# Improving the Quality of Degraded Document Images

Ergina Kavallieratou and Efstathios Stamatatos

*Dept. of Information and Communication Systems Engineering.*
*University of the Aegean*
*83200 – Karlovassi, Greece*
*{ergina,stamatatos}@aegean.gr*

## Abstract

*It is common for libraries to provide public access to historical and ancient document image collections. It is common for such document images to require specialized processing in order to remove background noise and become more legible. In this paper, we propose a hybrid binarizatin approach for improving the quality of old documents using a combination of global and local thresholding. First, a global thresholding technique specifically designed for old document images is applied to the entire image. Then, the image areas that still contain background noise are detected and the same technique is re-applied to each area separately. Hence, we achieve better adaptability of the algorithm in cases where various kinds of noise coexist in different areas of the same image while avoiding the computational and time cost of applying a local thresholding in the entire image. Evaluation results based on a collection of historical document images indicate that the proposed approach is effective in removing background noise and improving the quality of degraded documents while documents already in good condition are not affected.*
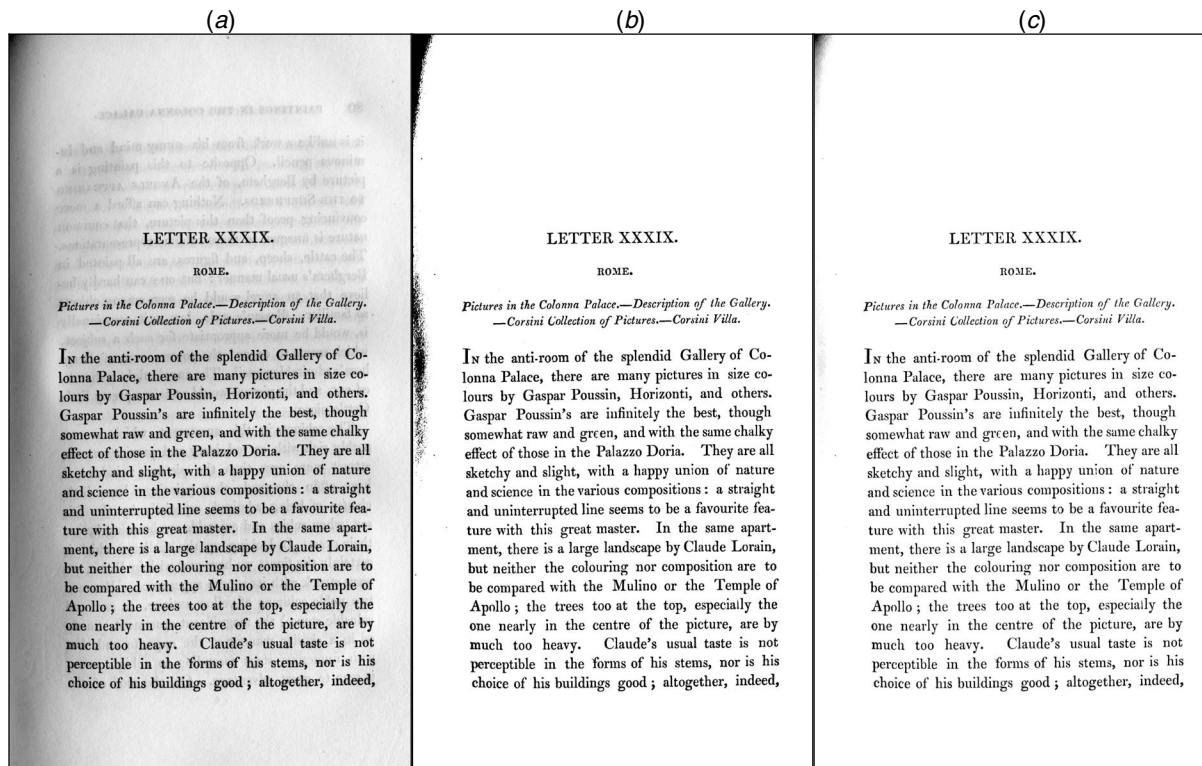
## 1. Introduction

Historical and ancient document collections available in libraries throughout the world are of great cultural and scientific importance [1-2]. The transformation of such documents into digital form is essential for maintaining the quality of the originals while provide scholars with full access to that information [3]. It is quite common for such documents to suffer from degradation problems [4]. Just to name a few, presence of smear, strains, background of big variations and uneven illumination, seepage of ink etc. are factors that impede (in many cases may disable) the legibility of the documents. Therefore, appropriate filtering methods should be developed in order to remove noise from historical document images and improve their quality before libraries expose them to public view. Within this framework, noise is considered anything that is irrelevant with the textual information (i.e., foreground) of the document image.

Image analysis systems use *binarization* as a standard procedure to convert a grey-scale image to binary form. An ideal binarization algorithm would be able to perfectly discriminate foreground from background, thus, removing any kind of noise that obstructs the legibility of the document image. The binary image is ideal for further processing [5] (e.g., discrimination of printed from handwritten text, recognition of the contents by applying OCR techniques etc). However, in the framework of a library collection of historical and ancient documents intended to be exposed to public view, the document images in many cases do not need further processing apart from removing the background noise and leave some "traces of time" behind. More importantly, given such a case, after the removal of background noise, it is possible for the document images to remain in grey-scale form. For instance, consider the images of figure 1. Figure 1*a* shows the original image, figure 1*b* the result of the binarization procedure, and figure 1*c* the corresponding grey-scale result after removing the background noise. In the latter case the remaining noise and the text characters are smoothed a little bit. This has the consequence of making the background noise practically invisible while the characters are more legible. Some binarization techniques support this option [6].

