

On the use of Lexeme Features for writer verification

Anurag Bhardwaj, Abhishek Singh, Harish Srinivasan and Sargur Srihari
Center of Excellence for Document Analysis and Recognition (CEDAR)
University at Buffalo, State University of New York
Amherst, New York 14228
{ab94,singh5,hs32,srihari}@cedar.buffalo.edu

Abstract

Document examiners use a variety of features to analyze a given handwritten document for writer verification. The challenge in the automatic classification of a pair of documents to belong to the same or different writer, are both (i)The task of proper selection and extraction of features from the handwritten document and (ii)The use of a proper model that is capable of utilizing the true discriminatory power of these features for classification. This paper describes the use of content specific skeleton based features for characters and pairs of characters(bigrams) and ascertains their discriminatory power. A triangulation skeletonisation procedure is first used to obtain the skeleton of the character(s), and features are computed from the skeleton. Experiments and results are conducted on content specific features extracted for two most frequently occurring bigrams(th , he), and characters(d and f). A neural network based on a Bayesian formulation was used to ascertain the discriminability power of these features. To combine these features with previously existing writer verification features, an alternative Naive Bayes model is also described and evaluated. From the results obtained, we conclude that bigram th has the highest discriminatory power followed by character d , f and bigram he . Also the paper highlights the significant increase in performance of writer verification($\sim 15\%$ more) with the use of Bayesian neural networks as against the Naive Bayes model.

1. Introduction

Writer verification is the problem of determining whether two handwriting samples were written by the same or different writers, which is of vital importance in Questioned Document Examination [1]. The task of proper selection and extraction of features from the handwritten document determines the performance of automatic writer verification. This paper describes the importance of various

features for bigrams(th,he) and unigrams d,f . We may define bigrams, unigrams and others as lexemes of handwriting. Features that which can be formally defined and easily extracted from handwritten documents, are termed as computational features [2]. The importance of computational features to the overall design of the writer verification is huge, since it enables the automation of writer verification. In this paper, we select and utilize such computational features for lexemes. The computational features are predominantly unique to the specific lexeme, and hence they may be termed as *content specific features*. We analyze the relative discriminatory power of these features using a bayesian formulation of a neural network. Additionally, an alternative method(Naive Bayes) that uses a gamma statistical model is also described in order to combine these features to previously existing feature sets[3]. The limitation of the simple Naive Bayes model in terms of writer verification accuracy is also seen when comparing it against the performance using the Bayesian neural network. The rest of the paper is organized as follows. Section 2 explains the selection and extraction of features followed by description of couple of writer verification models that utilize these features in section 3. Experiments and results are explained in section 4 followed by conclusion in section 5.

2. Features

Plenty of features for writer verification have been proposed[4][2]. These include document examiners features, as well as those that can be easily computed automatically. We focus more on features that fall under both categories.

2.1 Feature selection

The feature selection is a non-trivial task for all lexemes. We select for our analysis two bigrams(th,he) and two unigrams(d,f). The bigrams selected are the most frequently occurring lexemes in natural English writing. The reasons

Table 1. Computational Features for Bigram ‘th’

Feature Name	Feature Description
f1	Height of Bigram
f2	Width of Bigram
f3	Height to Width Ratio
f4	Vertical Distance between ‘t’-top and ‘h’-top
f5	Horizontal distance between ‘t’-top and ‘h’-top
f6	Euclidean distance between ‘t’-top and ‘h’-top
f7	Angle between ‘t’-top and ‘h’-top
f8	‘th’ Connected / Disconnected
f9	Average stroke-width of Bigram
f10	Slant of ‘t’
f11	Slant of ‘h’
f12	Presence of loop in ‘t’
f13	Presence of loop in ‘h’

Table 2. Computational Features for Bigram ‘he’

Feature Name	Feature Description
f1	Height of Bigram
f2	Width of Bigram
f3	Height to Width Ratio
f4	Vertical Distance between ‘h’-top and ‘e’-top
f5	Horizontal distance between ‘h’-top and ‘e’-top
f6	Euclidean distance between ‘h’-top and ‘e’-top
f7	Angle between ‘h’-top and ‘e’-top
f8	‘he’ Connected / Disconnected
f9	Average stroke-width of Bigram
f10	Slant of ‘h’
f11	Slant of ‘e’

for the specific selection of the unigrams relates to the type of features that are extracted. Features are in general, geometric(height,width etc.), angular(stroke angle, slant angle, etc.) and loop features(area/length/presence of loop)[2]. These features were intuitively found to be most varying in the unigrams *d,f*.

