

Writer Verification of Historical Documents Among Cohort Writers

Gregory R. Ball and Sargur N. Srihari
CEDAR
University at Buffalo
520 Lee Entrance, Suite 202
Amherst, NY 14228, USA
{grball, srihari}@cedar.buffalo.edu

Roger Stritmatter
Department of Humanities
Coppin State University
2500 West North Avenue
Baltimore, MD 21216-3698, USA
rstritmatter@coppin.edu

Abstract

Over the last century forensic document science has developed progressively more sophisticated pattern recognition methodologies for ascertaining the authorship of disputed documents. We present a writer verification method and an evaluation of its performance on historical documents with known and unknown writers. The questioned document is compared against handwriting samples of Herman Melville, a 19th century American author who has been hypothesized to be the writer as well as against samples crafted by several writers from the same time period. The comparison led to a high confidence result to the questioned document's writership, as well as gives evidence for the validity of the writer verification method in the context of historical documents. Such methodology can be applied to many such questioned historical documents, both in literary and legal fields.

1. Introduction

Over the last century forensic document science has developed progressively more sophisticated pattern recognition methodologies for ascertaining the authorship of disputed documents. These include advances not only in computer assisted stylometrics, but forensic handwriting analysis such as that used in CEDAR-FOX[1], a forensic examination tool that captures “the uniqueness of an individual writer by generating a probability distribution of distances between instances of character classes in the known document and the general population, as represented by a set of many samples” and employs a log-likelihood ratio “to determine the strength of confidence for the opinion.” In earlier work [2], we applied this method to a set of historical documents. In this paper, we establish the generalized

validity of the writer verification method for historical documents by performing analysis of historical documents created by writers matching similar demographics to the writers previously compared (similar historical era, gender, etc.).

Because documents typically studied by literary historians are usually printed rather than handwritten, the application of handwriting analysis in the humanities has been more limited than stylometrics [3, 6], primarily concentrating on a limited number of still unresolved questions such as the alleged appearance of William Shakespeare’s hand in the manuscript of Sir Thomas More (see Howard-Hill [4] for a survey). However, in cases where both known and unknown samples are handwritten, and contain sufficient sample size to generate reproducible results, forensic handwriting analysis may provide a gold standard for identifying authorship, for literary as well as legal documents.

2 Writer Verification

2.1 CEDAR-FOX

The CEDAR-FOX system has a number of interactive features that make it a useful as an experimental tool: feature extractors for handwriting, confidence value computation, an interface to scan handwritten document images, obtain line and word segments, and automatically extract features for handwriting matching after performing character recognition and/or word recognition (with or without interactive human assistance in the recognition process, e.g., by providing word truth).

One objective of CEDARFOX is “writer verification,” which provides a measure of confidence of whether two samples (questioned document and exemplar document) belong to the same individual or to two

different individuals. Verification is the task of determining whether a pair of handwriting samples was written by the same individual. Given a new pair of documents, verification is performed as follows: (i) writing element extraction, (ii) similarity computation, (iii) determining the log-likelihood ratio (LLR) from the estimated conditional probability density estimates.

2.2 Document Features

The system computes three types of features – macro-features at the document level, micro-features at the character level, and style features from bigrams and words. Each of these features contributes to the final result to provide a confidence measure of whether two documents under consideration are from the same or different writers.

Macro-features capture the global characteristics of the writer’s individual writing habit and style, which includes entropy of gray values, binarization threshold, number of black pixels, number of interior contours, number of exterior contours, contour slope components consisting of horizontal (0 degree or flat stroke), positive (45 or 225 degree), vertical (90 or 270 degree) and negative (135 or 315 degree), average height and average slant per line. In the current study, 11 of 13 macro-features (all except entropy and number of black pixels due to their variability depending on page size/scanning) were used in the experiments.

