - P_{block}(h^0|\mathbf{k}, \mathbf{q}) = 4.7 \cdot 10^{-10} \quad (19)$$

Table 2. Posterior scores for hypothesis h^0 in comparison of two words (Figure 3) and address blocks (Figure 4) obtained using trivial weights and different discount functions.

Discount function:	$P_{word}^s(h^0 \mathbf{k}, \mathbf{q})$	$P_{block}^s(h^0 \mathbf{k}, \mathbf{q})$	$ P_{block}^s(h^0 \mathbf{k}, \mathbf{q}) - P_{word}^s(h^0 \mathbf{k}, \mathbf{q}) $
$d_0 (\sum_i^n w_i) = \frac{n}{\sum_i^n w_i}$	$1 - 6 \cdot 10^{-11}$	$1 - 5 \cdot 10^{-10}$	$4.7 \cdot 10^{-10}$
$d_1 (\sum_i^n w_i) = \frac{n}{\max(m, \sum_i^n w_i)}$	0.964	0.991	0.027
$d_2 (\sum_i^n w_i) = \beta \cdot \log (\sum_i^n w_i)$	0.957	0.965	0.008

Note that the comparison of words (Figure 3) is based on 14 features and only one word of 4 letters. Whereas the second comparison (Figure 4) is based on 22 features and several words in address block.

Posterior scores for h^0 in the first (Figure 3) and the second comparisons (Figure 4) obtained using DWLR with trivial weights and discount functions d_0 , d_1 with threshold $m = 100$, d_2 with magnitude parameter $\beta = 0.05$ are shown in Table 2. As one can see from Table 2 posterior scores with trivial weights and trivial discount function d_0 coincide with posterior probabilities.

Posterior scores $P_{block}^s(h^0|\mathbf{k}, \mathbf{q})$, $P_{word}^s(h^0|\mathbf{k}, \mathbf{q})$ obtained using trivial weights and discount functions other than the trivial one are smaller in comparison to posterior probabilities $P_{block}(h^0|\mathbf{k}, \mathbf{q})$, $P_{word}(h^0|\mathbf{k}, \mathbf{q})$ which reflects that small number of used features (small number of characters in compared documents) leads to weaker results. Larger difference between posterior scores (see last column in Table 2) reflects that address block comparison with several words turns out to be more discriminative than one word comparison although original posterior probabilities difference is as small as $4.7 \cdot 10^{-10}$.

7. SUMMARY AND CONCLUSION

The probabilistic approach to the comparison of known and questioned handwriting samples has been considered. In the formulation the posterior score of the identification hypothesis is determined based on prior beliefs and an evidence set (observed features). A common way of using posterior probability as the score has a drawback. Nearly the same result can be based on many features and on few of them, which can be important to distinguish. We consider methods for distinguishing between such situations. One of these is to provide confidence intervals for the obtained result and another is to combine features using weights and a discount function.

The usage of intervals has the following problem: there is no common way to combine interval boundaries and the corresponding level into one score. The designed method has the following advantages: i) flexibility (it can be tuned for every particular application); ii) accuracy of the designed approach at least the same as of Bayesian approach; iii) there are no assumptions on likelihood distribution (in comparison to approaches based on credible or confidence intervals). Although the scoring method does not in general provide probabilities as output (which is a disadvantage), it distinguishes the relative strengths of having different numbers of features. The advantage of the method has been illustrated using several data sets including synthesized data and handwriting comparison cases.

Acknowledgement

The work was supported in part by Award No. 2010-DN-BX-K037 of the NIJ, Office of Justice Programs, U.S. DoJ. The opinions expressed are those of the author(s) and do not necessarily reflect those of the DoJ.

REFERENCES

- [1] I. W. Evett, J. A. Lambert, and J. S. Buckleton, "A Bayesian approach to interpreting footwear marks in forensic casework," *Science and Justice* **38**(4), pp. 241 – 247, 1998.
- [2] C. Aitken and F. Taroni, *Statistics and the Evaluation of Evidence for Forensic Scientists*, Wiley, second ed., 2004.
- [3] F. Taroni, C. G. G. Aitken, P. Garbolino, and A. Biedermann, *Bayesian networks and probabilistic inference in forensic science*, Wiley, 2006.

- [4] C. Neumann, C. Champod, R. Puch-Solis, N. Egli, A. Anthonioz, and A. Bromage-Griffiths, "Computation of likelihood ratios in fingerprint identification for configurations of any number of minutiae," *Journal of Forensic Sciences* **52**(1), pp. 54–64, 2007.
- [5] S. N. Srihari, "Computational methods for handwritten questioned document examination," *National Criminal Justice Research Report* , 2010.
- [6] Y. Tang and S. N. Srihari, "On evaluating the probability of identification in forensic science," *Submitted* , 2012.
- [7] S. N. Srihari and H. Srinivasan, "Comparison of ROC and likelihood decision methods in automatic fingerprint verification," *International Journal of Pattern Recognition and Artificial Intelligence (IJPRAI)* **22**(3), pp. 535–553, 2008.
- [8] A. Bolck, C. Weyermann, L. Dujourdy, P. Esseiva, and J. Berg, "Different likelihood ratio approaches to evaluate the strength of evidence of MDMA tablet comparisons," *Forensic Science International* **191**(1), pp. 42–51, 2009.
- [9] S. N. Srihari, C. Huang, and H. Srinivasan, "On the discriminability of the handwriting of twins," *Journal of Forensic Sciences* **53**(2), pp. 430–446, 2008.
- [10] P. Congdon, *Bayesian Statistical Modelling*, Wiley, USA, 2007.
- [11] D. Cox and D. Hinkley, *Theoretical Statistics*, Chapman and Hall/CRC, USA, 1979.
- [12] D. Kovalenko and S. N. Srihari, "On methods for incorporating evidences into posterior scoring of hypotheses," in *Proc. International Conference on Pattern Recognition*, 2012.
- [13] G. Antelman, A. Madansky, and R. McCulloch, *Elementary Bayesian Statistics*, Cheltenham, UK, 1997.
- [14] X. Wang, "Asymptotic properties of maximum weighted likelihood estimators," *Journal of Statistical Planning and Inference* **119**(1), pp. 37–54, 2004.
- [15] K. Murphy, *Conjugate Bayesian analysis of the Gaussian distribution*, UBC, 2007.