# Restoration of Ancient Document Images Using Phase Based Binarization

V. Supaja[1], Saudagar Nikhath Afreen[2], P Thanmai[3], P Chaitanya Lahari[4], S. Sri Varsha[5]

[1]Assistant Professor, Department of ECE, Ravindra College of Engineering for Women, Kurnool, Andhra Pradesh, India

[2,3,4,5]Department of ECE, Ravindra College of Engineering for Women, Kurnool, Andhra Pradesh, India

## ABSTRACT

The main defects present in historical documents are darkness, non-uniform clarification, bleed-through and faded characters. To remove these defects binarization method is used. In this paper a phase based binarization method is studied in which phase of ancient document images is preserved. This method is derived in to three steps: preprocessing, main binarization and post processing. In preprocessing phase preserved denoised image is derived. In main binarization two phase feature maps are derived are maximum moment of phase congruency covariance and a locally weighted mean phase angle. At last in post processing Gaussian and median filter is use for enhancement of image. It is also improve the performance of binarization methodologies.

Index Terms : - Historical document binarization, phase-derived features, document enhancement.

## I. INTRODUCTION

Historical documents go through different degradations due to getting old, complete use, some attempts of acquisition and ecological situation. The main defects present in historical documents are darkness, non-uniform clarification, smear, strain, bleed-through and faded characters. Those defects are problematic for document image analysis methods which presume even background and consistent quality of writing. In handwritten documents, the writer may use different amount of ink and pressure and make characters of dissimilar intensity or thickness, as well as faint characters. The same writer may write in different ways even within the same document .Similar problems, such as faint characters and non-uniform appearance of characters of the same font, are also encountered in historical machine-printed documents.

Today, there is a strong progress in the direction of digitization of these historical documents to save their content for future generations. The huge quantity of digital data created requires automatic processing, enhancement, and recognition. A key step in all document image processing workflows is binarization, but this is not a very complicated process, which is unfortunate, as its performance has a significant control on the quality of OCR results. lots of research studies includes the problem arising due to binarization of old document images characterized by many types of degradation and they are finding solution, including faded ink, bleed-through, show-through, uneven illumination, variations in image

contrast. There are also variations in patterns of hand-written and machine-printed documents, which add to the difficulties related with the binarization of old document images.

A phase-based binarization method is proposed for the binarization and enhancement of historical documents and manuscripts. The three main steps in the proposed method
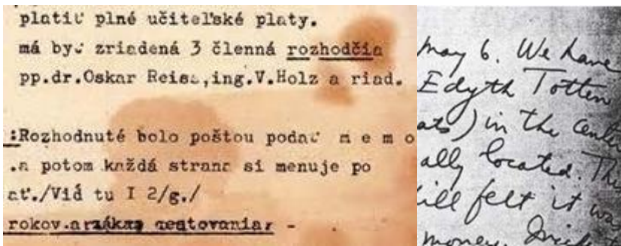


Fig. 1. Sample document images

A phase-based binarization method is proposed for the binarization and enhancement of historical documents and manuscripts. The three main steps in the proposed method are: preprocessing, main binarization, and post-processing. Phase of the image does not corrupt so phase information is mainly preserved. The preprocessing step involves image denoising. In which phase information is mainly preserved, followed by some morphological operations. Then edge map image is obtained by using canny edge operator. These two images are combined to obtained binarization image in rough form.

Then, for the main binarization step we use the phase congruency features. Phase congruency is dimensionless quantity that is invariant to changes in image brightness and contrast. The foreground of ancient documents can be modeled by phase congruency After completing the three binarization steps on the input images using phase congruency features and a denoised image the enhancement processes are applied. A median filter and a phase congruency feature are used to construct an object exclusion map image. This map is then used to remove unwanted lines and interfering patterns. The effect of each step on the binarized output image is discussed in each associated section.

The proposed binarization method is stable and robust to various types of degradation and to different datasets, thanks to its purpose-designed steps, and we provide comprehensive experimental results to demonstrate this robustness. The method outperforms most of the algorithms entered in the DIBCO'09 , H- DIBCO'10 , DIBCO'11 , competitions, based on various evaluation measures, including the F-measure, NRM, PSNR, DRD, and MPM.

