

Individuality of Numerals

Sargur N. Srihari, Catalin Tomai, Sangjik Lee and Bin Zhang
 Center of Excellence for Document Analysis and Recognition (CEDAR)
 Department of Computer Science and Engineering
 520 Lee Entrance, Suite 202, Amherst, NY 14228-2567
 {srihari,catalin,sjl1,binzhang}@cedar.buffalo.edu

Abstract

The analysis of handwritten documents from the viewpoint of determining their authorship has great bearing on the criminal justice system. In many cases, only a limited amount of handwriting is available and sometimes it consists of only numerals. Using a large number of handwritten numeral images extracted from about 3000 samples written by 1000 writers, a study of the individuality of numerals for identification/verification purposes was conducted. We have studied the individuality of numerals using clusterization and measured the numerals discriminability for verification purposes. The study shows that some numerals present a higher discriminatory power and that their performances for the verification/identification tasks are very different.

1. Introduction

The Forensic Document Examiner(FDE) often examines documents that contain a limited amount of handwriting. Sometimes the available handwriting includes only numerals and very few letters (e.g. records of drug dealing transactions, bank checks, tax forms, etc). The identification of numerals is therefore very important in many investigations and legal actions.

The individuality of handwriting in general was studied and reported upon by us in [1]. Audrey Giles ([2]) carried out the study of the degree of variability that may be observed in handwritten figures of a single individual. He also grouped different forms of numerals to illustrate the variability of handwriting (see Figure 1).

Ordway Hilton ([3]) suggested that a writer can be identified by his numerals simply by comparison with representative samples.

Edwin Alford in [4] observed that although an individual often takes great pains to disguise or distort his writing, he may completely overlook this in the formation of numerals

and therefore identification through comparison of numbers is important for FDE's. Jan Seaman Kelly ([5]) also observed that if a writer disguises the writing on a requested exemplar, the numbers are usually written naturally, free of disguise. In her study on 200 writers she suggested that naturally written numbers tend to be individualized and discussed the methods used by a writer asked to disguise the numbers in a requested exemplar.

Ellen Mulcrone Schuetzner ([6]) cataloged variation of numbers in both hand printing(manuscript) and handwriting(cursive) form.

James M. Davidson and Jon A. Keckler in [7] did a study of disguised handwriting, handprinting and numerals and showed the most common forms of disguise used.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0 |
|---|----|---|----|----|----|---|---|----|---|---|
| a | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0 |
| b | 11 | 2 | 3 | 44 | 55 | 6 | 7 | 8 | 9 | 0 |
| c | 11 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 0 |
| d | | 2 | 3 | 4 | 55 | | 7 | 8 | 9 | 0 |
| e | | 2 | 3 | 4 | 5 | | 7 | 88 | 9 | 0 |
| f | | 2 | 33 | 4 | 5 | | 7 | | | |
| g | | 2 | | | | | | | | |
| h | | 2 | | | | | | | | |

Figure 1. Different forms of numerals (Table taken from [2])

All these studies have experimentally shown the usefulness of the study of handwriting individuality. The goals of this work therefore are:

- Evaluate the individuality of numerals using automatic and manual clusterization.
- Measure the discriminatory power of numerals for writer verification.

- Determine the differences in discriminatory power among different numerals for a writer verification/identification system.

The paper is structured as follows: In Section 2 we presents the results of automatic and manual clustering of numerals In Section 3 we use the Bhattacharya distance to measure the discriminability of numerals for writer verification. In Section 4 we present the experimental framework, the verification/identification results obtained for different numerals and the analysis of these results. Section 5 contains some conclusions and ideas for future work.

2. Clustering of numerals

In his 1999 paper, Giles ([2]) observe that on his set of 110 writers, a significant number of them use two forms of numerals and only a very small number use three or more forms. For example, over 90% of writers used only one form for numerals 1 and 7 while less that 46% used only one form for numeral 6.

Our goal here is to estimate the individuality of numerals written by a much larger group of writers by measuring how many times the samples belonging to the same writer are distributed to the same or different clusters.

2.1 Automatic Clustering

We visually explore and evaluate the individuality of numerals by means of clusterization.

Suppose we have a set of sample numerals written by several writers. We want to see if (i) a clustering algorithm is able to automatically distribute the samples so that the samples belonging to the same writer are grouped in the same cluster (ii) the number of forms (ways) a person writes a numeral differs from numeral to numeral.

