

## The research on features of large deformation offline handwritten Chinese character recognition

Yajun Zhang <sup>a</sup>, Xiangnian Huang <sup>a</sup>, Gang Liu <sup>a</sup>

School of Computer and Software Engineering, Xihua University, Chengdu 610039, China

<sup>a</sup>candice0719@126.com

### Abstract

The offline handwriting recognition is one of the difficulties of pattern recognition, feature extraction results directly affects the recognition rate. Firstly, various existing structural features and statistical features will be analyzed, and feature fusion and feature extraction methods will introduction briefly. Finally, a method for feature extraction is presented, which is developed based on the multi-resolution theory and feature fusion theory. It will use these theories generate different feature and use fusion methods to generate new character feature, thereby improving the recognition rate.

### Keywords

Chinese character recognition, feature extraction, feature transformation, multi-resolution analysis.

### 1. Introduction

Off-line handwritten Chinese character recognition has been one of the difficulties in the field of pattern recognition, and some features of handwritten Chinese characters themselves bring a lot of adverse effects on off-line handwritten Chinese character recognition, especially the deformation of handwritten Chinese characters. Handwritten Chinese character recognition can be divided into online and off-line handwritten Chinese character recognition. The serious deformation of handwritten Chinese characters is the main reasons that off-line handwritten Chinese character recognition system has not reached the practical. Currently, available commercial products on the market cannot be completely used in practical stage, so the real off-line handwriting recognition stays in the laboratory research stage<sup>[1]</sup>.

At present, yet to find a way to recognize all handwritten characters a set of features. This article aims to analyze the current existing off-line handwritten Chinese character recognition features, and uses some theories generate a new set of features for Chinese character recognition. The general handwritten Chinese character recognition process can be divided into five steps: handwritten sample acquisition, image pre-processing, feature extraction, classification and post processing. As shown in Figure 1:

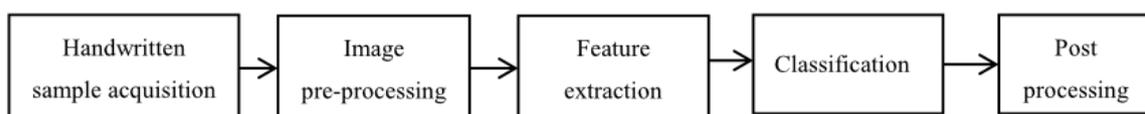


Figure 1 the general process of off-line handwritten Chinese character recognition

### 2. Handwritten sample collection

At present, the handwritten Chinese characters sample library includes Japan ETL-8, Chinese Academy of Sciences Institute of Automation, HCL200 other databases and so on, which provides test sets for the study of Chinese character recognition. To improve the quality of the collected samples, collection forms of special positioning mark are designed for the Chinese character recognition. Large deformation offline handwritten sample collection plan is designed as follows: Each A4 sheet of paper will be divided into 10lines \* 15rows and the positioning marks of black

rectangle are used at the top and left side of the form. The template word of Xing Kai will be written in the first column of each table, and the template word of large deformation cursive will be written in the sixth line. Figure 2 is a part of handwriting sample forms.

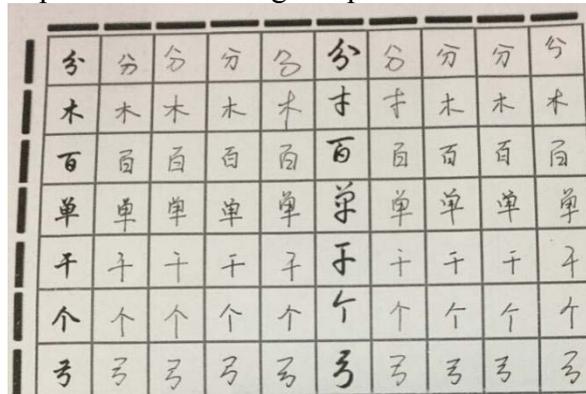


Figure 2. A part of Single handwriting sample forms

Depending on the different character structure, each structure separately collects 15 templates of Xing Kai and cursive. Then printing out the special collection forms on A4 paper as shown in figure 2, the people of different ages fill in these forms. In each sample forms, there are 30 sample templates and the number of samples in each sample template is 8, so the total number of collected handwriting samples will be 120.

