

SEPARATING TEXT AND BACKGROUND IN DEGRADED DOCUMENT IMAGES – A COMPARISON OF GLOBAL THRESHOLDING TECHNIQUES FOR MULTI-STAGE THRESHOLDING

GRAHAM LEEDHAM¹, SAKET VARMA², ANISH PATANKAR² and VENU GOVINDARAJU²

¹*School of Computer Engineering, Nanyang Technological University,
N4-2A-32 Nanyang Avenue, Singapore 839798
Email: asgledham@ntu.edu.sg*

²*Center of Excellence for Document Analysis and Recognition
UB Commons, 520 Lee Entrance, Suite 202, Amherst, New York 14228-2583, USA
Email: svarma | patankar | govind @cedar.buffalo.edu*

Before any processing of the textual content of a document image can be performed the text must be separated from the background of the image. Several thresholding algorithms have previously been proposed and are widely used in document processing. None have been shown effective at thresholding difficult documents where the background and foreground are non-uniform. In this paper we investigate the use of three global thresholding algorithms (Otsu's, Kapur's entropy and Solihin's quadratic integral ratio (QIR)) as the first stage in a multi-stage thresholding algorithm for use in degraded document images. It is concluded that Otsu's and Kapur's algorithms do not work well for difficult documents as they tend to over-threshold the image, thus losing much of the useful information. The QIR algorithm is more accurate in separating the foreground and background in these images, leaving a range of undecided, fuzzy, pixels for later processing in a subsequent stage.

1. Introduction

To process an image printed or written on any surface it is first necessary to separate the image pixels from the background pixels. In a document image this usually involves separating the pixels forming the printed text or diagrams (foreground) from the pixels representing the blank paper (background). The goal is to remove only the background, by setting it to white, and leave the foreground image unchanged. A very common way to perform this is to threshold the image.

The image after thresholding should:

1. retain all details of the originally created document including faint skate-on and skate-off penstrokes at

the beginning and end of handwritten pen-strokes in the foreground,

2. retain handwriting produced by a variety of writing instruments, and
3. remove the background image which may contain strong colours and/or patterns including extraneous markings and artefacts introduced before or after the original writing.

In this paper we investigate the use of three global thresholding algorithms for the first stage in a multistage thresholding method for separating the handwriting from the background in aged or otherwise degraded documents.

2. Image thresholding

2.1 Global thresholding

Numerous global image thresholding techniques have been proposed in the literature. Among those techniques, Otsu's [1] thresholding technique has been cited as an efficient and frequently used technique [2] followed, in order, by Kapur et al.'s Entropy technique [3], Abutaleb's entropy technique [4] and Kittler and Illingworth's minimum error technique [5]. The Integral Ratio method for global thresholding [6] has been shown better than Otsu's or the Entropy algorithm when thresholding images with high contrast background.

While the global thresholding approach is usually considered a one stage thresholding approach, the Integral Ratio method is a first stage in a multi-stage thresholding approach. In the first stage, the image is divided into three subimages (instead of two): foreground, background, and a fuzzy subimage where it is hard to determine whether a pixel actually belongs to the foreground or the background. Figure 1 shows the three subimages within

the intensity histogram. Two important parameters that separate the subimages are A , which separates the foreground and the fuzzy subimage, and C , which separate the fuzzy and the background subimage.

If a pixel's intensity is less than or equal to A , the pixel belongs to the foreground and if a pixel's intensity is greater than or equal to C , the pixel belongs to the background. If, however, a pixel has an intensity greater than A but less than C , it belongs to the fuzzy subimage and more information from the image is needed to decide whether it actually belongs to the foreground or the background.