To evaluate this scenario we considered a set of 1000 writers and have obtained 3 image samples of every numeral from each, thus having a total of 3000 image samples for each of the 10 numerals. After transforming the images into feature space using WMR features (explained below), we have applied a clusterization algorithm ([8]) on the resulting feature vectors. Figure 2 displays several numerals belonging to one of the clusters obtained for character "2".

The first plot of Figure 3 presents the results of running the clusterization algorithm on the samples for character "2". On the x axis we have the imposed number of clusters and on the y axis we have the percentage of writers that have their three samples grouped in 3 different clusters, 2 different clusters or 1 cluster. As expected, increasing the number of clusters also increases the number of writers that have their samples distributed to different clusters. However, the results confirm the findings in [2], for an imposed

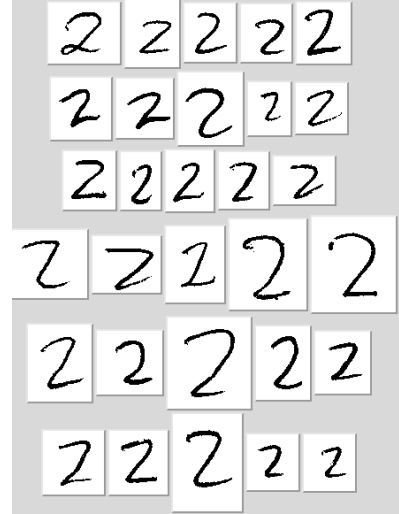


Figure 2. Image samples from one of the clusters for numeral "2"

number of 3 clusters, about 50% of writers have all their three samples grouped in one cluster (indicating that they use one single form), 40% into two clusters and less that 5% into three clusters.

The second plot of Figure 3 presents the percentages for different characters for maximum number of clusters fixed to 5. For digits "4" and "8" for example we have a higher percentage of writers with their samples distributed to 3 clusters, indicating a higher individuality. It is harder to group together image samples coming from the same writer for these numerals.

WMR features The WMR features are a set of 74 features composed of 2 global features: aspect ratio and stroke ratio of an entire character and 72 local features obtained as follows: each (numeral) image is divided into 9 sub-images; the distribution of the 8 directional slopes for each subimage form the set of $72(8 \times 3 \times 3)$ features. $F_{l_{i,j}} = s_{i,j}/N_i S_j$ for $i = 1, 2, \dots, 9$ and $j = 0, 1, \dots, 7$ where $s_{i,j}$ = number of components with slope j from subimage i , where N_i = number of components from subimage i and $S_j = \max_i(s_{i,j}/N_i)$

2.2 Manual Clustering

Here we have manually categorized the image samples into groups corresponding to the forms described in Figure 1. For example, for numeral '2' we have obtained the following sizes for each of the clusters:

Even if we don't obtain the exact correspondence between our clusters and Figure 1 groups, the differences between numeral forms were evident enough to warrant their clustering into groups of significant size, indicating a rather

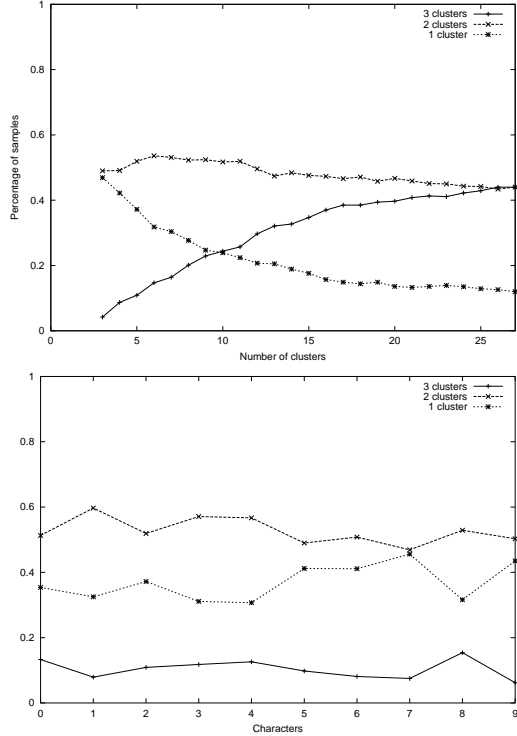


Figure 3.

| a | b | c | d | e | f | g | h |
|-----|-----|------|---|---|---|-----|---|
| 334 | 202 | 1411 | 6 | 7 | 3 | 413 | 0 |

Table 1. Sizes of manually obtained clusters that correspond to the forms from Figure 1 for numeral 2

high individuality for numeral 2 for example (larger number of clusters - higher individuality).

