

# An Improved Feature Extraction Method for Individual Offline Handwritten Digit Recognition

Wang Qinghui<sup>1,3</sup>, Yang Aiping<sup>2</sup>, and Dai Wenzhan<sup>1</sup>

1) Department of Automatic Control, Zhejiang Sci-Tech University, Hangzhou, 310018

2) Zhejiang University of Finance & Economics Hangzhou, 310018 (Email: Aiping\_yang@hotmail.com)

3) Department of Electron & Electric, Longyan University, Longyan, 364000 (Email: wqh - 126com@126.com)

**Abstract**—Offline handwritten digit recognition (OHDR) is considered as one of difficult problems in the field of pattern recognition. Because it is a challenging computational problem mainly due to the vast differences associated with the handwritten patterns of different individuals. In this paper, a novel method of feature extraction is presented based on structural feature for OHDR by simulating the process of human recognizing handwritten digit. Firstly state and state value are introduced, then the steps of how to determine the eigenvalue is explained in detail, last the method is applied in OHDR, and the result show its effectiveness.

**Index Terms**—Offline handwritten digit recognition, feature extraction

## I. INTRODUCTION

OHDR involves the recognition of handwritten character patterns in digital image. Feature extraction is an important step of OHDR and a core of pattern recognition. It is obvious that results of the recognition depends on feature extraction. If the feature extraction is done well, the classifier with high performance can be easily designed. Therefore feature selection directly affect the classifier's formation and its effectiveness.

At present, character has to be decomposed firstly for many approaches of OHDR and then feature extraction is carried out. However, decomposing OHDR probably brings on the error or rejection recognition for the simplicity of the strokes of digital. In order to improve recognition rates, in this method the shape of the handwritten digital strokes are analyzed and is made an approximate to the standard digital strokes in many digital character recognition methods, then is searched for the information that included in line segment, obtuse angle, acute angle, circle, etc. So a large number of eigenvalues are extracted in order to identify the distinction. At present, a large number of recognition methods are put

forward in the academe, such as global feature extraction method [1], pixel-by-feature extraction method [2] skeleton feature extraction method [2], the vertical direction data feature extraction method [2], 13 points feature extraction method [2] and so on [3-9]. Some of these methods have a lot of eigenvalues, or some method's algorithms are very complexity, or some method's recognition rates are not high.

The human eye is the most efficient recognizer in the world. When a person identify a handwritten digit, he will not care for the number of the line segments and the number of obtuse angle and the number of acute angles in the handwritten digit, he just concerns that whether the character has circle, the position of circle, and the number of circles. He looks at the overall left and right concave shapes in characters, then he can efficiently identify the handwritten digit.

In this paper, the process that human eye identify character of handwritten digit is simulated, a novel rapid handwritten recognition method is put forward based on the overall eigenvalues of the character. There are not involved complex thinning procedures on the digital character image in this method, so the possibility of error and rejection will be reduced. Because of thinning distortion and complex stroke feature analysis are not needed. So the speed of this method is quite fast, and recognition rate is high.

## II. STATE AND STATE VALUE

During the process of OHDR, in order to distinguish the handwritten digits effectively, feature extraction has to be extracted. Firstly feature extraction is sampled according to the principle as follows: (1) As much as useful information should be include; (2) The abstraction method of the eigenvalues should be simple and fast; (3) The correlation degree among each eigenvalues should be as little as possible; (4) The quantity of eigenvalues should be as little as possible; (5) the anti-jamming ability in the eigenvalues should be according to the above all principles. A kind of new digital character eigenvalues and feature extraction method have been put forward, namely 3 Points Feature Extraction.

