

## IMPROVING DEGRADED DOCUMENT IMAGES USING BINARIZATION TECHNIQUES

G. VARA LAKSHMI<sup>1</sup> & P. KAMALA<sup>2</sup>

<sup>1</sup>Research Scholar, Vignan Institute of Engineering for Women, Vishakhapatnam, A.P, India

<sup>2</sup>Assistant Professor, Department of ECE, Vignan Institute of Engineering for Women, Vishakhapatnam, A.P, India

### ABSTRACT

Restoring data from a degraded document image is a most important process in many applications. Those degraded document images are taken from the DIBCO (Document Image Binarization Contest)-2009 dataset. Segmentation of text from a degraded document images is very difficult task, because those document images are suffering from smear, ink-bleeding, smudge, etc. In this paper, we propose a document image binarization technique that addresses the above issues by using the adaptive local image contrast method. In the proposed technique, first we are performing the preprocessing step in order to eliminate the noise in the image by using the mean filter and then an adaptive contrast map is constructed by using the histogram equalization method for the preprocessed document image. Next step is to detect text stroke edge pixels of the document images by using the Canny's edge detector. The document text is then segmented by using the local threshold that is estimated based on the intensities of detected text stroke edge pixels. The proposed method is simple and it involves minimum parameter calculations. Performance of the proposed method is good compared to the previous by using some of the quality metrics like Peak Signal to Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM), etc.

**KEYWORDS:** Degradation, Adaptive Image Contrast, Histograms, Image Edge Detection, Pixel Classification, Image Segmentation

### INTRODUCTION

The Image segmentation is an essential task in the fields of image processing and computer vision. It is a process of partitioning the digital images and is used to locate the boundaries into a finite number of meaning full regions and easier to analyze. The Simplest method for image segmentation is thresholding. Thresholding is an important technique in image segmentation, enhancement and object detection. The output of the thresholding process is a binary image whose gray level value 0 (black) will indicate a pixel belonging to a print, legend, drawing, or target and a gray level value 1 (white) will indicate the background. The main complexity coupled with thresholding in document applications happen when the associated noise process is non-stationary. The factors that make difficult thresholding action are ambient illumination, variance of gray levels within the object and the background, insufficient contrast, object shape and size non-commensurate with the spectacle. The lack of objective measures to assess the performance of thresholding algorithms is another handicap. Many methods have been reported in the literature. It can extract the object from the background by grouping the intensity values according to the thresholding value.

Thresholding divides the image into patches, and each patch is thresholding by a threshold value that depends on the patch contents. In order to decrease the effects of noise, common practice is to first smooth a boundary prior to partitioning. The Binarization technique is aimed to be used as a primary phase in various manuscript analysis, processing

and retrieval tasks. So, the unique manuscript characteristics, like textual properties, graphics, line drawings and complex mixtures of the layout-semantics should be included in the requirements.

## RELATED WORK

Generally we are having so many thresholding techniques for the degraded document image binarization. But it is usually not a suitable approach for the degraded document image binarization. So we are using the Adaptive thresholding, which estimates a local threshold value for each document image pixel, is often a better approach to deal with different variations within the degraded document images. For example, the early window-based adaptive thresholding techniques estimate the local threshold value by using the mean and the standard deviation of the image pixels within a local neighborhood

The main drawback of these local thresholding techniques is that the thresholding performance depends heavily on the window size and the character stroke width. Other techniques also reported, including dynamic thresholding of gray level images, a threshold selection method from gray level histograms, adaptive document image binarization, document image binarization using background estimation and stroke edges, combination of document image binarization techniques, comparison of some thresholding algorithms for text/ background segmentation in difficult document images, survey over image thresholding techniques and quantitative performance evaluation. There are two features for segmenting text from the document background in a degraded document image. They are the local image contrast and the local image gradient. Because the document text usually has a certain image contrast to the neighboring document background. They are very useful and have been used in many document image binarization techniques. In Bernsen's paper [1], the local contrast is defined as follows:

$$C(i, j) = I_{\max}(i, j) - I_{\min}(i, j) \quad (1)$$

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 the minimum intensities within a local neighborhood windows of size  $(i, j)$ , 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  $I_{\min}(i, j)$  Bernsens's method is simple, but cannot work properly on degraded document images with a complex document background. We have earlier proposed a novel document image binarization method by using the local image contrast that is evaluated as follows:

$$C(i, j) = \frac{I_{\max}(i, j) - I_{\min}(i, j)}{I_{\max}(i, j) + I_{\min}(i, j) + \epsilon} \quad (2)$$

## IMPLEMENTATION

This section describes the various document image binarization techniques like optimized global thresholding using ostu's [2] method, sauv [3] and proposed method.

### Optimized Global Thresholding Using OTSU'S Method

In computer vision and image processing, Otsu's method, named after Nobuyuki Otsu is used to automatically perform clustering-based image thresholding, [4] or, the reduction of a gray level image to a binary image. i.e. This method assumes that the image contains two classes of pixels following bi-modal histogram (foreground pixels and background pixels), it then calculates the optimum threshold separating the two classes so that their combined spread (intra-class

variance) is minimal, or equivalently (because the sum of pair wise squared distances is constant), so that their inter-class variance is maximal.