## II.  LITERATURE REVIEW

In this section, some binarization methods are briefly described. Gatos et al propose an adaptive binarization method based on low-pass filtering, foreground estimation, background surface computation, and a combination of these. In an initial binary map is obtained using the multi-scale Sauvola's method , and then statistical methods are used to restore the missed strokes and sub- strokes. InValizadeh et al. map input images into a two- dimensional feature space in which the foreground and background regions can be distinguished. Then, they partition this feature space into several small regions, which are classified into text and background based on the results of applying Niblack'smethod.

Lu et al. propose a binarization method based mainly on background estimation and stroke width estimation. First, the background of the document is estimated by means of a one- dimensional iterative Gaussian smoothing procedure. Then, for accurate binarization of strokes and sub-strokes, an L1-norm gradient image is used. This method placed 1st of 43 algorithms submitted to the DIBCO'09 competition. In a local contrast image is combined with a Canny edge map to produce a more robust feature map.

Farrahi Moghaddam et al. propose a multi-scale binarization method in which the input document is binarized several times using different scales. Then, these output images are combined to form the final output image. This method uses different parameters for Sauvola's method to produce output images of the same size, but at different scales. In contrast, Lazzara and Gerard propose a multi-scale Sauvola's method which binarizes different scales of the input image with the same binarization parameters. Then, binary images with different scales are combined in some way to produce the final results.

Combination methods have also attracted a great deal of interest, and provided promising results. The goal of combining existing methods is to improve the output based on assumption that different methods complement one another. In several of these methods are combined based on a vote on the outputs of each. In a combination of global and local adaptive binarization methods applied on an inpianted image is used to binarize handwritten document images. The results show that this method performs extremely well; however, it is limited to binarizing handwritten document images only.

Learning-based methods have also been proposed in recent years. These methods are an attempt to improve the outputs of other binarization methods based on a feature map, or by determining the optimal parameters of binarization methods for each image In a self-training document binarization method is proposed. The input pixels, depending on the binarization method(s) used are divided into three categories: foreground, background, and uncertain, based on a priori knowledge about the behavior of every method used. Then, foreground and background pixels are clustered into different classes using the k-means algorithm or the random Markov field. Finally, uncertain pixels are classified with the label of their nearest neighboring cluster. The features used for the final decision are pixel intensity and local image contrast.

Another combined method based on a modified contrast feature is proposed. Lelore and Bouchara also classify pixels into three categories using a coarse thresholding method, where uncertain pixels are classified based on super resolution of likelihood of foreground. Howe proposes a method to optimize the global energy function based on a Laplacian image. In this method, a set of training images is used for optimization. Howe improved this method by tuning two key parameters for each image. In a learning framework is proposed to automatically determine the optimal parameters of any binarization method for each document image. After extracting the features and determining the optimal parameters, the relation between the features and the optimal parameters is learned. As we show in the Experimental Results and Discussion section, a problem associated with all these algorithms is that they are not reliable for all types of degradation and with different datasets.

## III. METHODOLOGY

The final binarized output image is obtained by processing the input image in three steps: preprocessing, main binarization, and postprocessing.

### 1. Pre-processing

In the preprocessing step, denoised image instead of the original image is used to obtain a binarized image in rough form. A number of parameters impact the quality of the denoised output image (ID), the key ones being the noise standard deviation threshold to be rejected (k), and the number of filter scales (N$\rho$) and the number of orientations (Nr) to be used. The N$\rho$ parameter controls the extent to which low frequencies are covered. The higher N$\rho$ is, the lower the frequencies, which means that the recall value remains optimal or near optimal. Based on our experiments, N$\rho$ = 5 is the appropriate choice in this

case. Therefore, to preserve all the foreground pixels, we set the parameters in the experiments as follows: k = 1, Nρ = 5 and Nr = 3.