Table 1 lists the *context specific features* for the bigram *th*[2]. Similarly table 2 lists those for bigram *he*. Features for unigram *d,f* are the same and are listed in table 3. The features are easy to correctly extract from the document and a sketch of the algorithm for the same follows.

Table 3. Computational Features for letter ‘d’ and ‘f’

Feature Name	Feature Description
f1	Height
f2	Width
f3	Height to Width Ratio
f4	Average Stroke Width
f5	Final Stroke Angle

2.2 Feature extraction

Once all features have been selected, the task of extraction becomes well directed. To begin with, the document is binarized using the well known Otsu’s algorithm [5] or using other adaptive binarization techniques. Segmentation of the document into lines is then performed by modeling the lines as bivariate Gaussian densities [6]. This is then followed by segmentation of each line into word images using a neural network based algorithm that classifies a gap between two connected components as a word gap or not. After this, the document is subjected to recognition. Various techniques for recognition include transcript mapping [7], word recognition using Lexicon[8] or automatic character recognition. The above ensures the availability of lexemes for feature extraction. The images that correspond to the specific lexemes in question in this paper are the inputs to the feature extraction algorithm. As discussed in[2], there are a number of algorithms for the same. Most of the features described in section 2.1(other than geometric features such as height, width etc.) are extracted from the skeleton of the image. The process of obtaining the skeleton of an image is generally done by a process called thinning. But thinning based skeleton extraction may not work well here since it introduces some artefact such as small connected components, spurious loops and extra branches, that will interfere with the extraction of accurate features. The skeletonization method that we use is based on a triangulation skeleton extraction scheme described in [9] which does a fast thinning of the binary images by decomposing a smoothed, polygon approximation of the boundaries of each word into a set of contiguous triangles.

Various algorithms for feature extraction from the skeleton, mentioned in [2] seemed to produce good results in tandem with our image preprocessing and skeleton extraction methods. To determine the co-ordinates(bounding box) of the bigrams, each individual character from the bigram is separated out using a character recognizer. After this, the co-ordinates are obtained by scanning the skeleton image. The feature corresponding to the distance between the top points(row 6 in table 1) and the angle formed(row 7 in table 1) can now be measured by using simple trigonometric operations. To determine whether ‘t’ and ‘h’ are connected to each other, a flood-fill algorithm is used with the initial start point as ‘t’-top(top left coordinate). If ‘h’-top is also marked during the flood-fill, it indicates that they are connected to each other.

The stroke width of the bigram is determined by the formula: $Strokewidth = \frac{2*N_s}{N_d}$. Where N_s is the number of foreground pixels and N_d is the number of edge pixels(calculated from chain code image). The slants are calculated using a hough transform[10] on the skeletonized im-

ages of the individual characters to check for angles upto 45 degrees on both sides of the vertical axis. The angle having the maximum number of black pixels is returned by the hough transform as the slant of the character. To determine presence of loops, we once again use the flood-fill algorithm.

