

A Digital Character Recognition Algorithm Based on the Template Weighted Match Degree

Mo Wenyong¹, Ding Zuchun²

¹Guangdong Zhongya Wireless Technology Co. Ltd.

²Guangdong Zhongya Wireless Technology Co. Ltd.

mowenyong@126.com, zuccding@gmail.com

Abstract. Template matching algorithm has the characteristics of high speed and real-time. This paper introduces template matching algorithm, and put forward an improved template matching algorithm which based on the weighted matching degree. After the completion of the pre-processing of input characters, the algorithm uses the moving match of the standard character template with respect to image character template. It uses a method of weighted matching degree. This algorithm can provide a higher matching rate of image character, and overcome the fallacious recognition produced by traditional calculation method. It not only guarantees the accurate recognition rate of general character, but also effectively avoids the influence of adherent noise and partial distortion, which greatly impacts the recognition rate of the character.

Keywords: character recognition, matching rate, weighted match

1 Introduction

The automatic meter reading system can take advantage of the camera to obtain the image information of the meters, and carry out a series of processing of the images. The processing contains several modules: character positioning, character segmentation, character normalization, and character recognition. During these modules, character recognition (OCR) plays a very important role in it and should be received our great attention [1-4]. In order to improve the recognition rate, currently, we have some common methods including template matching method, feature recognition method, peripheral contour method and so on. Methods used by most of the recognition system are based on template matching algorithm, and combined with feature recognition algorithm or other auxiliary algorithm, which have reached a high recognition level so far. However, we know that the quality of the images also depends on the external environment which is more or less different, so we can't guarantee the images received from the camera always keep in high quality. Because of uneven illumination, the phenomenon of partial image in different brightness is common. In this situation, in addition to Gaussian noise, the images contain irregular noise into a piece. After image pre-processing, most of them will appear adherent noise or characters of missing partial strokes. For this image, the recognition rate is still not very satisfactory, generally only reach 98%. This paper aims at this problem

and uses the method of template matching algorithm combined with algorithm of weighted matching rate to reduce the influence of adherent noise and characters of missing partial strokes. This method can strengthen the anti-interference ability of the system, and significantly improve the recognition rate of the template matching algorithm.

Section II introduces the basic knowledge of our algorithm. Section III will put forward the algorithm of digital character recognition based on the weighted matching degree of template matching method. In section IV, we will provide the actual experimental results which confirm the effectiveness and improvement of the algorithm..

2 The Basic Process of Template Match algorithm

The template matching method is one of the effective ways to achieve the classification of discrete input mode. It's essence is to measure a certain similarity between the input mode and the sample, then take maximum similarity as the category of the input mode, which is different from SVM and BP neural network classifier. It extracts features according to the intuitive form of the character and uses related matching principle to discern, that means the input character and standard character matching in a classifier. It is often very fast.

As an example, take account of the process of two-dimensional image as, the correlation matching algorithm is as following.

Let the input character to be presented by function $f(x,y)$ and the standard character to be presented by function $F(x,y)$. After contrast in the correlator, the output is $T(x,y)$. We use correlated variable quantity represent random variable quantity, and the correlated output will be

$$T(x_1-x_2, y_1-y_2) = \iint f(x, y) F[x+(x_1-x_2), y+(y_1-y_2)] dx dy \quad (2.1)$$

$$\text{When } x_1=x_2, y_1=y_2, T(0,0) = \iint f(x, y) F(x, y) dx dy \quad (2.2)$$

$$\text{When } f(x, y) = F(x, y), T(0,0) = \iint f(x, y) f(x, y) dx dy \quad (2.3)$$

(2.3) is the autocorrelation function of input character, and $T(0,0) \geq T(x, y)$. $T(0,0)$ is the main peak, and the secondary peak will be show in other standard character. As long as the secondary peak is not equal to the main peak, we can use the appropriate threshold to identify them.

Template matching method depends on different features obtained from the time of modeling, it can be divided into methods of image matching method, stroke analysis, geometric feature extraction method and so on. While modeling and matching comparison, image matching method getting matched based on graphic block itself, and the result of recognition is based on the degree of similarity. It is the method used in this paper.

