

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/221254210>

Converting Myanmar printed document image into machine understandable text format

Conference Paper · September 2011

DOI: 10.1109/ICDIM.2011.6093371 · Source: DBLP

CITATIONS

5

READS

3,687

3 authors, including:



[Khin Nweni Tun](#)

University of Computer Studies, Yangon

16 PUBLICATIONS 22 CITATIONS

SEE PROFILE

Converting Myanmar Printed Document Image into Machine Understandable Text Format

.....Htwe Pa Pa Win

.....University of Computer Studies

.....Yangon, Myanmar

.....hppwucsy@gmail.com

.....Mj lp'P y g'P K'Vwp

.....University of Computer Studies

.....Yangon, Myanmar

.....hppwp@gmail.com

Abstract — The large amount of Myanmar document images are getting archived by the Digital Libraries, an efficient strategy is needed to convert document image into machine understandable text format. The state of the art OCR systems can't do for Myanmar scripts as our language pose many challenges for document understanding. Therefore, this paper plans an OCR system for Myanmar Printed Document (OCRMPD) with several proposed methods that can automatically convert Myanmar printed text to machine understandable text. Firstly, the input image is enhanced by making some correction on noise variants. Then, the characters are segmented with a novel segmentation method. The features of the isolated characters are extracted with a hybrid feature extraction method to overcome the similarity problems of the Myanmar scripts. Finally, hierarchical mechanism is used for SVM classifier for recognition of the character image. The experiments are carried out on a variety of Myanmar printed documents and results show the efficiency of the proposed algorithms.

Keywords- Myanmar scripts, segmentation, feature extraction, OCRMPD, support vector machines.

I. INTRODUCTION

Optical Character Recognition is one of the most fascinating and challenging areas of pattern recognition with various practical applications. It can contribute immensely to the advancement of an automation process and can improve the interface between man and machine in many applications. And the World is witnessing a considerable transformation from print based-formats to electronic-based formats thanks to advanced computing technology, which has a profound impact on the dissemination of nearly all previous formats of publications into digital formats on computer networks. Then, one of the important tasks in machine learning becomes the electronic reading of documents. All various fields of the documents, magazines, reports and technical papers can be converted to electronic form using a high performance Optical Character Recognizer (OCR). And optical character recognition is a key enabling technology critical to creating indexed, digital library content, and it is especially valuable for scripts, for which there has been very little digital access [1], [2].

With the increasing demand for creating a paperless world, many OCR algorithms for English and other developed countries' languages have been developed over the years and

these can be available commercially or freely. But, development of an optical character recognition system for Myanmar languages is in little effort. This is because Myanmar (Burmese) scripts are rich in patterns while the combinations of such patterns makes the problem even more complex and hence the motivation to work further in this area. But properly utilized, OCR will help to make Burmese digital archives, practically accessible to local users and lay users alike by creating searchable indexes and machine-readable text repositories.

For an OCR system, segmentation phase is an important phase and accuracy of any OCR heavily depends upon segmentation phase. Due to the different nature of different scripts, the different segmentation mechanism is required for different languages [22]. Feature extraction is also a key step in the process of OCR, which in fact is a deciding factor of the accuracy of the system [23] and can efficient with the use of suitable approach depending on the nature of characters.

Nowadays, support vector machine is a useful technique for data classification and has been found to be successful when used for pattern classification problems. Fundamentally SVMs are binary classification algorithm with a strong theoretical foundation in statistical learning theory. The new pattern recognition SVM algorithms overcome some limitations with a strong underlying mathematical foundation [3].

In this paper, Optical Character Recognition System for Myanmar Printed Document (OCRMPD) is presented with a variety of proposing techniques, including a novel segmentation method to truly separate Myanmar characters, efficient Feature extraction method using zone and projection profile for isolated character data and the powerful SVM classifier to recognize Myanmar script features. These works are need to exert much effort to come up with better and workable OCR technologies for the local scripts in order to satisfy the need for digitized information processing.

The rest of the paper is organized as follow. Section 2 introduces the nature of Myanmar script. Section 3 presents the previous work as the background theories. Section 4 gives more details on the implementation of recognition system. Results are discussed in Section 5 and Section 6 is the conclusion.