3 Verification Model

A proper verification model needs to be implemented for the effective utilization of extracted features. We used two different models for this purpose - (i)Bayesian neural network and (ii)Naive Bayes

3.1 Bayesian Neural Network

The first model used a bayesian neural network with 3 hidden units. If y represents the output of the neural network and x the input(features), then the probabilistic framework of the objective function of the neural network is as given in equation 1.

$$G(w) = - \sum_{i=1}^N \log P(y|x;w) \quad (1)$$

where N is the total number of training data points and $P(y|x;w)$ is $\text{sigmoid}(\sum_{j=1}^H h_j \cdot w_j)$. h_j is the output at the j^{th} hidden neuron and w_j represents the weight for the edge between the j^{th} hidden neuron and the output neuron. The $\text{sigmoid}(\cdot)$ maps the input into $[0, 1]$ domain and hence can be considered as a probability. The two classes for classification correspond to $1 = \text{Same writer}$ and $0 = \text{Different writer}$. $G(w)$ in equation 1 represents the negative log probability of the data, and weights that are learnt minimizes $G(w)$ using back propagation (Gradient descent). Hence the probability of the data is given by $P(x|w) = e^{-G(w)}$. By Bayes rule we have $P(w|x) = \frac{P(x|w)P(w)}{P(x)}$ and hence, for different set of weights w , the posterior probability $P(w|x)$ can be evaluated as proportional to $P(x|w)P(w)$. $P(w)$ represents a prior over the weights, and the use of a regularizer of the form $E(w) = \frac{1}{2} \sum_i w_i^2$ results in a prior $P(w) = e^{-E(w)}$. When predicting the probability that a new data point x_{new} belongs to class 1, equation 2 is used.

$$P(x_{\text{new}} = 1) = \int P(y|x_{\text{new}};w)P(w|x_{\text{tr}})dw \quad (2)$$

where $P(w|x_{\text{tr}})$ represents the posterior probability of weights learnt during the training of the network. The above integral is replaced by a sample average estimate in our implementation.

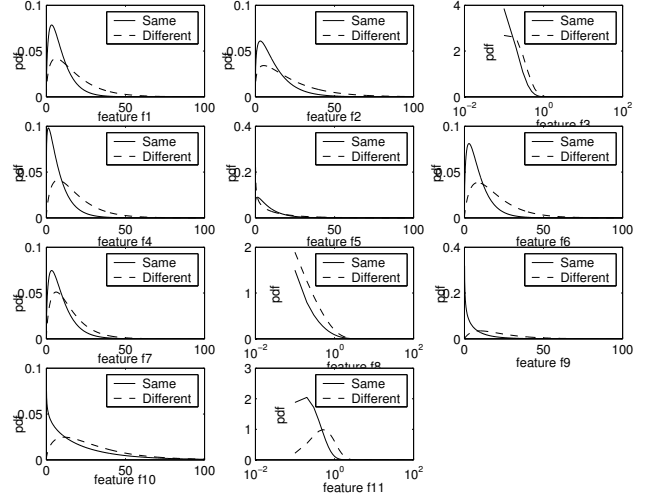


Figure 1. Modeling the histogram of features of lexeme 'he' using Gamma distribution for use by the Naive Bayes model

3.2 Naive Bayes

The second model used was a Naive bayes classifier which was implemented as a part of Cedar-Fox software[11]. This classifier used a gamma distribution to model the distance between feature vectors[3] corresponding to same writer and different writer pairs. The Gamma density function is as follows. $p(x) = \frac{x^{a-1} \exp(-x/b)}{(\Gamma(a)) \cdot b^a}$ 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 $b = \sigma^2/\mu$. The advantage of this classifier over the neural network classifier mentioned above is that it can be easily extended to be used with any other feature set available. Moreover, the performance of a particular feature can be interpreted using this naive bayes model as opposed to the neural network.