Traditional binarization approaches can be divided into two main categories:

1. Global thresholding methods: The pixels of the image are classified into text or background according to a global threshold. Usually, such methods are simple and fast. On the other hand, they cannot be easily adapted in case the background noise is unevenly distributed in the entire image (e.g., smear or strains) [7-8].

**Figure 1. Removal of background noise from historical document images: *a)* the original document image, *b)* the binarized image, *c)* the result in grey-scale form.**

2. Local thresholding methods: The pixels of the image are classified into text or background according to a local threshold determined by their neighboring pixels. Such methods are more adaptive and can deal with different kinds of noise existing in one image. On the other hand, they are significantly more time-consuming and computationally expensive [9-10].

From another point of view, binarization approaches can be divided as follows:

1. General-purpose methods: They are able to deal with any image. Therefore, they do not take into account specific characteristics of document images.

2. Document image-specific methods: They attempt to take advantage of document image characteristics (e.g., background pixels is the majority, foreground pixels are in similar grey-scale tones etc). In many cases, such methods are variations of general-purpose approaches [9].

Although it is reasonable the latter approaches should be more effective when dealing with historical document images, recent results show that general-purpose methods can be more reliable under certain conditions [11].

In previous work, we have presented an *Iterative Global Thresholding* (IGT) approach that is specifically designed for document images [6]. Apart from efficiency inherent in any global thresholding approach, this method has the additional advantage of providing the option to maintain the image in grey-scale after the removal of background noise, a more familiar form for human readers. In this paper, we propose a hybrid approach that attempts to combine the advantages of global and local thresholding. The main idea is that after the application of IGT to the document image, the areas that are more likely to still include significant amount of noise are selected and, then, IGT is re-applied to these areas separately. Hence, document images containing different kinds of background noise, unevenly distributed in the entire image, can be processed more effectively. Additionally, the time-cost of the approach remains on low level in comparison to original local thresholding techniques, since only a limited number of areas (instead of the entire image) need to be processed separately. To evaluate the proposed approach, we attempt to represent the degree in which the legibility of the image has been improved. To this end, a series of experiments based on both human user opinions as well as optical character recognition software is presented.

The rest of this paper is organized as follows: Section 2 includes related work and section 3 describes

our approach in detail. Evaluation results are presented in section 4 while the main conclusions drawn from this study and future work directions are summarized in section 5.

## 2. Related Work

Many methods have been proposed in local or global thresholding. One of the earlier methods in image binarization was proposed by Otsu [7] based on the variance of pixel intensity. Bernsen [10] calculates local thresholds using neighbours. Niblack [12] uses local mean and standard deviation. Sauvola [9] presents a method specialized on document images that applies two algorithms in order to calculate a different threshold for each pixel. As far as the problem of historical documents is concerned, Leedham [8] compares some of the traditional methods on degraded document images while Gatos [13] proposes a method using a combination of existing techniques. These are also the cases of Shi [14] and Yan [15] applied to some historical documents from the US library of Congress. Leydier [16] works with colored document images and implements a serialization of the *k*-means algorithm. Some of the above methods, apart from the basic binarization algorithm, have also used pre-processing or post-processing filters for improving the quality of the document image [13].

In our previous work, we presented a method for cleaning and enhancing historical document images [6]. Hereafter, this method will be called *Iterative Global Thresholding* (IGT). This method is both simple and effective. It selects a global threshold for a document image based on an iterative procedure. In each iteration, the following steps are performed:

(i)   The average pixel value is calculated.
(ii)  The average pixel value is subtracted from each pixel of the image.
(iii) The histogram is stretched so that the remaining pixels to be distributed in all the grey scale tones.

During the i-th iteration, document image $I_i(x, y)$ will be:

$$I_i(x, y) = 1 - \frac{T_i - I_{i-1}(x, y)}{1 - E_i} \qquad (1)$$

where $I_{i-1}(x, y)$ is the document image resulted in from the previous iteration ($I_0(x, y)$ is the original image), $T_i$ is the threshold calculated in the *i*-th iteration and $E_i$ is the minimum pixel value in the *i*-th repetition before the histogram stretching. After each iteration, an amount of pixels is moved from the foreground to the background. The iterations stop based on the following criterion:

$$|T_i - T_{i-1}| < 0.001 \qquad (2)$$

This approach works well on historical document images given that the foreground tone is darker than the background. It takes into account that the background corresponds to the great majority of the image pixels. Moreover, the foreground (text) will be roughly of the same grey scale tone and darker than the background. As a global thresholding method has relatively low time cost and it does not require complicated calculations. More importantly, it supports applications where the noise-free image should remain in grey-scale form.