Micro-features are designed to capture the finer details at the character level. Each micro-feature is a gradient-based binary feature vector, which consists of 512 bits corresponding to gradient (192 bits), structural (192 bits) and concavity (128 bits) features, known as GSC features (described fully in [7]). The use of micro-features depends on the availability of recognized characters, i.e. character images associated with truth. Four possible methods are available in CEDAR-FOX to get recognized characters: (i) manually crop characters and label each with its truth; this labor-intensive method has the highest accuracy, (ii) automatically recognize characters by using a built-in character recognizer; the method is error prone for cursive handwriting where there are few isolated characters, (iii) automatically recognize words using a word-lexicon from which segmented and recognized characters are obtained; error rates depend on the size of lexicon and can be as high as 40% for a page, and (iv) use transcript mapping to associate typed words in a transcript of the entire page with word images [5]; it involves typing the document content once which can then be reused with each scanned document. Since the full-page documents have the same content (CEDAR-letter), the transcript map-

ping approach was used. This method has an 85% accuracy among words recognized. Since most words are recognized, they are also useful for extracting features of letter pairs and whole words, as discussed next.

Style features are features based on the shapes of whole words and shapes of letter pairs, known as bigrams. A bigram is a partial word that contains two connected characters, such as “th,” “ed,” etc. Style features are similar to micro-features (i.e. the GSC features) where a binary feature vector is computed from the bigram or word image. While the bigram feature has the same size as the micro-feature (with 512 binary values), the word feature has 1024 bits corresponding to gradient (384 bits), structural (384 bits) and concavity (256 bits). As with micro-features, to use these two style features, word recognition needs to be done first. As mentioned above when describing micro-features, a transcript-mapping algorithm was used to do the word recognition automatically. The words being recognized are saved and used to compute the word-GSC features. Then characters are segmented from these words and the two consecutive segmented characters both being confirmed by a character recognizer are combined as one bi-gram component.

2.3 Description of Similarity Computation

For macro-features, since they are all real-valued features, the distance is simply the absolute value of their difference. For micro and style features, several methods have been recently evaluated, which has led to the choice of a so-called “correlation” measure being used to compute the dissimilarity between two binary feature vectors (described more fully in [8]).

The distributions of dissimilarities conditioned on being from the same or different writer are used to compute likelihood functions for a given pair of samples. Such distributions can be learned from a training data set and represented either as probability tables or as probability density functions. Probability distributions for discrete valued distances are represented as probability tables. In the case of continuous-valued distances parametric methods are used to represent probability density functions. Parametric methods are more compact representations than nonparametric methods. For example, with the k-nearest-neighbor algorithm, a well-known nonparametric method, we need to store some or all training samples. This compactness results in a faster run-time algorithm. Also parametric methods are more robust in a sense that nonparametric methods are more likely to over-fit and therefore generalize less accurately. Of course the challenge is to find the right parametric model for the problem at hand. While the

Gaussian density is appropriate for mean distance values that are high, the Gamma density is better for modeling distances since distance is never negative-valued.

2.4 Log Likelihood Ratio

When each document is characterized by more than one feature CEDAR-FOX makes the assumption that the writing elements or features are statistically independent, although strictly speaking this is incorrect. We investigated complex models, such as a Bayesian neural network, and got an improvement of 1-2% on overall accuracy, which is not significant. There are several other justifications for choosing naïve Bayes. First, its functioning is simple to explain and modify, e.g., since log-likelihood ratios are additive we can observe the effects of adding and removing features easily. Its simplicity goes back to the earliest QDE literature [1] where there is reference to multiplying the probabilities of handwriting elements. Second, as has been observed in other machine learning tasks, more complex models tend to overfit to the data, which can lead to poorer performance on large amounts of unseen data.

