

## An Efficient Segmentation and Binarization Method for Recognition of Degraded Devanagari script Scan document Images

Sushilkumar N.Holambe, Prof. Ulhas B.Shinde

*College of Engineering, Osmanabad. 413501*

*CSMSS,CSCOE, Aurangabad-431001, India*

*Corresponding authors: Sushilkumar N.Holambe*

**ABSTRACT:** Image segmentation is very helpful in degraded devanagari character recognition process. The Image segmentation is process of clustering the pixels depending on some property of the image, e.g., intensity gray levels, color, texture, depth, edge continuity. After segmentation, the whole image is partitioned into smaller regions, i.e., regions corresponding to individual surfaces, objects, or natural parts of objects. Segmentation can be used for object recognition, estimation of occlusion boundary within motion or stereosystems, image compression, image editing, and/or image database look-up systems, image compression, image editing, and/or image database look-up.

This paper, proposes technique to address the issues of degraded devanagari script scan images using adaptive image contrast. The degraded devanagari script scan adaptive image contrast technique is a combination of the local image contrast and the local image gradient. And they are tolerant to variation of text and background. Such variations are caused by number of document degradations. The proposed technique, constructs adaptive contrast map for degraded devanagari script scan image. The contrast map is combined with Canny's edge map, for the identification of text stroke edge pixels. Thresholding technique can be applied as global technique and local technique. Global thresholding is suitable for a degraded devanagari script scan document where there is uniform contrast delivery of background and foreground. However global thresholding fails to the applications where difference in contrast, Extensive background noise and difference in brightness exists. In such circumstances categorization of many pixels as a foreground or as a background is not so easy. Local thresholding plays significant role in such cases. Local thresholding technique uses local threshold  $t$ ; w.r.t. local window to segment the document image. This local threshold  $t$  is estimated based on the intensities of detected text stroke edge pixels. The proposed method is simple, robust, and involves minimum parameter tuning. It has been tested on three public datasets that are used in the recent document image binarization contest (DIBCO) 2010 & 2012 and handwritten-DIBCO 2011.

### General Terms

Segmentation of Color Thresholding, Degraded Devanagari Script

**Keywords:** Image segmentation • Binarization • Thresholding Document image binarization, Image thresholding • Adaptive local binarization

Date of Submission: 12-11-2018

Date of acceptance: 26-11-2018

## I. INTRODUCTION

The information contained in historical documents must be preserved therefore efforts are taken at National and international levels. This can be done by using a technique known as binarization. Binary image is such digital image that has just two values for every pixel. Two colors are used to represent these two values, i.e. black and white however any other colors can also be used. Therefore ground color is mainly representing the object whereas a simple way to binarize a degraded devanagari script scan image is through thresholding and separating the light and dark regions (background and foreground) according to the pixel intensities. In many image processing and pattern recognition applications, gray levels of pixels belonging to the object are substantially different from the gray levels of pixels belonging to the background. In such context, thresholding becomes a simple but effective tool for separating objects from the background. Thresholding creates binary images from a gray-level image by setting all the pixels below some threshold to zero and all pixels equal or above that threshold to one.

A simple way to binarize the degraded devanagari script scan image is through thresholding and separating the light and dark regions (background and foreground) according to the pixel intensities.

In many degraded devanagari script scan image processing and pattern recognition applications, gray levels of pixels belonging to the object are substantially different from the gray levels of pixels belonging to the

background. In such context, thresholding becomes a simple but effective tool for separating objects from the background. Thresholding creates binary images from a gray-level image by setting all the pixels below some threshold to zero and all pixels equal or above that threshold to one.

