

A Novel Segmentation Technique for Printed Malayalam Characters

Bindu Philip, R. D. Sudhaker Samuel and C. R. Venugopal, *Member, IACSIT*

Abstract— Segmentation of whole characters in Indian scripts is a rather tricky problem. This is because in Indian scripts the characters have modifiers which could be subscripts, superscripts, attached or nonattached vowel signs to the front or following the character forming a complex composite character. Often a character is composed of several sub-characters. Unlike other South Indian scripts like Kannada, Telugu and Tamil, in the Malayalam script the space between the sub-characters of the character is same as the space between the characters within a word rendering the character segmentation process quite complex as conventional profiling methods fail. This paper presents a novel segmentation algorithm for segmentation of character in such complex cases taking Malayalam as a typical example. The success of the algorithm is demonstrated by application of feature extraction on segmented characters and subjected to classification. Segmentation efficiency of 98.8 % is achieved which is very encouraging.

Index Terms — Segmentation, Feature Extraction, Classification.

I. INTRODUCTION

Segmentation is an integral part of any optical character recognition system. Proper segmentation assures efficiency of classification and subsequent recognition. For an optical character recognition (OCR) system, segmentation phase is an important phase and accuracy of any OCR heavily depends upon segmentation phase [1]. Incorrect segmentation leads to incorrect recognition. Segmentation phase includes line, word and character segmentation.

It is essential to obtain complete segmented character (a character along with its modifiers) without the presence of extraneous data to ensure robust feature extraction. Indian languages are uncharacteristic in the sense that unlike English and other European languages a single character would consist of several sub-characters sometimes to the extent of 3-4 segments.

It is further complex when the sub-characters of the characters are spaced to the same extent as that of the characters within a word. This is the motivation for addressing issues in segmentation of such complex

characters in a word.

Common classical segmentation approaches like horizontal and vertical profiling techniques fail to segment the complete characters in scripts like Malayalam; where it is observed that the distance of the space between the sub-characters of the character is same as the space between the characters within a word rendering the character segmentation process quite complex.

A major difficulty of developing a good OCR system in most Indian languages lies in the large set of basic, modified and compound characters in the text. Further, for some scripts the characters get topologically connected in a word. Hence, prior character segmentation is necessary to input the data in the character-based recognition engine. An Arabic OCR system using recognition based segmentation is discussed by Cheung, *et.al.* [2]. D. G. Elliman and I. T. Lancaster [5] presented a review of segmentation and contextual analysis for text recognition. K. S. Sesh Kumar, *et.al.*[6] proposed that spatial language models be used to segment the documents of a particular script. Kunte Sanjeev R and Sudhaker Samuel R.D [7] proposed a two-stage character segmentation scheme for Printed Kannada text which was designed to first check for the presence of subscripts in a word before segmentation. Ashwin T.V and Sastry P.S [8] proposed a font and size-independent OCR system for printed Kannada documents using support vector machines but the recognition rate obtained was 86-11%. Non-uniformity in the spacing of the characters within a word in most of the Indian scripts helps in segmentation. Consonants take modified shapes when attached with the vowels. Vowel modifiers can appear to the right, left, on the top or at the bottom of the base consonant. Such consonant-vowel combinations are called modified characters. In addition, two, three or four characters can combine to generate a new complex shapes called compound characters. These characters are very difficult for a machine to recognize. OCR preprocessing stage is the most important because it directly affects the reliability and efficiency in the segmentation and feature extraction process. Thus proper segmentation of a complete character plays a major role in increasing the recognition rate substantially. Character segmentation is a crucial phase in the character recognition engine. This is an intermediate step wherein the given input image is decomposed into a sequence of character segments of individual symbols. The segmentation module is a decisive step in a system as it evaluates whether a pattern isolated from the given image is meaningful or not and thus has a major contribution to improving the recognition rate of the system. Segmentation in Malayalam script is a challenge

Bindu Philip is a full time Research Scholar at JSS resaaech Foundation in the Electronics and communication Department at SJCE, Mysore, E-mail: binduthomas25@yahoo.co.in.

Dr. R. D. Sudhaker Samuel was with the Electronics and communication department, S. J. College of Engineering, Mysore(e-mail: sudhakersamuel@yahoo.com).

Dr. C. R. Venugopal is with the Electronics and communication department, S. J. College of Engineering, Mysore, E-mail: venu713@gmail.com

because of the fact that the connected components/ segments within a complete character and the segments of the following character are uniformly spaced between each other. Thus, segmenting a complete character within a word having a train of uniformly spaced segments has always been a challenging task. A novel character segmentation technique is proposed here which demonstrates 98.8 % efficiency.

