

On Computing Strength of Evidence for Writer Verification

Shrivardhan Kabra

Harish Srinivasan

Chen Huang

Sargur Srihari

Center of Excellence for Document Analysis and Recognition (CEDAR)

Department of Computer Science and Engineering, University at Buffalo

{skabra, hs32, chuang5, srihari}@cedar.buffalo.edu

Abstract

The problem of Writer Verification is to make a decision of whether or not two handwritten documents are written by the same person. Providing a strength of evidence for any such decision is an integral part of the writer verification problem. The strength of evidence should incorporate (i) The amount of information compared in each of the two documents (line/half page/full page etc.), (ii) The nature of content present in the document (same/different content), (iii) Features used for comparison and (iv) The error rate of the model used for making the decision. This paper describes the statistical model used for writer verification and also introduces a mathematical formulation to include the above four mentioned parameters, for calculating strength of evidence of same/different writer. The statistical model uses Gamma and Gaussian densities to parametrically model the distance space distribution arising from comparing ensemble of pairs of documents. The strength of evidence is mapped to a 9-point qualitative scale for the decision; one that is often used by questioned document examiners. Experiments and results show that with increase in information content from just a single word to a full page of document, the verification accuracy of the model increases.

1 Introduction

Writer verification is the task of determining whether two handwriting samples were written by the same or by different writers, a task of importance in Questioned Document Examination (QDE). A statistical model for this task was described in [1]. The model uses two distance space distributions resulting from comparing features from ensemble of pairs of documents written by (i) same writer and (ii) different writers. It was shown that a Gamma density is a better fit than Gaussian densities for modeling these distributions. Although a Gamma is a better fit (in terms

of Kullback-Leibler(KL) divergence) to one of the distribution(same), it is about the same for the other distribution(different). Firstly this paper shows that, modeling the same distribution with Gamma and the different distribution with Gaussian is better. Secondly, this paper focuses on a method to provide a strength of evidence of same/different that incorporates (i) The amount of information compared in each of the two documents(line/half page/full page etc.), (ii) The nature of content present in the document (same/different content), (iii) Features used for comparison and (iv) The error rate of the model used for making the decision. The first two of the parameters above were not included in framework of strength of evidence in [1]. The strength of evidence is then calibrated into a standard nine-point scale used by questioned document examiners.

2 Writer Verification Model

The three major components of the verification model include (i) Identifying discriminating elements (ii) Mapping from feature to distance space by using similarity measure and (iii) Parametric modeling of the distance space distribution. Each of the three is briefly described below.

2.1 Features - Discriminating elements

Features for writer verification have been split into Macro (global) and Micro (Local) features in [2] and features are termed as discriminating elements (DE [3]). The statistical model to be described can be used with any other set of features as well. Macro features(13 of them) are gray-scale based (entropy, threshold, no. of black pixels), contour based (external and internal contours), slope-based (horizontal, positive, vertical and negative), stroke-width, slant and height. These are real valued features. Micro features are for characters, bi-grams(two characters) and words. A set of 512 binary-valued micro-features corresponding to gradient (192 bits), structural (192 bits), and

concavity (128 bits) which respectively capture the finest variations in the contour, intermediate stroke information and larger concavities and enclosed regions, e.g., [4] are used for characters. For bi-grams and words, 1024 such binary valued GSC features [5] are used. The number of local features are dependent on the recognition of characters in the document. To assist recognition and increase the amount of local features, lexicon-based word recognition and transcript mapping approaches [6] can be used. The number of local features also depends on the amount of information in the document(full page/half page etc.).

2.2 Distance space distribution

Once a set of features described above in Section 2.1 are available for two documents, they can be compared. The comparison results in mapping from feature space to distance space. The macro features are real valued and hence the mapping to distance space is just the absolute difference between the two feature values. On the other hand, the similarity between binary valued feature vectors for local features can be calculated using a number of different measures such as Hamming distance, Euclidean distance etc. A detailed description of similarity measures for binary features is described in [7]. After much experimentation, the correlation similarity was decided as the best similarity measure.