In Otsu's method we are first finding the threshold value. In order to find the threshold we need to find intra-class variance (the variance within the class), defined as a weighted sum of variances of the two classes:

$$\sigma_w^2(t) = w_1(t)\sigma_1^2(t) + w_2(t)\sigma_2^2(t) \quad (3)$$

Weights  $w_i$  are the probabilities of the two classes separated by a threshold 't' and  $\sigma_i^2$  are the variances of those two classes. It is difficult to calculate two class variances. So we are going for inter-class variance.

Otsu says that minimizing the intra-class variance is the same as maximizing inter-class variance (variance between the classes):

Inter-class variance is represented as

$$\sigma_b^2(t) = \sigma^2 - \sigma_w^2(t) = w_1(t) w_2(t) [\mu_1(t) - \mu_2(t)]^2 \quad (4)$$

Which is expressed in terms of class probabilities  $w_i$  and class means  $\mu_i$ .

The class probability  $w_1(t)$  is computed from the histogram as 't'

$$w_1(t) = \sum_{i=0}^t p(i) \quad (5)$$

While the class mean  $\mu_1(t)$  is:

$$\mu_1(t) = \frac{[\sum_{i=0}^t p(i)x(i)]}{w_1} \quad (6)$$

Where  $x(i)$  is the value at the center of the  $i^{\text{th}}$  histogram bin. Similarly, you can compute  $w_2(t)$  and  $\mu_2$  on the right-hand side of the histogram for bins greater than 't'. The class probabilities and class means can be computed iteratively. This will give an effective algorithm.

### Algorithm

- Compute histogram and probabilities of each intensity level
- Set up initial  $w_i(0)$  and  $\mu_i(0)$
- Step through all possible thresholds  $t=1 \dots$  maximum intensity
  - Update  $w_i$  and  $\mu_i$
  - Compute  $\sigma_b^2(t)$
- Desired threshold corresponds to the maximum  $\sigma_b^2(t)$
- You can compute two maxima (and two corresponding thresholds).  $\sigma_{b1}^2(t)$  is the greater max and  $\sigma_{b2}^2(t)$  is the greater or equal maximum
- Desired threshold =  $\frac{\text{threshold}_1 + \text{threshold}_2}{2}$

From the desired threshold value we can segment the text from the document background. This method is having

computational redundancy because we are calculating many parameters and thus resultant image is having noise compared to the original document image.

### Sauvola Method

This method gives more performance than previous methods under some conditions like light variation on document image, light texture etc. In the Sauvola modification, the binarization is given by

$$T_{sauvola} = m * (1 - k * (1 - \frac{S}{R})) \quad (7)$$

Where  $m$  is the mean of pixels under window area,  $S$  is the dynamic range of variance and the value of  $k$  parameter may be in the range of 0.2-0.5.  $R$  is the maximum value of the standard deviation. The author fix the values  $k=0.5$  and  $R=128$ .

### Adaptive Image Contrast

This section describes the proposed document image binarization technique. Given a degraded document image, first we are converting the original document image into a binary image for clarity purpose then we are having preprocessed step, in this we are removing noise using mean filter and then an adaptive contrast map is constructed by using histograms and the text stroke edges are then detected through the canny edge map. The text is then segmented based on the threshold that is estimated from the detected text stroke edge pixels mean value.

### Preprocessing:

We are taking the document images from the dataset DIBCO-2009. Those documents contain unwanted pixels i.e. noise in order to eliminate that noise we are using mean filter. Mean filter will remove the noise present in the image. Mean value will be calculated from the surrounding pixels and the defected pixel will be replaced by the mean value.

### Contrast Image Construction:

The image contrast and 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 non-stroke edges from the background of degraded document that often contains certain image variations due to noise, uneven lighting, bleed-through etc. another way to construct a contrast map is by using histograms.

A histogram can also describe the amount of contrast. Contrast is a measure of the difference in brightness between light and dark areas in a scene. Broad histograms reflect a scene with significant contrast, whereas narrow histograms reflect less contrast and may appear flat or dull. This can be caused by any combination of subject matter and lighting conditions. Photos taken in the fog will have low contrast, while those taken under strong daylight will have higher contrast. The high contrast water has deeper shadows and more pronounced highlights, creating texture which "pops" out at the viewer. Contrast can also vary for different regions within the same image due to both subject matter and lighting. We can partition the previous image of a boat into three separate regions—each with its own distinct histogram.

Proposed method uses histogram equalization method to construct the contrast map. In histogram equalization we are having steps, first we have to find the running sum of the histogram values, then normalize those values by dividing the running sum with the total number of pixels, multiply the resultant values with the maximum gray level value and round off those values, finally map those gray level values to the input image.