II. MALAYALAM SCRIPT

The Malayalam script is an abugida (segmental writing system) of the Brahmic family, used to write the Malayalam language. Both the language and its writing system are closely related to Tamil, although Malayalam has a larger phoneme inventory. The Malayalam script is used to write the language of Malayalam spoken in the Southern Indian state of Kerala. It is a descendent of the Grantha script.

The characters are classified into two categories: svarams (or vowels) and vyanjanams (or consonants). In the Malayalam Script, there are many ways to form words. In the most straightforward case, the svaram and the vyanjanam can be strung together to form words. Often, it is more complicated because svaram characters are used only when a word begins with a vowel and vyanjanam characters are used to signify a consonant and a vowel [13]. Like other South Asian scripts, a Malayalam letter has an inherent vowel sound /a/ which is used with each unmarked or basic consonant symbol.

അ	ആ	ഇ	ഈ
ah	aah	yi	yee
ഉ	ഊ	ഋ	ഠ
uh	ooh	er	eru
എ	ഏ	ഐ	ഓ
eh	aeh	ai	oh
ഔ	ഓ	അം	അഃ
ohoh	ow	am	aha

Figure1: Malayalam svarams (vowels)

The character set of the Malayalam script consists of 16 vowels shown in Figure 1, 36 consonants shown in Figure 2, vowel signs shown in Figure 3 besides 30 commonly used conjuncts, few examples are shown in Figure 4.

ക	ഖ	ഗ	ഘ	ങ
ka	kha	ga	gha	nga
ച	ഛ	ജ	ഝ	ഞ
cha	chha	ja	jha	nja
ട	ഠ	ഡ	ഢ	ണ
ta	tta	da	dda	nha
ത	ഥ	ദ	ധ	ന
tha	thha	dha	dhha	Na
പ	ഫ	ബ	ഭ	മ
pa	pha	ba	bha	Ma
യ	ര	ല	വ	ശ
ya	ra	la	va	Sha
ഷ	സ	ഹ	ള	റ
shha	sa	ha	lha	Rha
ഴ				
zha				

Figure 2: Malayalam Consonants (vyanjanams)

To change this consonant to another, extra strokes called matras are added to the basic letter, in the Figure 3 column 3.

Vowel Sign	Left/Right of the Consonant/ conjunct	Vowel sign attached to ക (ka) (example)	
ഠ	Right	കഠ	kA
ഡ	Right	കി	Ki
ഢ	Right	കീ	kI
ണ	Right	കു	Ku
ൠ	Right	കു	kU
ൡ	Right	കു	kRu
വ	Left	കെ	ke
ൣ	Left	കേ	kE
൤	Left	കൈ	kai
൦ ഠ	Left & Right	കൊ	ko
ൣ ഠ	Left & Right	കോ	kO
ൡ	Right	കൗ	kau
ൠ	Right	കം	kaM
ൡ	Right	കഃ	kaH

Figure 3: Vowel Sign Modifiers

ക ക്ക ച ക്ക ച്ച ട ട്ത ഡ ഡ്ദ ണ ണ്ണ

(KkA) (ChA) (Tta) (Tha) (nga) (ngya)

Figure 4: Consonant conjuncts

The dependent vowels appear in combination with a consonant or a consonant cluster. Explicit appearance of a dependent vowel in a syllable overrides the inherent vowel of a single consonant character [14]. The independent vowels are used to write syllables which start with a vowel. The positioning of the dependent vowel may be to the left, to the right, or both to the left and right of the consonant/ conjunct, depending on the vowel sign being attached. Malayalam has remarkably distinct lateral variations when compared to many other Indian languages with a number of curls and twists in the characters. Another interesting feature of this script is that the number of columns varies from 53 to a phenomenal 347 columns over the entire extended character set of the language for the resolution used here.

III. SEGMENTATION

The digitized image is segmented to extract individual characters. The projection techniques, (horizontal, for line segmentation and vertical for the word segmentation) were proposed by [3]. Other techniques such as connected component approach [4] can also be adopted. Each individual character can be of different size and hence, we convert them to images of fixed size by scale normalization technique by preserving the length of the character but normalizing the height to m=50, explained in Section IV. A Malayalam character could consist of several uniformly spaced unconnected components such as a vowel (only at the beginning of a word) or a consonant/ conjunct along with vowel signs. Conventional techniques like horizontal and vertical projection profile methods fail to segment the complete character correctly because of the equal space

representing vowel signs appearing to the right of the consonant/ conjunct. This sub-character level of classification and recognition is done efficiently and fast by using the reduced classification search space as shown in Figure 7. The sequence of character segments enter the flow graph they go through the different passes to return the complete character clusters at the end of each pass. Thus a novel multiple pass segmentation algorithm is implemented with 98.8 % segmentation accuracy. A typical example with different valid sub-character sequence within a single word is used to explain the segmentation process in detail.



