

Determining Cursive or Printed Nature of Handwritten Samples

Gregory BALL, Danjun PU, and Sargur SRIHARI
Center of Excellence for Document Analysis and Recognition
State University of New York at Buffalo
UB Commons, 520 Lee Entrance, Suite 202
Amherst, NY 14228, UNITED STATES
 {grball, srihari}@cedar.buffalo.edu

Abstract. Forensic document examiners tend to compare documents to one another only when they are both cursive or both hand printed documents. Thus, for an automated system to process writer identification or verification, it is desirable to determine whether a sample of handwritten text is of a hand printed or cursive nature. We identify a feature based upon the number of segmented words and number of connected components which are identified as single characters. We evaluate the performance of this feature by using all 4,701 documents present CEDAR letter dataset and find that using this feature identifies documents as hand printed correctly 95% of the time and identifies them as cursive 91% of the time.

1. Introduction

Forensic document examiners have a long history of providing critical information and opinions in a wide variety of legal cases. The science of forensic document analysis stretches back many decades. Recently, the interest has increased in the question of how the type of writing in the sample (hand-printed or cursive nature) affects the results. Specifically, some United States federal district court judges referred to a lack of information regarding proficiency in identifying writers of hand-printed documents (Kam & al., 2003). For example, in the case of *US v. Jeffrey H. Feingold* (9th Cir, April 2004), CR 02-0976-PHX-SMM, Kam's error rates on analyzing handwriting were specifically questioned regarding their applicability to hand printing. In that case, the court called a months-long break so that the results could be separated. Kam and Lin's study demonstrated the superior proficiency of document examiners compared with laypeople in identifying writers of hand-printed document. Nonetheless, this problem illustrates the need for forensic document analysis to be conducted with the status of writing style in mind.

However, hand printed writing poses some unique challenges and has been specifically explored by some for decades (Conway, 1955). As Armistead notes, many writers are "less habituated to printing than to cursive writing," and "in cursive script, connecting strokes within and between letters constitute a critical feature for identification," (Armistead, 1984). Furthermore, he describes the importance of this problem as "handprinting is often the script of choice for writers of anonymous or disguised messages, apparently because many believe that handprinting is less identifiable than cursive script."

CEDAR-FOX. CEDAR-FOX (Srihari & al, 2003; Srihari & al, 2007) is a system for handwritten matching and recognition which has recently implemented the use of this feature, but is simply using a threshold point set by experiment and scaling the feature value from 0-100 around this threshold point. Formal analysis of its performance has not yet been performed and this is the first time analysis has been performed with ground truth provided by document examiners.

2. Features

Discussions with document examiners revealed that they utilize the intuitive feature of character connection to differentiate cursive and hand printed documents. In other words, if most characters within each word are connected, the document is deemed "cursive" in nature. If most are disconnected, the document is determined to be "hand printed." Mixed documents can exist which the document examiners term "running hand print," but often group with cursive documents since they have sections of connected characters useful for their analysis. To capture the idea of whether or not most characters are connected in the document, we identified several features:

Discreteness (f_1). We first examine the quality of having disconnected characters present in a document—the ratio of Isolated Character Count (ICC) and Word Count (WC). ICC is the cardinality of the set of groups of connected components recognized as individual characters. WC is the cardinality of the set of words in a given document. We use WC simply as a way to normalize ICC; thus, the feature we found to best capture the intuitive document examiner feature is the ratio $f_1 = ICC/WC$.

Loopiness (f_2). As cursive writing tends to have more loops present within each connected component, a second feature we considered is the "loopiness" of a document. We capture this with the ratio of interior to exterior contours present in a page of writing. Due to the complexity, this feature is not explored in depth in this paper.



Figure 1. Cursive (left) and hand printed (right) instances: (a) $f_1 = 0/1$, $f_2 = 5/2$; (b) $f_1 = 5/1$, $f_2=3/6$

Figure 1 shows two instances of the word “several,” one written in a cursive and one in a hand printed fashion. In the cursive instance, there are two large connected components; one containing ‘sev’ and one containing ‘eral’: ICC is 0 and WC is 1, yielding an f_1 of 0/1. The hand printed instance, however, contains six large connected components. Five are recognized as individual characters (‘s,’ ‘e,’ ‘v,’ ‘e,’ and ‘r’), giving a ratio of 5/1. Thus, our hypothesis was that cursive documents would have overall very low ratios and hand printed documents would have relatively high ratios. Even ideal cursive documents are likely to have some isolated characters (due to the words ‘a’ and ‘I’). To complicate matters, however, some spurious breaks often occur, even in predominantly cursive documents. An example is the instance of the word “enthusiasm” which occurred in a predominantly cursive document shown in Figure 2. In this case, there are six connected components (‘en,’ ‘t,’ ‘hus,’ ‘i,’ ‘as,’ ‘m’). Since three of these are individual characters in a single word, the f_1 is 3/1.

