

Ternary Entropy-based Binarization of Degraded Document Images Using Morphological Operators

T. Hoang Ngan Le

Dept. Computer Science & Software
Engineering, Concordia University
Montreal, Quebec, Canada
l_thihoa@encs.concordia.ca

Tien D. Bui

Dept. Computer Science & Software
Engineering, Concordia University
Montreal, Quebec, Canada
bui@encs.concordia.ca

Ching Y. Suen

Dept. Computer Science & Software
Engineering, Concordia University
Montreal, Quebec, Canada
suen@encs.concordia.ca

Abstract—A vast number of historical and badly degraded document images can be found in libraries, public, and national archives. Due to the complex nature of different artifacts, such poor quality documents are hard to read and to process. In this paper, a novel adaptive binarization algorithm using ternary entropy-based approach is proposed. Given an input image, the contrast of intensity is first estimated by a grayscale morphological closing operator. A double-threshold is generated by our Shannon entropy-based ternarizing method to classify pixels into text, near-text, and non-text regions. The pixels in the second region are relabeled by the local mean and the standard deviation. Our proposed method classifies noise into two categories which are processed by binary morphological operators, shrink and swell filters, and graph searching strategy. The method is tested with three databases that have been used in the Document Image Binarization Contest 2009 (DIBCO 2009), the Handwriting Document Image Binarization Contest 2010 (H-DBCIO 2010), and the International Conference on Frontier in Handwriting Recognition 2010 (ICFHR 2010). The evaluation is based upon nine distinct measures. Experimental results show that our proposed algorithm outperforms other state-of-the-art methods.

Keywords—binarization; degraded document image; ternary entropy-based; morphological operators; Shannon entropy

I. INTRODUCTION

Binarization means converting a multi-color image to a bi-color one by threshold selection techniques that classify the pixels of the image into either one (white) or zero (black). It has been studied for many years, and plays an important role in document analysis systems. As more and more documents are digitized, less time consuming and higher accuracy document image binarization has become increasingly necessary. However, thresholding badly degraded document image (DI) is still an unsolved problem due to features such as bleeding-through, large black border, ink-fading, uneven illumination, contrast variation, smear, various pattern background, and so on. Figs. 1 and 2 show some examples of such DI. Much research on binarization has been studied and reported in the literature [5][8][17]. Generally, binarization techniques can be divided into two categories, namely, global thresholding and local adaptive thresholding which are described below.

1) Global Thresholding Techniques

Global thresholding techniques partition an entire image based on a single threshold value which is obtained by the gray level histogram of the image. Many popular global thresholding techniques using different approaches such as: (1) fixed thresholding technique which binarizes with respect to a specified threshold value. This approach works well when the background and foreground intensities are clearly distinct and uniform over the image. This technique usually uses a mean value, for example 127, as default threshold value; (2) histogram shape-based binarization that makes use of the minimum (or the valley) between two peaks of the histogram as a threshold value [1]. This approach requires some knowledge of noise and the distance between peaks; thus, it is not suitable when peaks of the histogram are not clearly determined. To address this problem, some more pattern recognition methods are used to optimize foreground/background separation; (3) optimal thresholding binarization which applies a criterion function to each pixel in order to separate an image into two regions corresponding to white and black. The techniques of this group can be classified into sub-categories such as discriminant method [14], entropy thresholding [18][19], moment preservation [3][9], and minimum error thresholding [12].

2) Local Thresholding Techniques

Local thresholding techniques apply a separate threshold value to a single pixel or a particular region. The local threshold value can be calculated by different information of the input image. This may also be known as dynamic thresholding and can be divided into different approaches such as background subtraction [7], water flow model [16], illumination model [11], mean and standard derivation of pixel values [10], and local image contrast [6]. Some drawbacks of the local thresholding techniques are region size dependant, individual image characteristics, and time consuming. Therefore, some researchers use a hybrid approach that applies both global and local thresholding methods [4]. In general, local adaptive methods achieve better results than global ones aside from that they often rely on some specific parameters and have high computational cost.