Figure 8(a) complete word in English is Kaikolanae (means Accept) (b) Four full characters (c) sub-character sequence

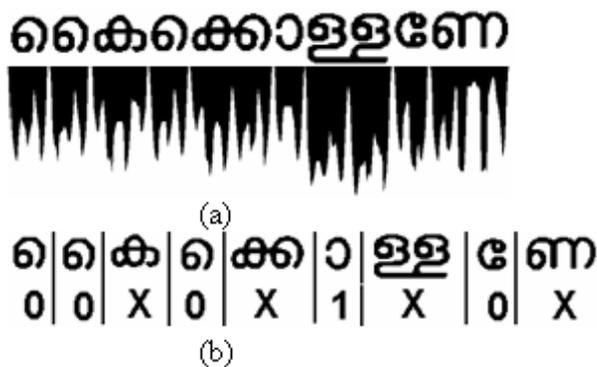


Figure 9(a) word and its vertical projection profile (b) all sub-character segments within the word

The word shown in the example requires four passes through our novel segmentation algorithm. Each pass is now explained. The first complete character segment **കൈ (0 0 X)** emerges from the first pass. The path taken for this is shown in Figure 10(a). The Second complete character segment **കൊ (0 X 1)** is obtained at the end of second pass. The path taken for this is shown in Figure 10(b). The third segment **ള്ള (X)** is returned at the end of third pass. The path taken for this is shown in Figure 10(c). The final and fourth segment **ണം (0X)** is returned at the end of fourth pass. The path taken for this is shown in Figure 10(d).

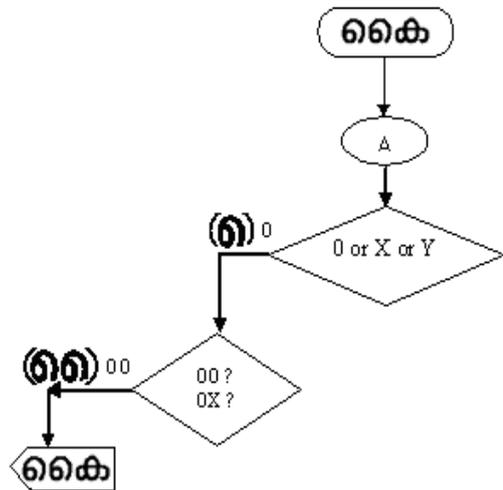


Figure 10 (a) Logical path to segment second complete character of കൈക്കൊള്ളണം pass1