On the other hand, there are some cases of degraded document images, IGT (and most of the existing methods) are unable to handle. First, in case there are stains or scratches of similar grey scale tone with that of the foreground, it is not possible to remove it without losing useful information. Second, in case the foreground is written in more than one main tones (e.g., presence of both printed and handwritten text) it is likely the lighter tone to be significantly attenuated (or even be removed). Unfortunately, such cases are not uncommon in historical document images. In this paper, we show how this method can be improved by separately processing areas where noise still remains.
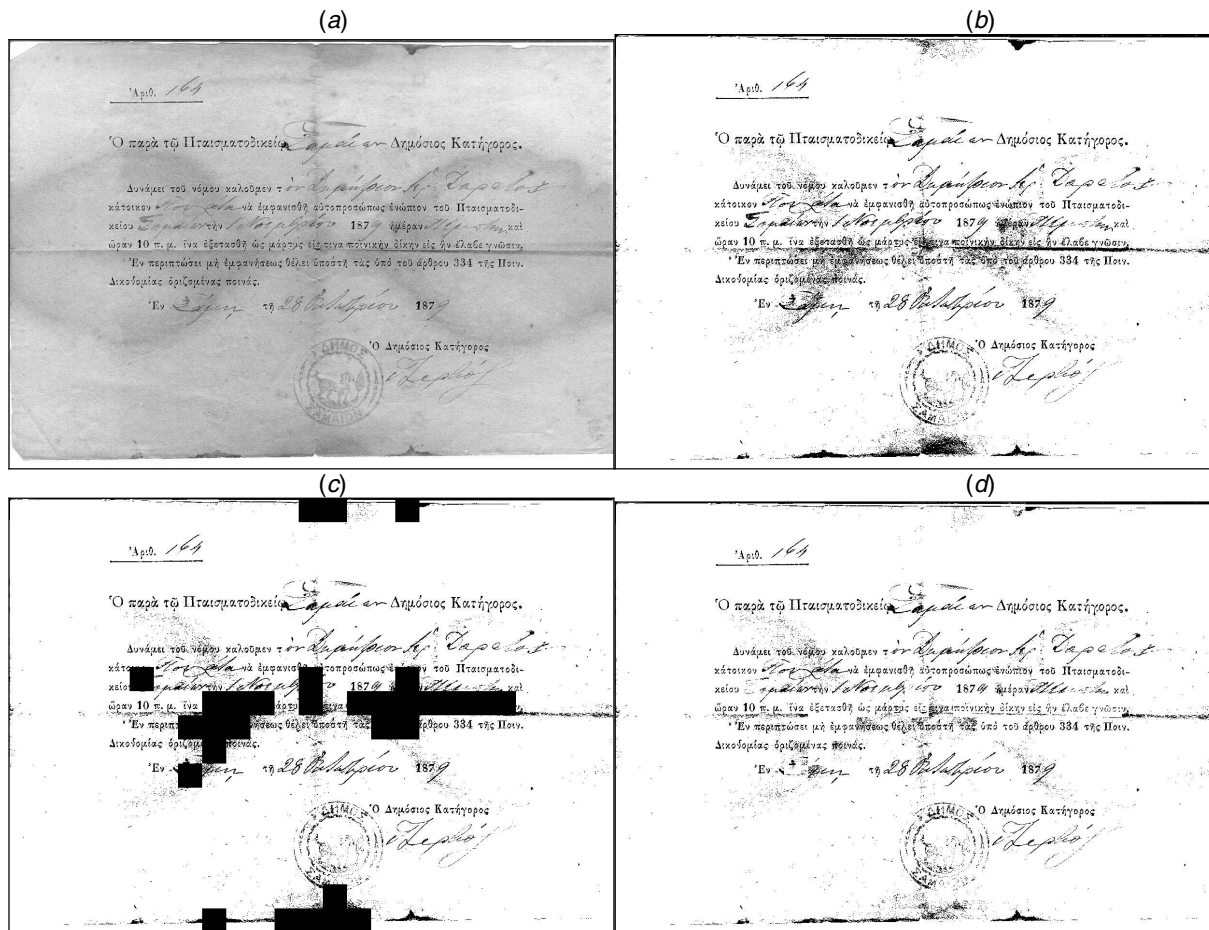
## 3. The Proposed Approach

We propose a hybrid approach for improving the quality of historical document images, that is, a combination of global and local thresholding. First, a global thresholding approach (IGT), is applied to the document image. Then, the areas that still contain noise are detected and re-processed separately. In more detail, the proposed algorithm consists of the following steps:

(i)   Apply IGT to the document image.
(ii)  Detect the areas in which it is more likely background noise to still remain.
(iii) Re-apply IGT to each detected area separately.

Figure 2 shows an example of the binarization of a historical image following the proposed approach. Figure 2*a*, 2*b*, 2*c*, and 2*d* show the original grey scale image, the result of applying IGT to the whole image, the detected areas with remaining noise, and the final result after applying IGT to each detected area, respectively. Note that the document image of figure 2 is a hard case since many kinds of noise coexist in the same image (uneven illumination, stains, and page crumples). Moreover, mixed text (both printed and handwritten) as well as stamps are additional obstacles for the noise removal procedure.

By selecting only specific areas of the image for processing based on local thresholding, we avoid the cost of applying local thresholding to the entire image.

IEEE
COMPUTER
SOCIETY

**Figure 2. Application of the proposed approach to a historical document image: *a*) the original image in grey-scale, *b*) the result of applying IGT to the entire image, *c*) the detected areas of remaining noise in black color, and *d*) the result after applying IGT to each area separately.**

In the next subsections we describe the procedures of detecting the areas with remaining noise and re-applying the IGT algorithm to these areas.
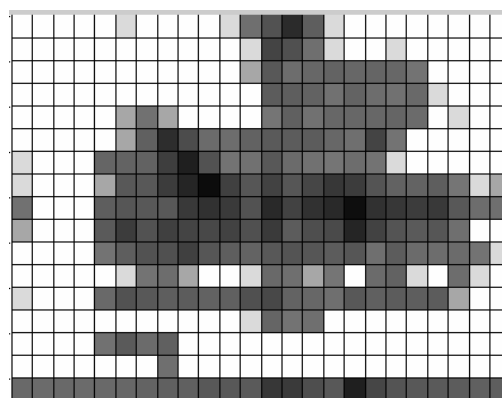
## 3.1 Detection of Areas with Remaining Noise

The detection of areas that need further processing is performed by using a simple method. The main idea is based on the fact that the areas that still contain background noise will include more black pixels on average in comparison with other areas. This is reasonable especially for document images that only include textual information.