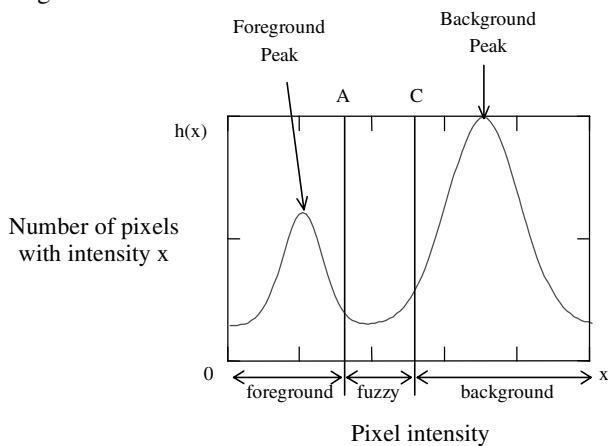


Figure 1: Integral Ratio (a two stage) thresholding approach

2.2 Local or adaptive thresholding

There have been several studies of applying local or adaptive thresholding algorithms to document images. One of the most recent studies [7] considered the problem of poor quality document images. The method is demonstrated to be effective for documents which are highly deteriorated because of variable illumination, shadow, smears and smudges. It is not likely to work well with handwriting or the italic style knib found in old documents.

2.3 Multi-stage thresholding

Multi-stage thresholding can be viewed as a process of reducing the search space of threshold candidate values stage by stage where each stage uses different information from the image until the final stage chooses the final threshold value.

A multi-stage thresholding algorithm is one that performs thresholding in n stages ($n > 1$). In each stage k , it produces

three subimages: one that is categorised as the foreground (F_k), one that is categorised as the background (B_k), and one that it is unable to categorise (D_k). Its input is the non-categorised subimage from the previous stage (D_{k-1}) and some specific information from the image. It passes the non-categorised subimage (D_k) to the next stage for further thresholding. This process continues and ends at stage n where D_n is empty, and all pixels are categorised as either the foreground or the background.

The multi-stage thresholding approach should not be confused with a traditional iterative technique. Although a traditional iterative technique produces a threshold value during each iteration, which is supposedly better than the previous threshold value, there are two main differences with a multi-stage technique:

1. Traditional iterative thresholding uses the same method of selecting a threshold value in each iteration, but multi-stage thresholding may use a different method for each stage using different information from the image.
2. The iterative technique produces only one threshold value for each iteration, while each stage in an n -stage thresholding technique may produce more than one threshold value.

Because traditional thresholding finds a threshold value in one stage, it can be considered as a special case of the multi-stage thresholding approach (i.e. one stage thresholding). In this respect, the multi-stage thresholding approach is a generalisation of traditional thresholding.

A unique characteristic of the multi-stage approach is that it does not follow the traditional global or local thresholding technique categorisation. Each stage may use either a global or a local thresholding technique. If global thresholding is used, the stage uses global threshold values to separate the images into the foreground, background, and the *don't know* subimages. If local thresholding is used, the stage uses a threshold value for each pixel to separate the image.

We are interested whether traditional global thresholding algorithms can be applied as part of a multi-stage thresholding approach.

3. Applying global thresholding to historical documents – an initial study

Old document images become degraded over time due to their storage and deterioration from moisture, infestation and many other factors. In addition, the conditions and materials used in their original production

were highly variable in consistency often resulting in a document in which the foreground and background are difficult to separate. The problem is particularly difficult because many of the document images have varying contrast, smudges, variable background intensity and the presence of the handwriting from the other side of the page (document) as shown in Figure 2.

Three global thresholding methods were evaluated for the initial segmentation of the foreground and background in grey-scale images of degraded documents. These were intended for the first stage in a multi-stage thresholding algorithm. The three methods were:

1. Otsu's method [1]
2. Entropy [3]
3. QIR – Quadratic Integral Ratio [6]
4. NIR – Native Integral Ratio [6]

A set of 20 degraded document images, similar to the one shown in Figure 2, were used for testing these algorithms and the results were compared by visual inspection. The visual inspection looked for errors in separation of background and foreground. A small sample was adequate to illustrate the general trend of the algorithm performance.

