

Online Text-independent Writer Identification Based on Temporal Sequence and Shape Codes

Bangy Li and Tieniu Tan

Center for Biometrics and Security Research, National Lab of Pattern Recognition,
Institute of Automation, Chinese Academy of Science, P.R. China
{byli, tnt}@nlpr.ia.ac.cn

Abstract

In this paper we present a novel method for online text-independent writer identification. Most of the existing writer identification techniques require the data to be from a specific text which is not applicable to cases where such text is not available, such as in criminal justice systems when text documents with different content need to be compared. Text-independent approaches often require a large amount of data to be confident of good results. We propose temporal sequence and shape codes to encode online handwriting. Temporal sequence codes (TSC) are to characterize trajectory in speed and pressure change in writing, and shape codes (SC) are to characterize direction of trajectory in writing handwriting. For TSC, we use two different codes to encode speed and pressure to code book: stroke temporal sequence codes (STSC) and neighbor temporal sequence codes (NTSC). At identification stage, we implement decision and fusion strategy to identify writer. Experimental results show that our proposed method can improve the identification accuracy with a small number of characters. Moreover, we find that the proposed method is even effective for cross-language (English & Chinese) writer identification.

1. Introduction

Pen-based interfaces have been recently popularized by the development and commercialization of mobile terminals such as PDA, electronic pen, electronic paper, or tablet PC notebooks[12]. Secure and automatic personal identification is becoming an important problem and writer identification using handwriting is based on the hypothesis that people write uniquely and can be characterized based on the information present in their handwriting. Automatic writer identification systems can be useful in a variety of applications including banks, criminal justice systems, de-

termining the authenticity of handwritten documents, etc. A writer identification system performs a one-to-many search in a large database with handwriting samples of known authorship and returns a likely list of candidates. The system must therefore be learnt from a set of handwriting samples of each individual candidate.

A comprehensive review covering the research work in automatic writer identification until 1989 is given in [10]. Recently, there have been efforts to automatically identify different writing styles. Scarce search results can be found along online text-independent handwriting. Individuality information can be present at various levels in handwritten data. Pitak et al. [13] propose point level information based on velocity of barycenter of pen-point movement. Primitive level information such as size, shape and style have been used by Namboodiri et al. [7] based on information about shape of curve. Yu et al. [15] build Gaussian models based on stroke level information such as dynamic features. Yasushi et al. [14] propose character level information that a different text including ordinary characters is used on every occasion of verification. Liwicki et al. [5] combine two levels of features about point and stroke level information extracted from a text line. In this paper, we propose a new method in which one can determine the identity of a writer based on a general stroke model and different combination of strokes types.

The paper is structured as follows. The next section presents motivation in text-independent handwriting. In Section 3 we describe three different codes as features. In Section 4, we present decision and fusion for identification. Experimental results and discussions are presented in Section 5. Finally, in Section 6 we conclude the paper.

2. Motivation

Our previous method builds Gaussian models each of which could be represented by its mean value and variance in 12 primary stroke types [15] and use the stroke's proba-

bility distribution function (SPDF) of four dynamic features to specify writer individuality, which aims to capture distribution of more dynamic aspects of the writing behavior of an individual [3]. However, this processing will lose information of dynamic attributes and the continuity of the writing. It also leads to inaccurate stroke segments due to the influence of emotions and physical state of writer.

The method proposed by Bulacu et al [6] do not consider dynamic features and the method proposed by Liwicki et al [5] need an amount of characters to build model. On-line handwriting has more information than offline handwriting in dynamic attributes. Firstly we consider TSC to characterize speed and pressure in temporal sequence relation and SC to characterize direction of trajectory in handwritten text. Secondly, we implement decision strategy in different scripts and fusion strategy in same scripts to identify writer. NTSC are to characterize relationship in the neighboring points and STSC are to characterize relationship in the strokes. Finally, SC are to exploit the relationship among three points containing more information in handwriting.

We consider dynamic features based on temporal sequence codes, which characterize dynamic relationship in temporal sequence, and SC is to encode shape information in orientation and slant to give a very stable characterization of individual handwriting style. This method effectively uses dynamic attributes of different temporal sequences and directions in the writing trajectory of handwritings.