We used Otsu's method on the normalized denoised image, where normalized denoised image is obtained by applying a linear image transform on the denoised image. This approach can also remove noisy and degraded parts of images, because the denoising method attempts to shrink the amplitude information of the noise component. The problem with this approach is that it misses weak strokes and sub-strokes, which means that we cannot rely on its output. To solve this problem, we combine this binarized image with an edge map obtained using the Canny operator. Canny combination those edges without any reference in the aforementioned binarized image are removed. We then compute a convex hull image of the combined image.
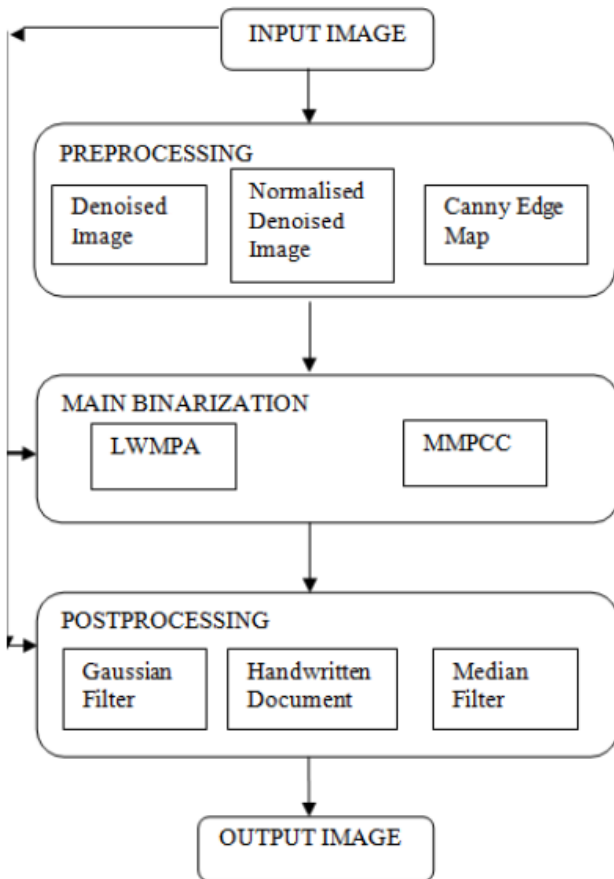


Fig.2. Flowchart of the proposed binarization method.

2. MainBinarization:

The next step is the main binarization, which is based on phase congruency features: i) the maximum moment of phase congruency covariance (IM); and ii) the locally weighted mean phase angle (IL ).

a) IM:

In this paper, IM is used to separate the background from potential foreground parts. This step performs very well, even in badly degraded documents, where it can reject a majority of badly degraded background pixels by means of a noise modeling method. Two-dimensional phase congruency is calculated by

$$PC_{2D,r}(x) = \frac{\sum_p W_r(x)[A_{pr}(x)\Delta\Phi_{pr}(x) - T_r]}{\sum_p A_{pr}(x)} \qquad (1)$$

Two-dimensional phase congruency is calculated by:

$$I_M = {}^{max}_r PC_{2D,r}(x) \qquad (2)$$

To achieve this, we set the number of two-dimensional log- Gabor filter scales ρ to 2, and use 10 orientations of two- dimensional log-Gabor filters r .In addition, the number of standard deviations k used to reject noises is estimated as follows:

$$k = 2 + \left[ a \times \left( \frac{\sum_{n,m} I_{Otsu,bw}(n,m)}{\sum_{n,m} I_{pre}(n,m)} \right) \right] \qquad (3)$$

where α is a constant (we are using α = 0.5); IOtsu,bw isthe binarization result of Otsu's method on the input image; and IPre is the output of the preprocessing step. Here, the minimum possible value for k is 2.

5.2.2 IL:

The two-dimensional locally weighted mean phase angle (IL) is obtained using the summation of all filter responses over all possible orientations and scales:

$$I_L(x) = arctan\left[ \sum_{p,r} e_{pr}(x), \sum_{p,r} o_{pr}(x) \right] \qquad (4)$$

We consider the following assumption in classifying foreground and background pixels using IL :

$$p(x) = \begin{cases} 1, & I_L(x) \leq 0 \\ 0, & I_L(x) > 0 \ and \ I_{Otsu,bw}(x) = 0 \end{cases} \quad (5)$$

where P(x) denotes one image pixel; and IOtsu,bw denotes the binarized image using Otsu's method. Because of the parameters used to obtain the IM and IL maps, IL produces some classification errors on the inner pixels of large foreground objects. Using more filter scales would solve this problem, but reduce the performance of IL on the strokes. Also, IL impacts the quality of the IM edge map, and of course requires more computational time. Nevertheless, the results of using Otsu's method to binarize the large foreground objects are of interest. Consequently, we used the IOtsu, bw image to overcome the problem.



Fig.3. Example of IM, IL, and ID maps.