In this paper, we propose a novel nonparametric local adaptive binarization method for poorly degraded DI. Our system firstly estimates the contrast image through a grayscale morphological closing operator. Secondly, two threshold values using the Shannon entropy-based ternarization classify the pixels into three regions corresponding to text, near-text, and non-text

regions. A double-threshold value is selected to maximize the entropy of the contrast image. Thirdly, mean and standard deviation-based local thresholding is applied on the near-text areas to best choose the text pixels. Finally, the binarized image is post-processed by two noise removal approaches corresponding to salt-pepper and block noise. In addition to binary morphological operators (MO), and shrink and swell filters, a graph searching strategy is used to get rid of undesired regions and bridge the breaks as well as to enhance quality of the image. Experiments show that our proposed method outperforms other state-of-the-art binarization schemes for DI from three testing databases of DIBCO 2009 [5], H-DIBCO 2010 [8], and ICFHR 2010 [21]. The performance is evaluated by nine distinct measures.

The rest of this paper is arranged as follows. We briefly review the basic MO on grayscale and binary images in Section 2. Shannon entropy binarization technique is given in Section 3. Section 4 presents the proposed scheme in detail. Experiment along with comparisons are described and discussed in Section 5. Finally, concluding remarks are summarized in Section 6.

II. MORPHOLOGICAL OPERATORS

Let f and se (a small window) be functions representing $N \times N$ image and $(2M+1) \times (2M+1)$ structuring element, respectively, where M is the width of the largest character stroke. Four basic MO, namely, dilation, erosion, opening, and closing, on both binary and gray-level images are in turn introduced in following sub-sections.

Dilation: The binary dilation of f by se is defined as $g(x,y) = f \oplus se = \bigcup_{i=-M}^M \bigcup_{j=-M}^M se(i,j) \cap f(x-i,y-j)$. That means if just one of the '1's in se matches with the input signal f then the output is '1'. In binary images, the dilation operator is to gradually enlarge the boundaries of foreground regions, character typically. Thus, the areas of foreground grow in size whereas the holes within those regions become smaller. The grayscale dilation operator is defined as a convolution-like operation and to increase the brightness of pixels which are surrounded by higher intensity. The gray-level extremum operation for grayscale dilation is given as $g(x,y) = f \oplus se = \text{Max}[f(x-i,y-j) + se(i,j)]$.

Erosion: The binary erosion of f by se is defined as $g(x,y) = f \ominus se = \bigcap_{i=-M}^M \bigcap_{j=-M}^M se(i,j) \cup f(x+i,y+j)$. That means if all '1's in se match with the input signal then the output is '1'. In binary images, the erosion operator is to get rid of irrelevant details from the images. That means the holes within those areas become bigger. The grayscale erosion has the opposite consequence of the dilation because the resulting erosion image tends to be darker whereas the other tends to be brighter. This operator is $g(x,y) = f \ominus se = \text{Min}[f(x+i,y+j) + se(i,j)]$.

Opening: Opening function consists of an erosion function followed by a dilation and is defined as $g = f \circ se = (f \ominus se) \oplus se$. In binary images, this operator is to remove gaps, spurs, and islands smaller than the structuring element se . In gray-level images, it removes bright spots isolated in dark regions and smoothes boundaries because bright borders reduced by the erosion are restored by the dilation.

Closing: consists of a dilation function followed by an erosion and is defined as $g = f \bullet se = (f \oplus se) \ominus se$. In binary images, this operator smoothes the original signal by filling gulfs, channels, and lakes smaller than structuring element se . In gray-level images it removes dark spots isolated in bright regions and smoothes boundaries.

