

Word Segmentation of Off-line Handwritten Documents

Chen Huang and Sargur N. Srihari

{chuang5, srihari}@cedar.buffalo.edu

Center of Excellence for Document Analysis and Recognition (CEDAR),
Department of Computer Science and Engineering,
University at Buffalo, The State University of New York, Buffalo NY, USA

ABSTRACT

Word segmentation is the most critical pre-processing step for any handwritten document recognition/retrieval system. This paper describes an approach to separate a line of unconstrained (written in a natural manner) handwritten text into words. When the writing style is unconstrained, recognition of individual components may be unreliable so they must be grouped together into word hypotheses, before recognition algorithms can be used. Our approach uses a set of both local and global features, which is motivated by the way that human beings perform this kind of task. In addition, in order to overcome the disadvantage of different distance measures, we propose an average distance computed using three different methods. The system is evaluated using an unconstrained handwriting database, which contains 50 pages (1026 line, 7562 words images) handwritten documents. The overall accuracy is 90.82%, which shows a better performance than a previous method.

1. INTRODUCTION

Line segmentation and word segmentation are the most critical pre-processing steps for any handwritten document recognition/retrieval task. The goal is to extract all the word images from a full page of handwritten document. It is very important because, first of all, in handwritten recognition, word recognition methods can be categorized into two categories: segmentation based and non-segment based, and both of them need to work on pre-extracted word images. Secondly, content-based image retrieval techniques, such as word spotting, also require all the word images in the documents to be pre-segmented properly. Wrongly segmented word images will fail most of the techniques in handwritten document recognition/retrieval system.

In the present paper we address the problem of separating a located text line into words. Separating handwritten text into words is challenging because handwritten text lacks the uniform spacing normally found in machine-printed text. Machine-printed text typically has inter-word gaps that are much larger than inter-character gaps (gaps between characters within one word). Therefore, there is little work on full page segmentation, with most of the previous work in handwriting focused on specialized domains like postal address and bank checks. For example, Seni and Cohen¹ evaluate eight different distance measures between pairs of connected components for word segmentation in handwritten postal addresses. Feldbach and Tonnie² present a system using constraints on the semantics to segment the date from church registers using a neural network. Marti and Bunke³ propose a full-page word segmentation algorithm and the evaluation is done by using the IAM database.⁴ The IAM database consists of text copied with care by a large number of writers and a ruler was used to ensure that the lines are straight and horizontal. Recently Manmatha and Rothfeder⁵ described a scale space approach for segmenting words from historical handwritten documents.

In this paper we propose a gap metrics based approach to perform the word segmentation task. The new approach has two main differences from previous methods. First of all, the gap metrics is computed by combining three different distance measures, which avoids the weakness of each of the individual one and thus provides a more reliable distance measure. Secondly, besides the local features, such as the current gap, a new set of global features are also extracted to help the classifier to make a better decision. The classification is done by using a three-layer neural network.

The remainder of this paper is organized as follows. Section 2 describes the method in detail including feature extraction and the neural network classifier. In section 3 we present some of the experimental results. Section 4 concludes the paper.

2. ALGORITHM DESCRIPTION

2.1. Preprocessing

The task of our algorithm is to segment a handwritten line image into words. Thus, the inputs of our algorithm are pre-segmented line images. This is done by using a statistical line segmentation approach proposed by Manivannan et al.⁶ This algorithm is robust to handle documents with skew and lines running into each other. It is based on modeling the lines as bi-variate gaussian densities that provide for accurate association of components to the respective lines. The use of piece-wise projection profiles to guide the lines drawn reduces the number of obstructing components.

2.2. Feature extraction

We form the word segmentation problem as a two-class classification problem, i.e., given a distance (or gap) between two components, classifying whether it is an inter-word gap. When a person makes a decision of whether a gap is an inter-word gap, he/she not only looks at the spatial separation between current pair of components, but also captures the some other local information, such as the size of the components, and the global features, such as whether the handwriting is cursive or hand-printed.

Therefore, inspired by human being, our algorithm computes 7 local feature as well as 4 global features from the entire line image. The local features are as follows.