2.3 Parametric model

The distribution in distance space can be modeled using known probability density functions (pdf). Assuming that similarity data can be acceptably represented by Gaussian or Gamma distributions, pdfs of distances conditioned upon the same-writer and different-writer categories for a single feature x have the parametric forms $p_s(x) \sim N(\mu_s, \sigma_s^2)$, $p_d(x) \sim N(\mu_d, \sigma_d^2)$ for the Gaussian case, and $p_s(x) \sim Gam(a_s, b_s)$, $p_d(x) \sim Gam(a_d, b_d)$ for the Gamma case. The Gaussian and Gamma density functions are as follows. Gaussian: $p(x) = \frac{1}{(2\pi)^{1/2}\sigma} \exp^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$. Gamma: $p(x) = \frac{x^{a-1} \exp(-x/b)}{(\Gamma(a)) \cdot b^a}$. Estimating μ and σ from samples using the usual maximum likelihood estimation, the parameters of the gamma distribution are calculated as $a = \mu^2/\sigma^2$ and $b = \sigma^2/\mu$. For macro features, we model both categories by Gamma distribution as $p_s(x) \sim Gam(\alpha_s, \beta_s)$ and $p_d(x) \sim Gam(\alpha_d, \beta_d)$. For micro features, while the “same-writer” category is modeled as $p_s(x) \sim Gam(\alpha_s, \beta_s)$ for Gamma distribution, the “different-writer” is modeled as $p_d(x) \sim N(\mu_d, \sigma_d^2)$ for Gaussian distribution. As an example, dissimilarity histograms corresponding to the same writer and different writer distributions for “slant” (macro feature) and for the letter “e” (micro feature) are shown in Figure 1 (a) and (b) respectively. Conditional parametric pdfs for “slant” and for “e” are shown in Figure 1 (c) and (d).

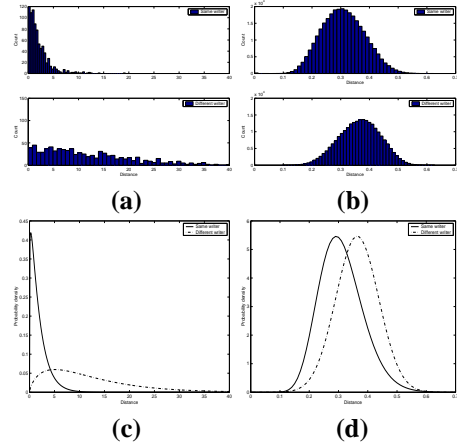


Figure 1. Histograms and parametric probability density functions

Once the distributions are modeled, the learning phase is complete. A new pair of unseen documents when compared results in a N distance values, 1 for each feature compared. For eg. if there were 10 common characters in the two documents, and 0 bigrams and 0 words common, the value of N would be 13 Macro + 10 Micro features = 23. It’s evident that the nature of the document affects the number of micro features but not the macro features. For one distance value $x_i, i \in \{1 \dots N\}$, the Likelihood ratio (LR) is given as $LR(x) = p_s(x)/p_d(x)$. Considering the features as independent, we can have the LR for N of them as $\prod_i^N LR(x_i) = \prod_i^N p_s(x_i)/p_d(x_i)$. The log likelihood ratio (LLR) is given as $LLR = \sum_i^N \log p_s(x_i) - \log p_d(x_i)$

2.4 Discussion on LLR

The understanding of the log likelihood ratio for the problem of writer verification is the key before the formulation of strength of evidence is presented. Firstly it is important to note that the magnitude of the LLR alone cannot provide a strength of evidence, though this is not obvious directly. Two reasons for this are the following

1. The reason can be illustrated with an example. Lets say that a LLR of +10 was obtained when only one word was available in the two documents indicating same writer pair. If the same LLR of +10 was obtained when a whole page of document was compared, then this is a weaker case. This is because, if a whole page of document were available, one would expect a larger positive LLR. This then a relatively small LLR value in magnitude can represent a strong strength of evidence.
2. The domain of LLR is $(-\infty, \infty)$ and hence it cannot be calibrated onto the required 9 qualitative scale (discussed later in Section 3. Even with the use of a

smoother like a sigmoid or a hyperbolic tangent, the resulting values lie in the extremity and does not take into account the amount of information compared between the documents, the features used and the nature of the content compared.

Hence there requires a new way to map from LLR to a qualitative assessment of strength of evidence. It is evident that the value of N (the number of features) directly affects the magnitude of the LLR. Below are two key factors that influence the value of N .

1. **The amount of information compared in the two documents (\mathcal{I}).** When two full page documents are compared, then it is more likely that more characters are common between them and hence greater is the value of N . If the model is perfect, then each of the terms in the summation is positive for same writer pairs and negative for different writer pairs. Hence, the magnitude of LLR is higher when documents with more information is compared.
2. **The nature of content(same/different) in the two documents (\mathcal{C}).** When two documents being compared have the same content(text) in them, then it is obvious that more characters/bi-grams/words recognized in each of them will be common. Hence, this also increases the value of N and results in a higher magnitude of LLR.

