

Machine Printed Character Segmentation Method using Side Profiles

Min-Chul Jung, Yong-Chul Shin and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition

State University of New York at Buffalo

Buffalo, New York 14260 U.S.A.

{mjung,ycshin,srihari}@cedar.buffalo.edu

Abstract

In this paper, a new segmentation method of machine printed character string with arbitrary length is proposed. It exploits recognition-based segmentation, combining with heuristic and holistic methods. The merged part of touching characters generates different shape of patterns from the primitive character patterns. However, far left side and far right side patterns in the touching characters will not be effected by the touching. The algorithm firstly constructs a line adjacency graph(LAG) from a word image. Blobs are found as connected components of the LAG and the small dot noises are removed. Secondly, as a word in English can be divided into three typographical zones such as the ascender, the x height and the descender, the location of the connected components among those zones are also examined. Thirdly, the rightward profile of the touching character is compared with that of the sample characters in the prototype and then the touching characters are segmented with the width of one of the candidates in the prototype. Finally, upward, downward and leftward profiles of the segmented pattern are compared with those of the candidate respectively. Third and final steps are continued until confirmed by successful matchings of the resulting character patterns. It has been tested with touching characters in Times and in Helvetica fonts that are proportional pitch fonts and found that the proposed method is promising.

Keywords : Machine Printed Character, Character Segmentation, Side Profile, Character Recognition, OCR

1 Introduction

In most existing OCR systems, word images are segmented into isolated characters so that character recognition can perform on individual characters and this process is called character segmentation. Character segmentation is an essential part of the system, as the

performance of the character segmentation can affect the performance of the overall system due to the fact that incorrectly segmented characters are not likely to be recognized properly [1]. As long as the characters are segmented correctly, degradation of individual character is not significantly affecting the overall system performance. If recognition is unacceptable, it is generally because the characters were difficult to segment [2]. Nevertheless, segmentation of touching characters in proportional pitch fonts is a challenging problem in OCR because character pitch information is not available for a document whose character pitch is proportional. Although it has been studied over a decade, a general solution to this non-trivial task is not readily available. Segmentation procedures are more heuristic in nature than recognition procedures and proper segmentation requires a priori knowledge of the patterns that form meaningful units, which implies recognition capability. Without understanding the symbols, there are no good criteria to avoid errors of segmentation. Therefore character segmentation should be closely coupled with character recognition. This paper presents a compromising solution of character segmentation with character recognition using side profile features – leftward, rightward, upward, downward – for machine printed touching characters.

2 Approach to Recognition-Based Character Segmentation

The merged part of touching characters generates different shape of patterns from the primitive character patterns. The motivation of recognition-based segmentation in this paper is that the merged part of touching characters generates different shape of patterns from the primitive character patterns, however far left side and far right side patterns in the touching characters will not be effected by the touching. The analysis based on those side views gives the solution of the character segmentation as well as the character recognition. Each character in English has a unique

side view, in other words, profile. For example, at the lower cases in Helvetica font, 'a, b, c, e, f, g, h, i, j, k, o, p, q, r, s, t, x, y and z' have unique rightward profiles and those characters can be recognized with a rightward profile alone. Note that 'd and l', 'm, n and u', and 'v and w' have the same pattern of rightward profiles respectively. In addition, 'a, d, f, g, i, j, p, q, s, t, x, y and z', 'b, h, k and l', 'c, e and o', 'm, n, r and u', and 'v and w' sets have also unique leftward profiles. Therefore, both a leftward profile and a rightward profile give the obvious feature of a character. Upward and downward profiles also give information on the character. Figure 1 shows some samples of four side profiles of characters and numerals in Helvetica font.

