Image Binarization and Lossless Compression of Ancient Documents using SPIHT Algorithm

Thumilvannan .P.S Assistant Professor, Department of CSE, AMACE, Vadamavandal-604410. Bhavani. S P.G. Student-M.E (CSE) Department of CSE, AMACE, Vadamavandal-604410. Sridevi. J P.G. Student-M.E (CSE) Department of CSE, AMACE, Vadamavandal-604410.

ABSTRACT

The ancient documents provides huge amount of information about the past, inventions and manuscripts which should be protected cautiously. These documents store a noteworthy amount of human heritage over time. However, many environmental factors which degrade the documents severely. A key step in all image processing is image binarization. Here the image is binarized using the gradient information including the phase feature. The three main steps are: preprocessing, image binarization, postprocessing. In the preprocessing step, both adaptive median filter and iterative wiener filter used to remove the noise from the image. In the postprocessing step, the BDND (Boundary Discrimination Noise Detection) filter is used to determine whether the current pixel is corrupted or not and the ABF (Adaptive Bilateral Filter) is used to sharpens the edge. For further image enhancement, the CLAHE (Contrast Limited Adaptive Histogram Equalization) is used to enhance the contrast of the input image. Finally, the binarized output need to be stored in effective manner but it consumes large amount of space and takes more number of bits to transfer from one place to another. So, the SPIHT (Set Partitioning in Hierarchical Tree) algorithm with Huffman Coding is used for lossless image compression, saves more number of bits to be transferred and provides enhanced image quality.

General Terms

Pre-processing, Post-processing, CLAHE, SPIHT, BDND Filter.

Keywords

Image binarization

1. INTRODUCTION

The preservation of ancient documents plays a major role in the task of image processing. Ancient documents is a valuable document which might be handwritten or printed documents in paper. The ancient documents provides more information about the historical events, manuscripts, nation's history, true knowledge of facts. They are not only the papers, books, old pictures, notebooks or civil acts. They are the living memory of our past. They are made up of fallible materials that often fade, rip or degrade over time due to many environmental factors, inadequate handling and poor quality of documents used in their formation cause them to suffer a huge amount of degradation. So, the various forms of degradation are bleed-through, faded ink, uneven illumination, high/low contrast, inappropriate temperature which affects the documents severely. These degradation which reduces the quality and richness of the document. For preserving the ancient documents we need to digitize the documents as images. So, there is a strong change towards digitization of the degraded ancient manuscripts. In the digitization, the ancient documents are scanned or captured

using camera and then the image is stored in the system. The major step in all scanned or captured degraded document image processing is binarization, but this is not a very stylish process, which is unsuccessful, as its performance has a noteworthy effect on the quality of OCR results.

The image binarization is a process of converting an image of up to 256 gray levels to a black and white image. The black and white image is referred as digital image and it takes only two possible values for each pixels. For white (full intensity) uses 1 and for black (no intensity) uses 0. The example of binarization is given below in Fig 1:



Fig 1 (a) Before binarization (b) After binarization

The basic method used for binarization is thresholding. In this method, there are two main types of thresholding, one is global and another one is local. In global thresholding, one threshold value is selected for the entire image according to the global or local information (for all pixels in the image it uses only one threshold value). In the local thresholding, the threshold values are determinedlocally i.e., pixel by pixel or region by region. Based on the threshold value chosen, it classifies the pixels into white (value above the chosen threshold) and black (value below the chosen threshold). If we choose optimal threshold value it will give correct and accurate results of binarization.

The major purpose of document image binarization is that it is used to easily identify the flaws in document image binarization only the image can be enhanced, restored, compressed, analysed, and recognized. After the binarization, the image is introduced to the filter where the unwanted artifacts and noise is removed using various types of filters.Each filter plays different kind of noise removing operation from the degraded ancient documents image.Finally we yield the accurate output of degraded ancient documents image and stored for later use. The amount of storage used for storing the binarized ancient documents image occupies more space to store the few number of images. The image takes lot of storage space, for example, 1024 x 1024 x 32 x bits images requires 4MB and also the desire to exchange the images over the internet, have led to a large interest in image compression algorithms. So, there is a necessity of storing the ancient document images to be effective and also it should consumes less amount space to store more binarized and filtered images.Here, if the images are compressed and stored means it takes less amount of space for huge amount of ancient document images.