Figure 2. An example of spurious breaks in a predominantly cursive : $f_1 = 3/1$, $f_2 = 8/6$

Approximating WC. In order to process arbitrary documents which lack ground truth, we first approximate the WC. To approximate WC, we first preprocess the document, segment it into lines and then segment the line into words. Finally, we count the segmented words. To segment lines, we utilize a stroke tracing method to intelligently segment the overlapping components and uses slope and curvature information of the stroke to disambiguate the path of the stroke at cross points. Once the overlapping components are segmented into strokes, a statistical method is used to associate the strokes with the appropriate lines. This method is capable of disentangling lines which collide by splitting and associating the correct character strokes to the appropriate lines (Kuzhinjedathu & al., 2008). To segment the lines from the previously mentioned method into words, we used a gap metric based approach (Huang & al., 2008) specifically designed to separate a line of unconstrained (written in a natural manner) handwritten text into words. This approach uses a set of both local and global features, motivated by the way that humans perform this type of task. The local features are distance between a pair of components, distance between neighboring components, width of left and right components, and height of the left and right components. The global features include ratio of the number of exterior contours and the number of interior contours, the average height of grouped components, average width of grouped components, and average distance between components. Two distance metrics are computed and averaged depending on the circumstances—the first metric is either the bounding box or minimum run-length distance and the second is the convex hull distance. The actual classification based on these features is done using a three-layer neural network.

Determining ICC. Determining ICC is performed by utilizing a character model recognizer (Favata, 1996). First, the full preprocessed page image is separated into its connected components. Each connected component is then passed to a character recognition algorithm which determines if the connected component is likely to consist of a single character. The features used for this determination are the Gradient, Structural, Concavity (GSC) features, a 512-dimensional feature vector, which are used in automatic character recognition for interpreting handwritten postal addresses. We extract local contour features based on the pixel gradients present in an image. A gradient map, i.e., gradient magnitude and direction at every pixel, is computed; this map is thresholded to avoid responses to noise and spurious artifacts, the map is then partitioned coarsely. Character contours are encoded by quantizing gradient directions into a few representative direction bins. This method does not require that images be normalized to a fixed grid-size before feature extraction. The features were used to train a 2-layer neural network classifier (Srikantan & al., 1993).

3. Dataset

The dataset used for our experiments is the entire CEDAR letter dataset. This dataset consists of 4701 documents. 1567 writers wrote the “CEDAR letter” exactly three times each. The “CEDAR letter” (Srihari & al., 2002) is a document which has content designed to contain 156 words which include all characters (letters and numerals), punctuation, and distinctive letter and numeral combinations (ff, tt, oo, 00). In particular, this

document set out to contain each letter of the alphabet in capital form at the initial position of a word and in lowercase form in the initial, middle, and terminal locations in a word (a minimum of 104 forms of each letter). The vocabulary size is 124 (that is, 32 of the 156 words are duplicate words, mostly stop words such as “the,” “she,” etc.). The general document structure allows extracting macro-document attributes such as word and line spacing, line skew, etc. The letters were written three times in each writer’s most natural handwriting using plain unlined sheets with a medium black ball-point pen. The repetition was to determine each writer’s variation of handwriting from one writing occasion to the next. The samples scanned at 300 dpi and 8-bit grayscale.

The documents in this data set were examined and tagged by two professional document examiners with a Boolean flag as either “hand printed” or not “hand printed.” The “not hand printed” group consists of mainly cursive documents with “a few questionable... running hand print” documents. This feature was created by the document examiners as it captured the sets information necessary for their work. An example of a predominantly cursive and a hand print document is shown in Figure 3. The document examiners tagged 621 documents as hand printed with the remaining 4080 being cursive with “a few” running hand print.