\*This work was supported by National High Technology Research and Development Program 863(No2009AA04Z139), Supported by Research Fund for the Doctoral Program of Higher Education of China (No: 20070338002), Natural Science Fund of Zhejiang Province (No: Y607556) Wen-zhan Dai: corresponding author and a professor of Zhejiang Sci-Tech University. Email: dwzhan@zstu.edu.cn

TABLE I  
THE STATE AND STATE VALUE

state	C-L	C-R	C-U	CC	S-C	W-C	NT
state value	1	2	3	4	5	6	0

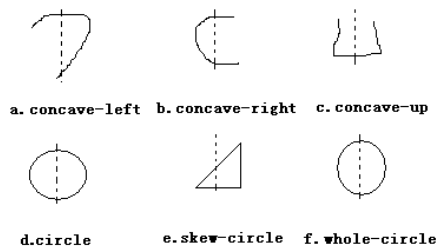


Fig. 1. The figure of states

According to our studying, the digital character can be decomposed into six basic states such as concave-left (C-L), concave-right (C-R), concave-up (C-U), circle (CC), skew-circle (S-C), whole-circle (W-C), nothing (NT), which is shown in Fig. 1. The corresponding values of states are shown in Table 1.

**C-L:** In this situation, there are two nodes when the broken line intersects through the real curve, the broken line is regarded as the boundary, the left area of boundary is not closed but the right area of the boundary is closed. The case is for that there is only two nodes, such as digit '7'. (The meaning of closed can be defined as follows: the broken line locates in the middle of the graph, the algorithm is designed to search line-by-line. For example, if it is found by the algorithm that there is no curve points in some line, it will be considered the part of area is not closed).

**C-R:** In this situation, there are two nodes when the broken line intersects through the real curve, the broken line is regarded as the boundary, the right area of the boundary is not closed but the left is closed.

**C-U:** In this situation, there is only one nodical which is located in the middle of the graph, it can be concluded from the algorithm both left and right area are closed.

**CC:** In this situation, there are two nodes, it can be concluded from the algorithm both left and right area is closed. The above node is located in the area of the top graph, and the below node is located in the area of the middle matrix rather than bottom of the matrix.

**S-C:** In this situation, there are two nodes. It can be concluded from the algorithm both left and right area is closed. The above node isn't located in the area of the top matrix, and the below node is located in the area of the middle matrix.

**W-C:** In this situation, there are two nodes. It can be concluded from the algorithm both left and right areas are

closed. The above node is located in the area of the top matrix, and the below node is located in the area of the bottom matrix (Top range is referred as 1/7 above digital picture matrix; the bottom range is referred as 1/7 under digital picture matrix).

**NT:** It is showed this there is nothing to be used for recognition.

Take some examples:

The state of digit '2' can be concluded from the algorithm, the above part of the character can be defined as 'C-L', the below part of the character can be defined as 'C-R'. The state of digit '4' can be concluded from the algorithm is 'S-C'. The state of digit '0' can be concluded from the algorithm is 'W-C'. The state of digit '5' can be concluded from the algorithm, the above part of the character can be defined as 'C-R', the below part of the character can be defined as 'C-L'.

### III. EXTRACT EIGENVALUE LAW IN HANDWRITTEN DIGIT

In this paper it is proposed that the digital eigenvalues information can be expressed by three eigenvalues which are respectively named as one letter. The 'A' is defined as the first eigenvalue, the 'B' is defined as the second eigenvalue, the 'C' is defined as the third eigenvalue. The eigenvalues are fixed according to the following steps:

Step 1:

The digital picture will be transformed into 1 or 0 in two-dimensional square matrix [10].

For examples: Fig. 2 (a) is the handwritten digit '4' before being dealt. Fig. 2 (b) is the picture after being dealt, the size of the dealt picture is reduced to the range that includes effective digital information. Fig. 3 is a matrix of handwritten digits '4'.

Step 2:

The number of the nodes which the boundary line intersects through the curves is calculated as the first eigenvalues A.

For examples: The number of node is one in handwritten digit '4' of Fig.4, so A=1; the number of nodes are 3 in handwritten digit '8', so A=3.