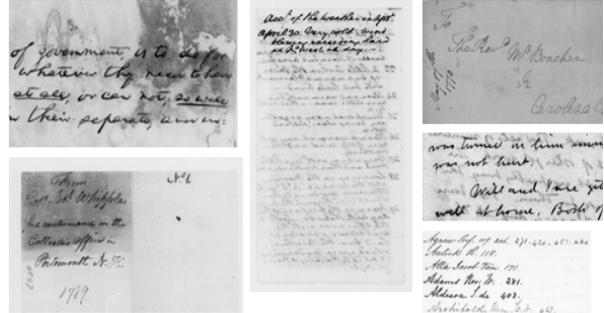


Figure 1. Representative images from DIBCO 2009 and H-DIBCO 2010



Figure 2. Representative example images from ICFHR 2010

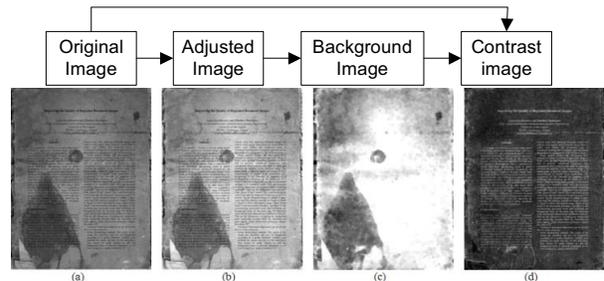


Figure 3. Contrast estimation

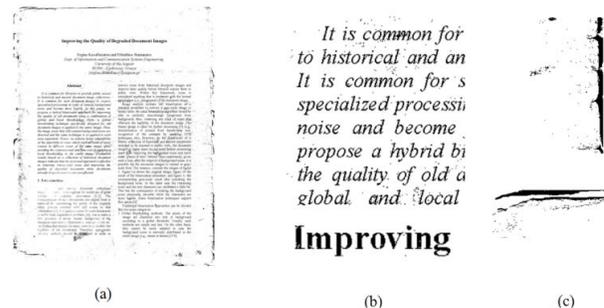


Figure 4. Binarization result without post-processing

```

Algorithm 1 Shannon entropy-based thresholding
Input: A grayscale image
Output: A threshold value Thres
Procedure:
    MAXIMUM=0
    Thres=0;
    For each t=MIN, t≤MAX, t++
         $p'(s) = \frac{p(s)}{\sum_{s=0}^{255} p(s)}$ ;  $p''(s) = \frac{p(s)}{\sum_{s=0}^{255} p'(s)}$ ;
         $H(S) = Hb + Hw = -\sum_{s=0}^{255} p'(s) \log(p'(s)) - \sum_{s=0}^{255} p''(s) \log(p''(s))$ ;
        if  $H(S) > \text{MAXIMUM}$  then
            MAXIMUM=H(S);
            Thres=t; \ end if \ end for

```

```

Algorithm 2 Ternarization system using Shannon entropy
Input: A grayscale image
Output: Two threshold values  $T_1$ , and  $T_2$ 
Procedure:
    MAXIMUM=0
     $T_1=0$ ;
     $T_2=0$ ;
    For each  $t_1=\text{MIN}$ ,  $t_1 \leq \text{MAX}$ ,  $t_1++$ 
        For each  $t_2=t_1+1$ ,  $t_2 \leq \text{MAX}$ ,  $t_2++$ 
             $p'(s) = \frac{p(s)}{\sum_{s=0}^{255} p(s)}$ ;  $p''(s) = \frac{p(s)}{\sum_{s=0}^{255} p'(s)}$ ;
             $H(S) = Hb + Hbw + Hw =$ 
             $-\sum_{s=0}^{255} p'(s) \log(p'(s)) - \sum_{s=0}^{255} p''(s) \log(p''(s)) - \sum_{s=0}^{255} p'(s) \log(p'(s))$ ;
            if  $H(S) > \text{MAXIMUM}$  then
                MAXIMUM=H(S);
                 $T_1=t_1$ ;
                 $T_2=t_2$ ; \ end if \ end for \ end for