3. Feature Extraction

The handwriting data is collected by Wacom Intuos2 tablet [2] which can be encoded as time-varying parameters such as x and y components concerning the pen-position at time t_i , the status of pen-down or pen-up s , the pen-pressure pre , the pen-altitude ϕ and pen-azimuth φ , and the raw handwritten data is represented as $\{x, y, t, s, pre, \phi, \varphi\}$ at each sampling point. Data collect can produce some spurious signal and noise with tablet PC, so preprocessing is included to remove spurious signals and noise in order to acquire effectively information in writing trajectory to enhance the input data.

3.1. Feature Extraction with NTSC

A stroke starts with a pen-down movement of the pen and ends with the next pen-up movement with the value preformed depending on the threshold value of $pr(t_i)$ whether the pen leaves (pen-up) or leaves (pen-down) the tablet surface. For a given stroke S consisting of point P_i to P_{i+n} , we compute the change of pressure in two points $PR_i = |pre_{i+1} - pre_i|$ and velocity $V_i = d(P_{i+1}, P_i)$ in two points as Fig. 1 shows.

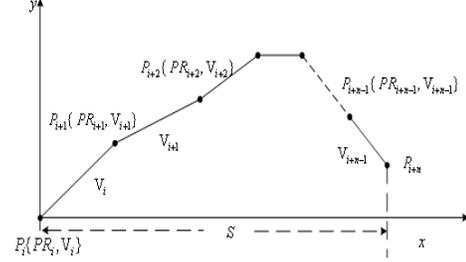


Figure 1. Dynamic attributes on a stroke.

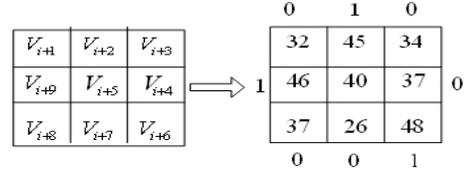


Figure 2. The velocity code of NTSC

For NTSC, a stroke includes sequence points $(P_{i+1}, P_{i+2}, \dots, P_{i+n})$. Firstly we compute the velocity code of NTSC (VCN) as follows:

$$VCN_i = \begin{cases} 0 & \text{if } V_i < V_{i+k_{nt}}, \\ 1 & \text{if } V_i \geq V_{i+k_{nt}}. \end{cases} \quad (1)$$

where k_{nt} controls code book sizes in handwriting ($i = 0, \dots, 2k_{nt}$). A NTSC code is encoded by $2k_{nt}$ bit binary value of temporal sequence. Fig. 2 shows a VCN is similar to the LBP operator [8] when $k_{nt} = 4$, a VCN becomes $\{01001001\}$ according to temporal sequence and the code book size is 2^{2*4} . Pressure code of NTSC (PCN) is computed similarly. We define two features (f_1, f_2) by distribution of VCN and PCN. f_1 and f_2 have 2^{2*4} dimensions.

3.2. Feature Extraction with STSC

For STSC, a stroke includes sequence points $(P_{i+1}, P_{i+2}, \dots, P_{i+n})$. Firstly we compute the velocity code of STSC(VCS) as follows:

$$VCS_i = \begin{cases} 0 & \text{if } V_i < \text{mean}(V(S)), \\ 1 & \text{if } V_i \geq \text{mean}(V(S)). \end{cases} \quad (2)$$

where $\text{mean}(V(S))$ is the average velocity in a stroke S , and k_{st} controls code book size in handwriting ($i = 0, \dots, 2k_{st} - 1$). A VCS is encoded by $2k_{st}$ bit binary value of temporal sequence. Fig. 3 shows velocity of temporal sequence encode a VCS when $k_{st} = 4$, a VCS becomes $\{01011011\}$ by temporal sequence and the code book size is

2^{2*4} . Pressure code of STSC (PCS) is computed similarly. We define two features ($f3, f4$) by distribution of VCS and PCS. $f3$ and $f4$ have 2^{2*4} dimensions.