II. NATURE OF MYANMAR SCRIPT

Myanmar (Burmese) script is recognized as Tibeto/Burman language group, developed from the Mon script and descended from the Brahmi script of ancient South India. It is the official language of Myanmar. The direction of writing is from left to right in horizontally. In Myanmar script, there is no distinction between Upper Case and Lower Case characters. The character set consists of 35 consonants (including ‘□’ and ‘□’), 8 vowels signs, 7 independent vowels, 5 combining marks, 6 symbols and punctuations, and 10 digits. Each word can be formed by combining consonants, vowels and various signs. It has its own specified composition rules for combining vowels, consonants and modifiers. There are total of above 1881 glyphs and has many similarity scripts in this language (e.g., □ and □, □ and □ and so on). When writing text, space is used after each phrase instead of each word or syllable. The shapes of Myanmar scripts are circular, consist of straight lines horizontally or vertically or slantways, and dots [11], [20].

III. RELATED WORK

Many researchers have proposed several ways to implement various OCR systems [4-10]. The authors of [24-26] are discussed for the feature extraction methods. But in [12-17], they stated that the SVM classifier can be used as the effective recognizer. Some of the existing techniques used in OCR for Myanmar scripts are presented here. A system of recognition for printed text in student application form and translated them to English words by using Hopfield Neural Network is proposed in [18], and they used all the pixel values as the features and has 97.56% accuracy rate. MICR based on statistical and semantic approaches for isolated handwritten character [19] is tested on 33 basic characters and 10 digits and gets 81 to 100% accuracy range. But they can only show the good accuracy for normal alphabets and not yet done for all the compound words. To the best of our knowledge, a comprehensive study on the success rate in terms of recognition accuracy for Myanmar printed text OCR system is yet to be reported.

IV. PROPOSED METHOD

As other traditional OCR systems, the proposed system also includes five processing steps as shown in Fig. 1. 6 different types of documents written in Zawgyi-One font and font size 12 are taken to test the system. These are scanned on a flatbed scanner at 300 dpi for digitization go for the preprocessing steps.

A. Preprocessing

Preprocessing step is the basic crucial part of the OCR system. The recognition accuracy of OCR systems greatly depends on the quality of the input text image. Firstly, we convert the raw input image into grayscale and then denoise it by removing noise using low pass Finite State Impulse Response (FIR) filter. Next, we binarize the clean image to a bi-level image by turning all pixels below some threshold to zero and all pixels about that threshold to one. We find this threshold value using Otsu method. Finally, we deskew the

binarized image with generalized Hough Transformed method. The detailed of the preprocessing steps are described in [21].

B. Segmentation

Segmentation is the process of the isolation of the individual character images from the refined image. It is considered as the main source of the recognition errors especially for small fonts. This is one of the most difficult pieces of the OCR system [4]. We use the X_Y cut method on the use of histogram or a projection profile technique for segmentation. It has been proven as a classical and more accurate method in Devnagari scripts such as Bangla and Hindi and some of the South East Asia scripts, English and some Greek OCR [7], [10]. The process of segmentation in our system mainly follows the following pattern:

- Line Detection and slicing
- Character Segmentation

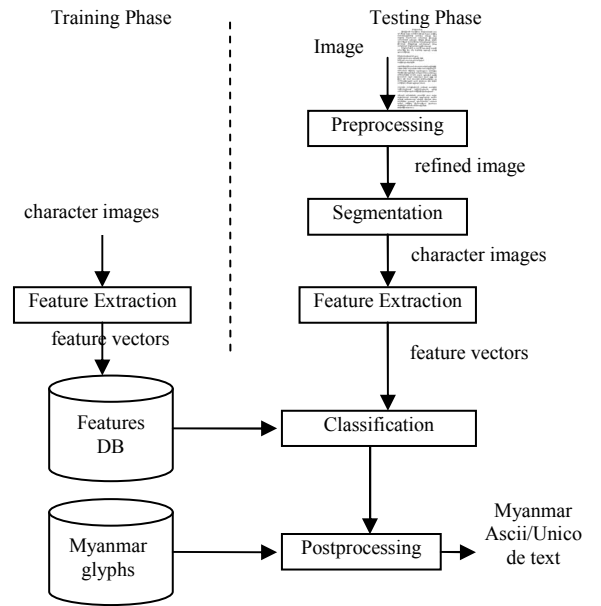


Figure 1. System Design of the Myanmar OCR system

1) *Line Detection and slicing* : To detect the lines, assume that the value of the element in the x th row and the y th column of the character matrix is given by a function f :

$$f(x, y) = a_{xy} \quad (1)$$

where, a_{xy} takes binary values (i.e., 0 for background white pixels and 1 for black pixels). The horizontal histogram H_h of the character matrix is calculated by the sum of black pixels in each row:

$$H_h(x) = \sum_y f(x, y) \quad (2)$$

And cut the lines depend on the $H_h(x)$ values. as shown in Fig. 2.