The process of reducing the amount of data required to represent a digital image is referred as image compression. Data compression algorithms can be broadly divided into two categories:

- 1) Lossless algorithms remove only redundancy present in the data. The reconstructed image is identical to the original, i.e., all of the information present in the input image has been preserved by compression.
- Lossy algorithms which creates more redundancy (by discarding some information) and then remove it. The lossy algorithms provide higher compression ratio but it losses some valuable information in the image.

The image compression should not lose any valuable information from the image and also it should give higher compression ratio. There are lots of image compression techniques are available in both algorithms such as lossy and lossless. The lossy and lossless compression ratio is differed by the coding of redundancy. Here the SPIHT algorithm is used to compress the ancient document image with the help of Huffman Coding. The SPIHT algorithm provides higher compression ratio, error protection and high quality. The Huffman Coding is a simple and robust coding to encode and decode the image perfectly and also it reconstruct the image with minimum redundancy codes.Using the proposed algorithm, it provides less space for huge amount of degraded document images and also it saves the number of bits to be transferred from one place to another place.

2. RELATED WORK

In this section, we deal with various kinds of related techniques which is interrelated to proposed method. In [1] document image binarization, the challenging task is to segment the text from the badly degraded document images they the text stroke edge pixels with the help of binarized adaptive contrast image with the canny edge map. After that, the local threshold value is determined from the text stroke edge pixels. In postprocessing, the quality of text stroke edge is improved. The proposed technique is more stable and easyto-use for document images. The document images as smear, smudge, bleed-through and low-contrast images only is considered for testing. It couldn't deal with the problems of uneven-illumination, brightness, and salt and pepper noise. In [2], the historical manuscripts and degraded document images are binarized using spatially adaptive statistical method with aid of the priori information and maximum likelihood classification and spatial relationship on the image domain. The Sauvola binarization method is used to process the independent pixels locally and also it preserves the weak connections and strokes in the image. The soft decision method [3] used for binarization doesn't need any parameter (parameterless) for thresholding. But the proposed method doesn't perform well in recovering very small loops in characters (highly degraded background) and doesn't preserve the connectivity between the broken parts of the text with both direction and intensity gradient of image pixels.

In [4], to segment the text from badly degraded historical document images they used the image contrast and it is defined in terms of the local image maximum and minimum. The document is segmented using the detected high contrast pixels (Otsu's global thresholding) and local threshold value. As in accordingly, it can introduce error if the background of the degraded document images contains a certain amount of pixels that are dense and also the proposed method heavily depends on the high contrast image pixels. In [5], for

binarization, it estimates the document background surface by using iterative polynomial smoothing procedure. Using L1-norm image gradient, text stroke edge is detected and the document is segmented using the local threshold value. The proposed method has different limitations such as it recovers the text from the document which suffered from bleedthrough only when the back-side text should be fairly brighter. The document images should not possess any weak slanting. The polynomial smoothing cannot handle the sharp variation of small size in the document folding.

Image enhancement [7] is a class of image processing operations whose goal is to produce an output digital image that is visually more suitable for human. There are two main categories of image enhancement techniques: Spatial and Frequency domain enhancement. In the spatial domain techniques, the manipulation is directly don't with image pixels. In the Frequency domain technique, the manipulation is done with the frequency or intensity of the image. From the spatial domain technique, the method CLAHE is used to enhance contrast of the image. It doesn't lead to the oversaturation in image and also it avoid the amplification of noise.

In [8], Kwame Osei Boateng had been used median filter to remove the impulse noise (salt and pepper noise) in the image and also it improves the degraded image quality. Here the process of median filter is as follows: first it calculates the median by sorting all the pixel values from the surrounding neighbourhood and then it replaces the noisy pixel with the median value. Here they proposed the spatial median filter which uses the overlapping window to filter the noise and also it calculates median for each window. Even though, it removes in effective manner it is not preferred one because it removes the fine details of information from the image.

