



# Segmentation of Degraded Document Text by Local Threshold Method

A.Anees Fathima<sup>1</sup>, T.N. Sudhashree<sup>2</sup>

Electronics and communication Engineering, Dhaanish Ahmed College of Engineering ,  
India<sup>1</sup>.

Asst.Professor, Dept.of Electronics and Communication Engineering, Dhaanish Ahmed  
College of Engineering, India<sup>2</sup>

**ABSTRACT**— Restoration plays a very important role in enhancing the degraded image. This paper proposes a novel document image binarization technique that addresses the issues by using adaptive image contrast. The adaptive image is the combination of local image contrast and local image gradient. For a degraded input image, adaptive image contrast is first constructed. The contrast map is then binarized and combined with canny's edge map to identify the text stroke edge pixels. The local threshold segments the document text that is estimated based on the intensities of detected text stroke edge pixels. Post processing techniques are applied to the image and the restored document is obtained.

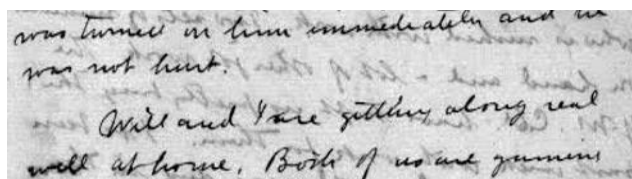
**Keywords**— Image Contrast, document analysis, image processing, degraded document, pixel classification.

## 1, INTRODUCTION

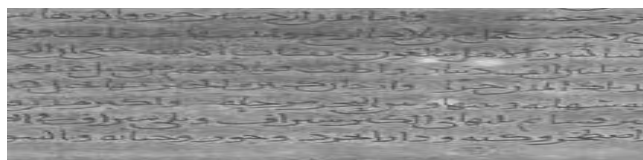
It is very important to preserve the historical document which reveals the information about the civilized past. The document posed much degradation due to weather conditions, method of preservation etc. The images so captured have major problems like broken letters, erased letters. Such degraded document consists of noise and blurred characters which are connected or split. One possibility to handle such documents is to improve the image quality. Document Image Binarization is performed in the pre processing stage for document analysis and it aims to segment foreground and background. A fast and accurate Binarization technique is important in ensuing document image processing task such as optical character recognition (OCR). Document image binarization is of great significance for the optical character recognition. The use of OCR as a means of evaluation of modern printed documents only, supported by contemporary engines. Most document analysis algorithm is built on taking advantage of the underlying binarized image data [1]. Document image requires logical and semantic content preservation during thresholding. Binarization has been a subject of intense research during last ten years. Most of the developed algorithms rely on statistical methods, not considering the special nature of document images. Binarization methodologies can be classified into global and local



adaptive. In this case if local adaptive, the adaption is based on parameters which are usually related to contrast and the gradient of the images. Fig 1(a) and (b) shows the Degraded document Image examples.



**Figure 1(a) Degraded Document Image**



**Figure. 1(b) Degraded Document image**

## 2, Existing System

This paper presents a document image binarization technique that segment the document text within it has a different intensity level compared with the surrounding document background. The text stroke edge is further detected from the compensated. Document image by using L1-norm image gradient. Finally, the document text is segmented by a local threshold that is estimated based on the detected text stroke edges. Local thresholding [6]-[9] for each document image pixel is estimated by adaptive thresholding is often a better approach to deal with different variations within degraded document images. The early window-based adaptive thresholding techniques [10], [11] estimate the local threshold by using the mean and standard variations of image pixels with in a local neighborhood window. Thresholding performance depends on the window-size set to 3 empirically and the character stroke width is the drawbacks of the window-based thresholding technique.

### 2.1 Local Image Contrast

Contrast is the difference in luminance and/or color that makes an object (or its representation in an image or display) distinguishable. The maximum contrast of an image is the contrast ratio or dynamic range. Document text has certain image contrast to the neighboring document background hence local image contrast and local image gradient are useful for segmenting the text from the document background. In Bersen's paper [12] the local contrast is defined as follows:

where  $C(i, j)$  denotes the contrast of an image pixel  $(i, j)$ ,  $I_{max}(i, j)$  and  $I_{min}(i, j)$  denotes the maximum and minimum intensities within a local neighborhood windows of  $(i, j)$  respectively.. If the local contrast  $C(i, j)$  is smaller than a threshold, the pixel is set as



background directly. Otherwise it will be classified into text or background by comparing with the mean of  $I_{max}(i, j)$  and  $I_{min}(i, j)$ . Bersen's method is simple, but cannot work properly on degraded document images with a complex document background.

A novel document image binarization method [5] by using the local image contrast that is evaluated as [13].