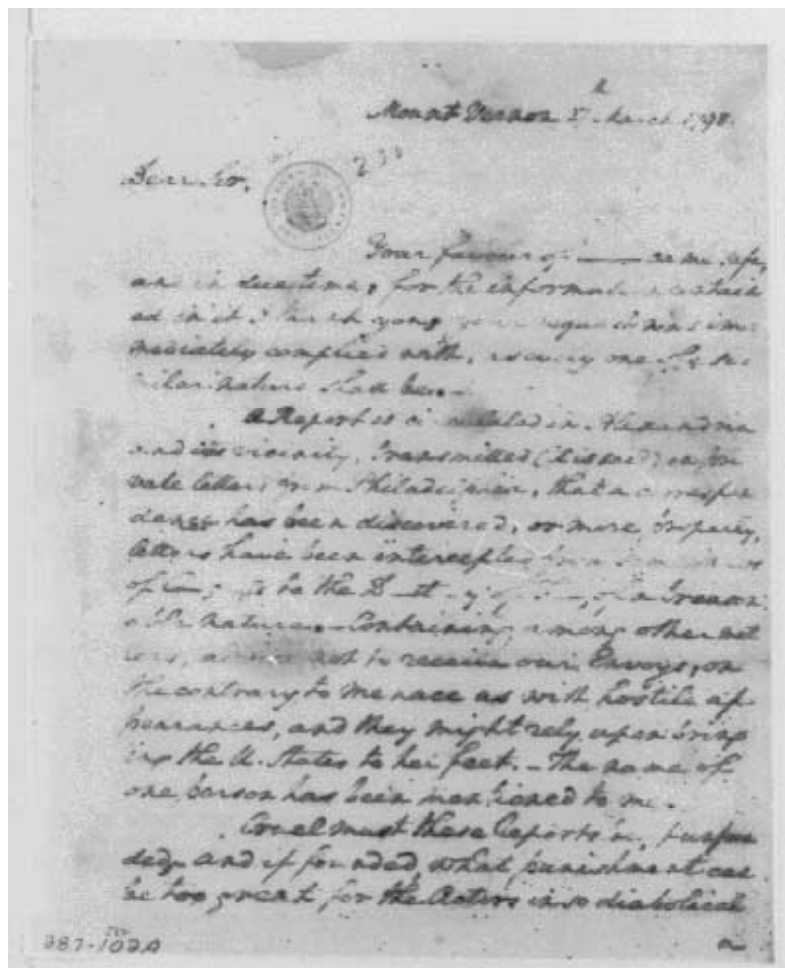


Figure 2. Example of an historic document written with quill pen and suffering from age degradation.

4 Results and Discussion

Based on visual inspection of the original and resulting image of each segmentation technique the following observations were made:

Otsu's algorithm - Always over-thresholded the image resulting in the foreground being lost. Generally, it does not produce good results when the background intensity is high and there is low contrast in the image. Figure 3(a) shows how this image of Figure 1 is very overthresholded such that very little of the original image remains.

Entropy - Performed well if the document images had good contrast. It achieved good smudge removal but tended to over threshold. An example is shown in Figure 3(b).

QIR - (Quadratic Integral ratio) – This algorithm performed well as it generally was able to separate definite foreground (dark) pixels and definite (background pixels) as shown in Figures 3(c) and 3(e) respectively. The uncertain or fuzzy pixels were clearly defined (see Figure 3(d) and required further processing to determine appropriate assignment to background or foreground. Better results were obtained for images with these characteristics:

- Contrast is not Good
- Foreground and background are non-uniform
- Can handle both dark and light foreground and background.
- Better results for document images where writing from the reverse side of the paper is visible through the paper.

NIR- (Native Integral Ratio) – Performed well if there was good contrast, uniform background and foreground.