Each of the two likelihoods that the given pair of documents were either written by the same individual or by different individuals can be expressed, assuming statistical independence of the features as follows. For each writing element $e_i, i = 1, \dots, c$, where c is the number of writing elements considered, we compute the distance $d_i(j, k)$ between the j th occurrence of e_i in the first document and the k th occurrence of e_i in the second document for that writing element. The likelihoods are extracted as

$$L_s = \prod_{i=1}^c \prod_j p_s(d_j(j, k))$$

$$L_d = \prod_{i=1}^c \prod_j p_d(d_i(i, k))$$

The LLR in this case has the form

$$LLR = \sum_{i=1}^c \sum_j [\ln p_s(d_i(j, k)) - \ln p_d(d_i(j, k))]$$

That is, the final LLR value is computed using all the features, considering each occurrence of each element in each document. The CEDAR-FOX verification system outputs the LLR of two documents being tested. When the likelihood ratio (LR) is above 1, the LLR value is positive and when the likelihood ratio is below 1, the LLR value is negative. Hence, if the final score

obtained is positive, the system concludes that the two documents were written by the same writer. Similarly, if the final LLR score is negative, the system concludes that the two documents were written by different writers.

3 Questioned and Known Documents

A recent literary discovery which might pose a test case for advanced forensic handwriting analysis is the Hydrachos manuscript (H.), an April 1846 satirical newspaper, The PHILADALPHIA GAZETTE - EXTRA, with seven pen and ink drawings accompanying a 437 word handwritten commentary on U.S. and world news.

Measuring 40 x 25 cm, the document contains satirical news content, primarily from the United States, Great Britain, Italy, and China, on both the recto (Figure 1(b)) and verso sides (Figure 1(a)).

Lexical, grammatical, thematic, visual, content, and situational analysis[9] all support the hypothesis that the manuscript's author is the New England author Herman Melville. Although Melville does not use the word Hydra<r>chos, many other phrases and allusions in the manuscript can be traced to his published writings. Qualitative stylistic analysis supports the attribution. In the 437 word document, Stritmatter[9] was able to trace 59 words and phrases, including Antidiluvian (sic); hoe cake (used in an identical context); John Bull and Jonathan (as sobriquets for England and America). Analysis of diction and content is also consistent with Melville's authorship.

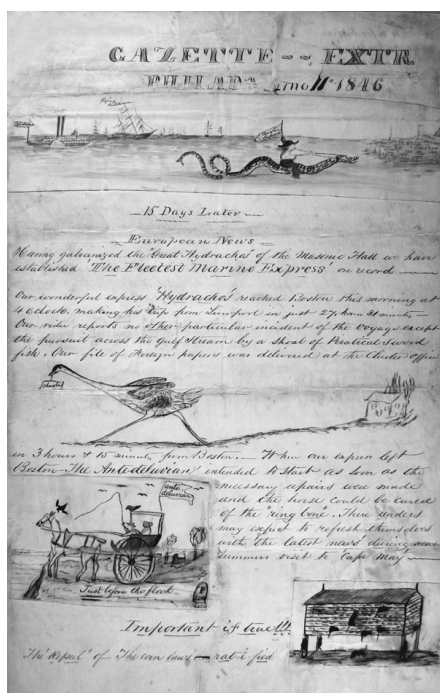
In addition, while the utility and accuracy of the writer verification system used in CEDAR-FOX is well-established on a large data set, it has not been tested extensively on documents of historic origin. To help establish that the system is era-independent, or at least appropriate for comparison to Melville's time period, letters written by two of Melville's male relatives (Thomas Melville, Jr. and Alan Melville) were obtained, both produced in 1818¹.

4 Experiments and Results

To evaluate the documents according to the writer verification method described in the previous section, we followed the following sequence of steps to prepare the documents for analysis:

Step 1: Manual preprocessing: We modified the documents to remove the hand-drawn non-text images, and remove the more major noise.