Where  $\epsilon$  is a positive but infinitely small number that is added in case the local maximum is equal to 0. Compared with the Bersen's contrast in equation 1, the local image contrast in equation 2 introduces a normalization factor to compensate the image variation within the document areas such as that in the sample document image in Fig 1(b) as an example. The small image contrast around the text stroke edges in equation 1 will be compensated by a small normalization factor (due to the dark document background) as defined in Equation 2.

## 2.2 Local Image Gradient

An image gradient is a directional change in the intensity or color in an image. Image gradients may be used to extract information from images.

The image gradient has been widely used for edge detection and it can be used to detect the text stroke edges of the document images effectively that have a uniform document background. On the other hand, it often detects many nonstroke edges from the background of degraded document that often contains certain image variations due to noise, uneven lighting, bleed-through, etc. To extract only the stroke edges properly, the image gradient needs to be normalized to compensate the image variation within the document background.  $(I_{max}(i, j) - I_{min}(i, j))$  is refers to the local image gradient that is normalized to  $[0, 1]$ .

## 3, PROPOSED TECHNIQUE

This paper is on the development of new approaches for restoration of degraded document images. In the proposed document image binarization techniques an adaptive contrast map [3] is first constructed for a given degraded document image and the text stroke edges are then detected through the combination of the binarized adaptive contrast map and the canny edge map. The text is then segmented based on the local threshold that is estimated from the detected text stroke edge pixels [6]. Some post-processing is further applied to improve the document binarization quality.



### 3.1 Construction of Adaptive Image Contrast

Adaptive image contrast is the combination of local image contrast and local image gradient.

is the weight between local contrast which is assigned high weight and gradient and produce good results. Otherwise the image gradient will be set high. The proposed system relies on the image gradient and avoid over normalization factor. where  $C(i, j)$  denotes the local contrast[3] in Equation 2 and  $(I_{max}(i, j) - I_{min}(i, j))$  refers to the local image gradient that is normalized to  $[0, 1]$ . The local windows size is set to 3 empirically.  $\alpha$  is the weight between local contrast and local gradient that is controlled based on the document image statistical information. Ideally, the image contrast [3] will be assigned with a high weight (i.e. large  $\alpha$ ) when the document image has significant intensity variation. relies more on image gradient and avoid the over normalization problem.

We model the mapping from document image intensity variation to  $\alpha$  by a power function as follows:

Where  $Std$  denotes the document image intensity standard deviation and  $\gamma$  is a pre-defined parameter. The power function has a nice property in that it monotonically and smoothly increases from 0 to 1 and its shape can be easily controlled by different  $\gamma$ .  $\gamma$  can be selected from  $[0, \infty]$ , where the power function becomes a linear function when  $\gamma = 1$ . the sample document has small intensity variation within the document background but large intensity variation within the text stroke, the use of the local image contrast removes many light strokes improperly in the contrast map, whereas the use of local image gradient is capable of preserving those light text strokes.

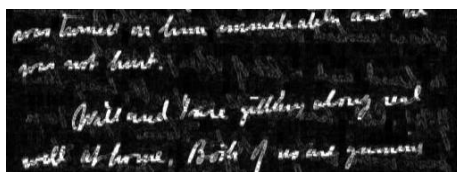


Figure 2(a)

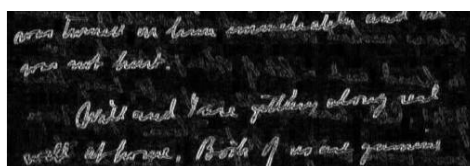


Figure 2(b)

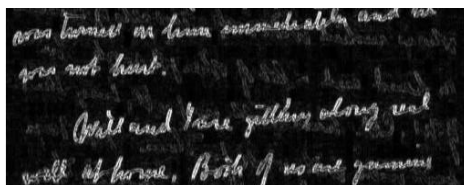


Figure 2 (c)

Figure 2. Contrast Image constructed using (a) local image gradient, (b) local image contrast, (c) adaptive image contrast. The local image maximum and minimum is used to suppress the background variation as described in Equation 2. The image contrast in Equation 2 has the limitation of not handling the bright text strokes. Fig 2 shows the contrast map of the document that is created by using local image gradient and local image contrast and our proposed method in Equation 3. The power function becomes linear function when  $\gamma = 1$ . Therefore the local image gradient plays a major role in Equation 3.

### 3.2 Binarization

Binary images are also called bi-level or two-level. This means that each pixel is stored as a single bit (0 or 1). We detect the text stroke edge pixel properly by using Otsu's global thresholding method. Binary image mapping can be further improved through the combination with the edges by canny edge detector, because canny edge detector has good localization property as it can mark all the edges to the real edge location in the detecting image.