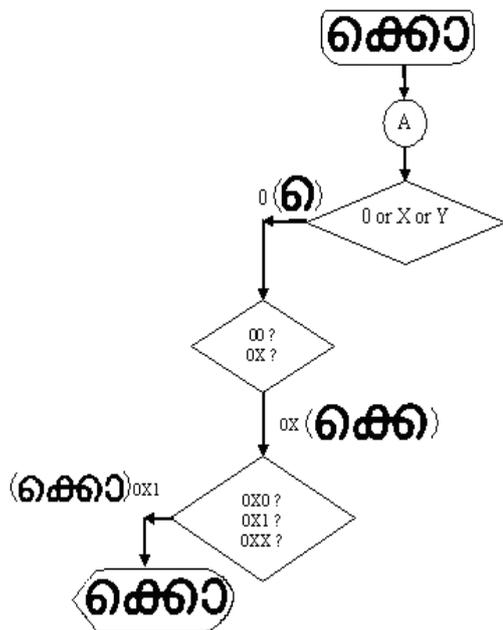


Figure 10(b) Logic path to segment second character of കൈക്കൊള്ളണം pass2

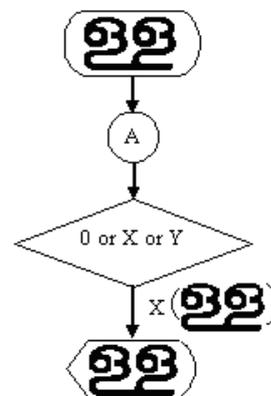


Figure 10(c) Logical path to segment second character of കൈക്കൊള്ളണം pass3

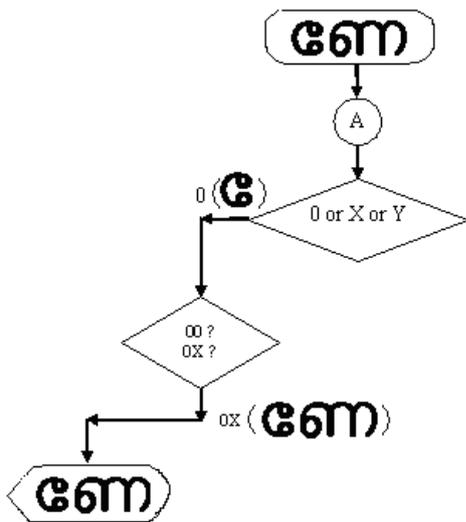


Figure 10(d) Logic path to segment second character of കൈക്കൊള്ളുണേ pass4

Thus a given word can be efficiently segmented to complete character segments with as many numbers of passes as that of the complete characters within the word.

Another word is considered again to illustrate the robustness of the novel segmentation algorithm to obtain segmented complete printed Malayalam characters. Figure 11 (a) which requires to be segmented as shown in Figure 11 (b). The valid character sequence for the word considered as per the representation used for sub-characters in our segmentation algorithm is shown in Figure 11 (c). The classical vertical projection profile segments the word as shown in Figure 12 (a) and the sequence of sub-character segments obtained is as shown in Figure 12 (b). This sequence of segments represented by either X, Y, 0 or 1 as in Figure 12 (b) now enters the start node of the segmentation flow chart as depicted by flowchart shown in Figure 6

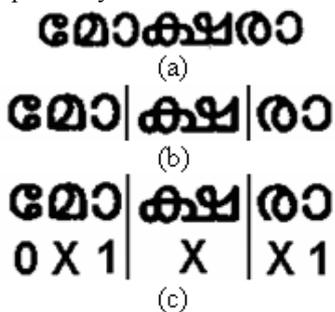


Figure 11(a) complete word- in English is Mokshara (means salvation)(b) three full characters (c) sub-character sequence

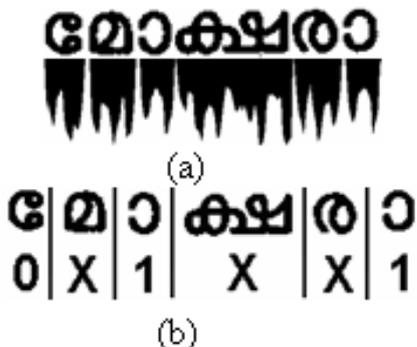


Figure 12(a) word and its vertical projection profile (b) all sub-character segments within the word

The word shown in the example 11 (a) requires three passes through our novel segmentation algorithm. Each pass is now explained. The first complete character segment മോ(0 X 1) emerges from the first pass. The path taken for this is shown in Figure 13(a). The Second complete character segment ക്ഷ(X) is obtained at the end of second pass. The path taken for this is shown in Figure 13(b). The third segment രാ(X 1) is returned at the end of third pass. The path taken for this is shown in Figure 13(c).

Thus here in this example the word has three full character segments and the given word is efficiently segmented to three complete character segments with three passes.