The likelihood ratio (LR), which summarizes the result, is given by $LR(x) = \frac{p_s(x)}{p_d(x)}$ [3], where $p_s(x)$ denotes the probability density function of the distances for feature x corresponding to the same pair of writers and $p_d(x)$ denotes the probability density function of the distances for feature 'x' corresponding to different pair of writers. The log-likelihood ratio (LLR), obtained by taking the natural logarithm of LR, is more useful since LR values tend to be very large (or small)[3]. The decision of same/different when two unknown documents are compared is based on the sign of the LLR value. An example of matching two 'th' bigram from documents written by the same writers is shown in table 4 and similarly for documents written by different writers is shown in table 5. As seen in the table, most LLR values for each feature in the same writer table is positive, and that for the different writers is neg-

Table 4. same writer

Feature	Questioned doc feature value	Known doc feature value	Euclidean distance	LLR
Height of Bigram	55.0	59.0	4.0	0.232
Width of Bigram	63.0	63.0	0.0	1.5
Height to Width Ratio	0.87	0.936	0.113	0.063
Vertical Distance between 't'-top and 'h'-top	10.0	8.0	2.0	0.238
Horizontal distance between 't'-top and 'h'-top	20.0	14.0	6.0	-0.02
Euclidean distance between 't'-top and 'h'-top	22.36	16.12	6.24	-0.04
Angle between 't'-top and 'h'-top	26.56	29.74	3.18	0.32
'th' Connected / Disconnected	1.0	1.0	0.0	0.11
Average stroke-width of Bigram	4.1	3.66	0.44	-0.02
Slant of 't'	6.0	3.0	3.0	0.56
Slant of 'h'	-4.0	3.0	7.0	0.02
Presence of loop in 't'	0.0	0.0	0.0	0.11
Presence of loop in 'h'	0.0	1.0	1.0	-0.16

Table 6. A comparison of writer verification accuracy with both verification models for all bigrams and unigrams

Lexeme	Accuracy with Naive Bayes Model	Accuracy with Neural Network
th	82%	92%
he	62%	86%
d	76%	88%
f	75%	88%

ative. Based on the cumulative LLR score, we determine whether the documents are written by the same writer or different writers. In the case where the document is characterized by more than one feature we assume that the writing elements are statistically independent[3]. The log-likelihood ratio (LLR) in this case has the form $LLR = \sum_{i=1} \sum_j \sum_k \ln p_s(d_i(j, k)) - \ln p_d(d_i(j, k))$.

4 Experiments and Results

For writer verification, the bigrams are extracted from both the documents which need to be verified. Feature extraction is then done for each bigram from both, the questioned and the known documents. The questioned and known documents are taken from varied sets. Training is done taking the questioned and known documents from the same writer sets to identify the distance values clusters for same writers, as well as from different writer sets to identify the clusters for different writers. When Euclidean distances are taken for the same writer documents, the distances are usually much smaller than those take for different writer documents. Thus the training can be completed and the results can be modeled as gamma distributions.

With the experiments conducted on bigrams **th**, **he** and characters **d** and **f**, the accuracies obtained for all the lexemes are shown in the table 6. we obtain an accuracy of 82% for **th**, 76% for **d**, 75% for **f**, 62% for **he** using a naive bayes verification model. Figure 3 and 3 shows the modeling of the distance space distribution(same and different writer pairs) of features for character **d**, as Gamma densi-

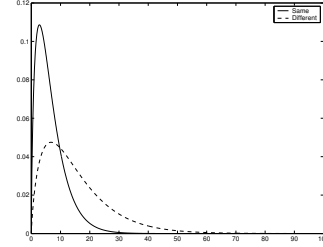


Figure 2. Feature - Height of character D. Modeling the distance space distributions as Gamma densities.

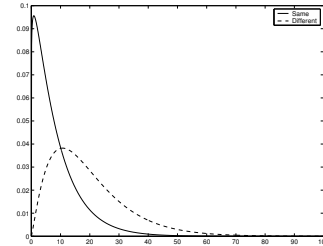


Figure 3. Feature - Stroke Width of character D. Modeling the distance space distributions as Gamma densities.

ties. Similarly modeling using Gamma for character **f**, are shown in Figure 4 and Figure 5 produced desired results.