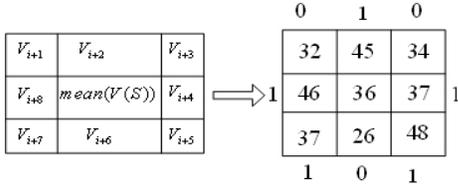


Figure 3. The velocity of code in LTSC

3.3. Feature Extraction with SC

We extract shape code based on strokes that mean a connected component from pen-down to pen-up. Thus a stroke is a sequence of points during a certain time interval between the time when the pen-tip touches the panel to the time it leaves. According to writing direction of each point in a stroke, we define $k_{sc} = 12$ angles as Fig. 4(a) shows, then we define SC as direction types in two points and angle computed by

$$\theta = \arctan((y_{k+1} - y_k)/(x_{k+1} - x_k)) \quad (3)$$

Therefore two directions (three points) form a SC and a large number of SC (as Fig. 4(b) shows) can be derived from each handwriting sample. The primitives can cover over all strokes according to the direction models and shape code is extracted from the whole handwriting data. The total number of SC consist of $(12 * 12)$ codes. Fig. 5 shows 144 shape codes.

We use the distribution of SC as feature $f5$ to characterize writer characteristics, which can also be interpreted as the transition probability between primitives in a simple Markov process.

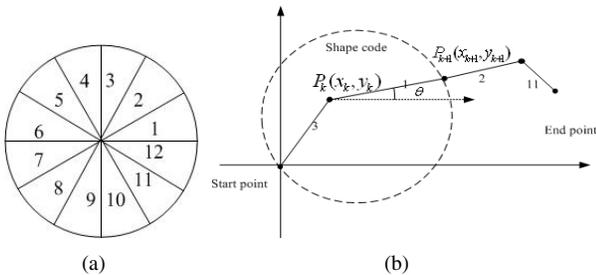


Figure 4. Shape code in handwriting (a) Model of direction. (b)A SC in a stroke.

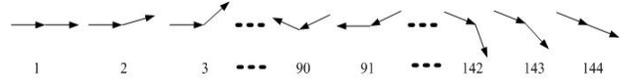


Figure 5. The SC of handwriting .

4 Decision and Fusion Strategy for Writer Identification

The identification of writers based on given feature vectors is a typical pattern recognition problem. Writer identification is performed using the the nearest neighbor classifier in a "leave-one-out" strategy. A large of number of distance measures are tested in our experiments: L1, L2, Cosine-Angle (CA), Chi-Square (CS) and Diffusion-Function (DF) distances [9] [11] [4]. Only the best performing ones (DF) are reported. For online independent handwriting of same scripts, firstly, we use distance of SC feature by a predefined decision threshold T to judge kinds of scripts and simple weighted distance to combine all features in same scripts, otherwise for the different scripts, we combine $f1, f2, f3$ and $f4$ to identify writer. Fig.6 shows decision and fusion

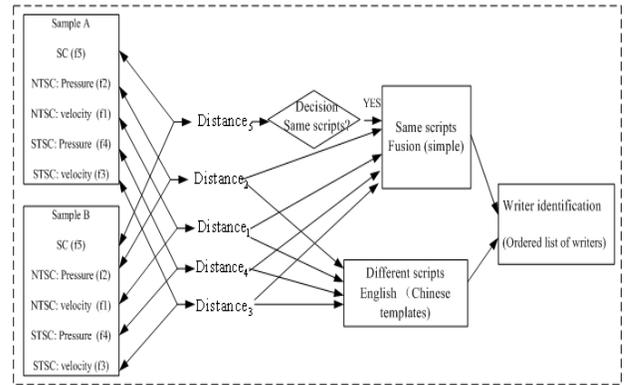


Figure 6. A decision and fusion strategy for writer identification

strategy for writer identification.

5 Experimental Results and Discussions

Our experiments are based on the NLPR handwriting database[1]. In this database, samples are collected using a Wacom Intuos2 tablet. The tablet has 1024 levels of pressure sensitivity, and its maximum data rate is 200 pps. This databases contains more than 1500 handwritten texts in on-line format from over 242 writers in two sessions. Each writer writes eight pages of texts which include four pages of Chinese texts and four pages of English texts respectively. In the first session, each subject is required to write

six pages of handwriting about 40 characters, i.e. one for fixed content in Chinese, two for free content in Chinese, one for fixed content in English, and two for free content in English. In the second session, each subject is required to freely write one page of Chinese and one page of English about 40 characters respectively. Therefore, our databases contain three sub-datasets: dataset I (DB I) includes two sessions in Chinese text. Dataset II (DB II) includes two sessions in English. we merged Chinese and English dataset to obtain dataset III (DB III).

