

Automatic Separation of Machine-Printed and Hand-Written Text Lines

U. Pal and B. B. Chaudhuri

Computer Vision and Pattern Recognition Unit

Indian Statistical Institute

203 B. T. Road, Calcutta - 700 035, INDIA

Email: {umapada,bbc}@isical.ac.in

Abstract

There are many types of documents where machine-printed and hand-written texts intermixedly appear. Since the optical character recognition (OCR) methodologies for machine-printed and hand-written texts are different, it is necessary to separate these two types of text before feeding them to the respective OCR systems. In this paper, we present such a scheme for both Bangla and Devnagari. The scheme is based on the structural and statistical features of the machine-printed and hand-written text lines. The classification scheme has an accuracy about 98.3%.

1. Introduction

Existence of hand-written and machine-printed text in a single document is common in many types of document. For example, question papers where answers are to be written by hand on the blank spaces given below the questions, business and personal letters, application forms etc. From the comprehensive survey papers [4,6,8,11] it can be understood that machine-printed and hand-written character recognition schemes are quite different from each other. So, if a document contains both machine-printed and hand-written portions, they should be separated and fed to the respective OCR systems [1].

There exist a few papers on the classification of machine-printed and hand-written text but they deal with English, Chinese and Japanese scripts. In 1993, Imade *et al.* [5] described a method to segment a Japanese document into machine-printed Kanji and Kana, hand-written Kanji and Kana, photograph and printed image. By extracting the gradient and luminance histogram of the document image, they use a layered feed forward neural network model in their system. Franke and Oberlander [3] reported a method to check whether a data field in a form is hand or machine printed. In 1995, using

directional and symmetrical features as the input of a neural network, Kuhnke *et al.* [8] developed a method to identify machine-printed and hand-written English characters. Recently, Fan *et al.* [2] described a method for the classification of machine-printed and hand-written text lines from English, Japanese and Chinese scripts. They used spatial features and character block layout variance as the prime features in their approach.

This paper deals with the separation of machine-printed and hand-printed text both in Bangla and Devnagari, the two most popular scripts in south Asia. The approach used here is based on the distinctive structural and statistical features of machine-printed and hand-written text lines in these scripts. To the best of our knowledge this is a pioneering work of its kind on Indian language scripts.

2. Bangla and Devnagari script properties

Hindi and Bangla are the most popular languages in Indian sub-continent, and the 4th and 5th most popular language in the world, respectively. The script form of Hindi is called Devnagari, while that of Bangla is called Bangla. Devnagari script is used to write Hindi, Nepali, Marathi and Sindhi languages while Bangla script is used to write Bangla, Assamese and Manipuri languages. Bangla and Devnagari script are originated from the ancient Brahmi script and because of their same origin, these two scripts have some structural features in common. These common features help us to build up the system.

The properties of Bangla and Devnagari scripts that are useful for the present work are given below.

1. There are 11 vowel and 39 consonant characters in modern Bangla alphabet. They are called *basic characters*. In Devnagari, there are 49 basic characters. Out of them 11 are vowels and 38 are consonants. The concept of upper/lower case character is absent in these scripts.

2. Many characters of Bangla and Devnagari alphabet have a horizontal line at the upper part. In Bangla, this line is called *matra*, and in Hindi it is called *Sirorekha*. However, we shall call them here as *head-line*. When two or more Bangla or Devnagari characters sit side by side in proper alignment to form a word, the matra or sirorekha portions touch one another and generate a long head-line, which is used as a feature to isolate machine-printed and hand-written text line.

3. In Bangla (Devnagari) some vowels following a consonant take a modified shape, which depending on the vowel is placed to the left, right or both (in case of Bangla), top or bottom of the consonant. They are called *modified characters*.

4. A Bangla or Devnagari text line may be partitioned into three zones. The *upper zone* denotes the portion above the head-line, the *middle zone* covers the portion of basic (and compound) characters below head-line and the *lower zone* is the portion where some of the modifiers can reside. The imaginary line separating middle and lower zone is called *base line*. A typical zoning is shown in Fig.1.