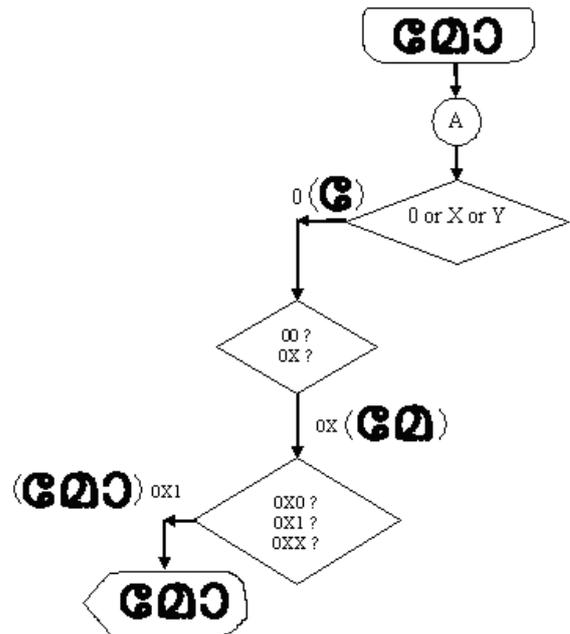


Figure 13 (a) Logical path to segment First complete character of മോക്ഷരാ pass1

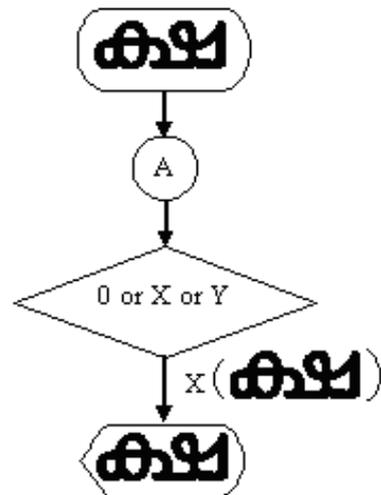


Figure 13 (b) Logical path to segment second complete character of മോക്ഷരാ pass2

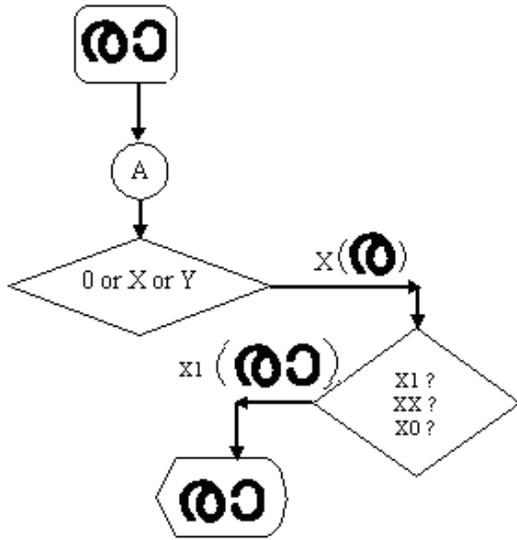


Figure 13 (c) Logical path to segment third complete character of മോക്ഷരോ pass3

IV. FEATURE EXTRACTION

The process of digitization followed by segmentation essentially renders the image in the form of an $m \times n$ matrix. These matrices are then generally normalized and then converted into a square matrix in order to apply the classical tools of linear algebra for characterization. It is easy to see that any form of characterization results in a reduction in dimensionality which essentially helps in the search process by classification over a large data base.

However, there are instances where rich information along rows would be lost in the process of reducing the image matrix to square. One good example is the segmented images of the characters of Malayalam language. It thus makes good sense to retain the number of columns. A rectangular, strictly black and white digital image preprocessed to remove any extraneous noise can be represented by a matrix A, where

$$A = (a_{ij}) \in \mathbb{R}^{m \times n} : a_{ij} = \{0, 1\}, \quad (1)$$

usually $n > m$.

In order to ensure practicality in classification and identification, for the matrix in Equation (1) reduction in dimensionality is performed to obtain a feature vector $x \in \mathbb{R}^m$, at the same time capturing the distinct information in all the n columns. The methods proposed in this paper essentially capture useful information along rows of matrixes and hence we call them lateral analysis. Selected features along rows are retained rather than losing them by compression, normalization or by looking only at the overall characteristic feature of a matrix. The matrix A as represented in Equation (1) has several lateral features. One of them is Frequency Capture. This process in principle captures the frequency of transitions along each row. The feature vector, $x \in \mathbb{R}^m$ of the matrix $A \in \mathbb{R}^{m \times n}$ in this

$$\text{approach is defined by } x_i = \sum_{j=1}^n |a_{i,j+1} - a_{i,j}| \quad (2)$$

This captures features of characters with multiple loops which is a distinct feature of Malayalam characters.

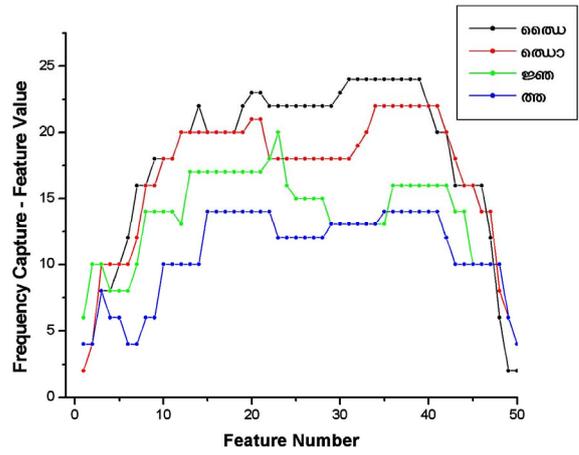


Figure 14 Frequency Capture Features versus Feature Number plot for Long Characters ജ്ജ (Jai), ജോ (Jo), ജ്ജ (jja), ത്ത (tta)

The extracted 50 features for four long characters ജ്ജ (Jai), ജോ (Jo), ജ്ജ (jja), ത്ത (tta) are plotted against their feature number and it is evident that the variation between class or interclass distance is large as shown in Figure 14. Similarly, the Frequency Capture Features versus Feature Number plot for Conjunct Characters ക്ക (kka), ത്ത (tma), ത്ത (ttha), ത്ത (dda) is shown in Figure 15 and the Frequency Capture Features versus Feature Number plot for General Characters ജ്ജ (jja), ഞ (zha), ത്ത (dda), ത്ത (eja) is shown in Figure 16.