The image is divided into segments of fixed size $n$x$n$. In each segment, the frequency of black pixels is calculated. The segments that satisfy the following criterion are, then, selected as:
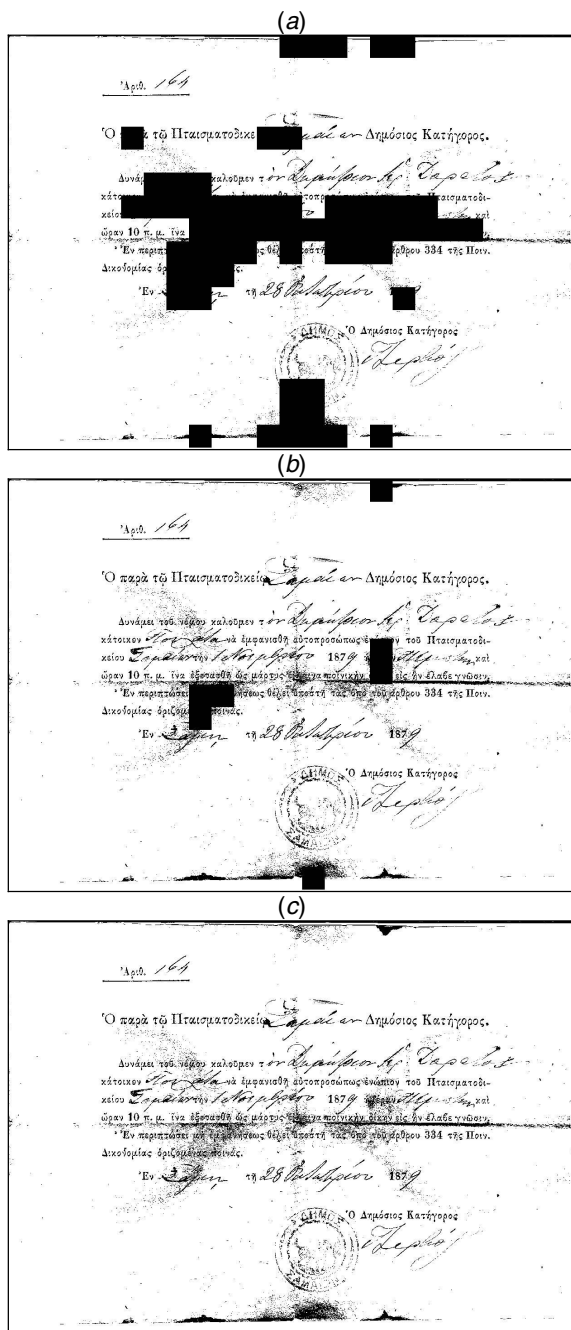
$$f(S) > m + ks \tag{3}$$

where $f(S)$ is the frequency of the black pixels in the segment $S$ while $m$ and $s$ are the mean and the standard deviation of the black pixel frequency considering the segments of the entire page, respectively. The selected segments form areas by connecting neighboring



**Figure 3. Detection of areas with remaining noise. Frequency of black pixels of the image of figure 2*b* (window size = 50x50, *k*=2). Darker fill color corresponds to higher average frequency of black pixels.**

segments in respect to their original position in the image. The row-by-row labeling algorithm is used [17]. Figure 3 depicts the distribution of black pixels in the document image of figure 2*b*.

(*a*)



(*b*)



(*c*)



**Figure 4. Detection of areas (indicated with black color) based on different values of parameter *k*: *a*) *k*=1, *b*) *k*=3, and *c*) *k*=5. In all cases window size is 50x50.**

The parameter *k* determines the sensitivity of the detection method. The higher the *k*, the less segments will be detected. This means that some of the areas that