3. Measure of discriminability of numerals for writer verification

We can assume that, in general, two 0's written by two different writers are more probable to be confused than let's say two 8's, simply because the numeral 8 is more complex than "0" (more strokes, curves, etc). Next we are going to verify this hypothesis experimentally.

Let's consider N characters c^i with $i = 1, N$ from an alphabet A of M characters, with $N \leq M$. We also consider W writers and from each writer' set of handwriting samples we extract K instances of a certain character c^i . The k 'th character image sample from the set of images extracted for writer w is $c_{w,k}^i$.

Once we convert each character image into a feature vector we compute sets of similarity distances of two types: $D_{w,w}^i = d(c_{w,l}^i, c_{w,m}^i)$, for $l \neq m$ - Distances between character images belonging to the same author $D_{w,x}^i = d(c_{w,l}^i, c_{x,m}^i)$, for $w \neq x$ - Distances between character image belonging to different authors

We approximate the histograms of $D_{w,w}^i$ and $D_{w,x}^i$ distances by normal distributions and we compute their mean and standard deviations. For each character c_i we have computed 6500 distances: 1500 $D_{w,w}^i$ and 5000 $D_{w,x}^i$.

Figure 4 displays the $D_{w,w}^i$ and $D_{w,x}^i$ distributions for characters 0 and 4. The confusion is measured by the amount of overlap between the two distributions. If we have two arbitrary distributions p_1 and p_2 , at any position x , the probability of classification error is given by

$$P(error) = \min(p_1(x), p_2(x)) / (p_1(x) + p_2(x))$$

An upper bound on $P(error)$ is given by the Bhattacharyya distance which measures the separation between two pdf's:

$$d_B = \int_{-\inf}^{\inf} p_1(x)^{1/2} p_2(x)^{1/2} \quad (1)$$

Assuming that the distributions are Gaussians (like the ones in the plots of Figure 4) the formula becomes:

$$d_B = \frac{1}{2} \sum_{i=1}^N \frac{\sigma_{i,1}^2 + \sigma_{i,2}^2}{2\sqrt{\sigma_{i,1}^2 \sigma_{i,2}^2}} + \frac{1}{2} \sum_{i=1}^N \frac{\mu_{i,1} - \mu_{i,2}}{\sigma_{i,1}^2 + \sigma_{i,2}^2} \quad (2)$$

The larger the overlap between the distributions for a certain character, the higher the uncertainty regarding the authorship of two images of that character. Therefore, numerals that are more "discriminatory" have a lower Bhattacharyya distance. Table 2 presents the ranking of the numerals based on these distance. As expected, numerals like 2 and 6 figure at the top of the discriminatory ability list, while others like 7 and 0 are at the bottom of it.

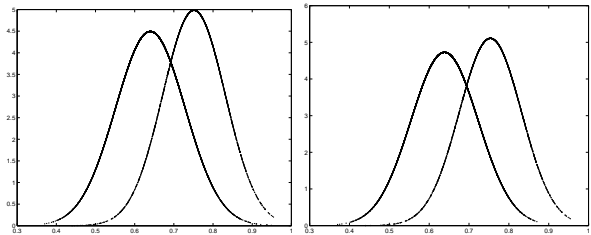


Figure 4. Same-writer and different-writer distances frequency distributions for the characters 0 and 2

| Digit | Distance |
|-------|----------|
| 2 | 0.7771 |
| 5 | 0.7785 |
| 3 | 0.7790 |
| 8 | 0.7794 |
| 6 | 0.7817 |
| 1 | 0.7947 |
| 4 | 0.7972 |
| 7 | 0.7972 |
| 9 | 0.8008 |
| 0 | 0.8027 |