In pre-processing stage, we applied the Sauvola binarization technique [14] to convert the input document into a binary image, which is robust to noise and blur. After binarization, complement the image to make the document text content as foreground and non-text content as background. Remove a few components from the binary image, whose heights and widths are less than the user defined value (UDV1) (5 pixels), which removes few noisy/non-text components from the document. Discard the non-text components such as figures, tables, and borders with the help of geometrical filters by comparing the component size to the document size. We have discarded those components, whose ratio of width to the number of columns of the image is more than UDV2 (20%) and similarly discarded those components, whose ratio of height to the number of rows of the image is more than UDV3 (10%). Such type of filters will remove the non-text components. All the UDVs stated in this paper are determined very carefully after so many experimentations on various kinds of Indian multilingual office document images to remove only noisy noncharacter components. Most of the tiny non-character components present in a scanned document will be smaller than 5 pixels in size, bigger non-character components occupies more than 20% of the width, and more than 10% of the height of the scanned document. Background color rest of the image. Separation foreground and background of documents images is the pre-processing step for the document analysis, carried out by Binarization. Gray-scale document image is converted into a binary document image i.e. with two values only. In document image processing like applications, binarization technique which is fast and accurate is very important to ensure correctly processing tasks of document images ,such as optical character recognition (OCR).calculation of thresholding for degraded document images is an unsolved problem due to, high intra/inter variation between the document background and text stroke across different document images. As illustrated in Figure 1(a), Historical documents are degraded due to bleed through. In addition, handwritten text documents are degraded due to a certain amount of variation in terms of the stroke width, stroke brightness, stroke connection, and document background as illustrated in Figure 1(b). In addition, historical documents are often degraded by different types of imaging artifacts as illustrated in Figure1(c). Thresholding errors are introduced due to such different types of degradations in document images. The recent Document Image Binarization Contest (DIBCO)[1] [2], [3] held under the framework of the International Conference on Document Analysis and Recognition (ICDAR) 2010 & 2012 and the Handwritten Document Image Binarization Contest (H-DIBCO) [4] held under the framework of the International Conference on Frontiers in Handwritten Recognition show recent efforts on this issue. We participated in the DIBCO 2010 and our background estimation method [5] performs the best among entries of 43 algorithms submitted from 35 international research groups. We also participated in the H-DIBCO 2011 and our local maximum-minimum method [6] was one of the top two winners among 17 submitted algorithms. In the latest DIBCO 2012, our proposed method achieved second best results among 18 submitted algorithms. This paper presents a document binarization technique that extends our previous local maximum-minimum method [6] and the method used in the latest DIBCO 2012. The proposed method is simple, robust and capable of handling different types of degraded document images with minimum parameter tuning.



**Fig a & b**



**(c)**

**Fig 1: Degraded document image examples (a), (b), (c).**

**(c) Is taken from Bickley diary dataset**

In particular, the proposed technique addresses the over-normalization problem of the local maximum minimum algorithm [6]. At the same time, the parameters used in the algorithm can be adaptively estimated.

## **II. EXISTING SYSTEM**

The Adaptive degraded devanagari script scan document image binarization uses early window-based adaptive thresholding technique. In this window-based technique estimation of the local threshold is done with the help of mean and the standard variation of image pixels within a local neighborhood window. The local contrast method proposed in Bernsen's "Dynamic thresholding of gray-level images," is simple and depends upon the maximum and minimum intensities within a local neighborhood windows of an image pixel (i, j) respectively

## **2.1 Degraded Devanagari Script Scan Document Efficient Thresholding for Binarization Algorithm Comparison**

These techniques are compared with some existing thresholding algorithms. Document images with background noise or illumination/contrast variations are used for the evaluation of algorithms. A threshold value for binarization can be calculated by following a simple method [4] as follows:

1. Select an initial estimate for T. (A suggested initial estimate is the average of the minimum and maximum intensity values in the image).
2. Segment the image using T. This will produce two groups of pixels: G1, consisting of pixels with intensity values  $\geq T$ , and G2, consisting of pixels with values  $< T$ .
3. Compute the average intensity values  $\mu_1$  and  $\mu_2$  for the pixels in regions G1 and G2.
4. Compute a new threshold value:

$$T = \frac{1}{2} * (\mu_1 + \mu_2)$$

5. Repeat steps 2 through 4 until the difference in T in successive iterations is smaller than a predefined parameter T0. The thresholding quality was assessed from resultant words in the background using Precision and Recall analysis of the same. Though all types of images are not handled by any single algorithm, each algorithm definitely works better than others for particular types of images. Appropriate algorithm(s) are combined to do task better.