may still need further improvement will not be selected. On the other hand, a low *k* guarantees that all the areas that still need improvement will be selected together with other areas in which the noise has been already removed. Moreover, the computational and time cost of applying IGT to more areas will be increased. Therefore, an appropriate value of *k* should be selected to deal with this trade off. Figures 4*a*, 4*b*, and 4*c* show the selected areas of a document image using *k*=1, 3, and 5, respectively. Note that for *k*=5 no area is selected.
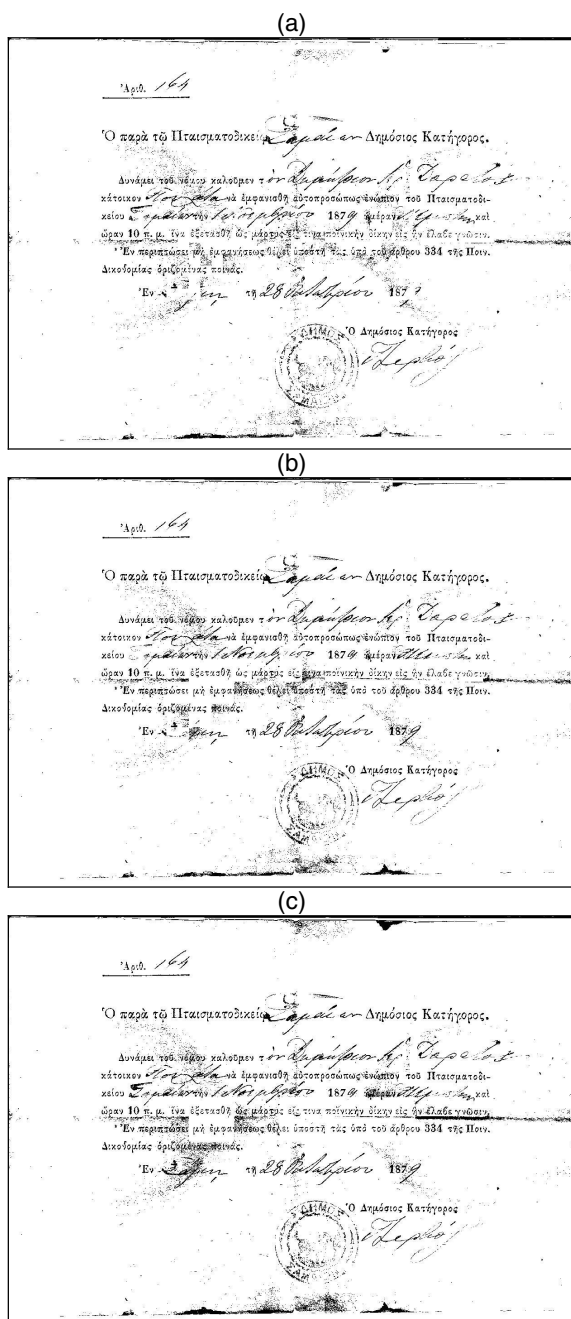
## 3.2 Local Thresholding on Selected Areas

The areas detected by the previously described procedure are separately re-processed based on local thresholding. That is, for a given area, the IGT method is applied to the corresponding area of the original image. The iterations stop when either the criterion of formula (2) is satisfied or the number of iterations exceeds the corresponding number of iterations required for the global thresholding on the entire image (from the first step of the proposed approach).

Given that the selected areas have relatively high average density of black pixels, IGT removes a lot of pixels during the first iterations. In comparison to the application of IGT to the entire page, the background noise in the selected areas is more likely to be removed since the area is likely to be more homogeneous than the entire image. In general, this procedure tends to move more pixels of the selected areas to the background in comparison with the previous application of IGT to the entire image. However, in case a selected area does not contain considerable amount of background noise, the foreground is not attenuated unless it consists of different grey-scale tones (e.g., presence of both printed and handwritten in the same area).

An important factor for the successful re-application of IGT in selected areas is the size of the window *n*x*n* used in the procedure of selecting the appropriate areas (described in the previous subsection). Figures 5*a*, 5*b*, and 5*c* depict the results of applying the proposed approach on a document image based on different window sizes (*n*=25, 50, and 100, respectively). In all cases, *k* was set to 2. As can be seen, a small window size forms more but smaller areas. On the one hand, there is the advantage of adapting area in more detail at the part of the image that still contains noise. On the other hand, the resulting areas are small and in many cases they do not provide enough information for successfully re-apply the IGT algorithm on them. In case of large window size, fewer but bigger areas are detected. On the one hand, they provide enough

information to the IGT algorithm in order to effectively remove background noise. On the other hand, these areas cannot be easily adapted to a specific part of the image that still contains noise. As a consequence, the final image may contain neighboring areas that have dissimilar amount of background noise.

(a)



(b)



(c)



**Figure 5. Results of re-applying IGT to the selected areas using different window size *nxn*: a) *n*=25, b) *n*=50, and c) *n*=100. In all cases. *k*=2.**

# 4. Evaluation

## 4.1 Document Collection

In order to evaluate the proposed approach, we collected an amount of historical document images. In more detail, we formed a collection of 183 document images taken from Korgialenios Library of Cephalonia (KLC) and the Daratos Private Archive (DPA) of old documents including the following:

- 129 document images taken from books of the 19th century (source: KLC). These books were written in English, Italian, or Greek.
- 19 document images taken from Greek newspapers of the early 20th century (source: KLC).
- 21 musical scores of handwritten Byzantine music notation (source: KLC).
- 14 document images of handwritten Greek text of the 18th century (source: DPA). The majority of these images were in especially bad condition.

**Table 1. Details about the condition of the document collection used in this study (total number of document images is 183).**

| | |
|---|---|
| In good condition | 7.65% |
| **Degradation problems** | |
| Uneven illumination | 55.19% |
| Holes on the page | 7.65% |
| Seepage ink | 42.07% |
| Stains, smearing | 22.95% |
| Page crumple | 25.14% |
| **Other problems** | |
| Mixed text (print./ handwrit.) | 10.38% |
| Presence of page lines | 3.83% |

Based on this collection, we attempted to cover a wide range of document types, that is, both printed and handwritten, variable degree of degradation as well as diversity on the form of included information (different languages, musical notation, etc.). Therefore, a representative amount of degradation problems is included in this collection, similar in any historical document collection. Table 1 shows details about the condition of the document images used in this study.