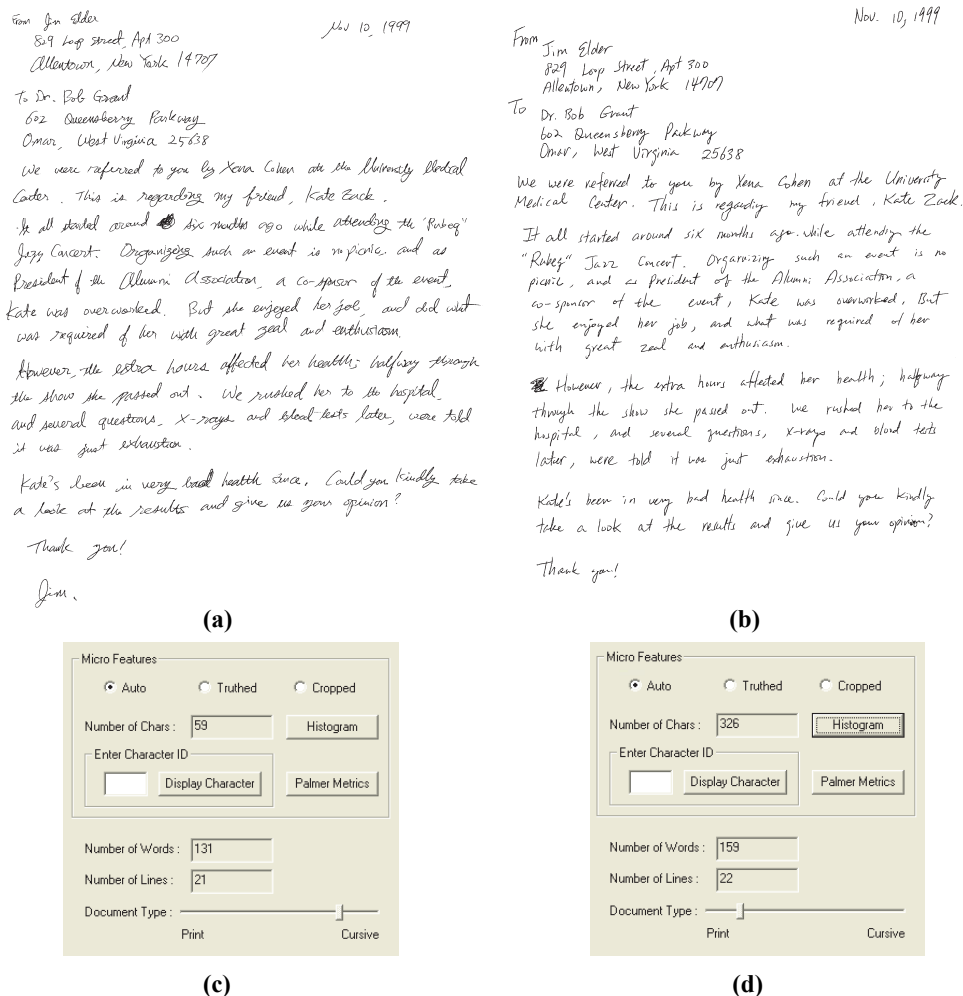


Figure 3. GUI for cursive/hand print recognition (a) predominantly cursive document, (b) predominantly hand print document; slider bar showing (c) predominantly cursive status, (d) predominantly hand print status

4. Experiments

We began by processing the CEDAR letter dataset in its entirety, generating our measurement on all documents to determine whether or not it had good predictive value. That is, we wanted to see how well the f_1 correlated with the Boolean feature value provided by the document examiners.

The ratio had a range of 0.10 to 6.86 over all 4701 documents. Of the hand printed documents, the range was 0.20 to 6.36 with a mean of 2.14. Cursive documents have values ranging from 0.10 to 3.05 with a mean of 0.57. Figure 4 shows a stacked histogram of f_1 (99% of the data is included; the few points with values > 3.2 are

omitted for space) as well as the data modeled with Gamma distributions. We determined a threshold by using the midpoint of the means of the two subsets (cursive and hand print) in the training set. The mean values of the cursive and hand print training sets were found to be 0.57 and 2.08 respectively. The threshold was set at 1.33. We performed a second experiment by modeling the two subsets with two Gamma distributions and choosing the class with higher probability; for the experiment, we approximated the distributions using Gaussian distributions.