- Distance between current pair of components.
- Distance between previous and next pair of components. This feature captures the neighbor information. If no previous or next pair (the first or last one), then the maximum distance is assigned as the feature value.
- Width of the left and right components.
- Height of the left and right components.

The global features are including:

- Ratio of the number of exterior contours and the number of interior contours. This feature captures the writing style (cursive or hand-printed) information.
- Average height of the grouped components.
- Average width of the grouped components.
- Average distance between components.

Before computing the distance between each pair of exterior connected components, the components will be clustered first such that the stray marks above and below the line will be grouped together with their primary components, i.e., if the horizontal range of a component spans over another component, these components should be put into the same group.

In order to overcome the weakness of different distance measures mentioned in,¹ we compute two distances measures and use the average of them as the final distance. The first one is measured using either the bounding box method or the minimum run-length method. The minimum run-length method is used only if the two bounding boxes are overlapping horizontally, as shown in Fig. 1 (a). Here a run-length is defined as the distance along a straight line between two connected components.

The second measure is the convex hull distance, which is computed as follows. For each grouped component, an approximate convex hull is first computed. Then the center of gravity (CG) will be computed for each hull. The line connected the CGs of two adjacent groups will be found. The intersections of the CG line with the two hulls also will be found. The distance of the two groups is defined as the distance of the two intersections, as shown in Fig. 1 (b).

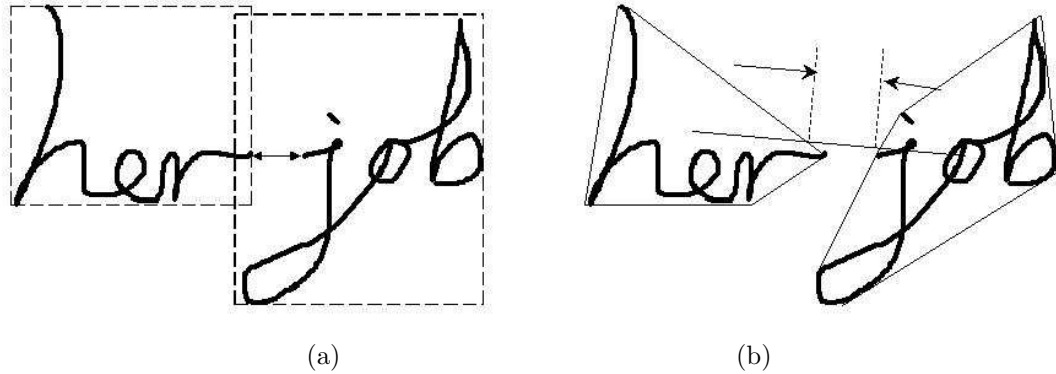


Figure 1. Examples of types of distance measures between a pair of connected components. The bounding box method and minimum run-length method are shown in (a), and the convex hull distance is shown in (b).

2.3. Classification

A three-layer neural network is used for the classification. At the input layer we have 11 features as mentioned above. While the hidden layer includes four hidden units, we have two units at the output layer, which usually gets a better performance than having one output unit for a two-class classification problem.

The training of the neural network is conducted by using a set of 600 line images, which are manually truthed. The line images are segmented from full page handwritten documents, which are a portion of a large collection of unconstrained handwriting documents. The dataset will be described in detail in the next section.

3. EXPERIMENTAL RESULTS

3.1. Dataset

The dataset used for experiment contains 50 pages (1026 line/7562 words images) of handwritten documents, which is a subset chosen from a large dataset created for forensic document examination studies.⁷ The content of the document is so called “CEDAR letter”, which was designed to contain 156 words including all characters (letters and numerals), punctuations and distinctive letter and numeral combinations (ff, tt, oo, 00). The vocabulary size is 124. That is, 32 out of 156 words are duplicate words, and most of them are the stop words, such as “the”, “she” and etc. About 1,500 individuals copied the CEDAR letter three times each in his/her most natural handwriting using plain unlined sheets, and a medium black ball-point pen. The samples were scanned using 300 dpi resolution and 8-bit grayscale. Figure 2 (a) shows a sample image and the content of the CEDAR letter.