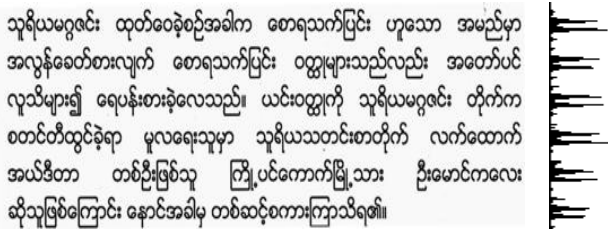


Figure 2. Example of line segmentation using projection

2) Character Segmentation

Similarly, the vertical histogram H_v of the character matrix is calculated by the sum of black pixels in each column of the line segment:

$$H_v(y) = \sum_x f(x, y) \quad (3)$$

Characters are segmented using these histogram values. However, this method alone is not enough for the Myanmar scripts. As for the small font, some character is not correctly segmented as shown in Fig. 3.



Figure 3. Example of wrong segmentation error with projection

And it may also be problem for some connected components. Moreover, the connected components can't extract earlier as other languages because it can appear not only in shorter segments but also in longer segments that of the line height. That's why the nature of Myanmar scripts cause over segmentation and under segmentation problems. To overcome overlaps and wrong segmentation cases, assume the points from (3) as the pre segment points and we need to add the following procedures to check the possible points according to line height:

Begin
 $CCs \leftarrow$ possible column points of connected components
 $mixcharwidth \leftarrow$ the minimum width of the character
 $densitythreshold \leftarrow$ the minimum density value for each column
 $bottomthreshold \leftarrow$ the threshold distance of the nearest pixel from the bottom

For each pre segmented point results from (3)

Begin
 Calculate *density* of the pixels vertically
 Calculate *bottomprojection* of each column
 If $density < densitythreshold$
 Begin
 Store the column point in $columnpoints[]$
 For each *column* in $columnpoints[]$
 Being

$remaininlength \leftarrow$ width of pre segment point - *column*
 If $column \in CCs$
 Begin
 If ($bottomprojection < bottomthreshold \ \&\& \ remaininlength > mixcharwidth$)
 Begin
 Denote final segment points
 End
 End
 End
 Else
 Denote pre segment points as the possible points.

End
 End
 End

C. Feature Extraction

Before the extraction of features we need to normalize the binary character images to have the standard width and height. We normalize all character images height into N and the equal amount is used for width with respecting the original aspect ratio.

Feature extraction involves extracting the attributes that best describe the segmented character image as a feature vectors. This process maximizes the recognition rate with the least amount of elements [5]. In our approach we employ two types of statistical features. The first one divides the character image into a set of zones and calculates the density of the character pixels in each zone as in [15]. The Myanmar characters are written into three main zones for horizontal and the minimum component for a truly segmented glyph is one and the maximum component may be four as shown in Fig 4. Therefore, we considered for the second type of features, the area that is formed from the projections of the top, middle and bottom as well as of the left, center and right character profiles is calculated.

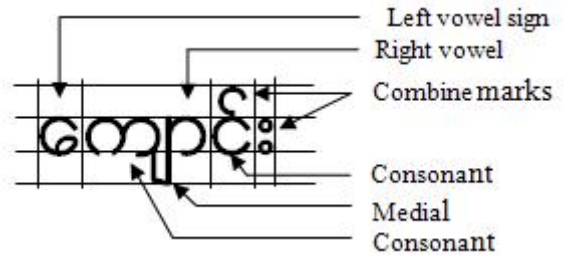


Figure 4. Sample of Myanmar Glyphs

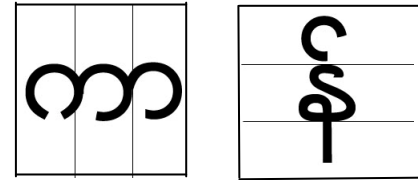


Figure 5. Division of each character depend on writing nature

Let $g(x, y)$ be the binary image array and w, h be the width and height of the segmented character. In the case of

features based on zones, the image is divided into equal zones. For each zones, we calculate the density of the character pixel as follow:

$$F_z(n) = \sum g(x, y), n = 0, \dots, Z_{\max} - 1 \quad (4)$$

Where, x, y be the pixel point in each zone.