### 3. Image pre-processing

Image pre-processing is a series of basic pre-processing operations which based on gray images, such as: binaryzation, normalization, de-noising, and character segmentation. The collected samples need to do a series of image pre-processing operations.

#### 3.1 Grayscale and binary

The scanned samples is a color image, for the final recognition, which be converted to grayscale image and then converted to a binary image. According to the formula which can achieve transformation from the color image to gray image, as follows:

$$Y = 0.299 \times R + 0.587 \times G + 0.114 \times B \tag{1}$$

Each pixel of grayscale image is presented by one of 256 different colors gradations, so gray image will be binaryzation for removing unnecessary color component. Setting the threshold value T can make the grayscale image into the binary image, which can be presented by two gray-scales that represent part of a foreground image and the background color. Existing binaryzation methods include global threshold method, local threshold and dynamic threshold method. After the gradation and binaryzation, the image is shown below:

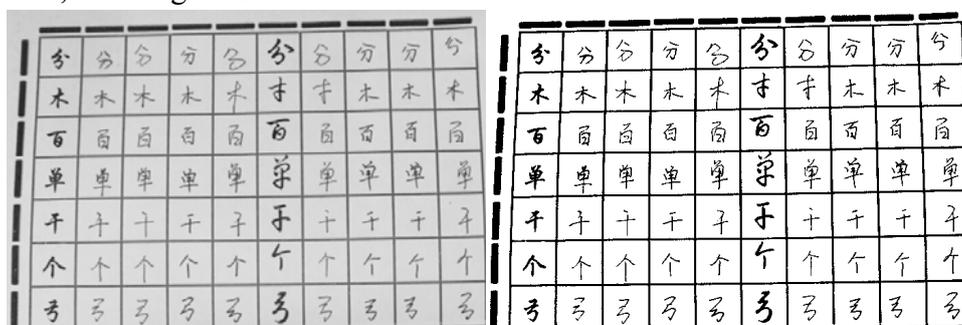


Figure 3 (a) gradation

Figure 3 (b) binaryzation

#### 3.2 Character segmentation

When designing the sample collection form, the positioning marks of black rectangle are already at the top and left side of each table, which can easily find the area of the sample characters (Form box). But the size and location of collected sample is not fixed, and cannot able to locate the characters

circumscribed rectangle. To achieve accurate segmentation of Chinese characters, it is necessary to find the handwritten Chinese characters outside rectangle accurately. According to the character segmentation algorithm of literature [2], this method can be used to scan the pixel of characters boundary accurately to identify the external rectangle.

### 3.3 Normalized

Because of different writing styles, the collected samples have different size and writing position so these samples are normalized for the post-extraction features. Normalization includes position normalization, size normalization, strokes normalization, etc.

## 4. Feature of offline Chinese handwritten character

Good features should have these characteristics: a strong classification ability and independent, easing to extract, high stability and anti-interference, and minimizing the number of features. According to the methods of extracting feature, features generally divided into two categories: statistical and structural features. The statistical features are divided into global statistics features and local statistical features.

### 4.1 Statistical features [3]

Statistical features are based on the binary or grayscale image, which is extracted by conversion of the character dot-matrix information. Statistical features have the following characteristics: extracting feature easily, tolerating characters distortion and noise and good anti-jamming capability.

#### Grid feature

Grid feature of characters are that the character dot matrix is divided into  $n * n$  evenly, and find the percentage of black pixels of strokes in each grid. Grid is usually divided into uniform grid and elastic grid. Elastic Grid is based on the distribution of character image pixel, which using a set of non-uniform lines divide the image. The principle is that adjacent pixel evenly distributed between the two lines. Depending on different partitioning elastic grid methods, elastic grid can be divided into global elastic grid and local elastic grid [4]. Double elastic grid divides vertical and horizontal direction and the diagonal direction of the image of Chinese characters into grids, which use a set of non-uniform grid lines. To a certain extent, double elastic grid can be well adapted to changes in the characters downwards-right direction, while addressing the lack of elastic grid. Uniform grid is that the acreage of each grid is equal [5]. Elastic grid can tolerate the stroke deformation of different writing styles in some extent.