Figure 3 shows the result of each algorithm on the original image given in Figure 2. Of all the global thresholding techniques studied, QIR performed the best. It can accurately separate background from foreground and identify uncertain pixels which lie in the fuzzy region. Subsequent thresholding stages can be used to separate the fuzzy pixels using a local adaptive thresholding technique.

The detailed greyscale threshold values can be seen in Table 1. The image depicted in Figures 2 and 3 is image number 4. As can be seen, only four of the images could be adequately segmented using a single global threshold (column 2 in Table 1). Of these (image numbers

2, 3, 7 and 9) the QIR, NIR and Entropy algorithms performed well. The remaining images had too much variation and contrast to be accurately. But the result of the QIR algorithm was generally good in that the fuzzy range between thresholds A and C were accurate.

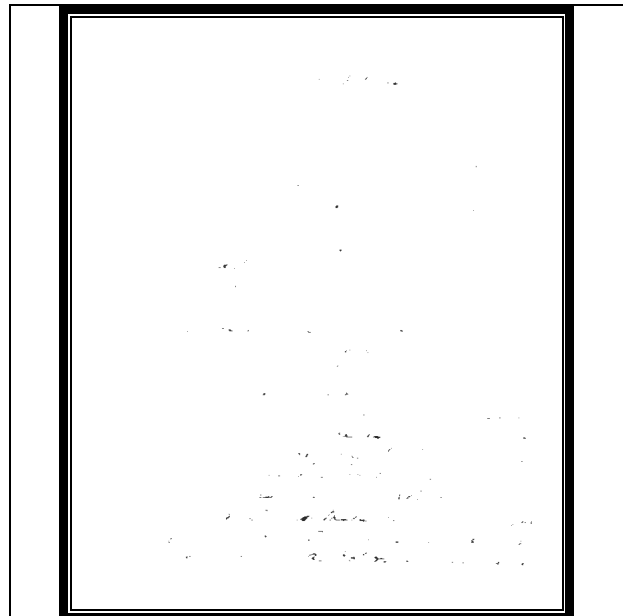


Figure 3(a) the result of OTSU's algorithm on the image in Figure 2

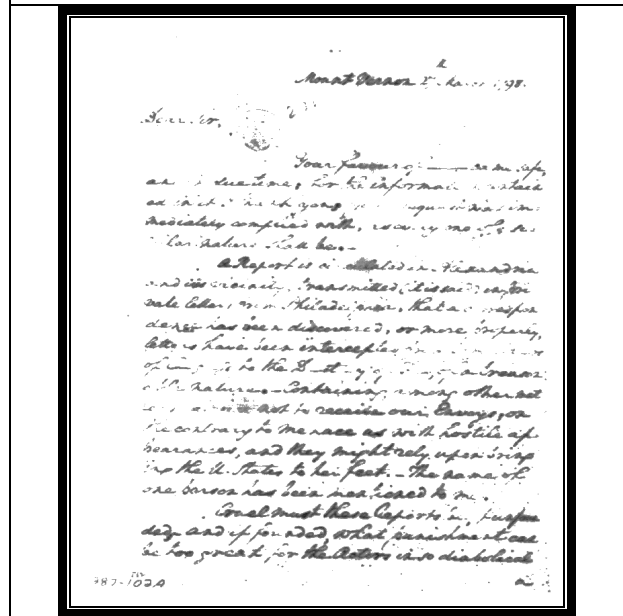


Figure 3(b) the result of Kapur's Entropy algorithm