```

III. ENTROPY BASED THRESHOLDING

Entropy is a measure of the uncertainty associated with a random variable. In Information Theory, it is assumed that there are N possible events, called s_i , which occur with probability $p(s_i)$, where $i=0,1, \dots, N-1$. Note that the entropy can be measured in bits/symbol. The Shannon entropy associated with the source S is defined as $H(S) = -\sum_{i=0}^{N-1} p(s_i) \log(p(s_i))$. In the entropy-based thresholding systems, the input image is considered as a source and the entropy is a summation of Hb and Hw corresponding to the entropy of black and white regions. The main point of these systems is to choose suitable parameters and an appropriate initial threshold. Some state-of-the-art schemes [2][18][19] recently used 2D-Tsallis entropy criteria to optimize the selection. A simple Shannon entropy-based thresholding is expressed in Algorithm 1.

1) Contrast Estimation

Contrast intensity is extracted by the grayscale morphological closing function. To ensure the contrast image contains all pixels that lie within the text area, the pixel intensity values in the input DI are first adjusted. Simply, all gray-levels are mapped to a new range from 0 to 255. Fig.3a shows an unevenly illuminated bleeding-through document image. Fig.3d illustrates the contrast image obtained by subtracting original image from the background (Fig.3c) which is generated by gray-level morphological closing operation on the adjusted image (Fig.3b).

2) Double Threshold Binarization

This section describes our method of generating the binarized image using the ternarization thresholding in combination with the local mean and the standard deviation. Let I and C be the original degraded DI and its contrast image. The binarizing procedure is illustrated in two phases. In Phase 1, an input image is classified into three regions corresponding to text, near-text, and non-text regions. In this phase, two thresholds T_1 and T_2 are generated by our improved Shannon entropy-based ternarization

as expressed in Algorithm 2. To relabel the pixels in the second region, we follow the method in Phase 2.

Phase 1: In our Shannon entropy-based ternarization, two threshold values are selected to maximize the entropy $H(S) = Hb + Hbw + Hw$. In this system, Hb , Hbw , and Hw are related to regions of text, near-text, and non-text pixels. Let T_1 and T_2 be the two threshold values see Algorithm 2, the classification is defined as follows

$$C(x, y) \in \begin{cases} Hw & \text{if } C(x, y) \leq T_1 \\ Hbw & \text{if } T_1 < C(x, y) < T_2 \\ Hb & \text{if } C(x, y) \geq T_2 \end{cases}$$

Phase 2: The pixels in the near-text areas are relabeled by applying the following method. In the $(2M+1) \times (2M+1)$ window C of the contrast image, $C(0,0) \in Hbw$ is the central pixel, Num is the number of high contrast pixels, namely, pixel intensities between T_1 and T_2 . The mean (I_{Mean}) and the standard deviation (I_{Std}) of these pixels are calculated:

$$I_{Mean} = \frac{\sum_{i=-M}^M \sum_{j=-M}^M I(i, j) \times C(i, j)}{Num}, \quad I_{Std} = \sqrt{\frac{\sum_{i=-M}^M \sum_{j=-M}^M (I(i, j) \times C(i, j) - I_{Mean})^2}{Num}}$$

The pixels can thus be relabeled by equation:

$$I(x, y) = \begin{cases} 0 & \text{if } \begin{cases} Num > 0 & \text{and } C(x, y) \in Hbw \\ \text{and } I(x, y) < \min(I_{Mean} + I_{Std}, T_2) \end{cases} \\ 1 & \text{otherwise} \end{cases}$$

Fig. 4a shows the binarized images which are generated by our ternary entropy-based binarization algorithm. However, like most DI thresholding, our proposed scheme introduces an amount of redundant pixels in the border and around the text (Figs. 4b and 4c).