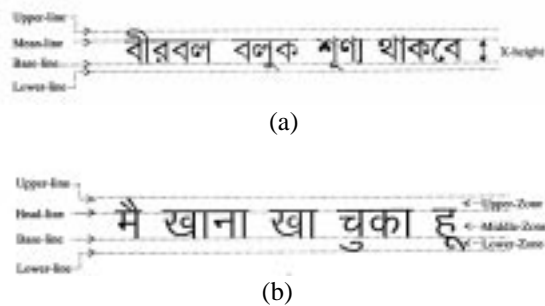


Fig.1 Different zones of (a) Bangla and (b) Devnagari script line.

3. Preprocessing

Text digitization for the experiment of the system has been done by a flatbed scanner (manufactured by HP, Model no ScanJet 4C). The digitized images are in gray tone and we have used a histogram based thresholding approach to convert them into two-tone images.

For accurate text classification, the system should properly detect individual text columns and should accurately segment the lines from each text column. Different columns of a text document are detected using the run length smoothing approach due to Wang et. al [11].

After detection of each text column, the mode of the text (portrait mode or landscape mode) of the document

is determined. A text column is called as portrait (landscape) mode if the text lines in that column are in horizontal (vertical) manner. The white space between characters is always much smaller than the white space between the lines. We use this geometric criterion for text mode determination.

The lines of a text block are segmented by noting the valleys of the projection profile. The position where profile height is least denotes one boundary line. A text line can be found between two consecutive boundary lines.

4. Classification of machine-printed and hand-written text line

Our separation scheme is a tree classifier where in the nodes of the tree we use some simple and easily detected features of machine-printed and hand-written texts. The flow-diagram of the classification scheme is shown in Fig.2. We shall discuss here the classification technique of portrait mode document. The classification technique for landscape documents can be done in a similar way. The features used in the scheme are as follows:

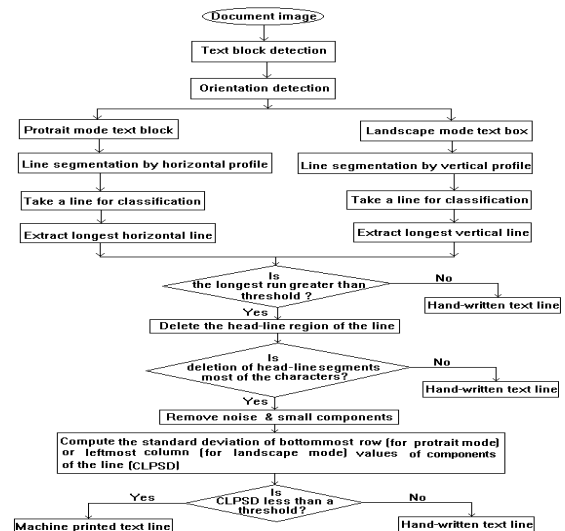


Fig.2: Flow diagram of the classification scheme.

(a) **First level feature :** Since characters of a word sit side by side in proper alignment in a machine-printed text line, the head-line portions of the characters in a word touch one another and generate a long head-line. At first, we use this feature for classification. The hand-written text lines can be separated from machine-printed text lines by computing the longest row-wise horizontal run. We

have noted that machine-printed text lines always generate a long horizontal run. The hand-written text lines may or may not generate such a long run. If the length of the longest run is less than a threshold T_1 then we classify that text line as hand-written line. The value of T_1 is set as twice the height of middle zone of a text line. Otherwise, it is either machine-printed or hand-written line. For illustration see Fig.3. In this figure, the first and third line are hand-written while the second line is machine-printed. For the second and third line, the longest horizontal run is greater than T_1 although the third text line is a hand-printed line, while for the first line this run is less than T_1 . Since the longest run for the first line is less than T_1 , we can classify this line as hand-written without using other features.