Table 2. The ranked list of numerals and their corresponding Bhattacharya distance

| Digit | Verification | Identification |
|-------|--------------|----------------|
| 0 | 74.8 | 6.8 |
| 1 | 70.6 | 1.8 |
| 2 | 79.6 | 12.0 |
| 3 | 77.5 | 11.4 |
| 4 | 79.6 | 12.9 |
| 5 | 79.8 | 13.8 |
| 6 | 78.4 | 10.3 |
| 7 | 76.7 | 11.2 |
| 8 | 76.8 | 10.9 |
| 9 | 78.9 | 9.9 |

Table 3. Identification and Verification performance for individual digits

4 Writer Verification and Identification

4.1 Experimental Framework

More than 1000 people wrote three samples of the CEDAR letter, each. For every writer we have extracted 3 samples of each numeral. Each image is then transformed into a micro-features (detailed in [9]) vector composed of binary values. For verification the distance between two numeral images is given by a real-valued distance vector computed using a similarity distance between the two binary vectors. For identification, each digit is described by its 512-binary vector.

In the following experiments, for verification, the training and testing sets consist of 1500 same-writer and 5000 different-writer distances (for each digit). The identification results were obtained on the whole set of 3000 documents.

We have estimated the verification/identification performance in two different scenarios: (i) individual performance of digits (ii) accumulated performance of digits - starting with digit 0, we add 1,2,...,9 and measure the performance on each set.

4.2 Results

The results obtained for the both verification and identification tasks are presented in Table 3. The variance in performance from digit to digit is evident. As expected the identification performance of digits is much smaller than the verification performance (many-class decision compared with two-class decision). Also, the ranking of digits based on their individual performance is similar, with digit 5 presenting the best performance for both identification and verification.

The results obtained when using a variable number of digits are presented in Table 4. The more digits we have the higher the performance in both cases.

| AD | V | I |
|-----|------|------|
| 0 | 74.8 | 6.8 |
| 0-1 | 70.6 | 8.8 |
| 0-2 | 76.2 | 26.1 |
| 0-3 | 74.7 | 37.1 |
| 0-4 | 79.2 | 50.6 |
| 0-5 | 78.7 | 60.6 |
| 0-6 | 78.9 | 68.7 |
| 0-7 | 79.3 | 73.2 |
| 0-8 | 79.0 | 73.6 |
| 0-9 | 79.6 | 78.4 |

Table 4. Verification and Identification performance for accumulated digits

5 Conclusions

We have studied the individuality of digits for handwriting verification/identification using different methods. Automatic clustering was used to evaluate the distribution of digit samples written by the same author into different clusters that may or may not represent different forms of writing a certain digit.

The Bhattacharya distance was used to measure the individuality of digits for verification purposes. Then, the individual and accumulated verification/identification performance of digits was evaluated on a large data set.

From these experiments we can draw several conclusions: (i) different digits have different discriminability power for authorship detection purposes (ii) this variance in individuality is generally given by the number of forms that a certain digit may possess (iii) digits perform rather poor for verification and extremely bad for identification (iv) accumulation of digits is essential for obtaining a reasonable decision. Even then, the addition of some digits may have a negative impact on the verification/identification rate.

6 Acknowledgment

This work was supported partly by the National Institute of Justice grant 2002-LT-BX-K007.

References

- [1] Sargur N. Srihari, Sung-Hyuk Cha, Hina Arora, and Sangjik Lee. Individuality of handwriting. *Journal of Forensic Sciences*, 47(4):1–17, July 2002.
- [2] Audrey Giles. Figuring it out. *Journal of the American Society of Questioned Document Examiners*, pages 62–64, 1999.
- [3] Ordway Hilton. Identification of numerals. *International Criminal Police Review*, pages 245–250, October 1970.
- [4] Edwin A. Alford. Identification through comparison of numbers. *Identification News*, 15(7):11–14, July 1965.
- [5] Jan Seaman-Kelly. The examination of disguised numbers. *Journal of Forensic Sciences*, pages 1027–1028, September 1999.
- [6] Ellen Mulcrone Schuetzner. Class characteristics of numbers. In *American Academy of Forensic Sciences*, Reno, Nevada, February 2000.
- [7] James M. Davidson and Jon A. Keckler. A study of disguised handwriting, handprinting and numerals. In *The 39th annual conference of the American Society of Questioned Document Examiners*, Houston, Texas, August 1981.
- [8] John A. Hartigan. *Clustering Algorithms*. John Wiley & Sons, 1975.
- [9] G. Srikantan. Gradient representation for handwritten character recognition. In *IFWHR III*, pages 318–323, 1993.