**Text Stroke Edge Pixel Detection:**

The purpose of the contrast image construction is to detect the stroke edge pixels of the document text properly. The constructed contrast image has a clear bi-modal pattern, where the adaptive image contrast computed at text stroke edges is obviously larger than that computed within the document background.

Text strokes are detected using edge detection techniques. We are using Canny's edge detector, because Canny's edge detector has a good localization property that it can mark the edges close to real edge locations in the detecting image. In addition, canny edge detector uses two thresholds and is more tolerant to different imaging artifacts such as shading, smearing, etc.

It should be noted that Canny's edge detector by itself often extracts a large amount of non-stroke edges. Further implementation

**Local Threshold Estimation:**

The text can then be extracted from the document background pixels once the high contrast stroke edge pixels are detected properly. Two characteristics can be observed from different kinds of document images: First, the text pixels are close to the detected text stroke edge pixels. Second, there is a distinct intensity difference between the high contrast stroke edge pixels and the surrounding background pixels. The document image text can thus be extracted based on the detected text stroke edge pixels. The algorithm calculates the mean value in a window and if the pixel's intensity is below the mean the pixel is set to white color, otherwise the pixel is set to black color. Finally text be segmented base on mean.

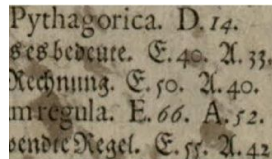
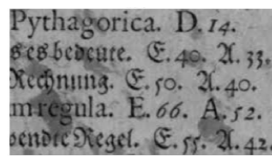
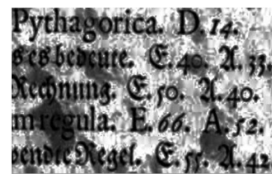
**RESULTS****Subjective Analysis of Proposed Method:****Figure 1: Original Image****Figure 2: Pre-Processing Image  
Using Mean Filter****Figure 3: Contrast Mapping**



Figure 4: Text-Stroke Detection

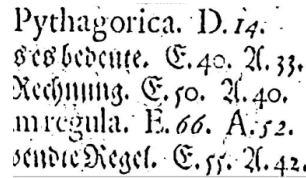
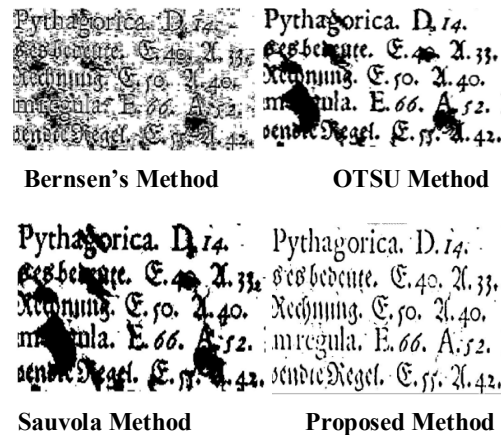


Figure 5: Resultant Image

Subjective Comparison of DIBCO-2009 Dataset Document Image



Objective Comparison:

Table 1: Evaluation Results of the Dataset of DIBCO 2009

Method	PSNR	SSIM	MPM	NRM
Bernsen	7.894	0.0037	0.070	0.414
Otsu	8.122	0.0039	56.60	0.410
Sauv	7.256	0.0038	0.073	0.594
Proposed	11.83	0.0050	0.033	0.532

CONCLUSIONS

In this work presents an adaptive image contrast based document image binarization technique that is tolerant to different types of document degradations like document smear, shading, smudge, ink-bleeding, blurring etc. The proposed technique is simple and robust, only few parameters are involved. Moreover, it works for different kinds of degraded document images. The proposed method has been tested on the various datasets. Experiments show that the proposed method outperforms most reported document binarization methods in terms of PSNR, SSIM, MPM, NRM. MPM and NRM values of our proposed method cannot give accurate results. Finally text in the degraded documents will be extracted clearly in proposed method than previous methods although MPM and NRM values are less.

**REFERENCES**

1. N. Otsu, "A threshold selection method from gray level histogram," *IEEE Trans. Syst., Man, Cybern.* Vol. smc-9, no. 1, January 1979.
2. J. Sauvola and M. Pietikainen, "Adaptive document image binarization," *Pattern Recognition* 33(2), pp. 225-236, 2000.
3. B. Su, S. Lu, and C. L. Tan, "Robust Document Image Binarization Technique for Degraded Document Images," *IEEE Transactions on Image Processing*, vol.22, no.4, April 2013.
4. B. Gatos, K. Ntirogiannis, and I. Pratikakis, "ICDAR 2009 document image binarization contest (DIBCO 2009)," in *Proc. Int. Conf. Document Anal. Recognit.* Jul. 2009, pp. 1375–1382.
5. S. Lu, B. Su, and C. L. Tan, "Document image binarization using back- ground estimation and stroke edges," *Int. J. Document Anal. Recognit.* vol. 13, no. 4, pp. 303–314, Dec. 2010.
6. 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.
7. M. Sezgin and B. Sankur (2004). "Survey over image thresholding techniques and quantitative performance evaluation". *Journal of Electronic Imaging* 13 (1): 146–165.

