

A System for Segmentation and Recognition of Totally Unconstrained Handwritten Numeral Strings

Z. Shi, S. N. Srihari, Y-C. Shin and V. Ramanaprasad
Center of Excellence for Document Analysis and Recognition
State University of New York at Buffalo, Buffalo, NY 14260, U.S.A.
E-mail: zshi@cedar.buffalo.edu

Abstract

In this paper we propose a system for segmentation and recognition of totally unconstrained handwritten numeral strings. The system is composed of several document analysis modules, namely, a preprocessing module, a segmentation module and a recognition module. The preprocessing module includes connected component analysis, identifying substring with touching digits and estimation of number of digits in the substring. The segmentation module is built with a new segmentation algorithm based on a thorough stroke analysis using contour representation of the strokes. In the recognition module, a high performance digit recognizer is used for the isolated digit images after segmentation and then a simple post-processing is called in the case when some punctuation marks or delimiters such as dash, comma and period are included in the numeral string. Due to the high performance of the segmentation module, the system is efficient and robust with high recognition performance.

1 Introduction

The recognition of numeral strings has gained much attention due to its pivotal role in large number of application such as the postal address interpretation, bank check recognition and reading of tax forms and of census forms. In all of these applications, with few exceptions, the handwritten numeral strings are unconstrained. The problems encountered in numeral string recognition include: estimation of digits in the string, grouping of broken components, splitting touching/overlapping digits, removing/recognizing various punctuation marks such as comma, period and hyphen.

There are several approaches to handwritten numeral string recognition and segmentation [10]. In *segmentation-based* approach, the string is segmented into isolated entities each a complete image or a representation of individual numerals. Then each entity is sent to a numeral recognizer for a recognition result. Methods using segmentation-based approach separate the segmentation from the recognition. The segmen-

tation for getting each individual numeral is usually based on heuristics and done before the recognition of each numerals.

Another approach is called *segmentation-free* in the sense that multiple segmentation candidates are validated by an isolated numeral recognizer. This approach, sometime also called *oversegmentation* is more rational than segmentation-based approach in the case when heuristics prevail, at the price of a higher computational complexity. Works using *segmentation-free* approach are found in [3, 8].

In this paper we apply both segmentation-based and segmentation-free methods. Works using this approach are found in [7, 4, 5].

We build a preprocessing module using the segmentation-base approach. We do not try to segment the numeral string to entities each representing a single numeral. Rather we try to cut the long numeral string down to smaller substrings – blobs which may be of a single numeral or several numerals. The main efforts are put in getting fragments of one numeral into the same blob. In most cases a numeral is in one connected component. In case of broken strokes caused by binarization process or two stroke numerals like 4 and 5, heuristic rules can be used to group them together. Hence we do not have to use recognition to help in this stage. Special punctuation marks and delimiters are also filtered out in preprocessing module.

The segmentation module is built with a new segmentation algorithm based on a thorough stroke analysis using contour representation. The segmentation of each connected or touching numerals is done by segmenting the numerals with splitting the contour representations. For each blob, it first detects the special contour points as splitting points. Then the split contour pieces are grouped by a grouping algorithm consisting of a set of heuristic rules. The grouping algorithm takes the information including the number of numerals in the blob, slant of the strokes and locations of each split contour pieces.

When we estimate the number of numerals in the preprocessing module, we try to get the blobs with

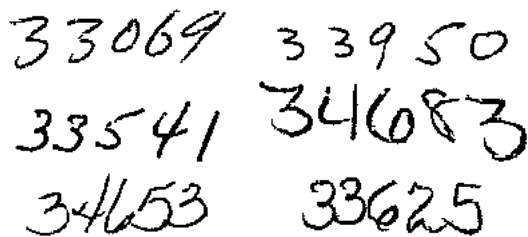


Figure 1: Numeral string samples

single numeral filtered out and send them to the recognizer directly. We take a very conservative approach to get the filtering with high success rate. In the segmentation module we go with an approach similar to the segmentation-free methods. But instead of going with a complicated search algorithm as presented in [3, 4] and [7], we just vary the grouping of the split contour pieces to get the numerals in the blob by varying the number of numerals around the estimated number. The optimal result is validated by a high accuracy numeral recognition. This design approach is based on the following facts. The accurate rate in estimation of the number of numerals in a blob is high so that the real number of numerals is within at most one away from the estimated number. The numeral recognizer, GSC used in the system, is not only highly accurate to accept valid numerals but also being able to reject with low confidence the invalid numerals. Therefore the number of calls to the recognizer is minimized to within at most three passes for each blob. The optimized recognitions will be the final results.

Examples in Figure 1 are images of some unconstrained numeral strings from postal zip codes.

In section 2 we shall describe preprocessing module, including connected component analysis, blob generation, punctuation marks and delimiters filtering and estimation of number of numerals in a string. In section 3 we present a novel approach to the problem of segmentation of touching characters. In section 4 we discuss the use of recognition to the segmentation in a loose fashion as segmentation-free approach.