3. Postprocessing:

In this step, we apply enhancement processes. First, a bleedthrough removal process is applied. Then, a Gaussian filter is used to further enhance the binarization output and to separate background from foreground, and an exclusion process is applied, based on a median filter and IM maps, to remove background noise and objects. Finally, a further enhancement process is applied to the denoised image. The individual steps are as follows:

1) Global bleed-through exclusion: Bleed-through degradation is a common interfering pattern and a significant problem in old and historical document images. In this paper, bleedthrough is categorized in two classes: i) local bleed- through; and ii) global bleed-through. Local bleed-through involves pixels located under or near foreground pixels, while global bleed-through involves pixels located far away the foreground text. Global bleed-through is one of most challenging forms of degradation, because there is no local to enable true text to be distinguished from bleed-through. At this stage, we investigate the possibility of the existence of global bleed-through. If it does exist, the parameters of the Canny edge detector are chosen to ensure that the output edge map contains only the edges of text regions which we expect to be located in a specific part, or parts, of the image. The existence of bleed- through is established by comparing the Otsu's result and the binary output obtained so far [19]. If there is a noticeable difference between these two binary images, we apply a global bleed-through exclusion method. Fig. 4 provides two examples of the global bleed-through exclusion process.
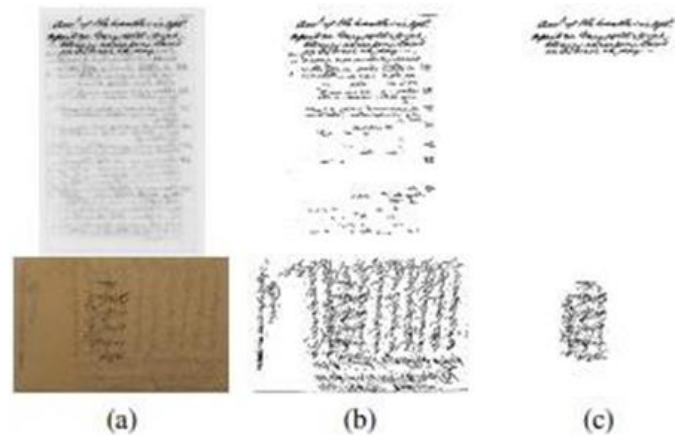


Fig.4 : Effect of using the proposed global bleed-through exclusion is shown in column (c). The left image (b) is the binarized image before the global bleed-through exclusion step has been applied.

2) Adaptive Gaussian filter: In this section, we take a similar approach to the one used in [47], except that a Gaussian smoothing filter is used to obtain a local weighted mean as the reference value for setting the threshold for each pixel. We use a rotationally symmetric Gaussian low-pass filter (G) of size S with σ value, estimated based on average stroke-width, where

σ is the standard deviation. This is a modification of the fixed S value used in [19]. The value for S is the most important parameter in this approach. Local thresholds can be computed using the following two-dimensional correlation:

$$T(x,y) = \sum_{i=-S}^{S} \sum_{j=-S}^{S} G(i,j) \times I(x+i, y+j) \,,$$

(6)

where I(x, y) is a gray-level input image. The result is a filtered image T(x, y) which stores local thresholds. A pixel is set to 0 (dark) if the value of that pixel in the input image is less than 95% of the corresponding threshold value T(x, y), and it is set to 1 (white) otherwise. We increased the value from 85% [47] to 95%, in order to obtain a near optimal recall value.

Some sub steps of the proposed binarization method work on objects rather than on individual pixels, and so it is important to separate foreground text from the background. The Gaussian filter described above is one of the methods used to achieve this. This filter is also applied to the equalized adaptive histogram image instead of the original image, in order to preserve weak strokes. The average stroke width is computed, in order to set S. There are various methods for computing stroke width [1], [3], [16]. In this paper, a very rapid approach, based on the reverse Euclidean distance transformation [48] of the rough binary image obtained so far, is used to estimate the average stroke width. This approach is dependent on the quality of the rough binary image, which has the potential to produce errors; however, it is a very fast way to calculate stroke width, and provides a good estimate of the average stroke width.

a) Document type detection: At this step, we need to determine the type of input document we are dealing with. We propose to apply the enhancement processes that are after this step to the handwritten documents only, and not to machine printed documents. The method we propose for detecting the type of document is straightforward and fast. We use the standard deviation of the orientation image that was produced during calculation of the phase congruency features. This image takes positive anticlockwise values between 0 and 180. A value of 90 corresponds to a horizontal edge, and a value of 0 indicates a vertical edge. By considering the foreground pixels of the output binary image obtained so far, we see that the standard deviation value of the orientations for these pixels is low for handwritten document images and higher for machine-printed documents. The reason for this is the different orientation values for interior pixels and edges. This approach works well for almost all the images we tested, including 21 machine-printed images and 60 handwritten document images, and only one classification error was found. It can be seen from the Fig.8 that the histogram of orientations of a handwritten document follows a U-shape behavior. Note that even if this approach fails to accurately detect the type of document, it nevertheless produces satisfactory output.