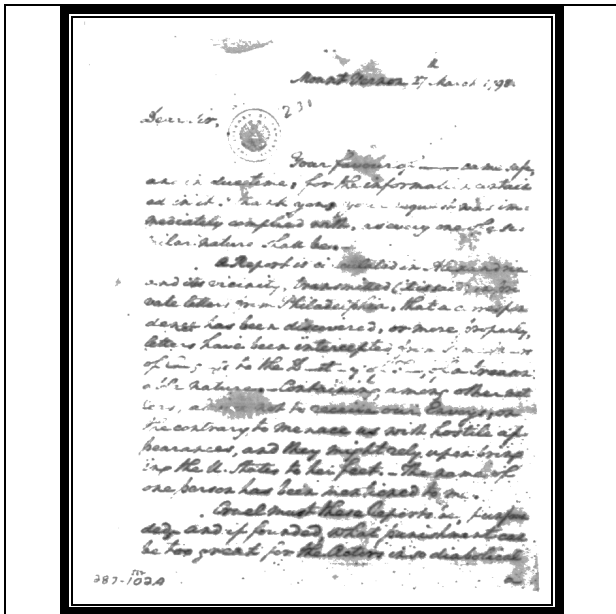


Figure 3(c) The result of the QIR algorithm – the Foreground pixels (threshold less than A).

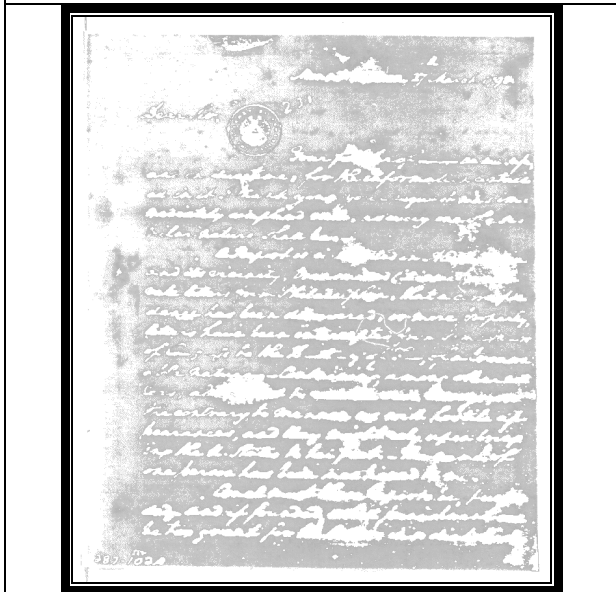


Figure 3(d) The result of the QIR algorithm - Fuzzy pixels (thresholds between A and C)

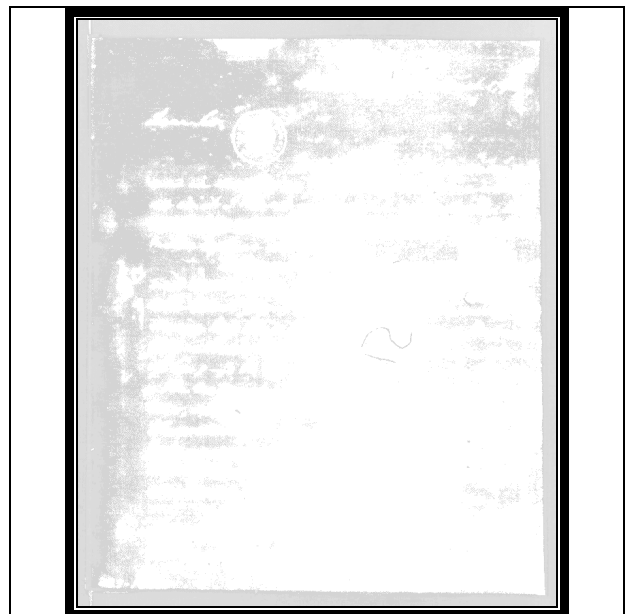


Figure 3(e) The result of the QIR algorithm - Background pixels (threshold greater than C)

Figure 3. Examples of the image shown in Figure 1 processed by various thresholding algorithms. (a) and (b) show the results of the Entropy and Qtsu algorithms on the original image shown in Figure 1. (c) (d) and (e) shown the three images resulting from the QIR algorithm.

5 Conclusion

No single global thresholding algorithm/scheme can work with degraded document images due to the peculiar characteristics of the image like smudges, varying background/foreground intensities, varying contrast and presence of the handwriting from the other side of the page. Local adaptive thresholding algorithms have not proven effective with this type of image. The alternative, currently under investigation, is the use of multi-stage thresholding.