#### Features of global transformation

Using various transformations for Chinese characters image generates features of global transformation, which usually use the transform coefficients as a feature. The common transformation includes K-L transformation, cosine transformation, fast transformation [6], Fourier transformation, wavelet transformation, Walsh transformation and Hough transformation [7]. The features of global transformation are not sensitive of the local transformation, and have strong anti-jamming capability. But, features of global transformation have the weakness of distinguishing between similar words.

#### Gabor feature

Gabor function can best balance resolution of the signal in the time domain and frequency domain. Image will be filtered by two-dimensional Gabor filter which formed by Gabor function, and then we can get a group of gabor feature vectors. First, the image will be constructed elastic grid, and then the midpoint of each elastic grid will be the sampling point  $(x_m, y_n)$ . In every sample points, filtering by using the Gabor filters (or filter group) can obtain different Gabor feature  $fgabor(x_m, y_n)$ , and finally the Gabor features of each sampling point will be spliced together to form feature vectors [8] [9].

#### Features of moment invariants

Features of moment invariants are linear feature, because of their stability in the scale, translation and rotation. So it is widely used in pattern recognition. Currently, features of moment invariants in the

feature extraction have a significant effect, such as: Hu moments, affine moments, Legendre moments and Zernike moments. Wherein, Zernike moments can be constructed any higher moments, which have an advantage over other moments in feature description.

Gradient feature

Gradient feature<sup>[10][11]</sup> describes the local direction feature of Chinese character stroke, which needs to be extracted in the grayscale or binary image. For Chinese characters images, using different image processing operators extracts gradient feature, such as Sobel operator, Roberts operator and Kirsh operator.

#### 4.2 Structural features

Structural features built on the strokes (usually the collection of strokes), which can be described by the strokes' spatial relationships and reflect the essential features of Chinese characters. Common structural features are: point feature, contour feature, stroke features, and component features. Currently, stroke feature extraction methods are stroke skeleton extraction method, stroke contour extraction method, stroke edge direction extraction, stroke contour direction angle and extraction based on wavelet transform. Structural features have both advantages and disadvantages, such as: insensitive of change in the shape of the Chinese character; strong ability in distinguishing similar words. So Structural features commonly used in distinguishing between similar characters. At the same time, structural features are difficult to extract and very sensitive to noise.

### 5. Feature Fusions and Feature Transform

#### 5.1 Feature Fusion

Feature fusion has serial and parallel feature fusion. Serial feature fusion makes two original feature vector end to end to form a new feature vector, which will lead to a sharp increase in the number of feature dimensions, so that the speed of feature extraction and recognition reduced greatly. Parallel feature fusion<sup>[12]</sup> makes two original vector features set ( $\alpha$ ,  $\beta$ ) to form complex vector space ( $\lambda$ ) by using a complex vector.

$$\lambda = \alpha + i\beta \quad (2)$$

Where  $i$  is an imaginary unit. If the two feature vectors measure is not uniform, in front of the feature vector integration, the two groups of feature vectors need to be normalized.

$$x_{norm} = \frac{x - \mu}{\sigma} \quad (3)$$

Wherein,  $x$  is the original feature vector,  $x_{norm}$  is a new feature vector which goes through feature normalized. During the parallel feature fusion, if the two sets of feature vector dimension are not equal, the low-dimensional feature vectors required zeros. The recognition speed of parallel feature fusion is better than the recognition speed of serial feature fusion. Therefore, for large deformation of off-line handwritten Chinese character, parallel features fusion method can be used for the recognition which combines local features with global features, in order to improve the overall recognition rate.