## **2.2 Degraded Devanagari Script Scan Document Segmentation of Color Thresholding**

The grey level thresholding algorithm with slight modifications is used for color thresholding. Multilevel thresholding has been conducted to the RGB color object. To study color information no of natural images are used. The color thresholding technique is being carried out based on the adaptation and slight modification of the grey level thresholding algorithm. Multilevel thresholding has been conducted to the RGB color information of the object extract it from the background and other objects. Different natural images have been used in the study of color information. The results showed that by using the selected threshold values, the image segmentation technique has been able to separate the object from the background. Image segmentation algorithms generally are based on one of the two basic properties of intensity values: discontinuity and similarity. Thresholding is a method of similarity category. It partitions an image into regions that are similar according to a set of predefined criteria. There are various thresholding techniques and it is also a fundamental approach to segmentation that enjoys a significant degree of popularity, especially in applications where speed is an important factor [17]. The result is supposed to be achieved if, by using the selected threshold, object is separated from the background. For document image binarization many thresholding techniques [7]–[10] have been reported. Global thresholding [12]–[14] is usually not a suitable approach as many degraded documents do not have a clear bimodal pattern. So the local thresholding estimator for each document image pixel is suitable approach. And this can be achieved through Adaptive thresholding [14], to deal with different variations within degraded document images.

## **III. DEGRADED DEVANAGARI SCRIPT SCAN DOCUMENT SEGMENTATION SUB-BLOCK CLASSIFICATION AND THRESHOLDING**

Degraded devanagari script scan document segmentation of sub-block classification and thresholding [15] consist of the below listed three feature vectors to test the local regions. These local regions are classified into three types:

- a) Heavy strokes
- b) Faint strokes
- c) Background (No stroke).

Content information is not included in the background. Lower values of edge strength and variance are covered in background. Small mean-gradient value is associated with a background though it is considered as noise-free. Strokes that are difficult to distinguish from the background are nothing but faint strokes. Heavy stroke areas have strong edge strength, more variance and larger mean-gradient value. The proposed weighted gradient thresholding method is applied to the different classes of degraded devanagari script scan sub block.

### **3.1 Degraded devanagari image**

#### **3.1.1 Enhancement**

To proceed further with degraded document images enhancement of faint strokes is necessary. Wiener filter was applied to avoid the noise enhancement. The stroke faint enhancement can be done as:

Step 1: Enhancing the degraded script image can be done by finding the Maximum and minimum grey value in the 3x3

Window.

Step 2: Mini =min (Total No. Of window elements)

Maxi = max (Total No. Of window elements)

Compare “pixel – mini” and “maxi – pixel”, where “pixel” is the pixel-value. If the former is greater, the “pixel” is closer to the highest grey value than the lowest value in this window; hence the value of “pixel” is set to the highest grey value (“pixel” = “maxi”). If the former is smaller, then the value of “pixel” is set to the lowest grey value (“pixel”= “mini”).

### 3.1.2 Thresholding

To Threshold faint stroke, weighted method is used. And this method is based on mean gradient. Various directional strokes are normally available with old devanagari scripts documents.

## IV. PROPOSED METHOD

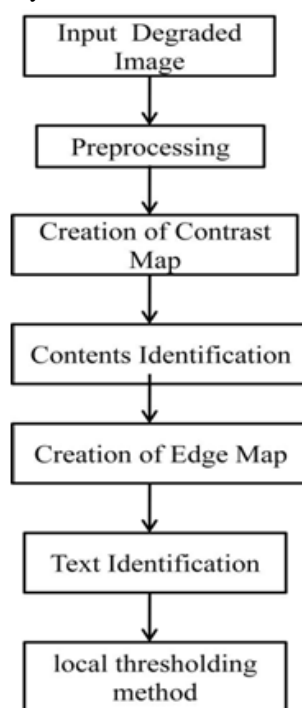
The proposed degraded devanagari script scan document image binarization technique works in no.of steps and described in this section.

- a. Degraded devanagari script scan contrast image construction.
- b. Degraded devanagari script scan text stroke edge pixel detection.

Local Threshold Estimation.

Post Processing.

The proposed method can be implemented as shown in fig.2 Contrast image construction is done in the preprocessing stage.canny’s edge detection method is used for the edge detection. Separation of text from the image is done through local thresholding method.post processing shows its importance in order to improve degraded devanagari script scan image quality.