¹The authors wish to acknowledge the New York Public Library for providing scans of the relevant documents



(a) Hydrachos manuscript, verso



(b) Hydrachos manuscript, recto

Figure 1. The unknown document—the Hydrachos manuscript

Step 2: Pre-processing with CEDAR-FOX: Binarization, line and word segmentation, document wide feature calculation, automatic character recognition.

Step 3: Word segmentation correction: Some manual correction to the word segmentation was performed since there were some slight word segmentation errors.

Step 4: Transcript mapping, manual transcript corrections: Using the “transcript mapping” function and correct transcripts, the known and unknown documents were manually verified to be correctly word truthed.

Some of the processed documents used are shown in Figure 2.



(a) Known document



(b) Unknown document

Figure 2. Pre-processed and truthed documents

Table 1. Overall LLR Results

	Total LLR	System Opinion
H. Melville Letter 1 vs. Hydrachos, verso	35.29	Identified As Same
H. Melville Letter 2 vs. Hydrachos, verso	129.42	Identified As Same
H. Melville Letter 1 vs. Hydrachos, recto	20.82	Highly Probable Same
H. Melville Letter 2 vs. Hydrachos, recto	60.12	Identified As Same
H. Melville Letter 1 vs. H. Melville Letter 2	333.52	Identified As Same
Hydrachos manuscript, verso vs. recto	42.87	Identified As Same
Full Hydrachos manuscript vs. H. Melville Letter 1	56.11	Identified As Same
Full Hydrachos manuscript vs. H. Melville Letter 2	189.54	Identified As Same
T. Melville Letter vs. H. Melville Letter 1	-14.84	Indications Did Not
T. Melville Letter vs. H. Melville Letter 2	12.38	Indications Did
T. Melville Letter vs. Hydrachos, recto	-68.18	Identified As Different
T. Melville Letter vs. Hydrachos, verso	-73.03	Identified As Different
A. Melville Letter vs. H. Melville Letter 1	-671.41	Identified As Different
A. Melville Letter vs. H. Melville Letter 2	-386.09	Identified As Different
A. Melville Letter vs. Hydrachos, recto	-318.96	Identified As Different
A. Melville Letter vs. Hydrachos, verso	-448.61	Identified As Different
A. Melville Letter vs. T. Melville Letter	-284.46	Identified As Different

5 Results

Figure 3. Comparison of “ea,” “ou,” “oo,” and “wh” between H. Melville’s writing and the Unknown document



We treated the problem as a case of four known (letters) and two unknown (verso and recto) documents and generated an LLR comparison value for each of the pairs of documents. Comparisons of some select bigram features are shown in Figure 3.

The comparison method was validated on a test set of 1,648 test cases[8]. When performing on the test cases using all features on documents containing different content, the overall error rate for the general population was 0-2.4%. These test documents were also used

to map the LLR score onto a nine point opinion scale, ranging in order as follows: identified as same, highly probable same, probably did, indications did, no conclusion, indications did not, probably did not, highly probable did not, identified as different. The overall LLR and opinion results for the documents are described in Table 1.

The LLR scores are compound scores consisting of several feature components: the global (macro), character (micro), bigram, word, and lexeme features.

The characters were automatically extracted based upon the word truth. All words were split into candidate characters and the automatically generated characters were validated by a character recognizer. Characters failing the automatic verification check were not used in the comparison so as to not introduce artifacts into the comparisons.

Each acceptable character was compared exhaustively to each acceptable corresponding character in the comparison document. When all LLR's of all characters are combined, overall trends are observed. In this case, all characters led to positive results, adding information suggesting same authorship.

To more characterize the strength of confidence in the LLR scores, we modeled the 1,648 validation cases and evaluated which cases had higher or lower LLR scores – that is, we tabulated the LLR values of all known same and different writer cases (with full page documents consisting of different content) and determined the percentage of such validation cases that had higher or lower LLR values than our current study's results. The results are presented in Table 2.