As can be seen, only a few documents are considered to be in good condition (e.g., homogeneous background noise not affecting the readability of the document). These documents considerably assist the evaluation of a noise removal algorithm since they should not be affected drastically. The rest of the documents have at least one problem due to either degradation or the content of the document. The most common problems are uneven background illumination and seepage of ink.
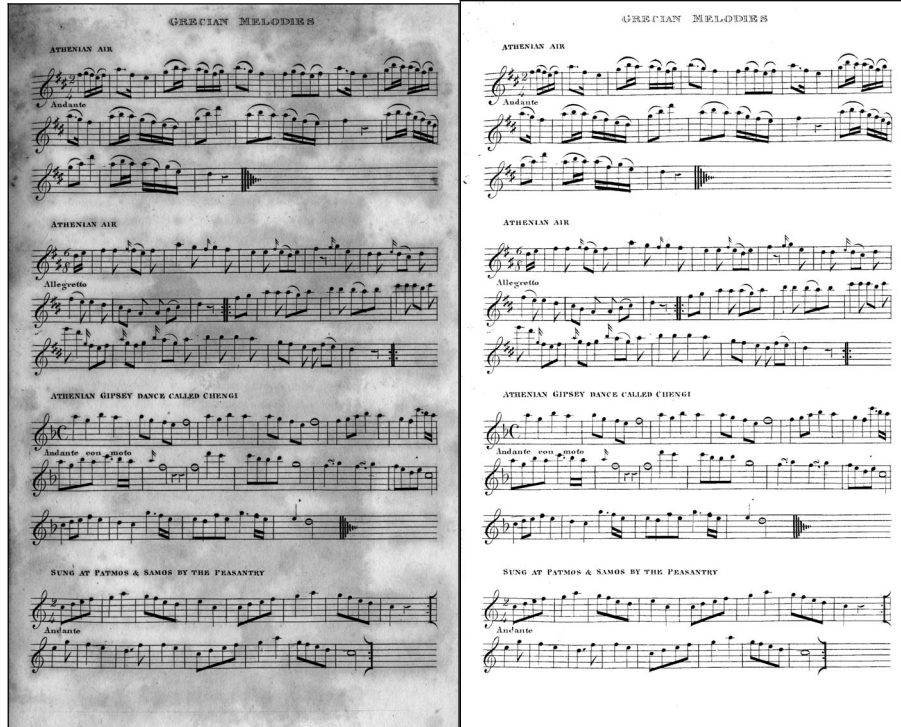
**Figure 6. A document image before and after the application of the IGT algorithm in the whole image with no remaining noise.**
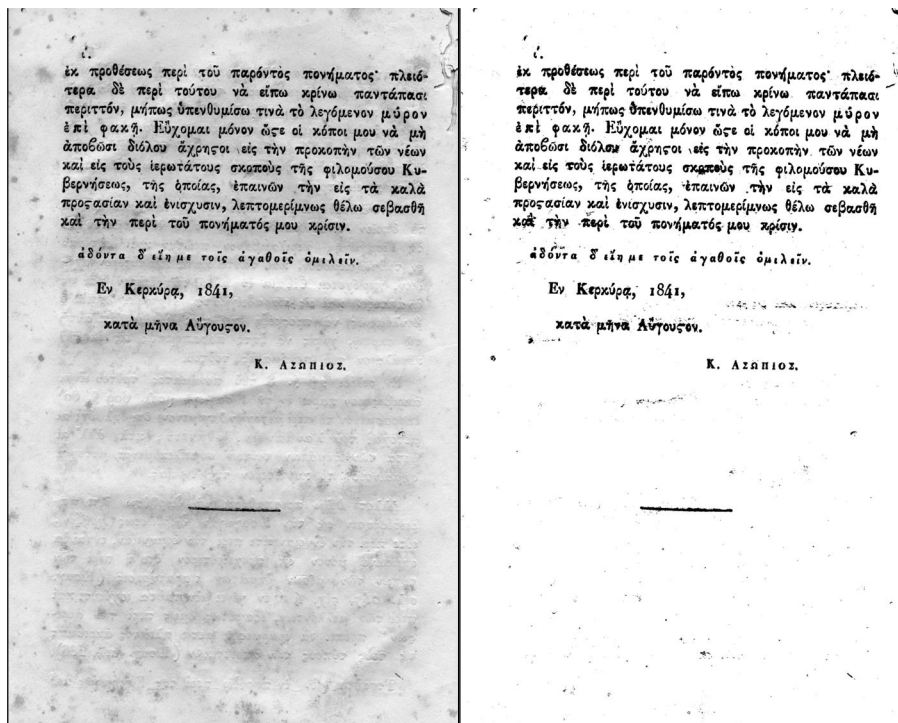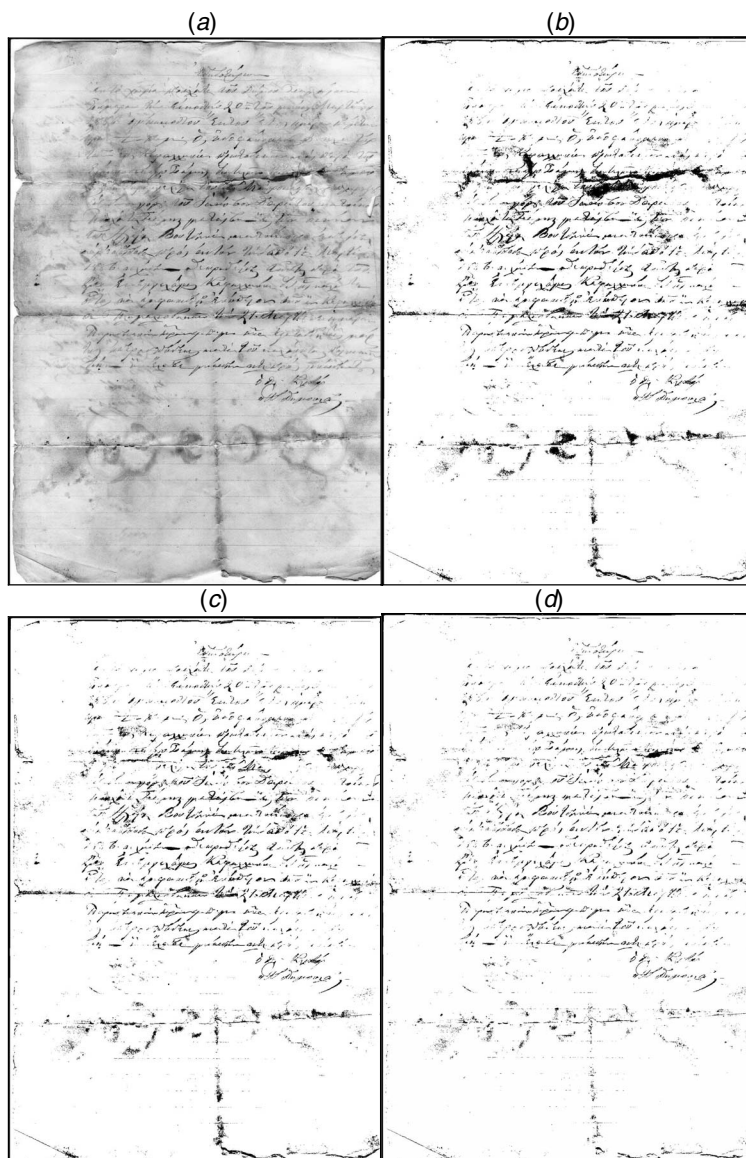