5.1 Performances of Same Scripts

For DB I (Chinese Database), we first use same Chinese texts as templates, different Chinese texts as testing samples. Then we also use different Chinese texts as templates and different Chinese texts as testing samples. Fig.7 gives average performance of writer identification in different Chinese templates based on Chinese text. For DB II (English Database), Fig.8 gives average performance of writer identification in different English templates based on English text.

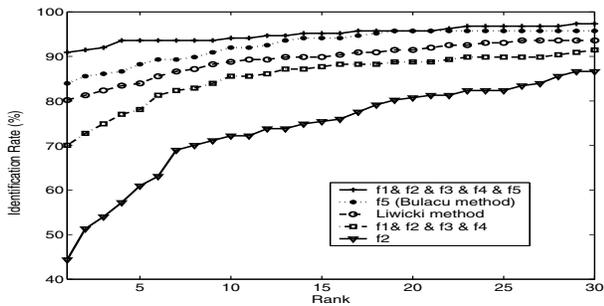


Figure 7. Average performance of writer identification in different Chinese templates as a function of rank size for Chinese text.

5.2 Performances of Different Scripts

For database III, we use different scripts as template and different scripts to test our method for writer identification. We use English texts as templates and Chinese texts as testing or Chinese texts as templates and English texts as testing. Fig.9 shows performance of writer identification in different scripts

5.3 Discussions

The code book size of each set of features can influence identification results. We choose $k_{nt} = 4$ in NTSC and $k_{st} = 4$ in STSC from experimental results. For

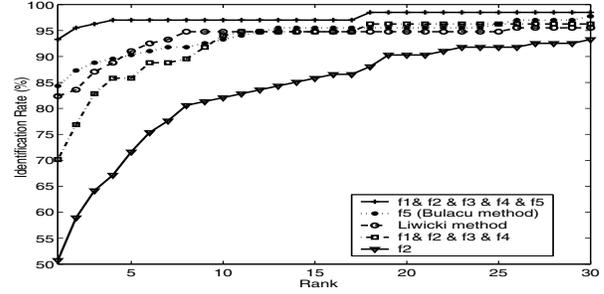


Figure 8. average performance of writer identification in different English templates as a function of rank size for English text. SC features with $k_{sc} = 18$.

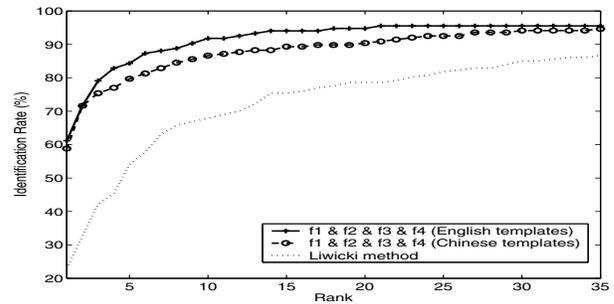


Figure 9. The identification results for using different scripts (SC features with $k_{sc} = 18$).

SC features, experimental results show SC features with $k_s = 12$ do not produce extra improvements over features with $k_s = 18$ in Chinese text. While we choose SC features with $k_s = 18$, the experimental results are better than those with $k_s = 12$ in English in that English scripts have more curves.

First, Fig.7 and Fig.8 show performance of combining all of features ($f1, f2, f3, f4$ and $f5$) is better than other combination of features. The important features ($f2$ and $f5$) are highly discriminative and give very satisfying results in writer identification. The method proposed by Bulacu et al [6] do not consider dynamic features that the experimental result are better than $f5$ ($f5$ feature is similar Contour-hinge PDF). The method proposed by Liwicki et al [5] need an amount of characters to train GMMS. Our method can improve performance and reduce the amount of characters required.

Second, as Fig.9 shows, we can see that the results of our methods are obviously better than other methods. The main reason is that the method proposed by Bulacu et al[6] did not consider writer identification of different scripts and the important feature (Contour-hinge PDF) is depend on at-

tribute of character that do not work in different scripts. The method proposed by Liwicki et al [5] need an amount of character and consider different structure of strokes in different scripts. For different scripts, Our method do not consider f_5 (SC) features and structure of strokes in different scripts, we only use dynamic feature to combine f_1 , f_2 , f_3 and f_4 in writer identification.

At last, we propose decision and fusion strategy for writer Identification based TSC and SC, which compute velocity and pressure in attribute of temporal sequence and shape of writing trajectory. While there are important differences in performance among different scripts and the best performer is DF distance in TSC.