When we consider features based on vertical profile projections, the character image is divided into S_v sections separated by the horizontal lines of y and calculated as follow:

$$y_i = i(h / S_v) - 1, i = 1, \dots, S_v - 1 \quad (5)$$

And for each section, we equally divide into blocks and calculate y_i , the distance between the base line and outermost pixel depending on the direction we considered as

$$y_s = \begin{cases} y_i - y_p, & \text{for bottom to top} \\ y_p - y_{i-1}, & \text{for top to bottom} \end{cases} \quad (6)$$

Where, y_p is the outermost pixel value of 1 and F_v be the total number of blocks to produce the vertical profiles and calculate the feature for each block as follow:

$$F_v(n) = \sum y_s(x), n = Z_{\max}, \dots, Z_{\max} + F_v - 1 \quad (7)$$

For the horizontal profile projections, the image is split into S_h sections separated by the vertical lines of x and calculated as follow:

$$x_i = i(w / S_h) - 1, i = 1, \dots, S_h - 1 \quad (8)$$

And for each section, we equally divide into blocks and calculate x_s , the distance between the base line and outermost pixel depending on the direction we considered as follow:

$$x_s = \begin{cases} x_i - x_p, & \text{for right to left} \\ x_p - x_{i-1}, & \text{for left to right} \end{cases} \quad (9)$$

Where, x_s is the outermost pixel value of 1 and F_h be the total number of blocks to produce the horizontal profiles and calculate the feature for each block as follow:

$$F_h(n) = \sum x_s(y), n = Z_{\max} + F_v, \dots, Z_{\max} + F_v + F_h - 1 \quad (10)$$

Therefore, the total feature for each character image is:

$$F_{\text{total}}(n) = F_z(n) + F_v(n) + F_h(n) \quad (11)$$

D. Classification

This process is responsible to match the test features of input images with the train features. SVM [27] is used as the recognizer for this OCR System.

The original form of SVM is the separating of hyperplane between two different classes. Because of the existence of a number of characters in any script, optical character recognition problem is inherently multi-class in nature. Every character in a language forms a class. The field of binary classification is mature, and provides a variety of approaches to solve the problem of multi-class classification [3], [12], [14].