Table 2. Confidence Score Breakdown

	% same writer having > LLR	% diff writer having > LLR
HM ₁ vs. Q _{verso}	14.32	0
HM ₂ vs. Q _{verso}	0	0
HM ₁ vs. Q _{recto}	41.99	0
HM ₂ vs. Q _{recto}	4.61	0
TM vs. HM ₁	0.97	74.51
TM vs. HM ₂	85.92	0.48
TM vs. Q _{recto}	0	1.21
TM vs. Q _{verso}	0	8.73
AM vs. HM ₁	0	0
AM vs. HM ₂	0	0
AM vs. Q _{recto}	0	0
AM vs. Q _{verso}	0	0
AM vs. TM	0	0

In short, none of the validation cases consisting of full page documents written by different writers with different content had LLR values that were higher than the H. Melville letters vs. Hydrachos sets presented, implying that the documents we compared are more similar than any of the validation pairings. In the case of Herman Melville Known Letter 2 vs. the average of the verso and recto sides of the questioned document, the LLR score implies similarity stronger than all but about 2% of the validation cases.

When comparing the T. Melville, Jr. letters to the Hydrachos manuscript, the system clearly identified the writers as different. When comparing the T. Melville, Jr. letters to the Herman Melville known letters, the results were essentially non-decisive, resulting in very low LLR scores when considering the large amount of text compared. One of the two scores was (incorrectly) slightly positive. When examining the incorrect result in more detail, one major factor for the score was that the letters being compared had very few large words in common, which meant that the word comparison had to be based on very short words, which yield less strong results.

When comparing A. Melville's letters to Herman Melville's letters and to the Hydrachos document, the system reported very strongly that the writers were different. When comparing the T. Melville, Jr. and A. Melville letters to one another, again the system reported a strong result that they were different. In fact, in all five cases, 0% of the validation test cases were the results so strong in either the same or different writer groups.

5.1 Conclusions

In all cases, the LLR numbers indicated the proposed writer likely penned the questioned document; in addition, comparisons between the questioned document and with known documents penned by various contemporary writers likely to have similar penmanship styles resulted in strong opinions that the documents were written by different writers. Further comparisons between the known different writers resulted in reasonable results indicating that the documents were indeed recognized as being created by different writers, with a single exception which resulted essentially in a non-conclusive result. This helps to establish that the writer verification method described returns valid results for the time period in question. Further results could be generated by comparing other known documents containing similar phrasing to the questioned document. Additional results could be gained by comparing these documents to other known documents.

References

- [1] G. R. Ball and S. N. Srihari. Comparison of statistical models for writer verification. In *Proceedings of Document Recognition and Retrieval XVI, SPIE*, January 2009.
- [2] G. R. Ball, R. Stittmeyer, and S. N. Srihari. Writer verification in historical documents. In *Proceedings Document Recognition and Retrieval XVII*, San Jose, CA, January 2010. SPIE.
- [3] D. Holmes. The analysis of literary style a review. *Journal of the Royal Statistical Society. Series A (General)*, 148(4):328–341, 1985.
- [4] T. H. Howard-Hill, ed., editor. *Shakespeare and Sir Thomas More: Essays on the Play and its Shakespearean Interest*. The University Press, Cambridge, 1989.
- [5] C. Huang and S. Srihari. Mapping transcripts to handwritten text. 2006.
- [6] R. D. Peng and N. Hengartner. Quantitative analysis of literary styles. *The American Statistician*, 56(3):175–185, 2002.
- [7] H. A. S. N. Srihari, S.-H. Cha and S. Lee. Individuality of handwriting. *Journal of Forensic Sciences*, 47(4):856872, July 2002.
- [8] S. N. Srihari, C. Huang, and H. Srinivasan. On the discriminability of the handwriting of twins. *Journal of Forensic Sciences*, 53(2):430–446, 2008.
- [9] R. Stritmatter. The great hydrachos satire of april 11, 1846: A lost literary manuscript attributed to herman melville. June 2009.