In the following sections, we briefly discuss about the proposed work. In this paper, we have used different kinds of filter in each and every stage to remove the different kind's noise and to improve the image quality.

3. PROPOSED SYSTEM

The main goal of the paper is that to enhance the image quality of the binarized image in effective manner and to improve the storage space efficiency. The proposed system architecture is given below in Fig 2.

The work of the system is described as: at first, the ancient document image is given as an input to the preprocessing stage which uses filters (AMF and IWF) to remove the unwanted halo artifacts from the image without creating any overshoot or undershoot. The filtered image is sent to the image binarization stage where the image pixels are clustered using phase information in the image and it track the edge by estimating gradient.

The image with proper edge information is sent to the postprocessing stage for further enhancement. In this stage, the ABF is used to sharpen the edge and the BDND is used to remove noise in the high frequencies. After the postprocessing stage, the enhanced image is further enhanced with certain features like contrast and brightness of the image using CLAHE method.

The enhanced image takes huge amount of space to store the image. So, here the lossless compression algorithm is used to provide efficient storage by using SPIHT algorithm with Huffman Coding.



Fig 2: System Architecture Diagram

3.1 Preprocessing

In the preprocessing phase, the Adaptive Median filter and Iterative Wiener filter is used and let us discuss about the filters below:

3.1.1 Adaptive Median Filter (AMF)

The AMF which performs spatial processing to determine which pixels in an image have been affected by impulse noise. It classifies pixels as noise by comparing each pixel to its surrounding neighbour pixels. It label the pixel as impulse noise such that the pixel that is differ from a majority of its neighbours and also the pixels that is not structurally aligned with those pixels to which it is similar. Finally it replaces the noise pixel with the median pixel value of the pixels in the neighbourhood. It smoothens the unwanted noise and reduce distortion. The AMF can changes the size of S_{xy} (the size of neighbourhood) during operation. The algorithm of AMF is explained as level A and level B.

Level A:
$$A1 = Z_{med} - Z_{min}$$

$$A2 = Z_{med} - Z_{max}$$

if A1 > 0 AND A2 < 0, go to level B

else increase the window size

if window size < Smax, repeat level A

Level B:
$$B1 = Z_{xy} - Z_{min}$$

$$B2 = Z_{xy} - Z_{max}$$

if B1 > 0 AND B2 < 0, output Z_{xy}

else output Z_{med}

The notations are Z_{min} -minimum gray level in S_{xy} , Z_{max} -maximum gray level value in S_{xy} , Z_{med} -median of gray levels in S_{xy} , Z_{xy} -gray level at coordinates(x, y), S_{max} -maximum allowed size of S_{xy} . Finally, the AMF removes the impulse noise from the ancient document image in greater extent.

3.1.2 b. Iterative Wiener Filter (IWF)

The original ancient document image is again given as an input to the IWF to remove the additive noise from the image and also it inverts the blurring simultaneously. The IWF is optimal in terms of mean square error. The basic iterative algorithm is work as follows:

- 1. The degraded image is used as an initial estimate of original image and a restored image is attained from the corresponding wiener filter.
- 2. The restored image is used as an updated estimate of the original image and it leads to a new restoration.
- 3. The iteration continue until the estimate converges.

The IWF outperforms well in case of motion blur, atmospheric turbulence and blurred signal-to-noise ratio. Iterative techniques can be applied in cases of spatially varying or nonlinear degradations or in cases where the type of degradation is completely unknown (restoration).



Fig. 3 (a)





Fig. 3 (c)

Figure 3: (a) Input Image (b) Output of Iterative Wiener Filter with PSNR=11.3 (c) Output of Adaptive Median Filter with PSNR=24.0

Finally, the filtered image from both the filters are compared with PSNR see Figure 3. The filtered output which has high PSNR value is taken as a better output and it is used as input to the image binarization stage.