3.2. Experimental results

Among these 1026 line images, 600 are used as the training set and the rest 426 lines (3273 word images) are used as the testing set.

In the testing set, 2907 out of 3273 words are extracted correctly. Therefore the system performance is about 90.82% on overall accuracy. A previous method designed for postal address application⁸ was also evaluated using the same testing dataset, and the overall accuracy is 87.36%. This indicates that the proposed new algorithm shows a better performance. An example of word segmentation for a full page document is shown in Fig. 2 (b).

Among those error segments, we observe that over-segmentation error rate is slightly higher than the under-segment error rate. While this is just a preliminary performance evaluation, in order to compare it with some other state-of-the-art algorithms, we are currently perform another testing based on the IAM database.

From Jim Elder
829 Loop Street, Apt 300
Allentown, New York 14707

Nov 10, 1999

To Dr. Bob Grant
602 Queensberry Parkway
Omar, West Virginia 25638

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Euck.
It all started around six months ago while attending the "Pubog" Jazz Concert. Organizing such an event is no picnic, and as President of the Alumni Association, a co-sponsor of the event, Kate was overworked. But she enjoyed her job, and did what was required of her with great zeal and enthusiasm.
However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions, X-rays and blood tests later, were told it was just exhaustion.

(a)

From Jim Elder
829 Loop Street, Apt 300
Allentown, New York 14707

Nov 10, 1999

To Dr. Bob Grant
602 Queensberry Parkway
Omar, West Virginia 25638

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Euck.
It all started around six months ago while attending the "Pubog" Jazz Concert. Organizing such an event is no picnic, and as President of the Alumni Association, a co-sponsor of the event, Kate was overworked. But she enjoyed her job, and did what was required of her with great zeal and enthusiasm.
However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions, X-rays and blood tests later, were told it was just exhaustion.

(b)

Figure 2. A handwriting sample from CEDAR letter dataset. (a) The original document. (b) The document after performing word segmentation (words are shown in different colors).

4. CONCLUSIONS

In the paper, we propose a new gap metrics based word segmentation algorithm. This method computes a new set of features including both local and global informations. In addition, a new distance measure is proposed to overcome the weakness of each individual distance measure method. The system was evaluated using an unconstrained handwriting database. The system performance is better than the previous method. A further evaluation is being conducted to obtain a formal comparison.

REFERENCES

1. G. Seni and E. Cohen, "External word segmentation of off-line handwritten text lines," *Pattern Recognition* **27(1)**, pp. 41–52, January 1994.
2. M. Feldbach and K. D. Tonnies, "Word segmentation of handwritten dates in historical documents by combining semantic a-priori-knowledge with local features," *Proc. of the 7th Int. Conference on Document Analysis and Pattern Recognition*, 2003.
3. U. V. Marti and H. Bunke, "Text line segmentation and word recognition in a system for general writer independent handwriting recognition," *Proc. of the 6th Int. Conference on Document Analysis and Pattern Recognition*, pp. 159–163, 2001.
4. U. V. Marti and H. Bunke, "A full english sentence database for off-line handwriting recognition," *Proc. of the 5th Int. Conference on Document Analysis and Pattern Recognition*, pp. 705–708, 1999.
5. R. Manmatha and J. L. Rothfeder, "A scale space approach for automatically segmenting words from historical handwritten documents," *IEEE Transactions on Pattern Analysis and Machine Intelligence* **27(8)**, pp. 1212–1225, August 2005.
6. M. Arivazhagan, H. Srinivasan, and S. Srihari, "A statistical approach to handwritten line segmentation," *Document Recognition and Retrieval XI, Proceedings of SPIE*, January 2007.
7. S. N. Srihari, S. Cha, H. Arora, and S. Lee, "Individuality of handwriting," *Journal Of Forensic Sciences*, pp. 856–872, 2002.
8. G. Kim, V. Govindaraju, and S. N. Srihari, "An architecture for handwritten text recognition systems," *International Journal on Document Analysis and Recognition*, pp. 37–44, 1999.