Also it is important to note that, when N is fixed (for eg. when only macro features are used), the range of the LLR values should not vary much. However, the accuracy of the model defined as percentage of correctly made decisions(same/different) is still influenced by the amount of information \mathcal{I} and the nature of the content \mathcal{C} . For example, when only macro features are used, the feature values are more accurate when more information is present in the document and as a result the LLR values classify the pair as same and different better.

3 Strength of Evidence

Having observed two important reasons that affect the LLR as mentioned in Section 2.4, the formulation of strength of evidence is parameterized based on two things (i) The amount of information compared(\mathcal{I}) and (ii) Whether the documents being compared have the same or different content(\mathcal{C}). In this paper, we discretize the values that \mathcal{I} can take on as one of *word, line, multiple lines, half page, full page*. Words of different lengths(short, medium and long) and words made of purely numbers were all considered as belonging to the type word. The two different values \mathcal{C} can take on is *Same Content, Different Content*. The value of \mathcal{I} can be automatically found during the

processing of the handwritten document which consists of line and word segmentation. The algorithm used for line segmentation is discussed in [8] and that for word segmentation is using a Neural Network to decide whether or not the gap between two connected components is a word gap or not. Hence using the number of lines L and the number of words in each line W_i $i \in \{1 \dots L\}$, the value of \mathcal{I} is decided as follows: $[L = 1, W_1 = 1 \Rightarrow \mathcal{I} = \text{word}]$, $[L = 1, W_1 > 1 \Rightarrow \mathcal{I} = \text{line}]$, $[L > 1 \ \& \ L < 4 \Rightarrow \mathcal{I} = \text{multiple-lines}]$, $[L \geq 4 \ \& \ L < 8 \Rightarrow \mathcal{I} = \text{half-page}]$, $[L \geq 8 \Rightarrow \mathcal{I} = \text{full-page}]$. Similarly for the value of \mathcal{C} , if the number of common words between the two documents is greater than 80% of the number of words in the smaller document(in terms of number of lines), then the value of \mathcal{C} is “Same Content”, or else it is different content.

3.1 Mathematical formulation

For each possible pairs of settings for \mathcal{I} and \mathcal{C} , the distribution of the LLR as observed on a validation set of ensemble of pairs is obtained. This ensemble of pairs consists of both, pairs from same as well as from different writers. The number of such pairs from the same and different writers are kept the same to avoid a bias in the distribution of LLR. Let D_{ic} represent the distribution of LLR for $\mathcal{I} = i$ and $\mathcal{C} = c$. Further, let D_{ic}^S represent the subset of D_{ic} where the samples truly belonged to the same writer and let D_{ic}^D represent the subset of D_{ic} where the samples truly belonged to different writers. It is clear that $D_{ic} = D_{ic}^S \cup D_{ic}^D$. Here it is important to note that the distribution D_{ic}^S and D_{ic}^D will be different for different sets of features used. For eg., the distribution can be further parameterized by a third variable that measures which feature set was used (macro only or macro+micro). We leave the discussion of inclusion of this third parameter to the experiments and results section. Using the distributions D_{ic}^S and D_{ic}^D , and for any given value of LLR \mathcal{L} , two percentages can now be obtained (i) P_{ic}^S : Percentage of samples in D_{ic}^S that had LLR values $> \mathcal{L}$ and (ii) P_{ic}^D : Percentage of samples in D_{ic}^D that had LLR values $> \mathcal{L}$. To be verbose, P_{ic}^S represents the percentage of same writer cases in the validation set that had LLR values even larger than the one for this. This implies that P_{ic}^S represents the percentage of same writer cases in the validation set that were *stronger* than the current case. Similarly, P_{ic}^D represents percentage of different writer cases that were *weaker* than the current case. Mathematically, they are defined as in equation 1.

$$\begin{aligned}
 P_{ic}^S &= \frac{|D_{ic}^S > \mathcal{L}|}{|D_{ic}^S|} \times 100 \\
 P_{ic}^D &= \frac{|D_{ic}^D > \mathcal{L}|}{|D_{ic}^D|} \times 100
 \end{aligned} \tag{1}$$

Scale	Opinions for same	P_{ic}^S
1	Identified as same	0.00 ~ 22.21
2	Highly probably same	22.22 ~ 44.43
3	Probably same	44.44 ~ 66.65
4	Indicating same	66.66 ~ 88.87
5	No conclusion	88.88 ~ 100.00

Table 1. Rules for obtaining the 9 point scale(1-5) when LLR \mathcal{L} is +ve

Scale	Opinions for different	$P_{ic}^{D'}$
5	No conclusion	88.88 ~ 100.00
6	Indicating different	66.66 ~ 88.87
7	Probably different	44.44 ~ 66.65
8	Highly probable different	22.22 ~ 44.43
9	Identified as different	0.00 ~ 22.21

Table 2. Rules for obtaining the 9 point scale(5-9) when LLR \mathcal{L} is -ve

where $|\cdot|$ represents cardinality. It is clear that $P_{ic}^{S'} = 100 - P_{ic}^S$ will represent that percentage of samples in D_{ic}^S that had LLR values $\leq \mathcal{L}$ and similarly we define $P_{ic}^{D'} = 100 - P_{ic}^D$. $P_{ic}^{S'}$ and $P_{ic}^{D'}$ represent the complement of P_{ic}^S and P_{ic}^D respectively.