3) Post-processing

The post-processing operations are introduced and applied on the binarized image in order to eliminate noise, fill the breaks, gaps or holes, and preserve stroke connectivity. According to the experimental results, two kinds of noise are considered in our work. The first type is salt and pepper noise which is defined as (1) white regions smaller than stroke width and usually inside character strokes, (2) black regions smaller than stroke width and settle around document text or scatter over the image as Fig. 4b. The second is called block noise which consists of black regions bigger than stroke width as Fig. 4c. According to the experiments, block noise is often located along the border of images [21]. To deal with the two different kinds of noise, two strategies are applied.

Strategy 1: To deal with salt and pepper. Firstly, the binary MO is applied to bridge unconnected pixels and remove isolated or spurious pixels. Secondly, to correct areas which are bigger than one pixel and smaller than character stroke width, shrink and swell filters are utilized. The filters remove the region noise from the background especially around the text as well as fill the regions in the foreground.

Strategy 2: To deal with block noise, a graph searching strategy is used. Each node in our graph is defined as a window bigger than twice the largest character stroke width. The root of the graph is a complete black node. Edge is detected if nodes are connected by black pixels.

IV. EXPERIMENTAL RESULTS

1) Testing datasets

The performance of our method as well as the comparison are evaluated qualitatively using three testing datasets of historical documents obtained from DIBCO 2009 [5], H-DIBCO 2010 [8], and ICHFR 2010 [21]. The DIBCO 2009 testing database contains five handwritten DI and five printed DI which are associated with ground truth images for the evaluation. The testing images of H-DIBCO 2010 include ten handwritten DI with the associated ground truth from the collections of the Library of Congress [20]. Some representative images of DIBCO 2009 and H-DIBCO 2010 are shown in Fig. 1. The ICHFR 2010 contest is related to quantitative assessment of historical documents image binarization algorithms with bleeding-through noise. Two different techniques for the blending are used, namely, the maximum intensity and the image averaging. ICHFR 2010 contest provides 10 ground truth images and 270 degraded DI [21]. Representative examples of the ICHFR 2010 database are shown in Fig. 2.

2) Evaluation Measures

In this section, criteria used to measure and evaluate a binarization algorithm in the last two contests DIBCO 2009, H-DIBCO 2010 are firstly described. Moreover, to make our comparison and evaluation more general, we consider binarizing as a bi-classification task which is detailed by Sokolova and Lapalme [13]. These criteria are briefly described as follows:

(a) F-Measure (FM): evaluates how well an algorithm can retrieve the desired pixels and is defined as

$$FM = \frac{2 \times pre \times rec}{pre + rec}, \text{ where } pre = \frac{TP}{TP + FP} \text{ and } rec = \frac{TP}{TP + FN}.$$

(b) pseudo F-Measure (p_FM): is based on the skeletonized ground truth image which is defined as $p_FM = \frac{2 \times pre \times p_rec}{pre + p_rec}$,

where p_rec is percentage of the skeletonized ground truth image SG that is detected in the resulting binary image B and defined as:

$$p_rec = \frac{\sum_{i=0}^{M-1} \sum_{j=0}^{N-1} SG(i, j) \times B(i, j)}{\sum_{i=0}^{M-1} \sum_{j=0}^{N-1} SG(i, j)}.$$

(c) Peak Signal to Noise Ratio (PSNR): is a measure of how close is an $M \times N$ image (I) to another (I') and defined as

$$PSNR = 10 \log \left(\frac{C^2}{MSE} \right), \text{ where } C \text{ is a constant that denotes the}$$

difference between foreground and background and is set to 1 in our case study and MSE is the mean square error between I and I' .

(d) Negative Rate Metric (NRM): is based on the pixel-wise mismatches between the ground truth image and prediction. It is defined as $NRM = \frac{1}{2} \left(\frac{FN}{TP + FN} + \frac{FP}{TN + FP} \right)$.