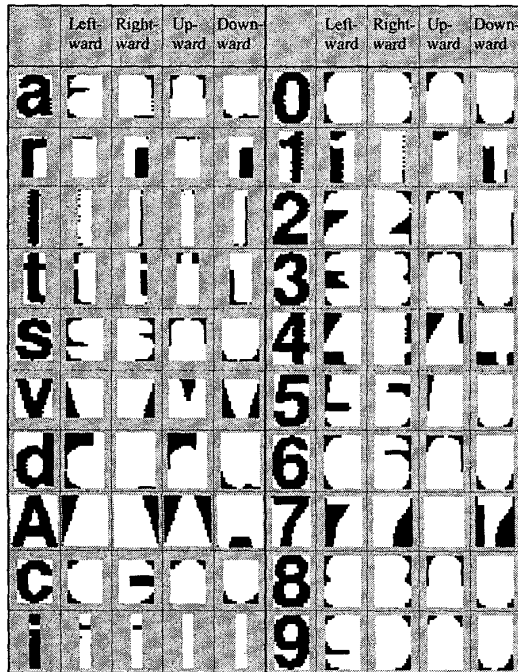


Figure 1: Some samples of four side profiles

When two characters are touching, it may be possible to recognize it without the character segmentation holistically. However, it is not known how many characters are in the touching characters before the character recognizer confirms. Since the profiles of the numerals are more noticeable than those of the characters and the number of numeral classes is 10 that is much less than that of character classes, the proposed method is more effective in numerals. The strategy in

this paper is that segmentation with recognition and then take it out from the touching character for the next segmentation. The algorithm proposed in this paper is illustrated in figure 2.

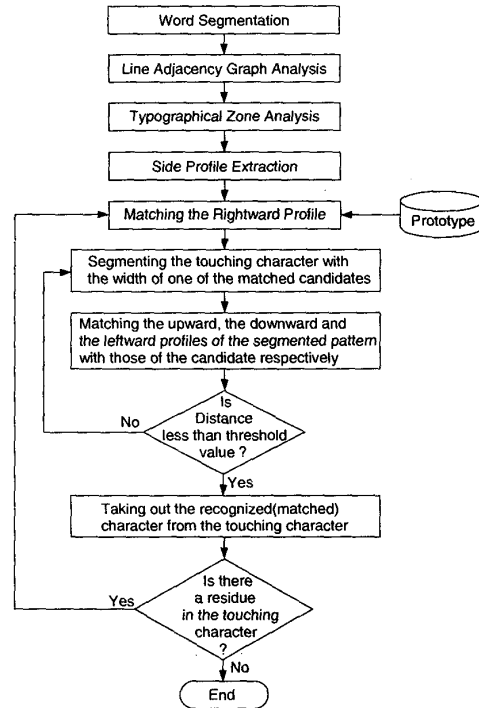


Figure 2: Flow Chart of the proposed method

3 Preprocessing

3.1 Line Adjacency Graph

From a word image, the algorithm first constructs a line adjacency graph(LAG). Blobs are found as connected components of the LAG and can be a single character or touching characters. The small dot noises, in other words, pepper noises are removed. The blobs are bound-boxed as shown in Figure 3c.

3.2 Typographical Zone Analysis

A word image is composed of three typographical zones: the ascender, the x height and the descender zones, which are delimited by four virtual horizontal lines, ascender, x height, base and descender lines [3]. From the vertical projection profile of a word image after the correction of skew angle, the typographical structure can be estimated. The analysis of the vertical projection profile shows that a word image has one of four types [4]. Figure 3b shows that the word is in the type II.

- Type I: The presence of all three zones,
- Type II: The presence of the ascender and the x height zones,
- Type III: The presence of the x height and the descender zones,
- Type IV: The presence of the x height zone only.

3.3 Font and Size

Font family can be identified by the analysis of the typographical attributes such as ascenders, descenders and serifs [5]. In this paper, Times and Helvetica fonts are chosen as proportional pitch fonts for serified and sans-serifed fonts, respectively. The profiles of serified font and sans-serifed font is different because of serif. Serif is a small decorative stroke at the ends of a vertical stroke and causes touching in most cases. If the scanning resolution is known, the font size can be extracted from the vertical pixel counts from the ascender, the x height and the descender zones [5]. The font size is 12 points and the scanning resolution is 200 dpi in the prototype. In that case, the vertical pixel numbers of the ascender zone, the x height zone and the descender zone are 6, 18 and 8 in Times font and 6, 19 and 7 in Helvetica font. As the input of the touching characters may have irregular size, the vertical pixel number is normalized according to the above values.

4 Recognition-based Character Segmentation

4.1 Four Side Profiles

The four profiles of one connected component encompassed with a bound box are obtained by counting the white pixels in the four directions, rightward, leftward, upward and downward respectively until the black pixels encountered. Figure 3d shows four directional profiles of the touching characters in the example. A histogram can be produced from a profile. The histogram appears as a graph with “character width or height” on the horizontal axis and “number of white pixels” on the vertical axis.