**Fig 2: Degraded devanagari script scan document Identification flow**

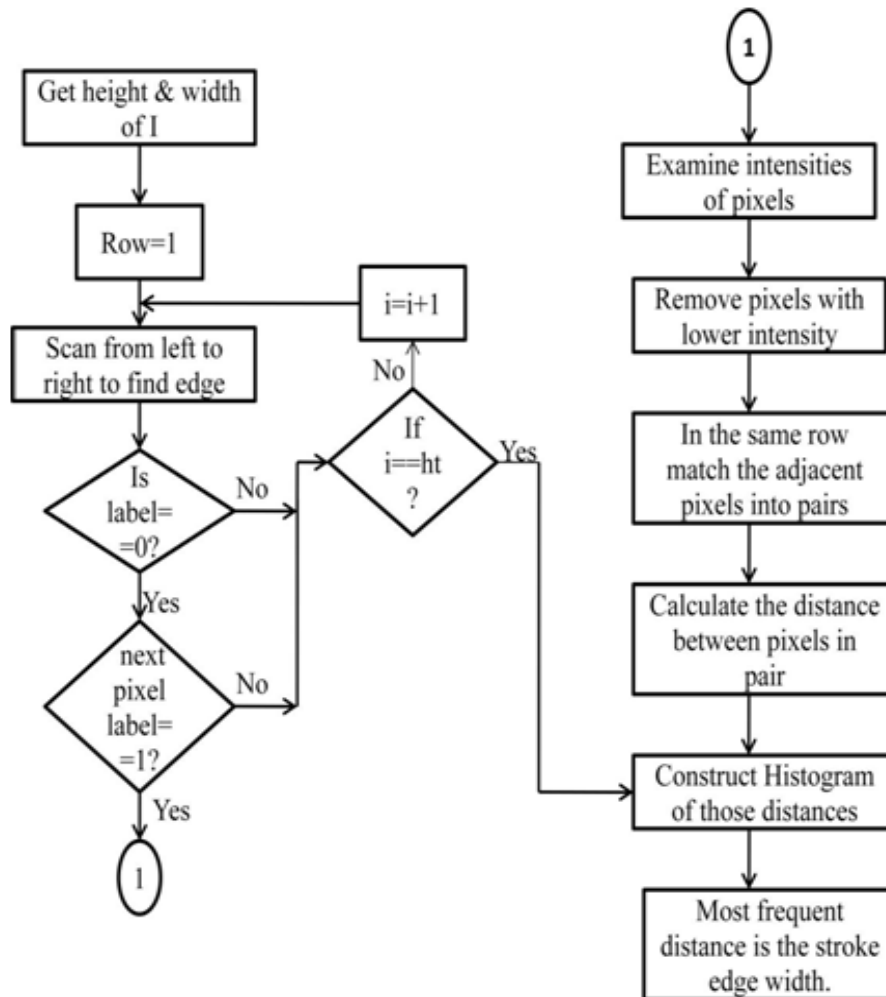
### 4.1 Degraded devanagari scrip Contrast Image Construction

Main aim of binarization is segmenting the degraded devanagri scan script document text from the document background. And this can be accomplished with the image features: local image contrast and local image gradient. Because the document text usually has certain image contrast to the neighboring document background. They are used in many document image binarization techniques [2] [3] and are very effective. When degraded devanagari script scan document image has noticeable intensity variations, need to work with

the image contrast with high weigh (i.e. Large  $\alpha$ ). To overcome over-normalization problem [1] we derive an adaptive local image contrast as:

$$. Ca(i,j)= \alpha C(I,j)+(1- \alpha)(I_{max}(I,j)-I_{min}(I,j)).(1)$$

The adaptive combination of the local image contrast and the image gradient in above equation can produce proper contrast maps for document images with different types of degradations.



**Fig 3** Degraded Devanagari Script Scan Document Thresholding Segmentation

#### 4.2 Degraded Devanagari Script Scan Document Text stroke edge pixel detection

The contrast image construction detects the stroke edge pixels of the degraded devanagari document text. At later stage edges are detected through canny edge detection algorithm. In this algorithm it smoothes the noise in the image then pixel at both sides of the text stroke will be selected as the high contrast pixel.