(b) Second level feature : We noted that characters in Bangla or Devnagari machine-printed word are connected through the head-line. Also, head-line of a word is properly aligned. Now, if we delete the head-line region from a text line then for machine-printed document all characters in that line get isolated whereas for hand-written document all characters may not be isolated because of the irregular alignment of the characters in words. If all characters are not isolated by the deletion of head-line region, we declare that line as hand-written. Else, it may be a machine-printed or hand-written line. See Fig.4 for illustration. Here, three text lines and their situation after deletion of head-line region are shown. From Fig.4(a) and Fig.4(b) it can be noted that the characters are topologically segmented due to the deletion of head-line region although Fig.4(a) is machine-printed and Fig.4(b) is hand-written text line. But for the hand-written text line shown in Fig.4(c), characters are not topologically segmented after deletion of head-line region. Thus, we can classify this line as hand-printed without using third level feature.

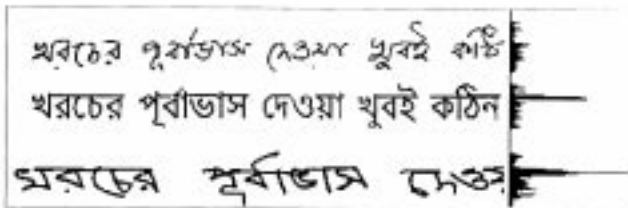


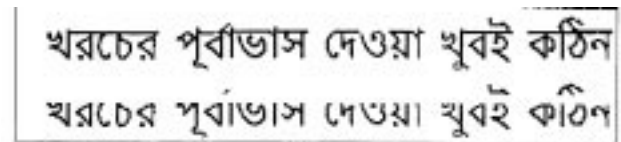
Fig.3: Example of longest run of three text lines. Here second line is machine-printed and others are hand-written.

The head-line region deletion is done as follows. From the horizontal projection of a line, we note the row corresponding to the peak of the profile. From the experiment we noted that the head-line region is about

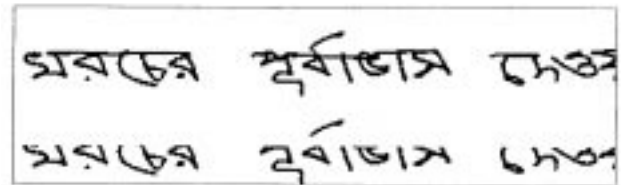
12% of a text line height. We delete 12% of the line height considering peak row as the center.

(c) Third level feature : Here, to identify a line we note the distribution of lowermost points of isolated components. Note that we use this feature when all characters of a line are properly segmented due to deletion of head-line regions. We note that the distribution of character lower most points is regular in machine-printed texts, and irregular in hand-written texts. This property is used at the third level feature for the identification of machine-printed and hand-written text lines.

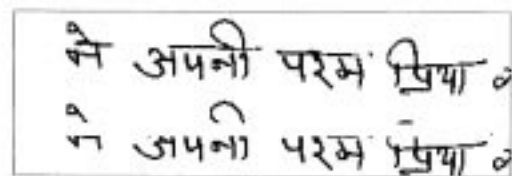
In machine-printed text, we note that the lowermost points of most of the characters of a text line lie only on two horizontal lines. For the characters to which a lower modifier is attached, the lowermost points lie on lower-line. Otherwise, the lowermost points lie on the base-line. For, example see the printed text lines shown in Fig.1. Here, the lowermost points of the characters lie either on base-line or lower-line. This is not true in hand-written text line.



(a)



(b)



(c)

Fig.4: Example of three text lines and their situation after deletion of head-line region are shown.

For a text line we compute two sets of lowermost points B and L corresponding to base-line and lower-line. If the lowermost point of a component does not lie on any one of these two lines then we include this point in one of the

two sets as follows. Let B_r and L_r be the row numbers corresponding to base-line and lower-line. Now, a component with lower-most row C_r belongs to the set B if $|B_r - C_r| \leq |L_r - C_r|$. Else, it belongs to L. Let b_1, b_2, \dots, b_m be m lowermost row values of m components belonging to set B and let l_1, l_2, \dots, l_p are p lowermost row values of p components that belong to set L. We noted that for machine-printed lines most of the elements of set B are equal i.e. they lie on the same row. This distribution is true for the set L also. But it is not true in hand-written text lines. Now, a spatial feature called *character lowermost point standard deviation* (CLPSD) is defined as

$$CLPSD = \sqrt{\frac{1}{m} \sum_{i=1}^m (b_i - \bar{b})^2} + \sqrt{\frac{1}{p} \sum_{i=1}^p (l_i - \bar{l})^2}$$

where

$$\bar{b} = \frac{1}{m} \sum_{i=1}^m b_i \quad \text{and} \quad \bar{l} = \frac{1}{p} \sum_{i=1}^p l_i$$