Step 3:

Handwritten digit is divided into two parts, upper-part and down-part, the boundary is in the middle of the matrix,

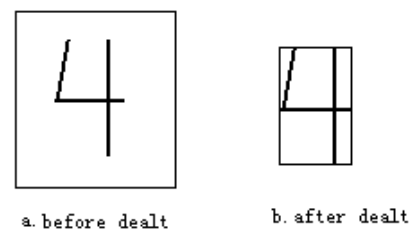


Fig. 2. the handwritten digit '4' is dealt with in earlier stage



Fig. 3. Two two-dimensional value matrix of the handwritten digit '4'



Fig. 4. Digit 4 and 8

and recognition process is carried on from top to bottom. Eigenvalues B represents the state of the upper-part of handwritten digits; and eigenvalues C represents the down-part of the handwritten digits.

Because eigenvalues B represents the upper-part of the handwritten digits, it is required that there is at least one node in the upper-part. In that case the recognition will be processed according to the above Table 1 to get the corresponding state value. Otherwise, it is shown that the part is no use for recognition, so let B=0. For example, the digit '4' in Fig. 5, B=5.

Eigenvalues C represents the down-part of the handwritten digits, it is demanded that there are at least two nodes in the down-part. In that case the recognition will be processed according to table 1 to gain the corresponding state value. Otherwise, it is shown that the part is no use for recognition so let C=0. For example, the digit '6' in Fig. 5, C=4.

For examples: The handwritten digit '4' in Fig. 2, because the node is one and it is located in the middle of the matrix, so B is valid and C is invalid. B=3, C=0; the handwritten digit '4' in Fig. 5, because the nodes are two, the upper-part characteristic state is 'S-C', and the down-part characteristic state is 'NT', so B=5, C=0; the handwritten digit '6' in Fig. 5, because the nodes are two, the upper-part characteristic

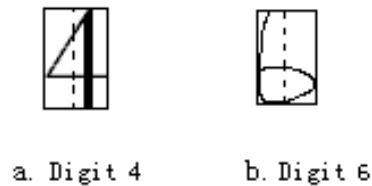


Fig. 5. Digit 4 and 6



Fig. 6. Digit 2 and 3

state is 'NT', and the down-part characteristic state is 'CC', so B=0, C=4; the handwritten digit '2' in Fig. 6, because the nodes are three, the upper-part characteristic state is 'C-L', and the down-part characteristic state is 'C-R', so B=1, C=2; the handwritten digit '3' in Fig. 6, because the nodes are three, the upper-part characteristic state is 'C-L', and the down-part characteristic state is 'C-L', so B=1, C=1.

#### IV. SIMULATION

M function in MATLAB is used to realize the extracting three eigenvalues. Eigenvalues ABC of common handwritten digits 0~9 of Fig. 7 can be extracted and shown as Table 2 and Table 3. It can be found that the nine digits have different corresponding values.

From Table 2 and 3, we can get the following several points:

TABLE II  
THREE CORRESPONDING EIGENVALUES(EV) OF DIGIT 0~4

Digit \ EV	0	1	2	3	4
A	2 2 2	1 1	3 3 3 3 3	3 3 3 3	3 1 1 2
B	6 6 6	0 0	1 1 1 1 1	1 1 1 1	5 3 3 5
C	0 0 0	0 0	2 2 2 2 4	1 1 1 4	1 0 0 0

TABLE III  
THREE CORRESPONDING EIGENVALUES(EV) OF DIGIT 5~9

Digit \ EV	5	6	7	8	9
A	3 3 3 3 3	3 2 3 3	2 2 1 2	3 4 4	3 3 2
B	2 2 2 2 2	2 0 2 2	1 1 1 1	4 4 4	4 4 4
C	1 1 1 1 1	4 4 4 4	0 0 0 0	4 4 4	1 1 0

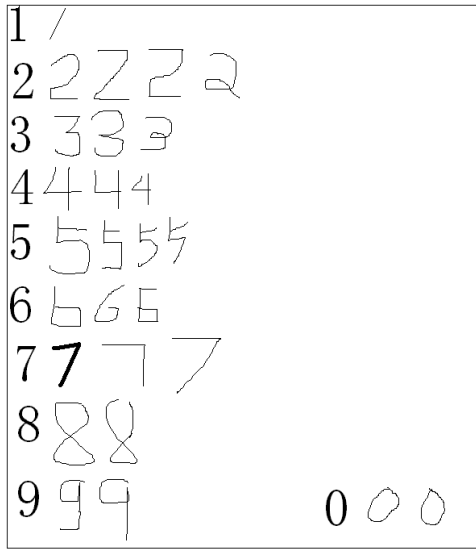


Fig. 7. Several hand-written forms of digit 0-9

(1) Though different writing style of digit 0~9 probably will results in different eigenvalues, the extracting value of ABC is strictly carried out in term of Table 1.