2 Preprocessing

In the preprocessing module, an input numeral string is divided to be substrings – blobs each containing complete numeral or numerals from touching. The special punctuation marks and delimiters are also detected and formed to separate blobs.

Since we do not require each blob contains only one numeral, i.e. there is no splitting is needed, the only difficulty will be grouping of fragments from broken numerals or two stroke numerals 4 and 5.

To solve the over fragmentation problem caused by binarization process, we blur the input image of a nu-



Figure 2: (a) There is a gap in the original image, (b) The vertically blurred image glues the gap.

meral string to get a masked image. The blurring is done by prolonging each vertical foreground runlength by one pixel. Therefore in the vertically blurred image, up to two pixel gaps in the original image will be glued together. We do not use horizontal blurring as it may cause close connected components to touch each other.

Then the connected components are generated from the masked image and sorted according to their horizontal positions. Average distance of gaps between components is estimated as well. The small components far away from the other components and long slim components are treated as candidates of noise which may be removed by a heuristic rule.

The base line of the numeral string is estimated with skew information included. Average height of the components is computed. We also compute an average width of numerals in the string using the average height as observation standard, i.e. only the components with horizontal size less than the average height are taken into the computation of the average width. The punctuation marks and delimiters are detected and the corresponding connected components are stored into blobs. The detection also classifies the other components as good big components, good small components and bad small components from noise.

A simple grouping algorithm is designed to group good small components to their closer neighboring components by measuring the distances from the center horizontal position of the good small components to their left neighboring component's right horizontal position and to their right neighboring component's left horizontal position. The grouped components are built into blobs for numerals.

For each numeral blob the number of numerals in the blob is estimated with a horizontal intersection method. For the blobs containing single numerals, the images are rebuilt from the components in the blobs to be images covered by the components in the original image. And the images each for a single numeral are sent to the recognition.

For the blobs with estimated number of numerals greater than one, the rebuilt images are sent to the segmentation module for further processing.

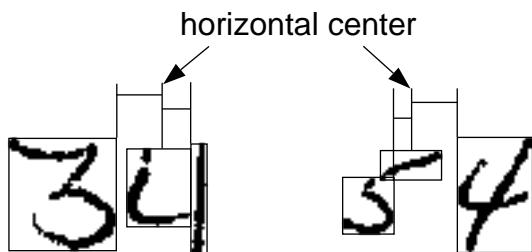


Figure 3: The detected small pieces are grouped to their neighboring good big pieces according to their locations.

3 Segmentation Algorithm

The most difficult part of the system is the segmentation of connected numeral strings. A survey of segmentation strategies is provided by Casey and Lecolinet [5]. Common approaches are often heuristic. While they have the advantage of efficiency, their accuracy is limited.

The method we present for segmentation of connected numerals is based on the contour representation of character strokes. Segmentation is done by splitting and grouping the contours.

3.1 Contour Classification and Segmentation

Contour is built of arrays of the contour points by tracing the boundary of the character stroke consecutively along counter clock wise. Stroke width is also estimated by averaging all the distances from each boundary points to their opposite boundary points.

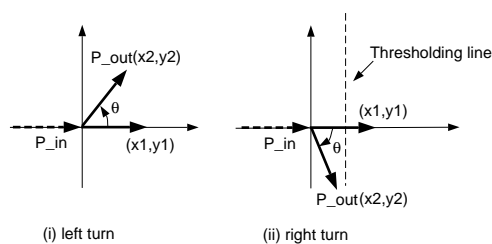
A close look at the touching points and ligatures between two digits reveals that at each touching point, the contour makes significant right turns. We use the significant right turning points to segment the contour to be classes of contour pieces.

To determine the *significant* right turning contour points, we first compute vectors P_{in} leading in to a contour point P from its several previous neighboring contour points and P_{out} going out of P to its next several contour points. These vectors are normalized and placed in a Cartesian coordinate system with P_{in} along the x -axis (Figure 4 (a)). A threshold THR is then selected such that any significant right turn satisfies the conditions:

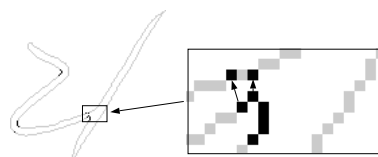
$$x_1y_2 - x_2y_1 < 0 \quad (1)$$

$$x_1y_1 + x_2y_2 < THR \quad (2)$$

where (1) indicates that the turn is to the right and (2) indicates that the turn is significant. Since the threshold THR is the x -coordinate of the thresholding line in Figure 4 (a), it can be determined experimentally to be a number close to zero. This ensures that



(a)



(b)

Figure 4: (a) The distance between the thresholding line and the y -axis gives a threshold for determining a significant right turn. (b) The significant right turning points and their opposite contour points are marked.

the angle θ made by P_{in} and P_{out} is close to 90° . Significant right turning points together with their opposite contour points located earlier divide the contour into contour pieces (Figure 4 (b)). Contour pieces may now be classified as belonging to one or the other digit.