Because traditional thresholding finds a threshold value in one stage, it can be thought of as a special case of the multi-stage thresholding approach (i.e. one stage thresholding). In this respect, the multi-stage thresholding approach is a generalisation of traditional thresholding.

From the results reported in this paper it is apparent that, while the widely used Entropy and Otsu algorithms work adequately with document images which consist of high contrast images with clear white background, they

do not work well with degraded images. They are therefore not suitable for use in a multi-stage thresholding approach. The QIR algorithm, generally performs well in separation of clearly background and foreground pixels.

Current research is investigating the use of the QIR algorithm as one stage in a multi-stage thresholding strategy with refinements to resolve the fuzzy pixels.

Acknowledgement

The authors would like to thank Yan Solihin for providing the source code to evaluate the QIR and NIR algorithms.

References

1. N. Otsu, A Threshold Selection Method from Gray-Level Histogram, *IEEE Trans. Systems, Man, and Cybernetics*, vol. 9, pp. 62-66, 1979.
2. Trier O.D. and Jain A.K. Goal-Directed Evaluation of Binarization Methods, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 17, no. 12, pp. 1191-1201, 1995.
3. J.N. Kapur, P.K. Sahoo, and A.K.C. Wong, A New Method for Gray-Level Picture Thresholding Using the Entropy of the Histogram, *Computer Vision, Graphics, and Image Processing*, vol. 29, pp.273-285, 1985.
4. Abutaleb, A.S. Automatic Thresholding of Gray-Level Pictures Using Two Dimensional Entropy, *Computer Vision, Graphics, and Image Processing*, vol. 47, pp. 22-32, 1989.
5. Kittler J. and Illingworth J. Minimum Error Thresholding, *Pattern Recognition*, vol. 19, no.1, pp. 41-47, 1986.
6. Yan Solihin and C.G. Leedham, Integral Ratio: A New Class of Global Thresholding Techniques for Handwriting Images, *IEEE Trans on Pattern Analysis and Machine Intelligence*, Vol. 21, No., pp.761-768, 1999.
7. Yibing Yang and Hong Yan, An adaptive logical method for binarization of degraded document images, *Pattern Recognition*, Vol 33, pp 787-807, 2000.

Image Number	Best threshold from visual inspection		Thresholds obtained by four global algorithms					Was the threshold result acceptable? If it was, which one?
	Global threshold possible? (Yes/No)	Best achievable global threshold (Grey scale value)	QIR		ENT	OTSU	NIR	
			A	C				
1	No	na	196	229	173	86	188	No
2	Yes	173	173	205	173	70	184	QIR/ENT
3	Yes	117	106	129	117	93	117	QIR/NIR/ENT
4	No	na	168	199	158	70	101	QIR
5	No	na	181	212	179	72	179	QIR/ENT
6	No	na	202	237	172	76	124	QIR
7	Yes	187	187	221	181	70	181	QIR/NIR/ENT
8	No	na	181	213	175	100	112	QIR
9	Yes	179	181	212	179	72	179	QIR/NIR/ENT
10	No	na	196	229	180	85	189	No
11	No	na	189	221	108	78	117	QIR
12	No	na	168	198	163	68	180	QIR
13	No	na	183	215	187	35	181	ENT
14	No	na	167	197	69	68	101	QIR
15	No	na	197	230	171	100	221	No
16	No	an	195	228	165	93	115	No
17	No	na	187	209	164	78	189	ENT
18	No	na	196	229	185	53	202	QIR
19	No	na	194	229	172	78	123	No
20	No	na	172	202	71	67	107	No

QIR = Quadratic Integral Ratio method, ENT = Entropy method, OTSU = Otsu's method, NIR = Native Integral Ratio method.

Table 1. Results of the thresholding algorithms on 20 typical historical document images.