

Offline Character Recognition of Printed Tamil Text using Template Matching Method of Bamini Tamil Font

D. Pugazhenthithi* and S. Arul Vallarasi

Department of Computer Science, Quaid-E-Millath Government College for Women (Autonomous), Anna Salai, Chennai – 600002, Tamil Nadu, India; pugazh006@gmail.com, arul.s.vallarasi@gmail.com

Abstract

Background/Objectives: Character recognition of the English alphabet using template matching method is a simple method to implement. This paper proposes Tamil character recognition of Bamini Tamil font using Template Matching method. **Materials/Methods:** The document image without skew is binarized at the preprocessing step. The preprocessed image is then segmented. Every line of text is segmented using horizontal projection analysis. Every character in a line is segmented using connected component processing. Then each character segmented is correlated with the preloaded templates of the system. The maximum correlation judges the character. In the same way, every segmented input is checked with the preloaded templates. These templates are mapped onto Tamil Unicode for recognition. The text is reconstructed using Unicode fonts and finally produces the Machine editable Unicode text in a text file. **Conclusion/Findings:** The system gives results considerably greater than 20 pixel base height of a character in the document image. **Applications/Improvements:** The possibilities of using other fonts character recognition is applied in future. Other methods are also considered for further implementation.

Keywords: Connected Component Labeling, OCR, Offline Character Recognition, Segmentation, Template Matching

1. Introduction

Optical Character Recognition (OCR) is a field of research in pattern recognition, artificial intelligence and computer vision. OCR is the mechanical or electronic conversion of typewritten or printed text images into editable text. It is a technology that converts different kinds of inputs, such as scanned documents, images captured by a digital camera into editable and searchable data. In other words, producing documents like a text file from a scanned image of a printed or handwritten page. OCR is a common method of digitizing printed texts so that it can be electronically edited, searched and used in machine processes such as machine translation (translate text or speech from one language to another), text-to-speech and text mining (process of deriving high-quality information from text).

The Optical Character Recognition can be widely categorized into two; online and offline¹⁻³. In offline character recognition scanned images of text is converted into machine editable text. Whereas in online character recognition the input of the OCR is received by means of digital pen, digital stylus like devices and the input is converted into machine editable text. There are many methods available for printed text recognition. In templates matching method, individual character templates are added into the system initially. Once the templates creation is over, the test image is fed into the system. The system segments the image into individual characters and the segmented image is correlated with the preloaded templates⁴. The maximum correlated character is declared as the character present in the image. The same process is repeated for every segmented character in the image.

* Author for correspondence

divided into three say 'n', 'g', 'h'. Likewise the compound letter 'gh' is divided into two say, 'g', 'h'. This is because the character has split parts in other words it has the white space between the three parts. And gets three individual labels.

Template matching is a technique in digital image processing for finding small parts of an image which match a template image¹³. The labeled image from the previous stage is split according to the label number. Every labeled image is correlated with the preloaded templates. The correlation coefficient is calculated using the below formula

$$\gamma = \frac{\sum_m \sum_n (A_{mn} - \bar{A})(B_{mn} - \bar{B})}{\sqrt{\left(\sum_m \sum_n (A_{mn} - \bar{A})^2\right) \left(\sum_m \sum_n (B_{mn} - \bar{B})^2\right)}} \quad (1)$$

Where m is row number and n is column number. Here A_{mn} is the pixel intensity or the gray scale value at a point (m,n) in the template image. B_{mn} is the gray scale value at a point (m,n) in the test image. \bar{A} = mean of (A) and \bar{B} = mean of (B)

The maximum correlated character is declared as the character present in the labeled image. The same process is repeated for every segmented character (labeled image) in the image. Then the text is reconstructed using Unicode font 'Arial Unicode MS'. Tamil characters range from '0B80' to '0BFF' in the Unicode consortium¹⁴. It finally produces the machine editable unicode text in a text file. The system is worked with Bamini Tamil True type font.

3. Results and Discussion

The system is implemented in Matlab R2013b. The

Sample Input Image of 40 pixel height (single character base height) is given below. The corresponding output for the input image is also given below. The sample output image in a text file. The document image is neither bold nor italic. The test image is passed through the system, the test document image is preprocessed, each character is separated and correlated with the preloaded templates, the maximum correlated template judges the character, and the output is given in a text file.

Table 1. Experimental Results

| Number | Base Height of a character in pixel | Percentage of Correctly judged characters | Percentage of Wrongly judged characters |
|--------|-------------------------------------|---|---|
| 1 | 10 | 48.14815 | 51.85185 |
| 2 | 11 | 43.51852 | 56.48148 |
| 3 | 13 | 41.66667 | 58.33333 |
| 4 | 15 | 26.85185 | 73.14815 |
| 5 | 17 | 56.48148 | 43.51852 |
| 6 | 19 | 82.40741 | 17.59259 |
| 7 | 21 | 88.88889 | 11.11111 |
| 8 | 23 | 93.51852 | 6.481481 |
| 9 | 25 | 90.74074 | 9.259259 |
| 10 | 27 | 90.74074 | 9.259259 |
| 11 | 29 | 92.59259 | 7.407407 |
| 12 | 31 | 91.66667 | 8.333333 |
| 13 | 33 | 92.59259 | 7.407407 |
| 14 | 35 | 92.59259 | 7.407407 |
| 15 | 37 | 91.66667 | 8.333333 |
| 16 | 39 | 90.74074 | 9.259259 |
| 17 | 41 | 93.51852 | 6.481481 |
| 18 | 43 | 96.2963 | 3.703704 |

Offline cursive handwritten character recognition and skew corrections¹⁵⁻¹⁷ are not added in the system. The