3.2 Grouping Contour Pieces Using Histogram Decision Lines

We assume that the number of numerals in the input image to the segmentation module is known. Then the simplest way of grouping the contour pieces to each group for a numeral is by using the vertical lines dividing the image to be boxes each contains one numeral as guide. A contour piece is classified as part of a numeral if the *center of mass* of the piece is in the box containing the numeral. In the ideal case, the numerals have similar widths and the touching strokes lie close to the center line. But this is not always the case. The method described fails, for example, for an image containing touching digits 5 and 1.

For general cases we first estimate a vertical slant using the consecutive contour structure. then we compute the histogram of extents along the slant direction where the extent of the image along each slant line is defined as the distance between the uppermost black pixel and the lower most black pixel in that slant line. The touching point may be expected to lie near the valley of the histogram. Valley areas are located to be intervals (Figure 5).

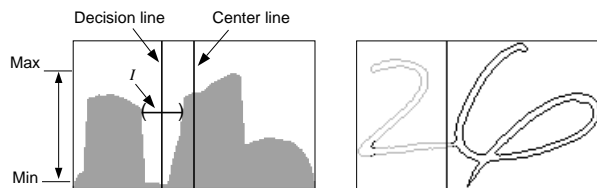


Figure 5: Histogram of vertical extents is used to find decision line, and the segmentation result.

A divide and conquer scheme is used to find the N decision lines for $N + 1$ numerals in the image. That is, the $N/2$ -th decision line is found first and the rest decision lines will be found in the two sub-fields from the image divided by the decision line.

For each decision line we find the intervals closest to the vertical bisecting line (the $N/2$ -th vertical line in the histogram) and find the minimal point within the valleys. The one which is closest to the bisecting line is taken as the position of the decision line (Figure 5).

We use the decision line as a guide line to output the contour segments as described earlier. Our segmentation algorithm allows us to classify entire strokes as belong to one or the other neighboring numerals, and results in a “cleaner” segmentation compared to traditional segmentation methods.

Images of segmented numerals may be recovered from the contour segments by drawing line segments from each boundary point and two neighboring boundary points to their opposite contour points, or the modified opposite boundary points which are built by using the estimated stroke width if any line segments to the opposite points are too long.

In fact, image recovery may be unnecessary since a contour based recognizer can directly use the separated contour segments for recognition.

4 Implementation and Performance Evaluation

A research system is implemented and tested on CEDAR CD-ROM database. A high performance numeral recognizer GSC developed in CEDAR is used in the system.

We tested the system on two sets of zipcode images from CEDAR CD-ROM and US Postal Database. In the first set there are 495 zipcode images for which a 85% field recognition correct rate and 95% segmentation correct rate are achieved. The segmentation correct rate is computed manually by checking all the field incorrect images. The errors are from incorrect grouping, incorrect splitting and incorrect recognition for digits in touching strings. In the second set there are 2579 US zipcode images. The field recognition cor-

rect rate is 85.7%. For both tests, zero rejection is assumed. Although we have not optimized this research system, the average speed per digit is 0.04 second on a Sparc-Ultra.

5 Conclusion

A system for segmentation and recognition of totally unconstrained handwritten numeral strings has been proposed. Preliminary experiments on real data from US postal zipcode showed very promising results in terms of segmentation and recognition rates when compared with other published systems. Moreover the proposed system can be used for segmentation and recognition of hand-printed characters with few modifications.

References

- [1] M. Shridhar and A. Badreldin, Context-directed segmentation algorithm for handwritten numeral strings, *Image and vision computing*, Vol. 5, No. 1, 3-9(1987)
- [2] R. G. Casey and E. Lecolinet, Strategies in character segmentation: A survey, *Proc. of the Int. Conference on Document Analysis and Recognition*, Montréal, Canada, 1028-1033(1995).
- [3] S. Seshadri and D. Sivakumar, A technique for segmenting handwritten digits, *Pre-Proc. of the Int. Workshop on Frontiers in Handwriting Recognition III*, Buffalo, New York, USA, May 25-27, 443-448(1993).
- [4] R. Fenrich, Segmentation of automatically located handwritten numeric strings, *From Pixels to Features III: Frontiers in Handwriting Recognition*, edited by S. Impedovo and J. C. Simon, Elsevier Science Publishers B. V. 47-59(1992).
- [5] G. Kim and V. Govindaraju, Handwritten Word Recognition for Real-Time Applications. *Proc. of the Int. Conference on Document Analysis and Recognition*, Montréal, Canada, 24-27(1995).
- [6] Z. Shi and V. Govindaraju, Segmentation and Recognition of Connected Handwritten Numeral Strings. *Pre-Proc. of the Int. Workshop on Frontiers in Handwriting Recognition V*,
- [7] T. M. Ha, D. Niggeler and H. Bunke, A system for segmenting and recognising totally unconstrained handwritten numeral strings, *Proc. of the Int. Conference on Document Analysis and Recognition*, Montréal, Canada, 1003-1009(1995).
- [8] H. Nishida and S. Mori, A model-based split-and-merge method for recognition and segmentation of character strings, *Advances in Structural and Syntactic Pattern Recognition*, H. Bunke(Ed.), World Scientific, 300-309(1992).