3) Object exclusion map image (IOEM): We construct an object exclusion map image (IOEM) based on a combination of a median filter and a binary map of IM (see Algorithm 1). Any object without a reference in this binary map will be removed from the final binarization results. This approach can remove noise, local bleed-through, and interfering patterns.

It is known that a median filter can reject salt-and-pepper noise in the presence of edges [49]. Like the method used in the previous section for the Gaussian filter, local thresholds are computed by applying an S × S symmetric median filter for each pixel in the input image. The value for S is estimated based on the average stroke width, instead of taking a fixed value as in . In turn, a filtered image equal in size to the input image is produced (Med), where its pixel values are local thresholds.
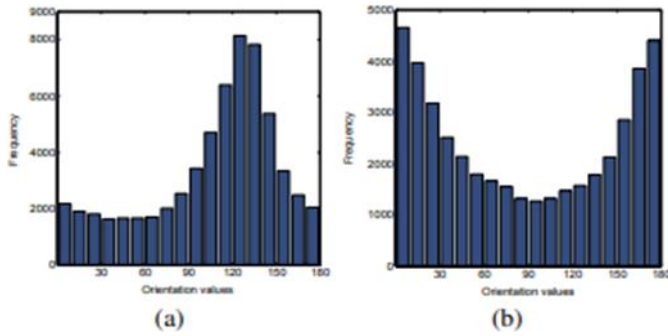
Fig. 5. Histogram of orientations of a handwritten document (a) and a machine-printed document (b).

A pixel is set to 0 (dark) if the value of that pixel in the input image is less than 90% of corresponding pixel value in Med, and it is set to 1 (white) otherwise. The output is called IMed.

4) Majority criterion: We propose a majority criterion based on a denoised image, ID. A majority criterion supposes that early binarization steps provide an optimal or near optimal recall value. Then, based on the fact that a foreground pixel should have a lower value than its adjacent background pixels, exclusion over the foreground pixels is performed. A majority criterion works as follows. For each foreground pixel in Ibwout, its $5 \times 5$ adjacent background pixels in ID are checked, and that pixel is removed from the foreground if its value in ID is less than that of any of the $5 \times 5$ background values in ID. This criterion works very well on noise, unwanted lines, and strokes and sub-strokes. Algorithm 1 provides the pseudo code of the proposed binarization method.

## IV. EXPERIMENTAL RESULTS AND DISCUSSION

The proposed binarization method is evaluated on a number of datasets. The following datasets were used: DIBCO'09 , H-DIBCO'10 , DIBCO'11 . These datasets provide a collection of images that have suffered different types of degradation, and which give enough information and are sufficiently challenging in terms

of evaluation setup to enable a meaningful examination of various algorithms.

First, we compare the subjective and objective performance of the proposed method with that of leading binarization methods in the literature. Then, we compared our proposed binarization method with state-of-the-art algorithms and the top ranking algorithm in each competition.

A. Subjective evaluation
In this section, we compare outputs of the proposed method with those of top-placing methods in each contest, whenever possible. Our proposed method performs a smooth binarization of the document images, thanks to the use of phase congruency measures and a denoised image. In Fig. 6, we compare the proposed method with three top-placing algorithms in DIBCO'11 the winning algorithm in DIBCO'09and the method proposed.

B. Objective evaluation
We used the well-known measures F-measure (FM), pseudo F- measure (p-FM), PSNR, distance reciprocal distortion (DRD) metric , the misclassification penalty metric (MPM), the negative rate metric (NRM) to evaluate various algorithms . The source code of the evaluation measures used in this paper is available . The results of the proposed method are compared with state-of- the-art binarization methods. PHIBD'12 is a dataset of historical Persian images consisting of 15 degraded document images. The results show that the proposed method achieved, on average, a 5% improvement over our earlier results .

It can be seen from these experimental results that other binarization methods produce different results for different datasets, whereas there is little difference between the results we obtained using the proposed method on different datasets, which shows the robustness of our method.In this paper, the standard deviation of the numerical results is considered to

measure the reliability of the various methods we compared.

C. Enhancement of other binarization methods

Preprocessing and main binarization in the proposed method are used as a mask to cross out false positive pixels on the output of other binarization methods, which resulted in an outstanding level of improvement. This mask has a high recall value with an acceptable precision value. Compared with previous works which were aimed at modifying other binarization methods, our proposed method shows even more improvement.