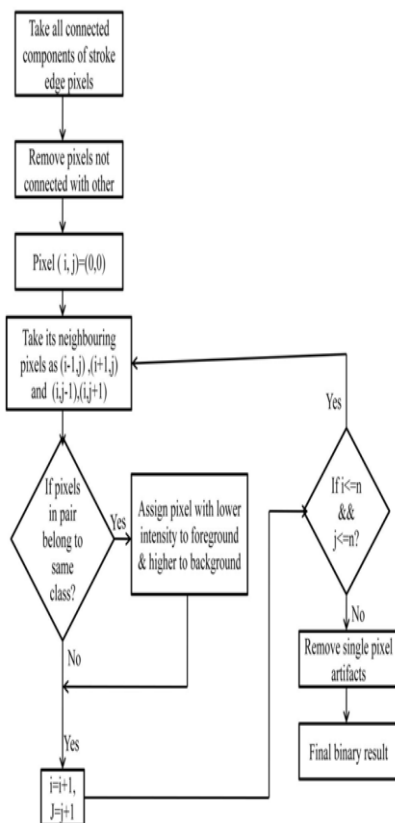
#### 4.3 Degraded Devanagari Script Scan Document Local Thresholding Segmentation

Once the text stroke edges are detected, we calculate the most frequent distance between two adjacent pixels that are on the edge. We perform it in the horizontal direction and use it as the estimated stroke width as shown in fig.3. degraded devanagari script scan document thresholding segmentation

#### 4.4 Degraded Devanagari Script Scan Document Post Processing Procedure

Binarization result can be improved by using Post- Processing method. By using the algorithm of post processing we remove single pixel artifacts along the text stroke boundaries after the degraded devanagari script scan document thresholding as shown in fig.4.





**Table 1. Evaluation Results Of The Dataset Of DIBCO 2014 For degraded devanagari scan script document**

Methods	F-Measure (%)	PSNR	DRD	MPM
OTSU	84.32	16.77	8.70	16.74
SAUV	84.64	17.78	8.10	9.35
NIBL	78.62	13.76	29.45	27.39
BERN	57.28	8.92	83.29	137.55
GATO	84.11	17.04	6.42	8.12
LMM	87.76	17.75	7.02	7.42
LELO	84.86	17.18	105.48	65.44
SNUS	86.2	18.17	16.67	10.08
HOWE	89.84	18.85	6.38	9.74
BOLAN	88.8	18.58	5.84	6.27
<b>Proposed</b>	<b>82.34</b>	<b>14.19</b>	<b>8.17</b>	<b>8.40</b>

**Fig 4: Degraded devanagari script scan document image post processing**

### V. CONCLUSION

The proposed system is tested throughly using variety of printed aged Indian Government office documents and we recorded 99% and 94% success rates for line extraction and degraded devanagari word segmentation, respectively. This system can be applicable in all the areas that are concerned with preserving the historical documents and managing the degraded documents. This also opens up the possibility of extending the work for a large variety of applications including city surveillance and medical image analysis where binarization is very important as a preprocessing step for subsequent identification of object recognition. The proposed method produces better average results for ten different performance evaluation metrics as compared to other widely used binarization methods as shown in table 1. The proposed method shows good noise immunity in the presence of salt-and-pepper and Gaussian noises. Moreover, the method is suitable for binarizing graphic and degrade devanagari script scan document images.

### REFERENCES

- [1]. George Nagy, "Twenty years of document image analysis in PAMI," IEEE Transactions on Pattern Analysis and Machine Intelligence 22.1, pp. 38-62, 2000.
- [2]. David L., et al. Milgram, Algorithms and hardware technology for image recognition.: MARYLAND UNIV COLLEGE PARK COMPUTER SCIENCE CENTER, 1978.
- [3]. Ioannis Pratikakis, and Stavros J. Perantonis. Gatos Basilios, "Improved document image binarization byusing a combination of multiple binarization techniques and adapted edge information," , 2008.
- [4]. Fatos T. Yarman-Vural Arica Nafiz, "An overview of character recognition focused on off-line handwriting," in Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on 31.2, , 2001, pp. 216-233.
- [5]. N. Otsu, "A threshold selection method from gray-level histograms," in IEEE Trans. Systems, Man, and Cybernetics, 1979, pp. 62-66.
- [6]. P.K. Sahoo, and A.K.C. Wong Kapur J., "A new method for gray-level picture. Thresholding using the Entropy of the Histogram," in Computer Vision Graphics and Image Processing vol. 29, 1985, pp. 273-285.
- [7]. Y. and C.G. Leedham Solihin, "Integral Ratio: A New Class of Global Thresholding Techniques for Handwriting images," in IEEE Trans. on PAMI, vol. 21, 1999, pp. 761-768.
- [8]. Yibing Yang and Hong Yan, "An adaptive logical method for binarization of degraded document images Pattern Recognition," in Pattern Recognition Society, Elsevier Science , vol. 33, no. 5, 2000, pp. 787-807.
- [9]. And Matti Pietikäinen Sauvola Jaakko, "Adaptive Document Image Binarization," in Pattern Recognition 33.2, 2000, pp. 225-236.