4.2 Matching of Profiles

As shown in Figure 2, the rightward profile of the touching character is compared with those of the sample characters in the prototype. The distance of the histogram from a character in the prototype (p) and the histogram from the touching character as a query (q) is

$$d[p, q] = \sum_{i=0}^N |x_{ip} - x_{iq}|$$



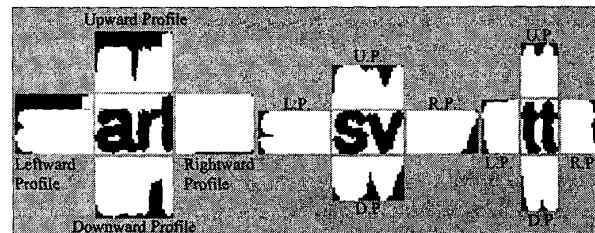
(a)



(b)



(c)



(d)

Figure 3: (a) A word image with touching characters. (b) The vertical projection profile of (a). (c) All bound boxed sub-images of the word. (d) Four side profiles of touching characters

where N is a normalized height or width of a character, x_{ip} is each value of the vertical axis of histogram, in other words, number of white pixels from a character in the prototype and x_{iq} is the same from the touching character. The characters whose distances are less than the threshold value are selected as the candidates of a last character in the touching characters. The touching characters can be segmented with the width of one of the candidates in the prototype. However, the character width is dynamically estimated during the following step.

4.3 Cutting Cost

Since a character width in the touching characters can be changed according to the scanning threshold, size normalization, noises and font style such as normal and bold even in the same font size, the width of one of

the candidate characters in the prototype gives a neighborhood of cutting path. Tsujimoto and Asada [6] introduced a break cost that is defined at each position between neighboring columns and is calculated by accumulating the black pixels of the image obtained by the AND operation between neighboring columns. The candidates for break position are obtained by finding local minima of a smoothed break cost function. In this paper, a few columns in its neighborhood give possible cutting paths and the best cutting path is decided after examining the cutting cost and the profiles of the segmented pattern. The cutting cost is how many black pixels should be fallen apart in order to make a cut of the touching characters. Figure 4 illustrates the cutting costs of an image with touching characters as shown in figure 3d. The image is size-normalized to the height of the candidate 'd' in the prototype. The aspect ratio of the image is kept. The cutting costs of every between columns in a neighborhood are investigated. In the figure, the cutting cost of line ① is 16, the one of line ② is 3 and the one of line ③ is 15. The cutting path that has a minimum of cutting cost is first considered to segment the touching characters.

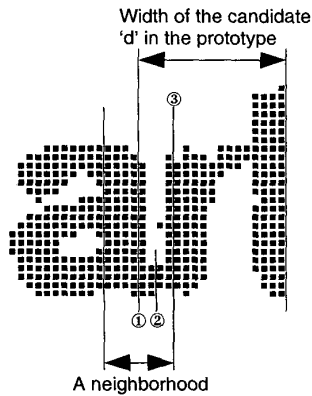


Figure 4: A sample image with touching characters of 'ar'. This illustrates the situation with the last character candidate being 'd'.

4.4 Segmentation with recognition

Upward, downward and leftward profiles of the segmented pattern are compared with those of the candidate using the above distance formula, respectively (See, figure 5).

A Normalized distance is needed in order to compare with distances from two candidates that have different sizes, in fact, different columns. The normalized distance is

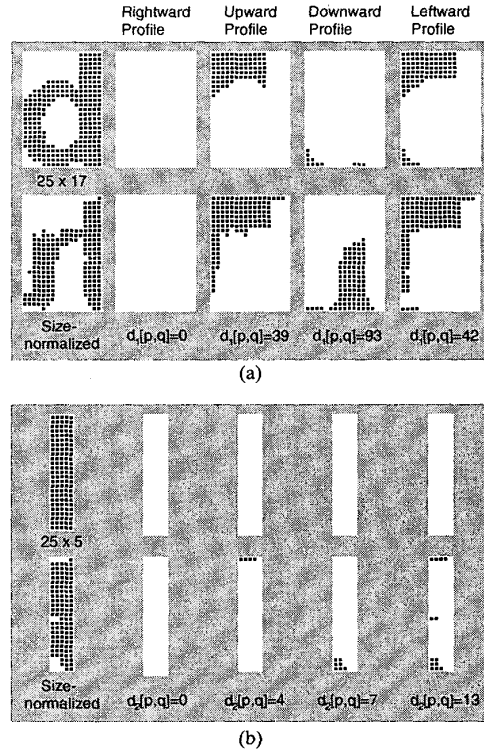


Figure 5: (a) Four side profiles of the candidate 'd' and the segmented pattern. (b) Four side profiles of the candidate 'l' and the segmented pattern.