#### 5.2 Feature Transform

Feature transform is that  $D$  feature by using appropriate transformation generate  $d$  ( $<D$ ) new features. Feature Transform reduces the feature dimensionality, which makes the classifier design achieved easily; and feature transformation can eliminate the correlation between features, which makes the new features classify more advantageously. Feature Transform including linear transformation and nonlinear transformation. The linear feature transformation is in common use, which includes principal component analysis (PCA), independent component analysis (ICA), and K-L transforms<sup>[13]</sup>. The basic idea of nonlinear transformation makes non-linear data which can be divided in the input space mapped to a high-dimensional feature space by using appropriate transformation, so that the original input space data becomes linearly separable in a high-dimensional space, and constructs an

optimal classification surface in a high-dimensional feature space which will be regarded as the original non-linear classification surface. Nonlinear transformation mainly has the follow methods: principal curves and manifolds, kernel principal component analysis (KPCA), multi-dimensional scaling, Isomap and locally linear embedding (LLE). Although the non-linear transformation method has obvious advantages in theory, each sample needs to use nonlinear transformation in practice, which will bring enormous computational burden, and also cause the disaster of dimensionality<sup>[14]</sup>.

## 6. Summaries

This paper discussed the current status of offline Chinese character recognition and various recognition features. This article finally proposed that characters can be extracted in different resolutions features in according to the multi-resolution analysis theory, and extracted feature will use feature fusion methods to generate a new set of features, which aims to improve recognition rate of large deformation handwritten Chinese characters.

## Reference:

- [1] Ding Xiaoqing: Chinese character recognition: a review [J]. *Acta Electronica Sinica*, 2002, 30(9): 1364-1368.
- [2] Xiao Bin: *Research on offline handwritten Chinese character recognition based on SVM*(MS. Xihua University, China 2009), p.28.
- [3] Guo Xiangdan: *Research on feature extraction technology applied to off-line handwritten Chinese character recognition* [D]. (MS. Hebei University of Technology, China 2011.), p.20.
- [4] Jin Lianwen, Xu Bingzheng: Directional cellular feature extraction with elastic meshing for handwritten Chinese character recognition [J]. *Journal of Circuits and Systems*, 1997, 2(3): 7-12.
- [5] LIU Wei, Zhu Ningbo et al.: A feature extracting method for handwritten Chinese character recognition based on elastic mesh and fuzzy feature of block [J]. *Journal of Chinese Information Processing*, 2007, 21(3): 117-121.
- [6] Reitboeck H, Brody T P: A transformation with invariance under cyclic permutation for applications in pattern recognition [J]. *Information and Control*, 1969, 15(2): 130-154.
- [7] Cheng F H, Hsu W H, Chen M Y: Recognition of handwritten Chinese characters by modified Hough transform techniques [J]. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 1989, 11(4): 429-439.
- [8] Su Y M, Wang J F: A novel stroke extraction method for Chinese characters using Gabor filters [J]. *Pattern Recognition*, 2003, 36(3): 635-647.
- [9] Wang X, Ding X, Liu C: Optimized Gabor filter based feature extraction for character recognition [C]// *Pattern Recognition, 2002. Proceedings. 16th International Conference on. IEEE*, 2002, 4: 223-226.
- [10] Liu H, Ding X: Handwritten character recognition using gradient feature and quadratic classifier with multiple discrimination schemes [C]// *Document Analysis and Recognition, 2005. Proceedings. Eighth International Conference on. IEEE*, 2005: 19-23.
- [11] Liu C L, Nakashima K, Sako H, et al.: Handwritten digit recognition: investigation of normalization and feature extraction techniques [J]. *Pattern Recognition*, 2004, 37(2): 265-279.
- [12] Yang Jian, Yang Jingyu, Gao Jianzhen: Handwritten Character Recognition Based on Parallel Feature Combination and Generalized K-L Expansion [J]. *Journal of Software*, 2003, 03: 490-495.
- [13] Yang Zhuqing, Li Yong, Hu Dewen: Independent component analysis: a survey [J]. *Acta Automatica Sinica*, 2002, 28(05): 762-772.
- [14] Li S Z, Lu J: Face recognition using the nearest feature line method [J]. *Neural Networks, IEEE Transactions on*, 1999, 10(2): 439-443.