- [10]. M Randolph T. Smith, "Enhancement of fax documents using a binary angular representation," in in Proceedings of, Int Symp on Intelligent Multimedia, Video and Speech Processing, Hong Kong,Cjina, 2001.
- [11]. Wu Sue Adnan Amin, "Automatic Thresholding of Gray-level Using Multi-stage Approach," in in Proceedings of the 7th International Conference on Document Analysis and Recognition (ICDAR 2003),IEEE, 2003.
- [12]. S. Perantonis Gatos B I. Pratikakis, "An Adaptive Binarization Technique for Low Quality Historical Documents," in Document Analysis Systems VI, vol. 3163, , 2004.
- [13]. Chen Y. and G. Leedham, "Decompose algorithm for thresholding degraded historical document images," in IEEE Proc.-Vis. Image Signal Process vol. 152, , December 2005.
- [14]. Chen Y. and G. Leedham, "Document binarization using Kohonen," in IET Image Process, 2007, pp. 67-85.
- [15]. Ioannis Pratikakis, and Stavros J. Perantonis. Gatos Basilios, "Improved document image binarization by using a combination of multiple binarization techniques and adapted edge information," in Pattern Recognition, vol. ICPR 2008. 19th International Conference on. IEEE, 2008, 2008.
- [16]. Konstantinos Ntirogiannis, and Ioannis Pratikakis Gatos Basilios, "ICDAR 2009 Document Image Binarization Contest (DIBCO 2009)," in ICDAR, vol. 9, 2009.
- [17]. Michael S. Brown, and Dong Xu. Huang Yi, "User-assisted ink-bleed reduction," in Image Processing, IEEE Transactions , oct 2010, pp. 2646-2658.
- [18]. Shelke,Apte,"Multistage handwritten marathi compound character recognition using neural networks"Journals of pattern recognition research 2, 253-268, 2011 .
- [19]. I.K.Sethi,"Machine recognition of Constrained Hand printed Devanagri" pattern recognition , vol.9,1977.
- [20]. Swapnil A. Vaidya, Balaji R. Bombade," A Novel Approach of Handwritten Character Recognition using Positional Feature Extraction", International Journal of Computer Science and Mobile Computing, IJCSMC, Volume 2, Issue.6, pg.179 – 186, June 2013.
- [21]. Sandhya Arora, Debotosh Bhattacharjee, Mita Nasipuri,L.Malik, M.Kundu and D.K.Basu, "Recognition of Non-Compound Handwritten Devnagari Characters using a Combination of MLP and Minimum Edit Distance".International Journal of Computer Science and Security (IJCSS), Volume (4), Issue( 1), 2014.
- [22]. P. Bhaskara Rao, D.Vara Prasad, Ch.Pavan Kumar, "Feature Extraction Using Zernike Moments" International Journal of Latest Trends in Engineering and Technology (IJLTET), ISSN: 2278-621X, Vol. 2, Issue 2 March 2013.
- [23]. Hamid Reza Boveiri,"On Pattern Classification Using Statistical Moments " International Journal of Signal Processing, Image Processing and Pattern Recognition processing and Pattern Recognition Vol. 3, No. 4, December, 2010.
- [24]. M.Hanmandalu and O.V. Ramana Murthy,"Fuzzy model based recognition for handwritte hindi numerals",International conference on recognition , pp.490-496, 2005.
- [25]. Ambadas B. Shinde, Yogesh H. Dandawate,"Shirorekha Extraction in Character Segmentation For Printed Devanagri Text In Document Image", IEEE 978-1-4799-5364-6/14 2014.
- [26]. Adawait Dixit, Ashwini Navghane, Yogesh Dandawate, "Handwritten Devanagari Character Recognition using Wavelet Based Feature Extraction and Classification Scheme", 11 th IEEE India conference INDICON, ISBN no.978-1-4799-5362-2,2014.
- [27]. Satish Kumar and Chandan Singh,"A Study of Zernike moments and it's use in devanagri Handwritten character recognition "Int. Conf. on Recognition, pp.514-520, 2005.