D. Time complexity

In this section, we evaluate the run time of our proposed method, performing our experiments on a Core i7 3.4 GHz CPU with 8 GB of RAM. The algorithm is implemented in MATLAB 2012a running on Windows 7. It takes 2.04seconds to operate it on a 0.3 megapixel image, and 20.28 seconds to produce output for a 3 megapixels image. It is worth mentioning that the proposed algorithm would run faster and would require much less memory if the phase congruency features are calculated using the alternative monogenic filters.



Fig. 6. Subjective comparison of six binarization methods. a) Lu's method b) Su's method c) The winner at DIBCO'11 d) Proposed method

## V. CONCLUSION

In this paper, we have introduced an image binarization method that uses the phase information of the input image, and robust phase-based features extracted from that image are used to build a model for the binarization of ancient manuscripts. Phase-preserving denoising followed by morphological operations are used to preprocess the input image. Then, two phase congruency features, the maximum moment of phase congruency covariance and the locally weighted mean phase angle, are used to perform the main binarization. For post-processing, we have proposed a few steps to filter various types of degradation, in particular, a median filter has been used to reject noise, unwanted lines, and interfering patterns. Because some binarization steps work with individual objects instead of pixels, a Gaussian filter was used to further separate foreground from background objects, and to improve the final binary output. Our experimental results demonstrate its promising performance, and also that of the postprocessing method proposed to improve other binarization algorithms.

We have also proposed a rapid method to determine the type of document image been studied, which will be of great interest. The behavior of ancient handwritten document images and machine-printed images shows differences in terms of binarization. The strokes and sub-strokes of handwritten images require accurate binarization, and the binarization of the interior pixels of the text of machine-printed images needs to be performed with care. Although the proposed binarization method works well on both handwritten and machine-printed documents, better results for both types of documents are achieved, when a priori information about the type of input document is available. In future work, we plan to expand the application of phase-derived features, which ensures the stable behavior of document images, to other cultural heritage fields, such as microfilm analysis and multispectral imaging.

## VI. REFERENCES

[1]. B. Su, S. Lu, and C. L. Tan, "Robust document image binarization technique for degraded document images," IEEE Trans. Image Process., vol. 22, no. 4, pp. 1408–1417, Apr.2013.

[2]. R. F. Moghaddam and M. Cheriet, "AdOtsu: An adaptive and parameterless generalization of Otsu's method for document image binarization," Pattern Recognit., vol. 45, no. 6, pp. 2419–2431,2012.

[3]. J. Sauvola and M. Pietikinen, "Adaptive document image binarization," Pattern Recognit., vol. 33, no. 2, pp. 225–236,2000.

[4]. B. Gatos, I. Pratikakis, and S. Perantonis, "Adaptive degraded document image binarization," Pattern Recognit., vol. 39, no. 3, pp. 317–327,2006.

[5]. R. Hedjam, R. F. Moghaddam, and M. Cheriet, "A spatially adaptive statistical method for the binarization of historical manuscripts and degraded document images," Pattern Recognit., vol. 44, no. 9, pp. 2184– 2196,2011.

[6]. K. Ntirogiannis, B. Gatos, and I. Pratikakis, "A combined approach for the binarization of handwritten document images," Pattern Recognit. Lett., vol. 35, pp. 3–15, Jan.2014.

[7]. B. Su, S. Lu, and C. Tan, "Binarization of historical document images using the local maximum and minimum," in Proc. 9th IAPR Int. Workshop DAS, 2010, pp.159–166.

[8]. B. Su, S. Lu, and C. L. Tan, "A self-training learning document binarization framework," in Proc. 20th ICPR, Aug. 2010, pp.3187–3190.

[9]. B. Su, S. Lu, and C. L. Tan, "A learning framework for degraded document image binarization using Markov random field," in Proc. 21st ICPR, Nov. 2012, pp.3200–3203.

[10]. P. Kovesi, "Phase preserving denoising of images," in Proc. Int. Conf. Digital Image Comput., Techn. Appl.,1999.

[11]. P. Kovesi, "Image features from phase congruency," Videre, J. Comput. Vis. Res., vol. 1, no. 3, pp. 1–26,1999.

[12]. K. Ntirogiannis, B. Gatos, and I. Pratikakis, "A performance evaluation methodology for historical document image binarization," IEEE Trans. Image Process., vol. 22, no. 2, pp. 595–609, Feb.2013.

## Cite this article as :