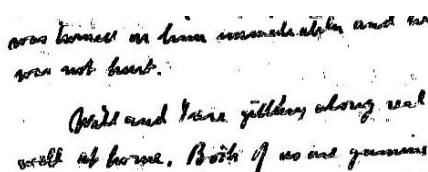


Figure 3(a)

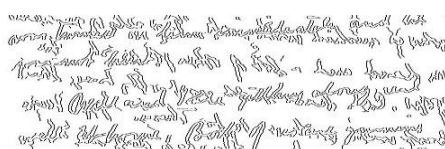


Figure 3(b)

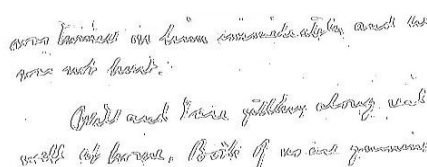


Figure 3(c)

Figure. 3. (a) Otsu thresholding (b) canny edge maps (c) combined output of edge maps of the sample document images



### 3.3 OTSU Threshold (Binarization)

**OTSU** is redirected from **Operational Test Support Unit**. It is automated choose threshold value in the image. This OTSU thresholding is based on the histogram of the input image.

Thresholding is a very basic operation in image processing. And, a good algorithm always begins with a good basis! Otsu thresholding is a simple yet effective global automatic thresholding method for binarizing grayscale images such as foregrounds and backgrounds.

### 3.4 Canny edge Detector

The canny edge detector is an edge detection operator that uses a multistage algorithm to detect a wide range of edges in images. It was developed by John.F.Canny in 1986.Canny edge detector uses two adaptive thresholds and is more tolerant to different imaging artifacts.

## IV. . CONCLUSION

This paper presents an adaptive image contrast based document image binarization technique [9] that is tolerant to different types of document degradation such as uneven illumination and document smear. The proposed technique is simple and robust, only few parameters are involved. Moreover, it works for different kinds of degraded document images. The proposed technique makes use of the local image contrast that is evaluated based on the local maximum and minimum.

## V. FUTURE ENHANCEMENT

In the combination of binarization and edge detection output we have some unwanted regions. The unwanted edges are removed using further process i.e. Edge Width Estimation, local Threshold of input image (binarization) and post processing. After post processing, removing unwanted edges of the image and clear output is obtained. Experiments show that the proposed method outperforms most reported document binarization methods in term of the Fmeasure, PSNR, etc

## REFERENCES

- [1] J.sauvola, M.pietikainen, page segmentation and classification using fast feature extraction and connectivity analysis, international conference on document analysis and recognition, ICDAR '95, Montreal, Canada, 1995,pp 1127-1131
- [2] B.Gatos, K. Ntirogiannis, and I.Pratikakis, 'ICDAR 2009 document image binarization contest(DIBCO 2009).'in proc .Int.Conf.Document Anal.Recogniy.,Jul.2009,pp. 1375-1382.



- [3] I.Pratikakis, B.Gatos, and K.Ntirogiannis, "ICDAR 2011 document image binarizationcontest (DIBCO 2011),' in proc,Int,Conf.Document Anal.Recognit., sep.2011,pp, 1506-1510
- [4] S.Lu,B.Su, and C.L. Tan, "Document Image binarization using background estimation and stroke edge,'Int.J.Document Anal. Recognit vol. 14, no. 4,pp. 303-314, Dec. 2010
- [5] B.Su, S.Lu, and C.L.Tan, "Binarization of historical handwritten document images using local maximum and minimum filter,' in proc.Int.Workshop Document Anal.Syst.,Jun.2010, pp. 159-166
- [6] G.Leedham, C.Ya, K.Taru, J.Hadi,N.Tan, and L.Mian, "comparison of some thresholding algorithm for text/background segmentation in difficult document images," in proc. Int. Conf.Document Anal.Recognit., vol. 13.2003, pp. 859-864.
- [7] M.Sezgin and B.Sankur, "survey over image thresholding techniques and quantitative performance evaluation," J.Electron.Imag.,vol13,no. 1,pp, 146-165, Jan. 2004.
- [8] O.D. Trier and A.K. Jain, "Goal-directed evaluation of binarization methods,"IEEE Trans.pattern Anal.Mach.Intell., vol. 17, no. 12, pp. 1191-1201, Dec. 1995.
- [9] O.D. Trier and T. Taxt, "Evaluation of binarizatoin methods for document images," IEEE Trans.Pattern Anal.Mach.Intell., vol 17, no 3,pp. 312-315,Mar.1995.
- [10] J.Sauvol and M.Pitikaien,'Adaptive document image binarization,'Pattern Recognit., vol.33, no. 2,pp. 225-236,2000.
- [11] W.Niblack, An Introducion to Digital Image Proessing, Englewood Cliffs,NJ: Prentice Hall, 1986
- [12] J.Bersen, 'Dynamic thresholding of gray-level images,' in Proc, Int.Conf.Pattern Recogniy., aoct 1986, pp.1251-1255
- [13] M.van Herk, "A fast algorithm for local minimum and maximum filters on rectangular and octagonal kernels," Pattern Recognit. Lett., vol.13, no. 7,pp.517-521, Jul.1992