3.2 Image Binarization

In this stage, it consists of two main process, one is phase feature extraction and another one is gradient estimation. Let we discuss about the two process below:

3.2.1 Phase Feature Extraction

In all the ancient document images, the most valuable information present in its phase feature. The phase in the image as tends to play different roles and in some situations many of the important features of a signal are preserved if only the phase is retained. The phase information is sufficient to reconstruct the degraded ancient document image to the high quality ancient document image. To estimate the phase feature, we first classify and cluster the pixel based on the assumption. The gray-scale values0 to 255 isclassify the pixels as 0 to 50, 51 to 100, 101 to 150, 151 to 200 and 201 to 255. After classifying, the pixels in the image are clustered to each group of classification of pixels.

After clustering, we use 2D convolution to estimate the phase feature from the image and then it estimates X-convolution and Y-convolution. From the estimated phase information, the gradient estimation is done in next process.

3.2.2 b. Gradient Estimation:

The image gradient is that the directional change in the intensity or color in an image. It is used to extract the information from the image. It calculates the global thresholding for gradient estimation. The algorithm used in this process is similar to Canny edge detection but the simple difference between these algorithm are, Canny edge detection algorithm which mainly detect the edge of the text in the image but the proposed algorithm which track (trace) the edge of the text perfectly. The algorithm runs in five steps such as,

- 1. Smoothing: Blurring of the image to remove noise.
- 2. Finding gradients: The edges should be marked where the gradients of the image has large magnitudes.
- 3. Non-maximum suppression: Only local maxima should be marked as edges.
- 4. Double thresholding: Potential edges are determined by thresholding.
- 5. Edge tracking by hysteresis: Final edges are determined by suppressing all edges that are not connected to a very certain (strong) edge.

3.2.3 Smoothing

It is inevitable that all images taken from a camera will contain some amount of noise. To prevent the edges, noise must be reduced. Therefore the image is first smoothed by applying a Gaussian filter. The kernel of a Gaussian filter with a standard deviation of Sigma= 1.4.

$$B = \frac{1}{159} \cdot \begin{bmatrix} 2 & 4 & 5 & 4 & 2 \\ 4 & 9 & 12 & 9 & 4 \\ 5 & 12 & 15 & 12 & 5 \\ 4 & 9 & 12 & 9 & 4 \\ 2 & 4 & 5 & 4 & 2 \end{bmatrix} - \text{Eq (1)}$$

3.2.4 Finding Gradient

The Canny algorithm basically finds edges where the grayscale intensity of the image changes the most. These areas are found by determining gradients of the image. Gradients at each pixel in the smoothed image are determined by applying what is known as the Sobel-operator. First step is to approximate the gradient in the x- and y-direction respectively by applying the kernels shown in Equation (2).

$$K_{\rm GX} = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} - {\rm Eq} (2)$$
$$K_{\rm GY} = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

The gradient magnitudes (also known as the edge strengths) can then be determined as a Euclidean distance measure by applying the law of Pythagoras as shown in Equation (3). It is sometimes simplified by applying Manhattan distance measure as shown in Equation (4) to reduce the computational complexity. The Euclidean distance measure has been applied to the test image.

$$|G| = \sqrt{G_{\rm x}^2 + G_{\rm y}^2} \tag{3}$$

$$|G| = |G_{\mathbf{x}}| + |G_{\mathbf{y}}| \tag{4}$$

where:

 G_x and G_y are the gradients in the x- and y-directions respectively. The image of the gradient magnitudes often indicate the edges quite clearly. However, the edges are typically broad and thus do not indicate exactly where the edges are. To make it possible to determine this, the direction of the edges must be determined and stored as shown in Equation (5).

$$\theta = \arctan\left(\frac{|G_{y}|}{|G_{x}|}\right)$$
 Eq (5)

3.2.5 Non-maximum Suppression

The purpose of this step is to convert the "blurred" edges in the image of the gradient magnitudes to "sharp" edges [12]. Basically this is done by preserving all local maxima in the gradient image, and deleting everything else. The algorithm is for each pixel in the gradient image:

- 1. Round the gradient direction to nearest 45°, corresponding to the use of an 8-connected neighbourhood.
- Compare the edge strength of the current pixel with the edge strength of the pixel in the positive and negative gradient direction. I.e. if the gradient direction is north (theta = 90°), compare with the pixels to the north and south.
- 3. If the edge strength of the current pixel is largest; preserve the value of the edge strength.
- 4. If not, suppress (i.e. remove) the value. A simple example of non-maximum suppression is shown below. Almost all pixels have gradient directions pointing north. They are therefore compared with the pixels above and below. The pixels that turn out to be maximal

in this comparison are marked with white borders. All other pixels will be suppressed.