Fig.10 shows our real-time online, text-independent writer identification system. When the system operates in the recognition mode, the rank-5 identities of the writer are listed on the status bar of the software interface when the user is writing, where circle indicates Chinese template and rectangle indicates English template. Writer lists will be change with the number of characters increasing. Our method is effective to identify writers in different scripts.

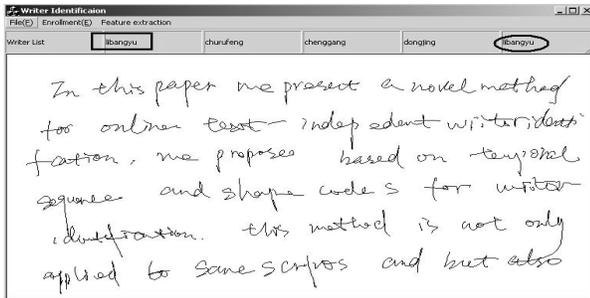


Figure 10. A real-time online writer identification system. Circle represents result of Chinese text as template, and rectangle represents result of English text as template.

6 Conclusions

In this paper, we have proposed a novel method based on temporal sequence and shape codes for online text-independent writer identification. It is an important attribute that dynamic features of relationship of velocity and pressure based on temporal sequence are to characterize writing rhythm and shape code is to characterize direction in writing trajectory. Experimental results show that these features based on TSC and SC are very stable personal characteristics and reveal individual writing style, and different combination of the dynamic features can apply in different scripts and improve the identification accuracy and reduce the amount of characters required.

Acknowledgments

This work is funded by research grants from the National Basic Research Program of China (Grant No. 2004CB318100), the National Natural Science Foundation of China (Grant No. 60736018, 60335010, 60702024, 60723005), the National Hi-Tech Research and Development Program of China (Grant No. 2006AA01Z193, 2007AA01Z162), and the Chinese Academy of Sciences.

References

- [1] NLPR database. <http://www.cbsr.ia.ac.cn/english/Databases.asp>.
- [2] Wacom website. <http://www.wacom.com/index.html>.
- [3] B.Li, Z.Sun, and T.N.Tan. Online text-independent writer identification based on stroke's probability distribution function. *Proc. of 2th ICB*, pages 201–210, 2007.
- [4] H.Ling and K.Okada. Diffusion distance for histogram comparison. *IEEE Conference on Computer Vision and Pattern Recognition*, pages 246–253, 2006.
- [5] L. Marcus, S. Andreas, B. Horst, B. Samy, M. Johnny, and R. Jonas. Writer identification for smart meeting room systems. In *Seventh IAPR Workshop on Document Analysis Systems, DAS*, number 70, 2006. IDIAP-RR 05-70.
- [6] M.Bulacu and L.Schomaker. Text-independent writer identification and verification using textural and allographic features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (4):701–717, April 2007.
- [7] A. Nambodiri and S. Gupta. Text independent writer identification from online handwriting. *Tenth International Workshop on Frontiers in Handwriting Recognition*, pages 131–147, October 2006.
- [8] T. Ojala, M. Pietikainen, and D. Harwood. A comparative study of texture measures with classification based on feature distributions. 29(1):51–59, January 1996.
- [9] R.O.Duda, P.E.Hart, and D.G.Stork. *Pattern Classification, Second Edition*. John Wiley & Sons, Inc., 2001.
- [10] R.Plamondon and G.Lorette. Automatic signature verification and writer identification -the state of the art. *Pattern Recognition*, pages 107–131, 1989.
- [11] T.Ahonen, A.Hadid, and M.Pietikainen. Face recognition with local binary patterns. *European Conference on Computer Vision*, pages 469–481, 2004.
- [12] S. M. Thierry Artieres and P. Gallinari. Online handwritten shape recognition using segmental hidden markov models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(2):205–217, 2007.
- [13] P. Thumwarin and T. Matsuura. On-line writer recognition for hta based on velocity of barycenter of pen-point movement. *Proc.of IEEE International Conference on Image Processing*, pages 889–892, 2004.
- [14] Y. Yamazaki, T. Nagao, and N. Komatsu. Text-indicated writer verification using hidden markov model. *Proc.of 7th ICDA*, pages 329–332, 2003.
- [15] K. Yu, Y. Wang, and T. Tan. Writer identification using dynamic features. *Pro.of ICBA*, pages 512–518, July 2004.