5. Summary and Conclusion

It was found that the 'th' bigram contains larger amount of information than the 'he' bigram and the single letter 'd' or 'f'. The discriminatory power is higher because both the characters of the bigram have an ascender. The writer verification task using the 'th' features alone gave an accuracy of 92%. This compares very closely to the ~ 95% accuracy reported by using all the features using the Naive Bayes model[3]. On an average the bayesian neural network outperforms the Naive Bayes model by approximately 15%. The incorrect assumption of independence of features by the Naive Bayes model becomes critical and it highlights the drawback of the model. On the other hand, the flexibility of adding/removing features on the fly is an important

Table 5. different writer

Feature	Questioned doc feature value	Known doc feature value	Euclidean distance	LLR
Height of Bigram	50.0	66.0	16.0	-0.35
Width of Bigram	57.0	108.0	51.0	-0.74
Height to Width Ratio	0.87	0.611	0.259	-0.22
Vertical Distance between 't'-top and 'h'-top	9.0	-7.0	16.0	-0.52
Horizontal distance between 't'-top and 'h'-top	8.0	16.0	18.0	-0.31
Euclidean distance between 't'-top and 'h'-top	12.04	26.92	14.88	-0.21
Angle between 't'-top and 'h'-top	48.36	-15.07	63.43	-0.75
'th' Connected / Disconnected	1.0	0.0	1.0	-0.11
Average stroke-width of Bigram	4.01	4.31	0.30	0.29
Slant of 't'	7.0	-8.0	15.0	-0.34
Slant of 'h'	-5.0	-14.0	9.0	-0.08
Presence of loop in 't'	0.0	0.0	0.0	0.11
Presence of loop in 'h'	0.0	0.0	0.0	0.16

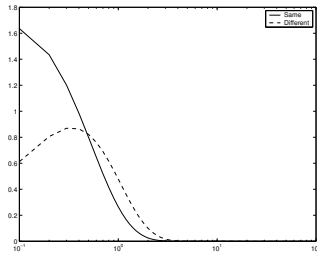


Figure 4. Feature - Height Width ratio. Gamma modeling of distance between features for character F.

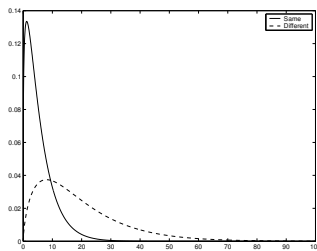


Figure 5. Feature - Stroke Width. Gamma modeling of distance between features for character F.

advantage of the Naive Bayes model. There thus is a need for a sophisticated model that in addition to superior performance can give the required flexibility. It suffices to say that the use of Bayesian neural network is the direction towards achieving this goal.

References

- [1] R. Huber and A. Headrick, *Handwriting Identification: Facts and Fundamentals*. CRC Press, 1999.
- [2] V. Pervouchine, *Ph.D. Thesis: Discriminative Power of Features Used by Forensic Document Examiners in the Analysis of Handwriting*, 2006.
- [3] K. B. V. S. P. K. S. N. Srihari, M. J. Beal, "A statistical model for writer verification," in *Proceedings of the eight International Conference on Document Analysis and Recognition (ICDAR 2005)*, vol. 2, 2005, pp. 1105–1109.
- [4] 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.
- [5] N. Otsu, "A threshold selection method from gray level histograms," *IEEE Trans. Systems, Man and Cybernetics*, vol. 9, pp. 62–66, Mar. 1979, minimize inter class variance.
- [6] M. Arivazhagan, H. Srinivasan, and S. Srihari, "A statistical approach to handwritten line segmentation," February 2007.
- [7] C. Huang and S. Srihari, "Mapping transcripts to handwritten text," October 2006, pp. pp. 15–20.
- [8] G. Kim and V. Govindaraju, "A lexicon driven approach to handwritten word recognition for real-time applications," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. Vol. 19(4), pp. pp. 366–379, April 1997.
- [9] W. B. M. Melhi, S. S. Ipson, "A novel triangulation procedure for thinning handwritten text," in *Pattern Recognition Letters*, vol. 22, 2001, pp. 1059–1071.
- [10] P. Hough, "Machine analysis of bubble chamber pictures," *International Conference on High Energy Accelerators and Instrumentation, CERN*, 1959.
- [11] 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.