$$D_1[p, q] = d_1[p, q] \times C_2$$

$$D_2[p, q] = d_2[p, q] \times C_1$$

where C_1 is a column of one candidate and C_2 is a column of the other. The character whose four normalized distances are all less than the threshold values is finally selected as a matched (recognized) character. As shown in figure 3d, for the 'ar' image, its rightward profile represents that its last character is either 'd' or 'l'. The algorithm segments 'ar' with the cutting path of a minimum of cutting cost in neighborhood of the width of the candidate 'd' and does the same with the candidate 'l' in the prototype. The segmented patterns are size-normalized to the sizes of the candidates. Upward, downward and leftward profiles of the segmented pattern are compared with those of the candidates 'l' and 'd'. Figure 5 shows those profiles and $d[p, q]$ distances. The normalized distances for a candidate 'd' are that $D_1[p, q]$ of rightward profiles is 0, $D_1[p, q]$ of upward

profiles is 195, $D_1[p, q]$ of downward profiles is 465 (It can be observed that two downward profiles are significantly different in the figure.), and $D_1[p, q]$ of leftward profiles is 210. The normalized distances $D_2[p, q]$ for a candidate 'l' are 0, 68, 119 and 221, respectively. Since the normalized distances show that 'l' is closer to the segmented pattern of the touching characters than 'd', 'l' is finally chosen as the last character of the touching characters. After taking the matched (recognized) character 'l' out from 'arl', the same procedure as the above with 'ar' repeats until there is no residue in the touching characters. In this example, the last character of the touching characters is either 'd' or 'l' at the step of 'Matching the Rightward Profile'. However, most of cases, the last character of the touching characters is recognized at the step of 'Matching the Rightward Profile' because of unique rightward profiles of characters. Unlike the other approaches based on hypothesis generation of potential cutting points, the approach presented here finds patterns that match with previously stored prototypes, and hence the error rate is low.

5 Experiments and Results

The proposed method was implemented in C language on SUN SPARC work station. In experiments, most touching characters are successfully separated, including special character pairs such as "ff", "tt", "fi", "fl", "rf", "rt" and "rp" those are narrow and often touching in proportional pitch font. The method is more effective when the last character of the touching characters has a unique and obvious rightward profile. The approach may decompose one character 'd' into two characters 'c' and 'l' because the rightward profiles of 'd' and 'l' are the same and 'd' includes 'l' pattern. For example, touching characters "ed" may be regarded as 'e', 'c' and 'l', in other words, the last character can be either 'd' or 'l' with the same confidence values. In that case, 'c', the leftover pattern of 'd' has a low confidence value. Beside, in this example, the pattern 'e' after segmenting with the width of 'c' also has a low confidence value. Therefore, the segmentation decisions should be tentative until confirmed by successful classification of the resulting character patterns. The combinations "nn", "mm" and "rm" are also successfully segmented, but the combination "rn" can be segmented into 'r' and 'n' and can be regarded as 'm' as well.

6 Conclusion

This paper proposed a new method to segment touching characters in machine printed documents using the four basic features of side profiles. The experi-

ment confirms that this method is effective in segmenting touching characters, and applicable to a variety of document analysis. Since the side profiles of serified and sans-serifed font is significantly different, font information helps the performance of the segmentation. As possible improvements in near future, instead of the matching of pixel histograms, the contour matching of character outlines may be considered.

References

- [1] Yi Lu "Machine Printed Character Segmentation - An Overview," *Pattern Recognition*, Vol. 28, No. 1, 1995, pp. 67-80.
- [2] Mindy Bokser "Omnidocument Technologies," *Proceedings of the IEEE*, Vol. 80, No. 7, Jul. 1992, pp. 1066-1078.
- [3] Abdelwahab Zramdini and Rolf Ingold, "Optical Font Recognition Using Typographical Features," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 20, No. 8, Aug. 1998, pp. 877-882.
- [4] Abdelwahab Zramdini, *Study of Optical Font Recognition based on Global Typographical Features*, Ph. D. Dissertation, University of Fribourg, Switzerland, 1995, pp. 68-69.
- [5] Min-Chul Jung, Yong-Chul Shin and Sargur N. Srihari, "Multifont Classification using Typographical Attributes," *Fifth International Conference on Document Analysis and Recognition*, 1999.
- [6] Shuichi Tsujimoto and Haruo Asada "Resolving Ambiguity in Segmenting Touching Characters," *First International Conference on Document Analysis and Recognition*, 1991, pp. 701-709.