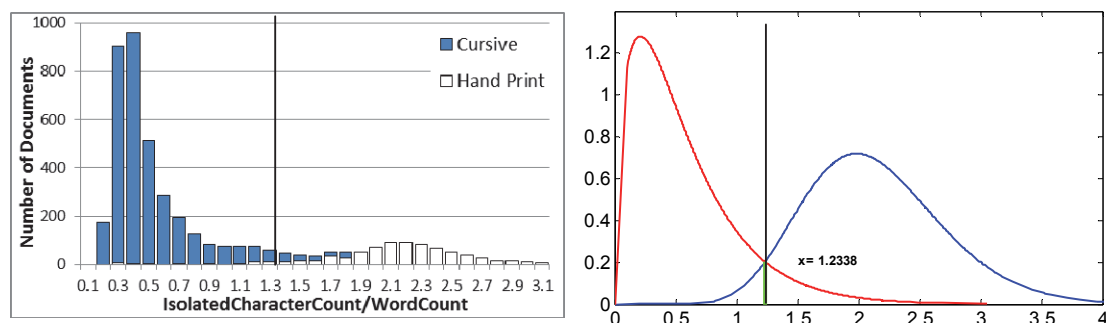


Figure 4. Histogram of binned feature f_1 (left) and representative Gamma distributions (right) with their thresholds

5. Results

1,775 documents in the cursive validation set were found to have values below the threshold, yielding the correct classification of 90.8%. 375 of the hand printed validation set were found to have values above the threshold, yielding correct classification of 92.0% of the documents. Overall, this led to the correct classification of an average of 94.5%. The Gaussian experiment yielded very similar results with performance improving on the hand printed set to 95.0% and decreasing slightly to 90.5% on the cursive set.

6. Conclusions and Future Work

Feature f_1 identifies the nature of documents correctly roughly 92.8% of the time, making a suitable feature for determining whether or not a document is written in a cursive or hand printed fashion.

In the future, we will present a detailed analysis of the above-mentioned f_2 feature as well as an analysis of a combination of these features. This will hopefully capture the cursive or hand printed nature of the handwriting by reflecting such things as the loops present in the writing. While we achieved success with a single feature, adding additional features affords us the opportunity to create more complex models and capitalize on multiple indications pointing to whether a document is cursive or hand printed.

Acknowledgements

This work was supported in part by the Department of Justice under 2010-DN-BSK037. The authors wish to thank Kirsten Singer and Traci Moran for providing valuable assistance in ground-truthing the data.

References

- Armistead, T. (1984). Issues in the Identification of Handprinting—A Case Study in Anonymous Death Threats, *Journal of Police Science and Administration*, 12(1), 81-98.
- Conway, J. (1955). The Identification of Handprinting, *The Journal of Criminal Law, Criminology, and Police Science*, 45(5), 605-612.
- Favata, J. T (1996). Character model word recognition. Proc. Fifth International Workshop on Frontiers in Handwriting Recognition, Essex, England, September, pp. 437-440.
- Huang, C. & Srihari, S.N. (2008). Word Segmentation of Off-line Handwritten Documents. Proc. Document Recognition and Retrieval (DRR) XV, IST/SPIE Annual Symposium, San Jose, CA (USA). January. SPIE Vol. 6815.
- Kam, Moshe & Lin, Erwei. (2003). Writer identification using hand-printed and non-hand-printed questioned documents. *Journal of forensic sciences*, 48.
- Kuzhinjedathu, K., Srinivasan, H., Srihari, S.N. (2008). Robust line segmentation for handwritten documents. Proc. Document Recognition and Retrieval (DRR) XV, IST/SPIE Annual Symposium, San Jose, CA (USA). January, SPIE Vol. 6815.
- Srihari, S.N., Cha, S., Arora, H., Lee, S. (2002). Individuality of Handwriting. *Journal Of Forensic Sciences*, pp. 856–872.
- Srihari, S.N., Srinivasan, H., Desai, K. (2007). Questioned Document Examination using CEDAR-FOX. *Journal of Forensic Document Examination*, 18(2), 1-20.
- Srihari, S.N., Zhang, B., Tomai, C., Lee, S., Shin, Y.-C. (2003). A System for Handwriting Matching and Recognition. Proc. Symp. Document Image Understanding Technology, Greenbelt, MD, April, pp. 67-75.
- Srikantan, G., Lam, S., Srihari, S.N. (1993). Gradient-based contour encoding for gray-scale character recognition. Proc. SPIE, San Jose, CA, 1906(2).