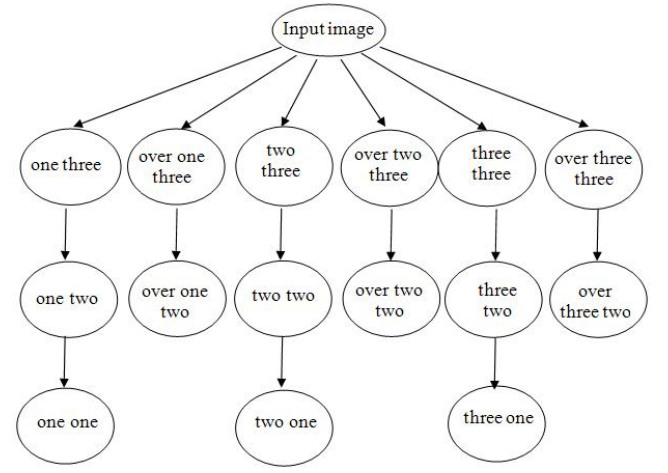


Figure 6. Hierarchical mechanism for Myanmar characters

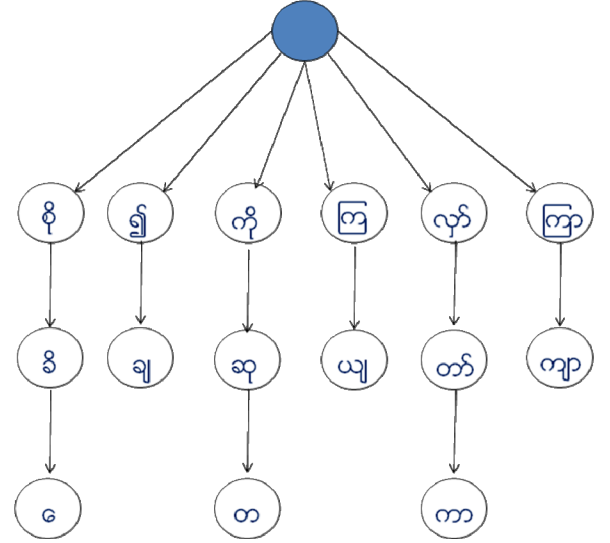


Figure 7. Example of Hierarchical mechanism for Myanmar characters

The Hierarchical mechanism is used for Multi-class SVM classification to reduce search space as there are a large number of characters in Myanmar scripts and there is the similarity between them. Firstly, the similar characters are clustered based on the nature of the writing style of the characters. As a result of this, all characters of 1881 classes can be reduced into 15 classes. And then perform the classification to extract the right class. The hierarchical group of characters is shown in Figure 6. The first level is for the characters for one column, over one column, two columns, over two columns, three columns and over three columns width written in three zones. The second level is for the characters written in two zones with the above column widths and the last level is for the characters for one column, two columns and three columns width that are written in one zone. The example mechanism is shown in Figure 7.

E. Postprocessing

This process is to produce the relevant text from the recognition results. This stage is also called the converting process because it converts the recognized character image or classified character image into related ASCII or Unicode text. The final result of this system, the output text can be modified and saved into any format.

V. EXPERIMENTAL RESULTS

The implementation is based on Java Environment using open source tool Eclipse and MySql Database. For experiment, 6 Myanmar Printed Documents that are written in Zawgyi-One font with size 12 and scanned on a flatbed scanner at 300 dpi are taken to test the system. These documents have some noise variations. The experiments are carried out for comparing segmentation accuracy, the effects of feature extraction on the accuracy and recognition accuracy. Table 1 shows the segmentation results of the proposed mechanism. Figure 8 compare the effectiveness of hybrid feature extraction method on accuracy rate and Figure 9 reveal the recognition rate of the proposed OCR system.

Table 1. Segmentation Accuracy for Printed Document

Document	Contained Characters	Truly Segmented Characters		Accuracy (%)	
		Projection only	OCRMPD	Projection only	OCRMPD
1	89	87	89	97.75	100
2	95	91	92	95.79	96.84
3	193	184	192	95.34	99.48
4	303	285	301	94.06	99.34
5	364	342	359	93.96	98.63
6	1048	1006	1038	95.99	99.05
Average				95.48	98.89

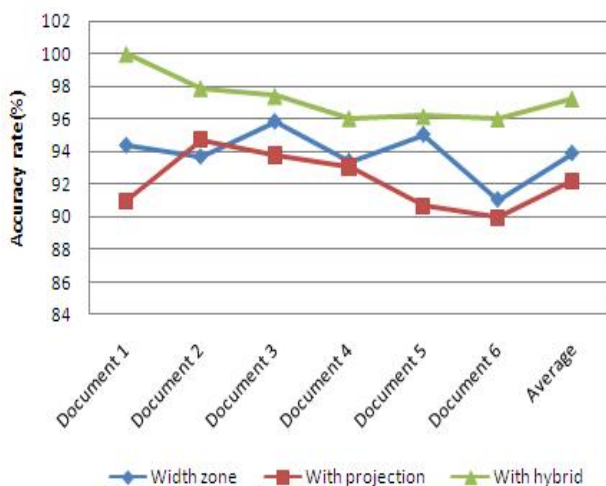


Figure 8. Accuracy Results with various Feature Extraction Methods

The accuracy of the OCR system is directly proportional with the accuracy of segmentation. The higher the accuracy rate of character segmentation can be got, the better the accuracy rate of the OCR system can be obtained.

The character image is normalized into 30x30 and 25 features are used for zoning method and 60 features are for projection profile method.

Figure 10 shows the average execution time per each document. These times are computed on a PC with 2.4GHz CPU and 2GB of RAM.