Fig 3. Non-maximum suppression

3.2.6 Illustration of non-maximum suppression: In Fig 3, the edge strengths are indicated both as colors and numbers, while the gradient directions are shown as arrows. The resulting edge pixels are marked with white borders.

3.3 Postprocessing

In this stage, there are two main filters are used such as, BDND (Boundary Discriminative Noise Detection) Filter and ABF (Adaptive Bilateral Filter). Let we briefly discuss about the filters used below:

3.3.1 BDND Filter

The BDND filter is applied to each pixel in the noisy image to determine whether it is corrupted or uncorrupted [6]. In this process, the artifacts in the high frequency edges is removed and also it provides window selection to predict the uncorrupted pixels in the image. After the filtering process, the uncorrupted pixels is indicated as 0's and corrupted pixels with 1's. The BDND filter algorithm works as:

- 1. At first, it selects the window size.
- 2. By sorting the pixels it classifies the pixels of a localized window, centering on the current pixel, into three classes- lower intensity impulse noise, uncorrupted noise and higher intensity impulse noise. The center pixel is considered as uncorrupted pixel.
- 3. Based on the center pixel, all other two classes are compared and then it detects the lower and higher intensity impulse noise.
- 4. Suppose, the pixel is found as uncorrupted pixel means, the iteration stops otherwise the iteration starts from the 1st step by selecting small window size.
- 5. Stop the iteration after the middle pixels is found as corrupted or uncorrupted.

Finally the corrupted pixels which has high and low density frequency is easily removed from the binarized ancient document image.

3.3.2 Adaptive Bilateral Filter (ABF)

The ABF is used to smoothen and sharpens the edges of the noisy image and also it provides clean edges. ABF works well for both gray scale and color images using domain and range filters. When comparing to normal bilateral filter, the ABF performs better in removing noise. It is mainly based on the weight assigned to each pixel in the image. It is also used to restore the denoised image to the original noise less image. It consists of two filters such as Domain and Range filter. In the domain filter, it gives higher weight to pixels that are spatially close to the center pixel and it is manly used for noise removal. In the range filter, it gives higher weight to pixels that are similar to the center pixel in gray value. The ABF has two modifications one is the offset introduced to the range filter and another one is offset and width of the range filter are locally adaptive. The ABF outperforms the bilateral filter in noise removal and at the same time, it renders ore sharper images than the bilateral filter.

3.4 Image Enhancement

In this stage, the binarized and less noisy image is given as an input to the image enhancement stage for further enhancement in its color and brightness. The enhancement of image will provides high visual impact on the image and provides high quality. There are two methods are used such as CLAHE and PCA. The brief description of the methods are discussed below:

3.4.1 Contrast Limited Adaptive Histogram Equalization (CLAHE)

It is a simple and effective method to enhance the contrast of the binarized image and also it preserves the input brightness of the image and it provides better enhancement. The basic idea of HE (Histogram Equalization) is to re-map the gray levels of an image from the histogram map. It is a kind of spatial domain technique which directly manipulates the pixels in the image. In the histogram equalization, it transforms the values in an intensity image, so that the histogram of the output image approximately matches a specified histogram.

The function *adapthiseq* [9]performs enhancement of contrast in image. The CLAHE operates on small data regions (tiles). Each tile's contrast is enhanced so that the histogram of each output region approximately matches the specified histogram. The contrast enhancement can be limited in order to avoid amplifying the noise which might be in the image. And also it limit on the level of enhancement can also be set, thus preventing the over-saturation caused by *histeq* method. The simple steps of CLAHE is: first load the image, resize the image, enhance the grayscale image and then enhance the color image.