In order to make data processing more rapid and accurate, we should do image preprocessing. The image preprocessing is a series of conversion processing on the input image, so the image meet the requirement of character module of recognition. The first step of the image preprocessing is that isolating the character from the input

A Digital Character Recognition Algorithm Based on the Template Weighted Match Degree

grayscale image. The most common methods is that binarization, denoising, positioning, segmentation. Image binarization means that to convert the grayscale into the image that containing only two type of gray value, which needs an appropriate threshold to distinguish between background and character. The mathematical expression is as follows:

$$f(x, y) = \begin{cases} 0, & f(x, y) < T \\ 1, & f(x, y) \geq T \end{cases} \quad (2.4)$$

T is the threshold value of binarization.

Then the next step is character positioning. Finding out the location where the character located in the picture, removing the background area, and extracting character area. The third step is character segmentation. Divide the characters into single ones. The final step is character normalization. Normalize each character into the same size of template characters.

3 Moving template matching method

Because of the different shooting environment and the illumination brightness changed over time, even the same characters show differences more or less after the processing of binarization. Moving template matching algorithm mainly aims at the characters effected by the noise and the missing strokes which existed normalized positioning error. It can capture the best position where the character is by moving standard template, thereby carry out the exact match. Moving template matching method is based on the template of target character, using the template of standard character to match the target character from eight directions of up, down, left, right, upper left, lower left, upper right, lower right. Its mathematical expressions are as follows:

Let the target character template to be presented by function $f(x,y)$, and the standard character template to be presented by function $F(x,y)$. Then:

Template upward moving match:

$$T(x, y) = f(x, y) \&\& F(x, y + \Delta y) \quad (3.1)$$

Template downward moving match:

$$T(x, y) = f(x, y - \Delta y) \&\& F(x, y) \quad (3.2)$$

Template leftward moving match:

$$T(x, y) = f(x, y) \&\& F(x - \Delta x, y) \quad (3.3)$$

Template rightward moving match:

$$T(x, y) = f(x + \Delta x, y) \&\& F(x, y) \quad (3.4)$$

Template upper leftward moving match:

$$T(x, y) = f(x, y) \&\& F(x + \Delta x, y + \Delta y) \quad (3.5)$$

Template lower leftward moving match:

$$T(x, y) = f(x, y - \Delta y) \&\& F(x + \Delta x, y) \quad (3.6)$$

Template upper rightward moving match:

$$T(x, y) = f(x + \Delta x, y) \& \& F(x, y + \Delta y) \quad (3.7)$$

Template lower rightward moving match:

$$T(x, y) = f(x + \Delta x, y + \Delta y) \& \& F(x, y) \quad (3.8)$$

As picture 1 showing below, graph (a) is target character and subgraph (b) is standard character. The number 1 presents foreground pixels, and the number 0 presents background pixels. The area enclosed by red lines is the correlated matching part between target character template and standard character template. Here, the standard character template rightward moving four pixels to match the target character template.

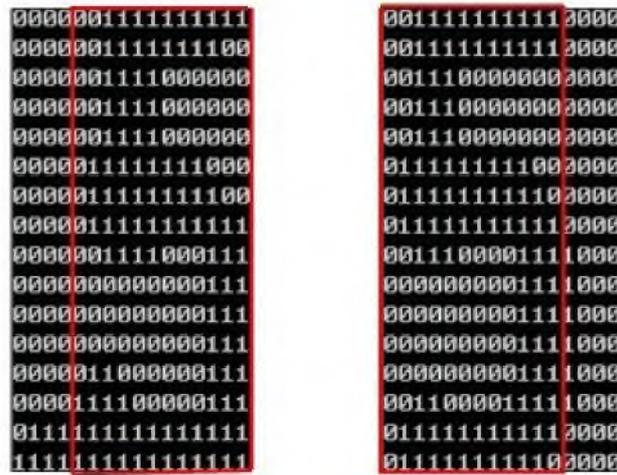


Figure1 (a)target character; (b)standard character

From the figure above, we can see that the influence of adherent noise have been avoided by using moving template matching method, and the right area where the target character is have been captured accurately. Standard character template moved one pix column each time, and recalculated the C value, M value and N value of the overlapping area(that means the area enclosed by red lines). C value, M value and N value are used to calculate the matching rate of combinational value. The meaning and the mathematical expressions of C value, M value and N value are as follows:

C value means: the number of the overlapping foreground pixels in the overlapping area.

$$C = \sum f(x,y) \& \& F(x',y') \quad (3.9)$$

M value means: the number of the foreground pixels of the standard character in the overlapping area.

$$M = \sum F(x,y) \quad (3.10)$$

A Digital Character Recognition Algorithm Based on the Template Weighted Match Degree

N value means: the number of the foreground pixels of the target character in the overlapping area.

$$N = \sum f(x, y)$$

Moving template matching method can greatly avoid the influence of adherent noise and missing partial strokes. Flexible move of the template made standard character and target character achieved a high degree of overlap, which has improved the matching degree of the character.