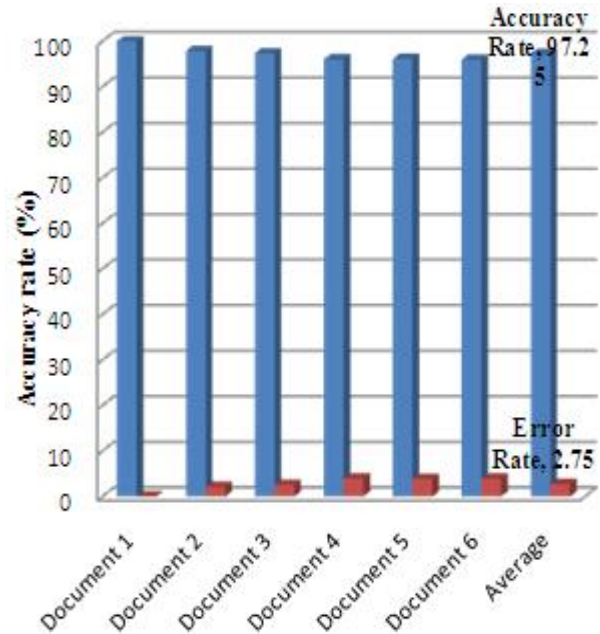


Figure 9. Recognition Accuracy for Myanmar Printed Documents of OCRMPD

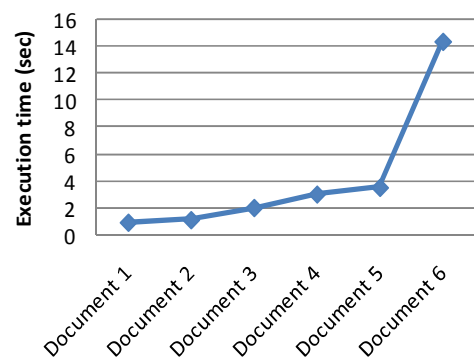


Figure 10. Execution time for each document

VI. CONCLUSION AND FUTURE WORK

This paper proposes a novel segmentation method to truly separate characters, an efficient feature extraction method and hierarchical classification mechanism for Myanmar Printed

document recognition system, OCRMPD, and shows the good result for the system. This result proved the advantages of the innovations. The segmentation scheme can be used for all Myanmar printed documents without user intervention. The combination of feature extraction methods can produce good results but it takes a more time than the normal zoning method. The hierarchical classification scheme can improve accuracy and save the processing time of classifier. The advancement of the system to recognize bilingual documents and historic documents are future works for the Digital Library Requirement.