(e) Misclassification Penalty Metric (MPM): is a measure of how well the resulting image representing the contour of ground

truth image and defined as $MPM = \frac{1}{2D} \left(\sum_{i=1}^{FN} d_{FN}^i + \sum_{j=1}^{FP} d_{FP}^j \right)$, where d_{FN}^i and d_{FP}^j denote the distance of the i^{th} false negative and the j^{th}

false positive pixel from the contour of the text in the ground truth image. The factor D is the sum all the pixel-to-contour distances of the ground truth object.

(f) Sensitivity (Sens): proportion of actual positives which are predicted positive: $Sens = \frac{TP}{TP + FN}$.

(g) Specificity (Spec): proportion of actual negative which are predicted negative: $Spec = \frac{TN}{TN + FP}$.

(h) Balanced Classification Rate (BCR): gives balanced assessments on the two classes which have to be adopted such as $BCR = (Sens + Spec) / 2$.

(i) β -measure (β_FM): weighted harmonic mean between sensitivity and specificity $\beta_FM = \frac{2 \times Sens \times Spec}{Sens + Spec}$.

In contrast to FM, p_FM, β_FM , PSNR, BCR, the binarization quality is better for lower NRM, MPM.

3) Results and Comparison

The experimental results are conducted on a PC with an Intel(R) Core i7, 2.67 GHz, and a 3-GB RAM. The operating system is Windows 7 32-bit, and our algorithms were programmed by Matlab 2010. We compare our results with [6][7]. Evaluation of DIBCO 2009 and H-DIBCO 2010 showed that Tan et al.'s scheme outperforms 43 document thresholding algorithms submitted to the DIBCO 2009 and 17 submitted methods to H-DIBCO 2010. The top results of Tan et al.'s algorithm are directly downloaded from [22] and evaluated in the distinct nine measures. Obviously, the algorithm of Tan et al. gives very good quality binarized images when applying to DIBCO 2009 and H-DIBCO 2010 testing databases. The experiments show that our results generally are quite comparable to Tan et al.'s. Although our F-Measure is less (0.10714 on the average), our proposed algorithm is more stable due to the difference between maximum and minimum is smaller than that in Tan's method. Taking DIBCO 2009 as an example, the comparison is presented in Table I.

In comparison with testing databases of DIBCO 2009 and H-DIBCO 2010, the database given by ICFHR 2010 contest is more practical and difficult because of various problems in the historical documents. The performances of our scheme and Tan et al.'s algorithm on some representative images from ICFHR 2010 database are shown in Table II. Herein, the first column shows the results of Tan et al.'s algorithm. Our results are presented in the second column. In each column, the whole page is on the left and some details are provided on the right. As Table III shows, our proposed method achieves much better binarization results with higher scores in FM, p_FM, PSNR, BCR, β -FM while keeping lower NRM, MPM on the testing database of ICFHR 2010 contest.

As further proof of the practicality and effectiveness of our proposed method, the time complexity is taken into account. For

TABLE I. COMPARISON ON DIBCO 2009

Algorithm		FM	p_FM	PSNR	NRM(%)	MPM(%)	Sens	Spec	BCR	β_FM
Tan et al's	Maximum	96.47303	99.44131	24.29	12.79824	0.857	0.93927	0.9992	0.96882	96.79155
	Minimum	84.9803	90.3902	14.14	0.636156	0.06	0.74561	0.98775	0.87202	85.36924
	Average	91.10432	96.854918	18.615	0.8361557	0.363	0.87719	0.995578	0.936385	93.16681
Our	Maximum	96.38722	99.66753	23.48	7.09048	4.367	0.93852	0.99902	0.96685	96.57843
	Minimum	88.44938	89.0464	14.75	3.31456	0.066	0.86903	0.97702	0.9291	92.6082
	Average	90.99718	94.2475	17.86689	5.152341	1.070333	0.90762	0.989333	0.94846	94.6484

TABLE II. RESULTING IMAGES ON ICFHR 2010

B. Su, S. Lu, C. L. Tan's algorithm		Our algorithm	
Whole page	Some details	Whole page	Some details