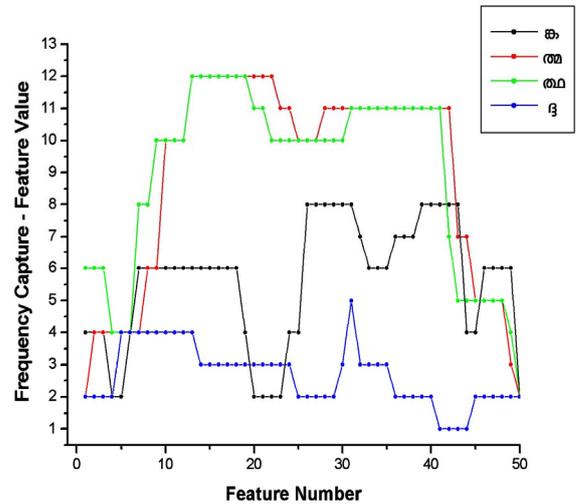


Figure 15 Frequency Capture Features versus Feature Number plot for Conjunct Characters ക്ക (kka), ത്ത (tma), ത്ത (ttha), ത്ത (dda)

In the same way the Frequency Capture Features for Similar Characters are plotted and as shown in Figure 17 (a) plot for ഞ (a) and ഞ (aa) (b) plot for ത്ത (Na) and ത്ത (nna) (b) (c) plot for ക്ക (ka) and ക്ക (kka), (d) plot for ത്ത (e) and ത്ത (ee). Thus the frequency capture features can be effectively used to classify segmented complete printed Malayalam characters efficiently.

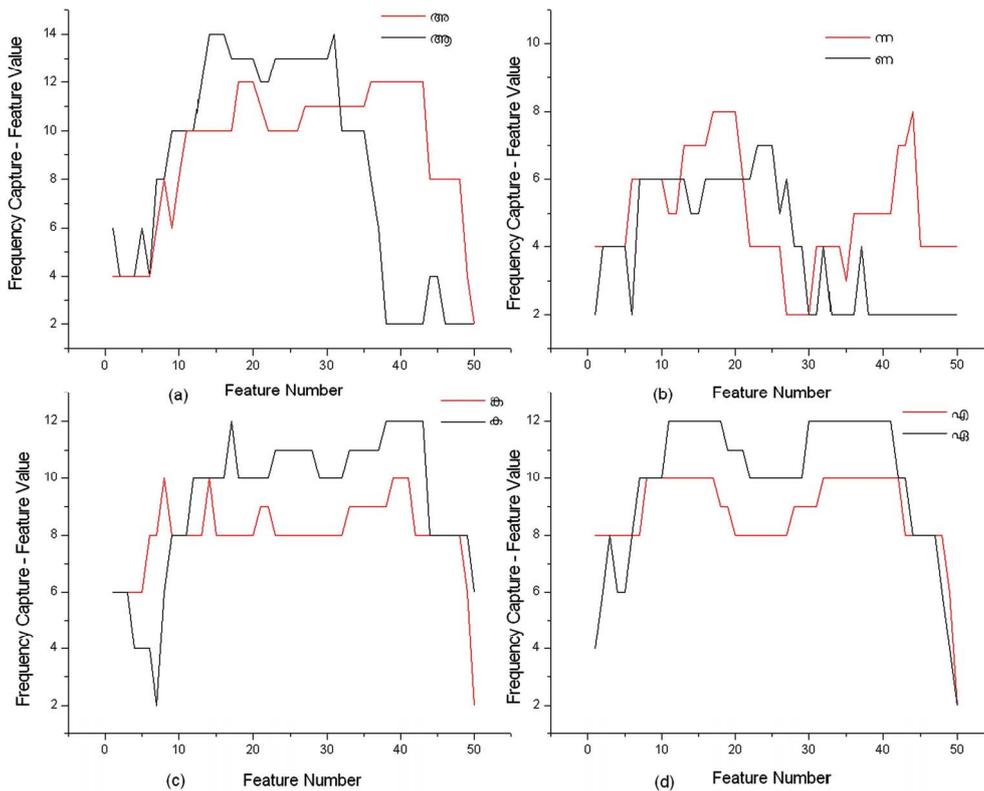


Figure 16 Frequency Capture Features for Similar Characters (a) plot for A (a) and B (aa) (b) plot for E (Na) and Y (nna) (c) plot for K (ka) and kka), (d) plot for F (e) and G (ee)