REFERENCES

- [1] V. Govindaraju and S. Setlur, "Guide to OCR for Indic Scripts: Document Recognition and Retrieval", 2009
- [2] "General guidelines for designing bilingual low cost digital library services suitable for special library users in developing countries and the Arabic speaking world", World Library and Information Congress: 75th IFLA General Conference and Council, 23-27 August 2009, Milan, Italy.
- [3] K. Shivsubramani, R. Loganathan, C. J. Srinivasan, V. Ajay and K. P. Soman, "Multiclass Hierarchical SVM for Recognition of Printed Tamil Characters", Centre for Excellence in Computational Engineering, Amrita Vishwa Vidyapeetham, Tamilnadu, India, 2007.
- [4] N. S. Sarhan and L. Al-Zobaidy, "Recognition of Printed Assyrian Character Based on Neocognitron Artificial Neural Network", The International Arab Journal of Information Technology, Vol 4, No.1, January 2007.
- [5] R. Singh and M. Kaur, "OCR for Telugu Script Using Back-Propagation Based Classifier", International Journal of Information Technology and Knowledge Management, July-December 2010, Vol. 2, No. 2, pp. 639-643.
- [6] R. Singh, C. S. Yadav, P. Verma and V. Yadav, "Optical Character Recognition (OCR) for Printed Devnagari Script Using Artificial Neural Network", International Journal of Computer Science & Communication Vol.1, No. 1, January-June 2010, pp. 91-95.
- [7] D. Achaya U, N. V. S. Reddy and Krishnamoorthi, "Hierarchical Recognition System for Machine Printed Kannada Characters", IJCSNS International Journal of Computer Science and Network Security, Vol. 8 No.11, November 2008.
- [8] H. Guo and J. Zhao, "A Chinese Minority Script Recognition Method Based on Wavelet Feature and Modified KNN", Journal of Software, Vol. 5, No. 2, February 2010.
- [9] H. A. Al-Muhtaseb, S. A. Mahmoud and R. S. Qahwaji, "Recognition of Off-line Printed Arabic Text Using Hidden Markov Models", Information and Computer Science Department, King Fahd University of Petroleum & Minerals, Dhahran 31261, Saudi Arabia and Electronic Imaging and media communications department, University of Bradford, Bradford, UK, 2008.
- [10] B. Chaulagain, B. B. Rai and S. K. Raya, "Final Report on Nepali Optical Character Recognition, NepaliOCR", July 29, 2009.
- [11] "Myanmar Orthography". Department of the Myanmar Language Commission, Ministry of Education, Union of Myanmar, June, 2003.
- [12] J. Dong, A. Krzyzak and C. Y. Suen, "An improved handwritten Chinese character recognition system using support vector machine", Pattern Recognition Letters, Vol. 26, 2005, pg- 1849-1856.
- [13] S. Rawat et al., "A Semi-automatic Adaptive OCR for Digital Libraries", Centre for Visual Information Technology, International Institute of Information Technology, Hyderabad - 500032, India, 2006.
- [14] M. Meshesha and C. V. Jawahar, "Optical Character Recognition of Amharic Documents", Center for Visual Information Technology, International Institute of Information Technology, Hyderabad - 500 032, India, 2007.
- [15] G. Vamvakas, B. Gatos, N. Stamatopoulos, and S. J. Perantonis, "A Complete Optical Character Recognition Methodology for Historical Documents", The Eighth IAPR Workshop on Document Analysis Systems, 2008.
- [16] B. Philip and R. D. Sudhaker Samuel, "Preferred Computational Approaches for the Recognition of different Classes of Printed Malayalam Characters using Hierarchical SVM Classifiers", International Journal of Computer Applications (0975 - 8887) Vol. 1, No. 16, 2010.
- [17] G. G. Rajput, R. Horakeri and S. Chandrakant, "Printed and Handwritten Mixed Kannada Numerals Recognition Using SVM", (IJCSSE) International Journal on Computer Science and Engineering, Vol. 02, No. 05, 2010, pg-1622-1626.
- [18] T. Swe and P. Tin, "Recognition and Translation of the Myanmar Printed Text Based on Hopfield Neural Network", Asia-Pacific Symposium on Information and Telecommunication Technologies (APSITT), pp 99-104, Myanmar, November 9-10, 2005.
- [19] Y. Thein and M. M. Sein, "Myanmar Intelligent Character Recognition for Handwritten", University of Computer Studies, Yangon, Myanmar, 2006.
- [20] S. Hussain, N. Durrani and S. Gul, "Survey of Language Computing in Asia 2005", Center for Research in Urdu Language Processing, National University of Computer and Emerging Sciences, 2005.
- [21] H. P. P. Win and K. N. N. Tun, "Image Enhancement Processes for Myanmar Printed Documents", the fifth Conference on Parallel & Soft Computing, University of Computer Studies, Yangon, Myanmar, December 16, 2010.
- [22] M. Agrawal and D. Doermann, "Re-targetable OCR with Intelligent Character Segmentation", The Eight IAPR Workshop on Document Analysis Systems, 2008.
- [23] R. Ramanathan et. al., "Robust Feature Extraction Technique for Optical Character Recognition", International Conference on Advances in Computing, Control, and Telecommunication Technologies, 2009.
- [24] S. V. Rajashekaradhy and Dr. P. V. Ranjan, "Efficient Zone Based Feature Extraction Algorithm for Handwritten Numeral Recognition of Four Popular South Indian Scripts", Journal of Theoretical and Applied Information Technology, 2008.
- [25] G. Vamvakas, B. Gatos and S. J. Perantonis, "A Novel Feature Extraction and Classification Methodology for the Recognition of Historical Documents", 10th International Conference on Document Analysis and Recognition, 2009.
- [26] Ngodrup et al., "Study on Printed Tibetan Character Recognition", International Conference on Artificial Intelligence and Computational Intelligence, 2010.
- [27] C. W. Hsu, C. C. Chang, and C. J. Lin, "A Practical Guide to Support Vector Classification", April 15, 2010.