**Figure 7. A document image before (*a*) and after (*b*) the application of the IGT algorithm in the whole image with small amount of remaining noise.**

**Figure 8. A document image before (*a*) and after (*b*) the application of the IGT algorithm in the whole image with considerable amount of remaining noise. In addition, the result of applying the proposed hybrid approach with different settings: c) *k*=2 and *n*=50, d) *k*=2 and *n*=25.**

## 4.2 Results

The IGT algorithm as described in [6] was first applied to all the document images of the collection (including those considered to be in good condition). The examination of the resulting images proved that this process increased the percentage of the document images considered to be in good condition. Now, 23.50% of the documents do not need further improvement. The rest of the new document images still contain some level of background noise. Figures 6, 7, and 8 show examples of documents with removed noise, documents with small amount of remaining noise, and documents with considerable amount of remaining noise.

We, then, applied the proposed approach in order to further improve the quality of the produced images. The detection of areas that still contain background noise was performed using window size of different size (*n*=25, 50, and 100) and parameter *k* value (*k*=1, 2, 3, 4, and 5). Finally, the IGT algorithm was re-applied to the selected areas separately. A human expert has manually checked all the resulting document images and was asked to compare the produced image with the

**Table 2. Results of the evaluation of our approach for different window sizes and parameter k values based on the decisions of a human expert. Total number of document images is 183. Precision is calculated based on the 43 images already in good condition.**

| n | k | Better | Same | Worse | Score | Precision |
|---|---|--------|------|-------|-------|-----------|
|    | 1 | 35.52% | 42.08% | 22.40% | 24 | 97.67% |
|    | 2 | 38.25% | 42.62% | 19.13% | 35 | 97.67% |
| 25 | 3 | 42.08% | 40.98% | 16.94% | 46 | 95.35% |
|    | 4 | 37.70% | 58.47% | 3.83% | 62 | 100% |
|    | 5 | 32.79% | 66.67% | 0.55% | 59 | 100% |
|    | 1 | 43.72% | 45.36% | 10.93% | 60 | 97.67% |
|    | 2 | 60.11% | 39.34% | 0.55% | 109 | 100% |
| 50 | 3 | 40.98% | 58.47% | 0.55% | 74 | 100% |
|    | 4 | 34.43% | 65.57% | 0.00% | 63 | 100% |
|    | 5 | 28.42% | 71.58% | 0.00% | 52 | 100% |
|    | 1 | 34.97% | 57.38% | 7.65% | 50 | 97.67% |
|    | 2 | 31.15% | 66.12% | 2.73% | 52 | 97.67% |
| 100 | 3 | 25.14% | 74.32% | 0.55% | 45 | 100% |
|    | 4 | 19.67% | 80.33% | 0.00% | 36 | 100% |
|    | 5 | 10.38% | 89.62% | 0.00% | 19 | 100% |

image produced by the initial application of IGT algorithm to the entire image and classify each image in one of the following classes:

- Better: The produced image has been improved.
- Same: The produced image and the simple IGT image are practically the same.
- Worse: The produced image contains more noise than before.

Figure 8*c*/8*d* shows an example of producing a better/worse version of the simple IGT output (given in figure 8*b*). Given that $d(i)$ is the decision of the human expert for the *i*-th image, let $d(i)$ be 1, 0, or -1 in case the decision is better, same, or worse, respectively. Then, a performance *score* is calculated for the document collection of *N* document images as follows:

$$score = \sum_{i=1}^{N} d(i) \qquad (4)$$

The evaluation results of this procedure are given in Table 2. Note, the last column of this table, labeled *precision*, refers to the images that did not need further improvement after the application of the simple IGT algorithm to the entire page. Recall that those images were 23.50% of the document collection. Precision is, then, defined as the percentage of this set of images (already in good condition) that remain on the same level of quality. This kind of measure is especially important since any approach aiming at the removal of background noise should not affect drastically those images. As can be seen, the proposed approach achieves high level of precision in all cases.