V. CLASSIFICATION

At any instance the search is restricted to the possible sub-characters that can occur. For example at decision level 'A', search is restricted to only all possible left vowel signs and consonants/ conjuncts. This takes away considerable computational load and reduces possibilities of mis-classification. Thus the first segment will be searched for in a space consisting of vowels (13 in all), consonants (36 in all) and the conjunct characters (30 in all), resulting in 73 possible classes. If the identified sub-character is a vowel segmentation takes place and the algorithm returns the classified Vowel (Vowels are not dependent and can appear only as the first character of a word). If the identified sub-character is a Left Vowel Sign (which is dependent and should be followed by a consonant or conjunct) segmentation does not take place but the next sub-character is identified at the next level of search among only consonants (36 in all) and the conjunct characters (30 in all) (a total of 66 classes within this reduced search space). After this sub-character is identified there is a possibility of yet a third sub-character segment with (X) being a possibility. The other possibility is that this sub-character segment could be either a vowel sign to the left (2 in all), consonants (36 in all) or the conjunct characters (30 in all). This constitutes the third level of search. If the identified sub-character segment is a consonant/ conjunct then the next segment is searched among right vowel signs, consonants (36 in all) / conjunct characters (30 in all).

Y (Vowels)
അ അതു ഇ ഇതു ഉ ഉതു ഋ ണി ണി ണി
ഒ ഓ ഔ അം അഃ
X (consonants)
ക ഖ ഗ ഘ ങ ച ഛ ജ ഝ ഞ
ട ഡ ള ഴ റ ത മ ഡ ണ
പ ഫ ബ ഭ ഢ ഡ റ റ ല വ ശ
ഷ സ ഹ ള ഴ
Or X (Conjuncts)
ക ച ള ത പ ള ള ജ ള ള ള ള ള ള
ന മ ള ള ള ള ള ള ള ള ള ള
ജ ള ള ള ള ള ള ള ള ള ള
ന മ ള ള ള ള ള ള ള
0 (left Vowel Sign)
ഒ , ഓ , ഔ
1 (Right Vowel Sign)
ഓ , ഴ , ഴ , ഴ , ഴ , ഴ , ഴ , ഴ , ഴ , ഴ
Vowel Sign appearing both to left and right of the consonant /conjunct
ഒ ഴ , ഴ ഴ

Figure 17: Different classes of Y, X, 0 and 1

If a right vowel sign (9 in all) is identified, segmentation takes place with the segmented character having a consonant/conjunct followed by the identified right vowel sign (X 1). But if a consonant/ conjunct is identified, segmentation takes place at the previous level to return a consonant/ conjunct. This process continues until all the segments of the characters of a given word are returned in the

proper sequence. The different classes of X, 0 and 1 are shown in Figure 14.

VI. RESULTS AND CONCLUSION

Malayalam was found most appropriate to evaluate the proposed novel segmentation technique. Malayalam characters could consist of several uniformly spaced unconnected components such as a vowel (only at the beginning of a word) or a consonant/ conjunct along with vowel signs. The equal space between the characters of a word and the sub-characters of a character within a word makes segmentation even more complex. Thus conventional techniques like horizontal and vertical projection profile methods fail to segment the complete character segments correctly. A database is created with complete set of Malayalam characters including vowels, consonants conjunct characters, vowel signs appearing to the left and right. A training set is formed for all the characters in the database and the feature vectors of each of these characters are extracted using frequency capture method. Feature vectors of the 94 characters form the training set.

An efficient segmentation algorithm for a complicated situation is presented with Malayalam script as example. The flow chart in Figure 6 can be modified for similar situations where most commonly used for segmentation, the profiling techniques fail. The inherent advantages of the proposed segmentation algorithm are that

- 1) The characters get classified during the process of sequencing their sub-characters.
- 2) At any decision level the search is over successively reduced spaces

These inherent features lead to a robust OCR in terms of segmentation and classification algorithm. The algorithm was tested over all possible sub-character sequence and Malayalam was found most appropriate to evaluate the proposed novel segmentation technique.

These inherent features lead to a robust OCR in terms of segmentation and classification algorithm. The algorithm was tested over all possible sub-character sequence and character combination, each time resulting in correct segmentation and classification. The proposed algorithm was verified with 1250 words.