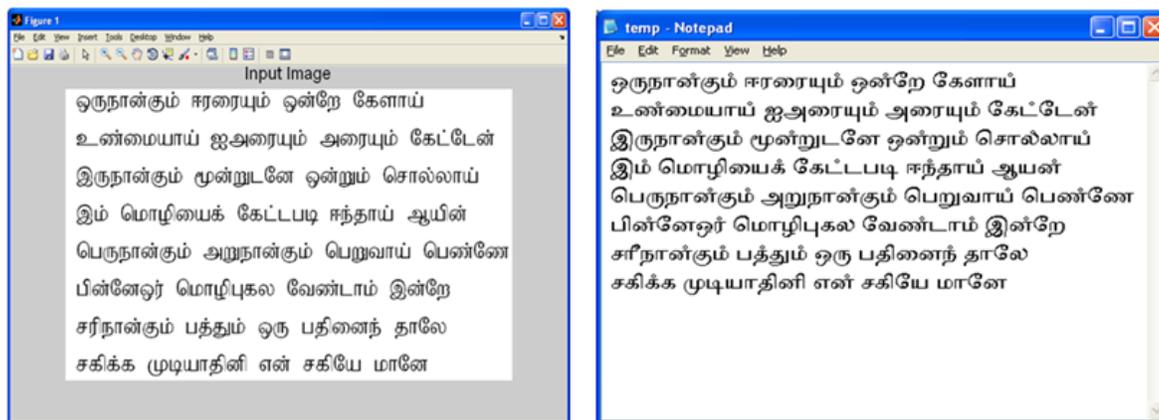


Figure 2. Input document image and its corresponding output in notepad.

Tamil character document image of 108 characters long is recognized in the system and the result is listed in below Table 1. For every character base height in pixel, correctly judged and wrongly judged character count is listed. For small font images the system gives more wrong judgment. For big font images the system gives better results. It gives reliable results greater than 20 pixel height (single character base height). If the pixel height of the character in the document image is greater than 20 pixels, then the output text is readable (with minimum error).

4. Conclusion

This paper presents Tamil Optical Character Recognition system which uses image binarization, Line segmentation, Character segmentation, Template matching and Tamil character reconstruction. The advantage of this method, Template matching makes Tamil character recognition simple. The system gives considerable results greater than 20 pixel base height of a character in the document image. The discovered result gives good accuracy for big fonts. The system loses its accuracy when the font size in the document image is small. But in normal practice while scanning the document, the image with high character pixel height is not possible. So the template matching method is less desirable for the Tamil character recognition system. But mostly the documents have small fonts. The system has not considered the skewed image. The future direction of the work is to apply skew correction in the scanned document image, and the Tamil character recognition system for all size of fonts.

5. Reference

- Chopra SA, Ghadge AA, Padwal OA, Punjabi KS, Gurjar GS. Optical character recognition. *International Journal of Advanced Research in Computer and Communication Engineering*. 2014; 3(1):4956–8.
- Dedgaonkar SG, handavale AA, Sapkal AM. Survey of methods for character recognition. *International Journal of Engineering and Innovative Technology*. 2012; 1(5):180–9.
- Kannan RJ, Prabhakar R. A comparative study of optical character recognition for Tamil script. *Euro J Scientific Res*. 2009; 35(4):570–82.
- Charles PK, Harish V, Swathi M, Deepthi CH. A review on the various techniques used for optical character recognition. *International Journal of Engineering Research and Applications*. 2012; 2(1):659–62.
- Shrivastava V, Sharma N. Artificial neural network based optical character recognition. *arXiv preprint arXiv:1211.4385*. 2012; 3(5):73–80.
- Alon J, Athitsos V, Sclaroff S. Online and offline character recognition using alignment to prototypes. *IEEE Proceedings 8th International Conference on Document Analysis and Recognition*; 2005. p. 839–43.
- Yasser A. *Preprocessing techniques in character recognition*. INTECH Open Access Publisher; 2010. p. 1–20.
- Su B, Lu S, Tan SL. Combination of document image binarization techniques. *Document Analysis and Recognition (ICDAR)*. International Conference on IEEE; 2011. p. 22–6.
- Richard GC, Lecolinet E. A survey of methods and strategies in character segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 1996; 18(7):690–706.
- Karmakar P, Nayak B, Bhoi N. Line and word segmentation of a printed text document. *International Journal of Computer Science and Information Technologies*. 2014; 5(1):157–60.
- Aparna KG, Ramakrishnan AG. A complete Tamil optical character recognition system. *Document Analysis Systems V*. Springer Berlin Heidelberg; 2002. p. 53–7.
- Kanimozhi VM, Muthumani I. Optical character recognition for English and Tamil script. *IJCSIT*. 2014; 5(2):1008–10.
- Due TO, Jain AK, Taxt T. Feature extraction methods for character recognition-A survey. *Pattern Recognition*. 1996; 29(4):641–62.
- Seethalakshmi R, Sreeranjani TR, Balachandar T. Optical character recognition for printed Tamil text using Unicode. *Journal of Zhejiang University Science A*. 2005; 6(11):1297–305.
- Kannan RJ, Prabhakar R, Suresh RM. Off-line cursive handwritten Tamil character recognition. *IEEE International Conference on Security Technology, SECTECH'08*; 2008. p. 159–64.
- Dan SB, Kopec GE, Dasari L. Measuring document image skew and orientation. *IS&T/SPIE's Symposium on Electronic Imaging: Science and Technology*. International Society for Optics and Photonics; 1995. p. 302–16.
- Adams M. Algorithm for text document de-skewing. *EECS 490-Digital Image Processing Session-1*; 2004.