The results of the evaluation process indicate that a window of medium size (50x50) achieves the best results as concerns both improving the quality of the majority of the images and maintaining precision at high level. As concerns the parameter *k*, its best value depends on the window size. For small window size, *k* needs to be high. That is, given many small areas, the algorithm should select only the more important ones. On the other hand, when the window size is large, the *k* should be low. That is, given a few big areas, the algorithm should select most of them. However, this also affects precision. For medium window size, *k* should be of medium value as well (2 and 3 seem to be the most effective). As can be seen, for those values a significant percentage of the document images is improved while only a few images are affected negatively.

## 5. Conclusions

In this paper we presented a hybrid binarization approach aiming at removal of background noise from historical and ancient documents. This way, we attempt to combine the advantages of global and local thresholding, that is, better adaptability of various kinds of noise at different areas of the same image based on low computational and time cost. The evaluation results using a historical document collection indicates that the proposed approach is able to deal with hard cases (where different kinds of background noise coexist) while keeping precision on high level (i.e., the document images already in good condition are not affected). Thus, it can be used in the framework of libraries willing to provide public access to their historical document collections as well as a

preprocessing step in document image analysis systems.

As concerns the parameters of the proposed algorithm, a medium size window ($n$=50) provides the most appropriate solution for detecting the areas with remaining noise. In addition, the best results were obtained by using $k$=2. However, in the framework of cost-sensitive application where the cost of worsening at least one document image is too high, $k$ should be higher to ensure that only the most important areas would be selected. Note that in the presented experiments, for $n$=50 and $k$>3 no document image is affected negatively by the proposed approach.

The evaluation of the proposed approach was based on a human expert. Although subjective, this methodology provides direct evidence on the quality improvement of the input document images. The human expert was asked to range the produced images using a minimal quality scale (better/same/worse). Therefore, the provided results can only be viewed as a rough estimation of the quality of document images. A more detailed decision scale would provide a better view on the amount of remaining background noise. For example, a multipoint scale or a Likert scale could be used towards this direction. In addition, opinions from multiple human experts would strengthen the evaluation procedure.

On the other hand, an objective evaluation approach would involve off-the-shelf OCR tools for testing the readability of the noise-free images. This methodology enables the quantitative comparison with other binarization approaches. However this evaluation method could only be applied to printed documents. We plan to perform such alternative evaluation procedures in the near future.

## Acknowledgment

## References

[1] Antonacopoulos, A.; Karatzas, D. "Document image analysis for World War II personal records" First International Workshop on Document Image Analysis for Libraries (DIAL'04), p. 336-341, 2004.

[2] Marinai, S.; Marino, E.; Cesarini, F.; Soda, G. "A general system for the retrieval of document images from digital libraries " DIAL'04, p. 150-173, 2004.

[3] Venu Govindaraju; Xue, H. "Fast handwriting recognition for indexing historical documents" DIAL'04, p. 314-320, 2004.

[4] H. S. Baird. "Difficult and Urgent Open Problems in Document Image Analysis for Libraries" DIAL'04, p. 25-32, 2004.

[5] Couasnon, B.; Camillerapp, J.; Leplumey, I.. "Making handwritten archives documents accessible to public with a generic system of document image analysis" DIAL'04, pp. 270-277, 2004.

[6] E.Kavallieratou, "A Binarization Algorithm Specialized on Document Images and Photos", Eighth International Conference on Document Analysis and Recognition (ICDAR'05), p.463-467, 2005.

[7] Otsu, N. "A threshold selection method from gray-level histograms". IEEE Trans. Systems Man Cybernet. pp. 62-66, 9 (1), 1979.

[8] Leedham, G., S. Varma, A. Patankar, V. Govindaraju "Separating Text and Background in Degraded Document Images" Proceedings Eighth InternationalWorkshop on Frontiers of Handwriting Recognition, pp. 244-249, September, 2002.

[9] Sauvola, J., Pietikainen, M., "Adaptive Document Image Binarization", Pattern Recognition, pp. 225-236, 33 (2000).

[10] Bernsen, J."Dynamic thresholding of grey-level images", ICPR86, pp 1251-1255, Paris, France, October 1986.

[11] J.He, Q.D.M.Do, A.C.Downton, J.H.Kim, "A Comparison of Binarization Methods for Historical Archive Documents", ICDAR'05, p.538-542, 2005.

[12] Niblack, W. "An Introduction to Digital image processing", pp 115-116, Prentice Hall, 1986.

[13] Gatos B., Pratikakis I. and Perantonis S.J. "An adaptive binarisation technique for low quality historical documents". IAPR Workshop on Document Analysis systems (DAS'2004), Lecture Notes in Computer Science (3163), Florence, Italy, pp. 102-113.

[14] Shi, Z., V. Govindaraju, "Historical Document Image Segmentation Using Background Light Intensity Normalization", SPIE Document Recognition and Retrieval XII, 16-20 January 2005, San Jose, California, USA.

[15] Yan, C., G. Leedham, " Decompose-Threshold Approach to Handwriting Extraction in Degraded Historical Document Images" IWFHR'04 pp. 239-244, Kokubunji, Tokyo, Japan, October, 2004.

[16] Leydier Y., LeBourgeois F., Emptoz H., Serialized Unsupervised Classifier for Adaptative Color Image Segmentation: Application to Digitized Ancient Manuscripts, ICPR, pp 494-497, Cambridge, 23-26, 2004

[17] L. Shapiro, G.Stockman "Computer Vision," Prentice Hall, 2001.