ACKNOWLEDGEMENT

This work was supported in part by research grants from UGC for Major Research Project in Science and Technology, F.No. 32-113/2006

REFERENCES

- [1] K.H.Aparna, Sumanth Jaganathan, P.Krishnan, V.S.Chakravarthy, "Document Image Analysis: with specific Application to Tamil Newsprint", Department of Electrical engineering, IIT, Madras.
- [2] A Cheung, M Bennamoun, N. W. Bergmann, An Arabic OCR system using recognition based segmentation, Pattern Recognition, Vol 34, pp 215- 233, 2001.
- [3] M. B. Sukhaswami, P. Seetharamulu, Arun K Pujari. Recognition of Telugu characters using neural networks. Int. J. of Neural Systems, 6(3):317-357, 1995.
- [4] Atul Negi, Chakravarthy Bhagvati, B.Krishna, An OCR system for Telugu, 6th ICDAR, Seattle, USA, Sept. 2001.
- [5] D. G. Elliman and I. T. Lancaster. A review of segmentation and contextual analysis for text recognition, Pattern Recognition, vol. 23, no. 3/4, pp. 337-346, 1990.
- [6] K. S. Sesh Kumar, Anoop M. Namboodiri, and C. V. Jawahar, Learning segmentation of documents with complex scripts. In ICVGIP, pp 749-760, 2006.
- [7] Kunte Sanjeev R, Sudhaker Samuel R D 2006 A two-stage character segmentation scheme for Printed Kannada text. J. Graphics, Vision and Image Processing 6: 1-8
- [8] Ashwin TV, Sastry P S, Afont and size-independent OCR system for printed Kannada documents using support vector machines. Sadhana 27: 35-58, 2002.
- [9] Y.He, A.Kundu: 2-D shape classification Using Hidden Markov Model, IEEE Trans. On PAMI, vol.13, 1991, pp.1172-1184.
- [10] Casey RG, Lecolinet E (1996) A survey of methods and strategies in character segmentation. IEEE T Pattern Anal 18:690 - 706
- [11] S. Harikumar, K. Jithesh, K. G. Sulochana, R. Ravindra Kumar, "Script based line & character segmentation techniques for Malayalam document images", In Proceedings of the International Symposium on Machine Translation (iSTRANS 2004) New Delhi, India, pp. 122-127, 2004.
- [12] U. Pal, S. Sinha, B. B. Chaudhuri, "Multi-Script Line identification from Indian Documents", In Proceedings of the ICDAR (ICDAR'03), pp. 880-884, 2003.
- [13] Janardhanan P. S. Issues in the development of OCR systems for Dravidian languages - proceedings of Akshara 94, BPB Publications, New Delhi, India 1994.
- [14] Malayalam standardization report May 2001.



Bindu Philip: Bindu Philip is a research scholar in the JSS research Center, S. J. College of Engineering, Mysore, Karnataka. She received her B.E degree in Electronics and Communication from Kuvempu University and M.Tech degree in Industrial Electronics from Visvesvaraya Technological University. Currently she is pursuing her doctoral studies in Electronics from the University of Mysore. She has contributed Thirty

papers in National and International conferences in the area of image processing and Four in International Journals. Her areas of interests include Image processing, Pattern Recognition, Computer vision and special systems for the differently enabled.



R.D Sudhaker Samuel: Dr Sudhaker Samuel was the Professor and Head of the Department of Electronics and Communication at Sri Jayachamarajendra College of Engineering, a Government-aided and autonomous institution at Mysore, India. He received Ph.D. from Indian Institute of Science, Bangalore and his M.Tech and B.E degrees from the University of Mysore, India. He has published over 150

research papers, three books and has produced Four PhDs, presently supervising five research scholars working towards their PhD. He has completed five sponsored research projects besides the four presently under execution including an international joint-research project. He has completed 22 industrial consultancy projects. His research interest includes Industrial Automation, Image Processing, Robotics, Embedded Systems and Biometrics.



C.R.Venugopal: Dr C.R.Venugopal is currently the Professor and Head of the Department of Electronics and Communication at Sri Jayachamarajendra College of Engineering, at Mysore, India. He received Ph.D. from IIT Bombay. He received his M.E. degree from Indian Institute of Science, Bangalore, India. He has published over 50

research papers, produced one PhD, presently supervising five research scholars working towards their PhD. He has six sponsored research projects presently under execution. His research interest includes Real Time Operating Systems and Algorithms, VLSI Architecture, and Content and Information Retrieval.