3.2 9 point scale

A standard 9 point scale is used by questioned document examiners to verify whether two documents are written by the same person or not. This scale is as follows: [1-Identified as same, 2-Highly Probable same, 3-Probably same, 4-Indicating same, 5-No conclusion, 6-Indicating different, 7-Probably different, 8-Highly probably different and 9-Identified as different]. The sign(+ve,-ve) of the LLR \mathcal{L} between a pair of documents makes a decision of same or different writer. The strength of evidence is based on this decision. The scale $1 \dots 9$ for a particular pair of document can be obtained using either P_{ic}^S (if $\mathcal{L} > 0$) or $P_{ic}^{D'}$ (if $\mathcal{L} < 0$). (In both cases, we are evaluating the percentage of samples that were stronger than the current case.) These two values can be calculated using Equation 1 provided i and c are known. The beginning of this section described the method to calculate these i and c . If the LLR \mathcal{L} is +ve, then the opinion scale is in the range 1-5 and in the range 5-9 if it is -ve. Note that, in either cases, the scale “5-No conclusion” can be obtained. Table 1 and 2 summarizes the rules for obtaining the 9 point scale for +ve and -ve LLR values.

4 Experiments

All experiments were conducted using the CEDAR-FOX system [9] that is an application designed for writer verification. A total of 824 writers, with each writer having written 3 samples of full page document is the size of the

Information \mathcal{I}	Content \mathcal{C}	Sys. Accuracy	Min.LLR	Max.LLR
Line	Same	80.46%	-82.47	16.95
	Different	62.98%	-72.14	11.05
Multiple lines	Same	88.71%	-64.03	16.84
	Different			
Half page	Same	93.56%	-63.82	22.83
	Different	94.41%	-76.4	22.38
Full page	Same	93.93	-73.87	24.76

Table 3. System accuracy and range of LLR when only Macro features are used

total learning/validation set to obtain the four distributions $P_{ic}^S, P_{ic}^D, P_{ic}^{S'}$ and $P_{ic}^{D'}$. These documents were full page documents, and hence smaller content half-page, line, word etc. were derived from this set. Another set of the same size was used as the testing set. The experiments were conducted primarily to show that the system accuracy and the range of LLR values obtained do vary with the amount of information content \mathcal{I} and the nature of the content \mathcal{C} . The strength of evidence presented above is just a mathematical formulation to provide for a qualitative assessment of the models decision. A third parameter other than \mathcal{I} and \mathcal{C} is introduced here as was told in Section 3.1. This parameter is introduced due to the interactive nature of the system that has the capability of interactively turning on/off the micro features and hence Tippett plots(distributions) inclusive/exclusive of micro features can also be obtained. The results are presented for each of setting of this parameter.

4.1 Results with macro features only

Without the micro features (only macro available), the value of N (number of features) to calculate the LLR is a constant and hence as expected, the distribution of LLR values does not vary much. The minimum and maximum values of the LLR are shown in Table 3. A row for word level is not included in the table since for one word, the macro features are extremely noisy. This is because, global document features cannot be represented well with just one word. But, this row is included for experiments when macro as well as micro is used (see Section 4.2). Also the row corresponding to full page different content is also not available due to lack of such data. The system accuracy column in the Table 3 is the percentage of times the system decided on same/different writer correctly. Although the range of LLR values is the same, it is evident from the table that as the information present in the document increases, the system accuracy also increases and saturates when half page of document is present. The mapping from LLR to 9 point scale strength of evidence is the same as discussed in Section 3.1. The LLR distributions corresponding to each row of the table is stored from performing experiments on the validation set. Once this is done, equation 1, Table 1 and Table 2 are used to calculate the 9 point scale.