- [28]. U.Bhattacharya, B.B.Chaudhari, R.Ghosh, M.Ghosh,"On recognition of handwritten Devnagri Numerals ",In Proc. of the workshop on learning algorithm for pattern recognition , Sydney, pp.1-7,2005.
- [29]. N.Sharma,U.Pal,F.Kimura and S.Pal,"Recognition of offline handwritten devnagri characters using quadratic classifiers", ICVGIP, LNCS4338, pp.805-816 ,2006.
- [30]. Bikash Shaw, Swapan Kumar Parui, Malayappan Shridhar, "Offline Handwritten Devanagari Word Recognition: A Holistic Approach Based on Directional Chain Code Feature and HMM,"pp. 203-208, International Conference on Information Technology, 2008.
- [31]. A.A.Ghatol, R.B.Ghongade,"A brief performance evaluation of ECG feature extraction techniques for artificial network based classification", TENCON 2007-2007 IEEE region 10 conference, 10/2007.
- [32]. Utpal Garain Et.Al. In "Segmentation Of Touching Characters In Printed Devnagari And Bangla Scripts Using Fuzzy Multifactorial Analysis",IEEE, transactions on systems, man, and cybernetics,vol. 32, no. 4, pp. 1094-6977/02,nov 2002
- [33]. R. Jayadevan et.al. in "Offline Recognition of Devanagari Script: A Survey",IEEE, transactions on systems, man, and cybernetics, vol. 41, no. 6,pp. 1094-6977,nov 2011.
- [34]. Naveen Sankaran et.al. in "Recognition of Printed Devanagari Text Using BLSTM Neural Network",ICPR, International Conference on Pattern Recognition,Report No: IIIT/TR/2012/-1,Nov 2012
- [35]. Manas Yetirajam et.al. in "Recognition and Classification of Broken Characters using Feed Forward Neural Network to Enhance an OCR Solution",IJARCET, Volume 1, Issue 8, ISSN: 2278 – 1323,Oct 2012
- [36]. Saiprakash Palakollu et.al. in "Handwritten Hindi Text Segmentation Techniques for Lines and Characters",WCECS, Proceedings of the World Congress on Engineering and Computer Science, Vol I, ISBN:978-988-19251-6-9,Oct2011.
- [37]. Er. Binny Thakral et.al in " Devanagri Handwritten Text Character Segmentation Techniques and Related Issues - A Review"IJARCSSE, Volume 4, Issue 7, pp. 1030-1034,July 2014.
- [38]. Priyanka U. Barve et.al. in "A Survey: Problems of Overlapped Handwritten Characters in Recognition process for Devanagari Script",IJCTA , Vol 5, p.p 941-946, June 2014
- [39]. Rajiv Kumar et.al. in "Character Segmentation in Gurumukhi Handwritten Text using Hybrid Approach",IJCTE, Vol. 3, No. 4, August 2011
- [40]. C. Halder et.al in "Word & Character Segmentation for Bangla Handwriting Analysis & Recognition",IEEE, Third National Conference on Computer Vision, Pattern Recognition, Image Processing and Graphics,p.p 978-0-7695-4599-8, 2011
- [41]. ZhongHua Cao et.al in "A New Drop-Falling Algorithms Segmentation of Touching Character"IEEE,p.p 978-1-4244-6055-7, 2010
- [42]. M Swamy Das et.al in "SEGMENTATION OF OVERLAPPING TEXT LINES, CHARACTERS IN PRINTED TELUGU TEXT DOCUMENT IMAGES",IJEST, Vol. 2(11), P.P 6606-6610, 2010
- [43]. Atallah. M. AL-Shatnawi et.al in "Skeleton Extraction: Comparison of Five Methods on the Arabic IFN/ENIT Database",IEEE, 6th International Conference on CSIT, p.p 978-1-4799-3999-2, 2014
- [44]. Munish Kumar et.al in " Segmentation of Isolated and Touching Characters in Offline Handwritten Gurmukhi Script Recognition",IJITCS , VOL 02, P.P 58-63,Jan 2014
- [45]. Dr. Jenila Livingston L.M. in " Text Detection From Documented Image Using Image Segmentation",IJTEEE, VOL 1, ISSUE 4, p.p 2347-4289 , 2013