A line is classified as machine-printed if the value of CLPSD is smaller than a value r_1 . Otherwise, it is called hand-printed. The value of r_1 is computed as

$r_1 = (\text{Average height of the components whose lower points are considered})/10$

Due to the dots of some characters in Bangla and Devnagari, or due to some punctuation marks like comma, or due to salt and pepper noise sometimes we may get high CLPSD value in a machine-printed text line and hence it may be wrongly identified as hand-written line. To tackle this situation we find lowermost points only of those components whose bounding box widths are greater than half of the average bounding box width of all components in the line. Hence, small and irrelevant components like dots of the characters as well as noise and punctuation marks are mostly filtered out.

5. Results and Discussion

To demonstrate the feasibility and validity of our proposed approach a wide variety of document images were tested. We applied our separation scheme on 100 different document images. The images were scanned from question papers, money order form, application

form, letter etc. We have noted that accuracy of the system is about 98.3%. We also noted that most of the identification errors are obtained from very short lines only.

Because of the use of simple features, which are easy to compute, our separation scheme is very fast. Also, the scheme does not depend on size and font of the characters in the text line.

India is a multi-lingual country, where a single document page may contain two or more language scripts. So, multi-script OCR is useful and important for such a country. To make a successful multi-script OCR, it is necessary to separate different script lines before feeding to OCR system of individual script. Thus, the work presented here has a strong direct application potential. This work is the first of its kind and can be applied for the separation of other similar scripts.

References:

1. B. B. Chaudhuri and U. Pal, "An OCR system to read two Indian language scripts: Bangla and Devnagari", *In Proc. 4th ICDAR*, pp. 1011-1015, 1997.
2. K. C. Fan, L. S. Wang and Y. T. Tu, "Classification of machine-printed and hand-written texts using character block layout variance", *Pattern Recognition*, Vol. 31, pp. 1275-1284, 1998.
3. J. Franke and M. Oberlander, "Writing style detection by statistical combination of classifiers in form reader application", *In Proc. 2nd ICDAR*, pp. 581-584, 1993.
4. V. K. Govindan and A. P. Shivaprasad, "Character recognition -- a review", *Pattern Recognition*, Vol. 23, pp. 671-683, 1990.
5. S. Imade, S. Tatsuta and T. Wada, "Segmentation and Classification for Mixed Text/Image Document Using Neural Network", *In Proc. 2nd ICDAR*, pp. 930-934, 1993.
6. S. Impedovo, L. Ottaviano and S. Occhinegro, "Optical character recognition -- a survey", *Int. J. Pattern Recognition Artif. Intell.*, Vol. 5, pp. 1-24, 1992.
7. S. Kahan, T. Pavlidis and H. S. Baird, "On the recognition of Printed Character of any font and size", *IEEE Trans. on Patt. Anal. and Mach. Intell.*, Vol. 9, pp. 274-288, 1987.
8. K. Kuhnke, L. Simoncini and Z. M. Kovacs-V, "A. system for machine-written and hand-written character distinction". *In Proc 3rd ICDAR*, pp. 811-814, 1995.
9. S. Mori, C. Y. Suen and K. Yamamoto, "Historical review of OCR research and development", *Proceedings of the IEEE*, Vol. 80, pp. 1029-1058, 1992.
10. G. Nagy, "Chinese character recognition : a twenty-five year retrospective", *In Proc. 9th ICPR*, pp. 163-167, 1988.
11. K. Y. Wang, R. G. Casey and F. M. Wahl, "Document analysis system", *IBM J. Res. Development*, Vol. 26, pp. 647-656, 1982.