4 Weighted Matching Degree Algorithm

Traditionally, the common methods related with template match of calculating matching rate are follows:

(1) Matching rate:

$$P_1 = c/m \quad (4.1)$$

(2) Matching rate:

$$P_2 = c/n \quad (4.2)$$

It is one-sided and inaccurate to use method (1) or (2) to calculate the matching rate between standard character and target character. As method (1), P_1 is equal to the result of the number of the foreground pixels of the standard character in the overlapping area divided by the number of the foreground pixels of the target character in the overlapping area, which is not conducive to the characters of complicated structure. For example, the character 4, as the area enclosed by the red lines of the following picture 2(a) shown that the right vertical stroke of the character 4 is very similar to the standard character 1. When the standard character 1 matches with it, the matching rate could be 100%. On the contrary, when the standard character 4 matching with it, the matching rate could be 95%. Actually, the matching rate of 95% has been very high, and because of the final result is determined by the standard characters, so that, the final result is just unilaterally reflect the matching degree. That is easy to cause the wrong judgment. In addition, the disadvantage of method (2) is in contrast with the disadvantage of method (1). Using method (2) is not conducive to the characters of simple structure. As it is shown in the picture 2(b), when the standard character 4 matching with the target character 1, the matching rate could be 100%. However, when the standard character 1 matching with it, the matching rate is 98%. That's why it is one-sided and inaccurate to use method (1) or (2) to calculate the matching rate between standard character and target character.

Actually, the matching rate between standard character and target character means the similar degree of them. Only by making an combination of similar degree between standard character and graphic character can we obtain an result that can totally and accurately reflect the similar degree. All in all, the weighted matching degree algorithm that put forward by this paper is to make an combination of P_1 and P_2 with the weight of 0.6:0.4 to calculate the matching rate between standard character and target character. Only by making an combination of them with an suitable weight, can we overcome the disadvantage of each other and retain the advantage of each other. The mathematical expression of weighted matching degree algorithm is follow:

$$W = 0.6 \times P_1 + 0.4 \times P_2 \quad (4.3)$$

W is the matching rate of weighted matching degree algorithm.

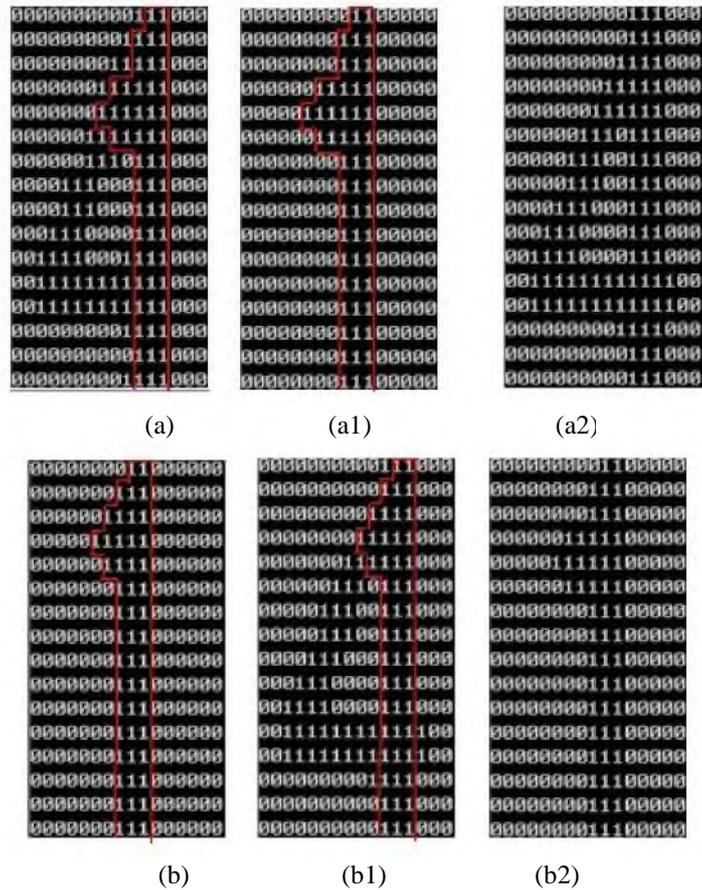


Figure2 (a)target character 4; (a1)standard character 1; (a2)standard character 4;
 (b)target character 1; (b1)standard character 4; (b2)standard character 1.