3.5 Compression

After binarizing the ancient document images, it takes huge amount of space to store the document images. For that purpose, the compression method is used and also here the lossless compression is used. One of the algorithm used for lossless compression is SPIHT (Set Partitioning In Hierarchical Tree) and for the encoding and decoding Huffman coding is used to code the image with minimum redundancy codes. The SPIHT algorithm is a advanced version of EZW (Embedded Zero tree Wavelet) coder and also it provides successive approximation and bit-plane processing. The Huffman Coding is an optimal compression algorithm where the frequency of individual letters are used to compress the data.

It employs spatial orientation trees and uses set partitioning sorting algorithm. Coefficients corresponding to the same spatial location in different subbands in the pyramid structure display self-similarity characteristics. SPIHT defines parentchild relationships between these self-similar subbands to establish spatial orientation trees. It classifies the pixels mainly into two groups such as List of significant pixels (LSP) and list of insignificant sets (LIS). The algorithm works as:

Step 1: Converts the image from spatial domain to frequency domain using wavelet transform.

Step 2: Initialize the partition of image transform X by using hierarchical trees into two sets: S=root and I=X-S (S-list of significant pixels and I-list of insignificant set)

Step 3: Sort the pixels in increasing order and converts it to column matrix.

Step 4: Process S(S) for checking the set S for significance and Process I() for checking the set I for significance.

Step 5: Huffman coding finds the repeated pixels and compress the insignificant pixels.

Step 6: Reconstruct the matrix to the form nxn and to produce the output with highest MSB at first.

Step 7: Quantization

Finally, the binarized ancient document image is stored effectively in less amount of space and it provides error protection, good image quality and high PSNR.

4. CONCLUSION AND FUTURE WORK

The median filter performs well as long as the spatial density of the impulse noise is not large. However the adaptive median filtering can handle impulse noise with probabilities even larger than these. An additional benefit of the adaptive median filter is that it seeks to preserve detail while smoothing non impulse noise. Considering the high level of noise, the adaptive algorithm performed quite well. The iterative wiener filter is used to eliminate the random noise throughout the document image. The document image under test is attempted to binarize with the help of clustering approach while estimating most likely background information using iterative algorithm. In each iteration the average intensity of the document image is adopted as midpoint between the clusters. After, removing the noise from the image the gradient is estimated and the edge is tracked by using the non-maximum suppression information. Compression increases the capacity of the communication channel due to the redundant information. Considering the important role played by digital imaging and video in today's world, it is necessary to develop a system that produces high degree of compression while preserving critical image or video information. The traditional image coding technology uses the redundant data in an image to compress it. But these methods have been replaced by digital wavelet transform based compression method as these methods have high speed, low memory requirements and complete reversibility. Now in this work we are considering SPIHT as a placement for wavelet compression methods. We are comparing it with wavelet encoding scheme and comparing the final results in

terms of bit error rate, PSNR and MSE. The combination of SPIHT with Huffman Coding provides high quality result in the image compression. In future, we are planning to use different compression algorithms to improve the storage efficiency and to provide the high quality image. In future, the different compression algorithms will be used to produce the effective results of compression and to store the image without noise.

5. REFERENCES

- 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. 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
- [3] J. Sauvola and M. Pietikinen, "Adaptive document image binarization," Pattern Recognit., vol. 33, no. 2, pp. 225– 236, 2000
- [4] 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
- [5] S. Lu, B. Su, and C. Tan, "Document image binarization using background estimation and stroke edges," Int. J. Document Anal. Recognit. vol. 13, no. 4, pp. 303–314, 2010
- [6] P. H. Sangave, Prof. G.P. Jain, "Modified Boundary Discriminative Noise Detection and Removal Technique for Salt and Pepper Noise Removal", Int. J. of Emerging Technology and Advanced Engg. Vol. 2, Issue 4, April 2012.
- [7] S.S. Bedi1, Rati Khandelwal,"Various Image Enhancement Techniques- A Critical Review," Int. J. of Advanced Research in Computer and Communication Engg. Vol. 2, Issue 3, march 2013
- [8] Kwame Osei Boateng, Benjamin Weyori Asubam and David Sanka Laar, "Improving the Effectiveness of the Median Filter" Int. J. Electronics & Comm. Engg, vol. 5, no. 1(2012), pp. 85-97
- [9] http://in.mathworks.com/help/images/ref/adapthisteq.htm