(2) The value of ABC depends on digits.

If 13 eigenvalues extracting method[2] is used to identify the handwritten digits in Fig.7, the eigenvalues to be extracted are so many, it will be fluctuated very much because character needs to be divided into 8 parts equally. If the method of the overall characteristic of digit[1] is adopt to recognize the digit, according to its algorithm, it is necessary to emit 8 radials from the directions of up, down, left, right, left, left-up, left-down, right-up, right-down, and then judges whether the ray cross character and the number of the rays which crosses characters.

If it had adopt the method of the overall characteristic of digit to recognize the digit, according to its algorithm, it has been required from all background points to emitting 8 radials from the directions of up, down, left, right, left, left-up, left-down, right-up, right-down, it needs to be judged the ray whether cross character, the number of the rays cross characters, the value of this background is set, so the calculating amount is very big, can not be identify reliably some handwritten digits that are offered in Fig. 8, for example, handwritten digits '4' with the characteristic of 'up open' can not be identified or the common number '8', it is unable to recognize.

Compared with the above methods, handwritten digital character ABC eigenvalues algorithm proposed in this paper is simple and with higher distinguishing rates.

## V. CONCLUSION

There are some weaknesses of low-recognition rates and complex-recognition process and low-execution rates in the

existing method of identifying the hand-written digits, this paper is based on the characteristic of OHDR, it is put forward a improved OHDR method, and it has been given the concrete procedures of the 3 Points Feature Extraction of OHDR, because this method hasn't been involved any complexity processing, only it needs to be extracted three eigenvalues of the effective digits, so with higher recognizing speed, and the simulating result has been confirmed the validity and swiftness of the algorithm.

## REFERENCES

- [1] F. Ye and F. Li, "Fast hand-written digital character recognition based on global feature", China: Computer Engineering and Design, 2006.2722.
- [2] L. H. Zhong and W. Hu, "A new Feature Extraction Method on Hand-written Digits Recognition System", China: Journal of university of Sichuan, 2007.44(5).
- [3] B. Liu, "Research on character Recognition Of License Plate Recognition System Based on Neural Networks", Wuhang, China: Academic dissertation of the master of Wuhan University, 2004.
- [4] S. Q. Shi and Y. J. Wang and X. Y. Tang, "Manuscript Chinese Ideograph Recognition Based on Extracting of Integral Form Character and Fuzzy Recognition", China: Computer Technology and Development, 2004.14(1).
- [5] R. Plamondon, S. Srihari, "On-line and off-line handwriting recognition: A comprehensive survey", IEEE Trans. Pattern Anal. Mach. Intell. 22(1)(2000)63-84.
- [6] J.P. Drouhard, R. Sabourin, M. Godbout, "Evaluation of a training method and of various rejection criteria for a neural network classifier used for off-line signature verification", in: IEEE Int. Conf. Neural Networks, 1994, pp. 294-299.
- [7] J. Richiardi, H. Ketabdard, A. Drygajlo, "Local and global feature selection for on-line signature verification document analysis and recognition", in: 8th International Conference on Document Analysis and Recognition, vol. 2, 2005, pp. 625-629.
- [8] M. Ammar, "Progress in verification of skillfully simulated handwritten signatures", Int. J. Pattern Recogn. Artif. Intell. 5(1)(1991)337-351.
- [9] N.A. Murshed, R. Sabourin, F. Bortolozzi, "A cognitive approach to off-line signature verification", Int. J. Pattern Recogn. Artif. Intell. 11(5)(1997)801-825.
- [10] M. Zhang and Z. Q. Yu and S. W. Yao, "Image pretreatment research in recognition of hand-written numerals", China: Micro-computer information, 2006, 22(6-1).