The rate of combination is getting from a large number of experience. I have used five rate of combination to take the experiment on more than 600 samples. Finally, the 0.6:0.4 is the rate which can reach a recognition rate of 100% on all the samples.

5 Analysis of Experimental Results

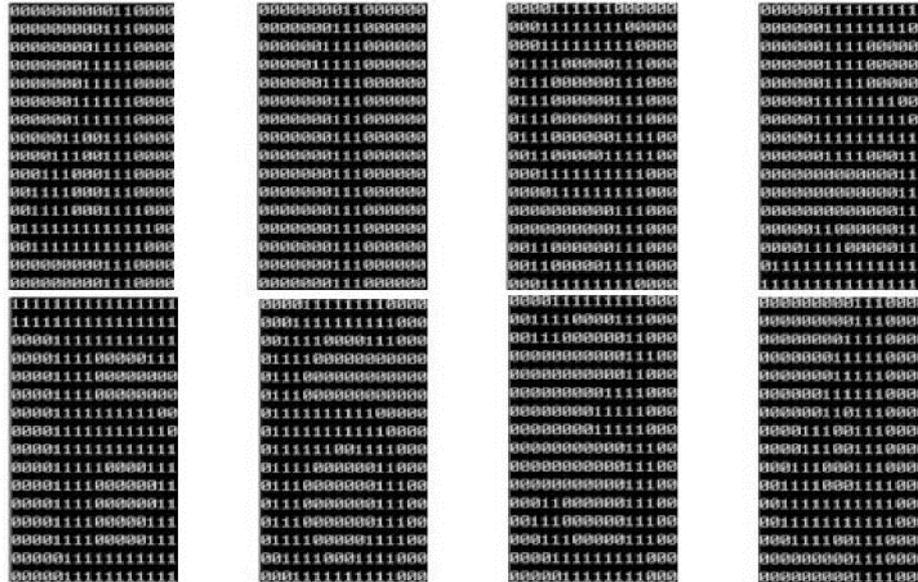


Figure3 Results. From left to right and up and down are subgraphs a b c d e f g h. They are target characters.

Table1 the experimental results

Target Charact ers	Matching Rate						
	P ₁	P ₂	0.2P ₁ +0.8P ₂	0.4P ₁ +0.6P ₂	0.5P ₁ +0.5P ₂	0.6P ₁ +0.4P ₂	0.8P ₁ +0.2P ₂
a	0	1	1	1	1	1	1
b	1	0	1	1	1	1	1
c	0	1	1	1	1	1	1
d	1	1	1	1	1	1	1
e	1	0	0	0	1	1	1
f	1	0	0	0	0	1	1
g	1	0	0	1	1	1	1
h	0	1	1	1	1	1	0

Note: "0" represents the wrong recognition, "1" represents the right recognition.

We can see the results from table 1 above, while used the first or the second method to calculate the matching rate between template character and target character, the program made many errors. Using other rate of combination, more or less, it made errors. Only using the rate of combination of 0.6:0.4, the program obtained the recognition rate of 100%

6 Conclusion

In the algorithm based on the weighted matching degree of template matching method, the template's moving matching can reduce the influence of the noise and the missing partial strokes, and make an accurate matching between standard character and target character; on the other hand, algorithm based on the weighted matching degree of template matching method has made the program to obtain an exhaustive and accurate matching rate between standard character and target character, which avoided the errors by the method of one-sided. The ability of recognition of the program has been greatly improved. Support by a lot of experience, the recognition rate of this program can reach 99.2%, which is a comparatively high level. In the further research, I wish to supplement the characteristic recognition algorithm of similar characters, and ensure the recognition rate of similar characters to improve the ability of recognition of the program.

References

1. Li Dan, Sui Chenghua, Tang Yijun: Research on Recognition of Dynamic Characters in Multi Digital Instruments. (2007).
2. Zhou Nina, Wang Min, Huang Xinhan, Lv Xuefeng, Wan Guohong: Pretreatment Algorithm of Auto Plate Characters' Recognition, Computer Engineering and Applications. (2003)
3. Gu Chenqin, Ge Wancheng. :Character Recognition Based on Template Matching Method. (2009).
4. Information on <http://www.ccs.neu.edu>