No of letters	Word/String	Sys. Accuracy	Min.LLR	Max.LLR
9	“Allentown”	81.43%	-71.95	42.51
6	“Street”	73.03%	-63.98	30.53
5	“Elder”	80.09%	-63.35	34.13
4*	“York,Loop”	74.63%	-73.98	22.54
3*	“Apt.New”	74.63%	-71.81	26.44
5	“14707”	79.85%	-74.31	31.45
3*	“829,300”	73.05%	-65.95	24.32
1	“1”	61.76%	-106.58	10.97

Table 4. System accuracy and range of LLR when only one word/string is used. It is seen that discriminability power of word is more than numerals and longer words give better accuracy. * denotes that the average accuracy was taken when two words of same length were considered.

Information \mathcal{I}	Content \mathcal{C}	Sys. Accuracy	Min.LLR	Max.LLR
Line	Same	86.40%	-93.82	115.57
	Different	62.98%	-72.14	11.05
Multiple lines	Same	93.81%	-105.02	96.83
	Different			
Half page	Same	93.08%	-322.59	698.64
	Different	94.78%	-111.83	172.28
Full page	Same	95.75%	-90.1	67.93

Table 5. System accuracy and range of LLR when only Macro and Micro features are used

4.2 Results with macro and micro features

With the inclusion of micro features, the range of LLR values as well as the system accuracy varies. Table 4 shows how the system accuracy varies with the number of characters available when only one word is available in both documents. It is seen that the system accuracy on longer words is better than those in shorter words. Also the discriminability power of words with characters is seen to be greater than strings of numbers. Table 5 shows the system accuracy and the LLR range when macro and micro features are used for other page types. The calculation of strength of evidence for a new unseen test case is as before.

4.3 Discussion on strength of evidence

From the above experiments it is conclusive that the LLR values and hence the model’s accuracy in decision making depends on the amount of information as well as the nature of content. When presented with a new pair of documents to be assessed, the following information is used to calculate the strength of evidence (i) Amount of information \mathcal{I} (word,line,multiple-lines, half-page,full-page), (ii) Nature of content compared \mathcal{C} (same content,different content), (iii) List of features compared (only Macro or Macro+Micro). and (iv) LLR value obtained by the model \mathcal{L} . The first three of the above decides on which LLR distribution(D_{ic} as in Section 3.1) to use for the assessment. The LLR value is then used to calculate the percentages and 9 point scale as in Equation 1 and Table 1 and 2. Note that the 9 point scale opinion is conditioned on the above four mentioned quanti-

ties and hence it suffices to say that the strength of evidence is reflective of the complete model/setup used.

5 Summary and Conclusion

A statistical model for writer verification using Gamma and Gaussian densities to model the distance space distribution was shown as a better fit than earlier used models[1]. A mathematical formulation to provide for a qualitative 9 point scale assessment of strength of evidence was introduced and it was shown that this formulation takes into account (i) Amount of information present in the documents, (ii) The nature of the content compared in the two documents and (iii) The set of features used(macro/micro) for comparison. Experiments and results on a large set indicated that the models accuracy and range of LLR values does indeed depend on all of the above mentioned factors and hence the correctness of the formulation of strength of evidence.

References

- [1] S. Srihari, M. Beal, K. Bandi, V. Shah, and P. Krishnamurthy, “A statistical model for writer verification,” August 2005, pp. 1105–1109.
- [2] S. Lee, S. Cha, and S. N. Srihari, “Combining macro and micro features for writer identification,” *SPIE, Document Recognition and Retrieval IX, San Jose, CA*, pp. 155–166, 2002.
- [3] R. A. Huber and A. M. Headrick, “Handwriting identification: facts and fundamentals,” *CRC Press*, 1999.
- [4] S. N. Srihari, S. H. Cha, and S. Lee, “Individuality of handwriting,” in *Journal of Forensic Sciences*, 2002, pp. 856–872.
- [5] B. Zhang, S. N. Srihari, and C. Huang, “Word image retrieval using binary features,” in *SPIE*, E. H. B. Smith, J. Hu, and J. Allan, Eds., vol. 5296, 2004, pp. 45–53.
- [6] C. Huang and S. Srihari, “Mapping transcripts to handwritten text,” October 2006, pp. 15–20.
- [7] B. Zhang and S. N. Srihari, “Binary vector dissimilarity measures for handwriting identification,” *Proc. Document Recognition and Retrieval X, Santa Clara, CA*, pp. pp 28–38, January 2003.
- [8] M. Arivazhagan, H. Srinivasan, and S. Srihari, “A statistical approach to handwritten line segmentation,” February 2007.
- [9] S. N. Srihari, B. Zhang, C. Tomai, S. Lee, Z. Shi, and Y.-C. Shin, “A system for handwriting matching and recognition.” *Proc. Symposium on Document Image Understanding Technology*, 2003, pp. 67–75.