

Document Image Binarization

Rushi B. Mogal¹, Shubham S. Kahale², Prof. Rajshri S. Rahane³

PG Student, Department of Computer Science^{1,2}

Assistant Professor, Department of Computer Science³

K. R. T. Arts, B.H. Commerce, A.M. Science (KTHM) College, Nashik

Affiliated to SPPU, Maharashtra, India

Abstract: Document image binarization is a difficult task, especially when trying to separate text in old or damaged documents. Binarization is an important first step in Optical Character Recognition (OCR), where the goal is to turn the document into black and white — separating the text (foreground) from the background. This makes it easier for computers to read and process the image. Since document quality can vary a lot, researchers have developed many different methods to handle different types of damage and degradation. This paper gives an overview of the main techniques used for binarization. We also discuss some of the main problems and challenges in document image binarization. In this paper, we evaluate different binarization methods to understand their weaknesses and suggest ideas that could help improve future research in this area.

Keywords: Degraded document images; binarization; threshold processing

I. INTRODUCTION

Overview:

Document image binarization is done in the preprocessing stage of document analysis. Its main goal is to separate the text (foreground) from the background. A fast and accurate binarization method is very important for later steps in document image processing.

Even though researchers have worked on this problem for many years, it is still difficult to choose the right threshold for binarizing degraded documents. This is because the text and background can vary a lot in different images. In our method, we detect the edges of the text by combining a grayscale image with a new image segmentation technique, which helps create a clear binarized image.

Document image binarization is a crucial preprocessing step in document analysis and recognition systems. Its primary goal is to separate the textual content (foreground) from the background, enabling better readability and more accurate downstream processing, such as optical character recognition (OCR) and information retrieval. A fast and reliable binarization technique significantly impacts the performance of subsequent stages in the document processing pipeline.

In this study, we explore a binarization approach that begins with grayscale conversion, followed by the application of two thresholding techniques: global thresholding using Otsu's method and local thresholding using Niblack's method. While global methods are effective for documents with clean and uniform backgrounds, they often fail in more complex cases. Local thresholding methods, on the other hand, adapt to pixel-level variations, making them more suitable for degraded and noisy images. This paper aims to analyse the effectiveness of these binarization techniques and processing strategies, and to highlight the strengths and limitations of each in dealing with various types of degraded documents.

Problem Statement:

Document image binarization is an important first step in document analysis. Its main goal is to separate the text (foreground) from the background. A fast and accurate binarization method helps improve the quality of further processing steps.



Even though this topic has been researched for many years, binarizing degraded document images is still a difficult problem. This is because the text and background can look very different in each image, making it hard to choose the right threshold. To solve this, we detect the edges of the text using a combination of a grayscale image and a new image segmentation method, which helps produce a clear and clean binarized image.

Objectives

1. To develop an effective binarization method that can accurately separate text from the background in degraded document images.
2. To improve the quality of document images for better readability and more accurate OCR (Optical Character Recognition) results.
3. To compare different binarization techniques and identify their strengths and weaknesses.
4. To implement serial and parallel processing approaches based on the size of the document image for faster execution.
5. To provide a reliable solution that can handle different types of degradation, such as stains, smudges, and uneven lighting.

II. LITERATURE SURVEY

Binarization is a technique used to separate the region of interest, such as text, from the background in an image, and it represents one of the fundamental methods of image segmentation. This process involves converting a grayscale image into a binary black and white image. Thresholding is a widely used tool in image segmentation, including global, local, and various automatic thresholding methods. In this paper, we have summarized some challenges and difficulties in the field of document image binarization [1].

In this paper, we propose a classification framework to combine different thresholding methods and produce better performance for document image binarization. Different binarization methods may create different corresponding binary image. Some binarization methods perform superior on certain kinds of document image, while others create better results for other kinds of document images. By combining different binarization techniques, better performance can be achieved with carefully analysis [2].

Image binarization is an important aspect of image analysis, such as scene text detection and medical image analysis. Especially in the field of document image processing, binarization has a wide range of applications as a basic method of digital image processing, including text recognition, document image segmentation, image morphological processing and feature extraction. This paper seeks to conduct a comprehensive review and synthesis of prevalent methods for document image binarization within an open research framework [4].

The main objective of this paper is to evaluate the different image binarization techniques to find the gaps in existing techniques. Document binarization is typically execute in the preprocessing phase of several document image processing associated fields such as optical character recognition (OCR) and document vision retrieval. Image binarization exchanges a picture up to 256 gray levels to a black and white picture [5].

III. METHODOLOGY

Description of Research Design:

Binarization is an active research area in the field of Document Image Processing. Binarization converts grey image into binaries image. Document image binarization is the most important step in preprocessing of scanned documents to save all or maximum sub components such as text, background and image [1].

In the 1960s, research on threshold segmentation primarily focused on global thresholding, local thresholding, and adaptive thresholding methods. However, these approaches encountered difficulties in effectively handling complex images with uneven pixel distributions and noise interference [4].

The main objective is to evaluating the short comings of algorithms for degraded image binarization [6].



Data Collection:

For this research, the datasets used to evaluate and compare various document image binarization techniques are sourced from previous research papers, specifically the DIBCO (Document Image Binarization Contest) and HDIBCO 2016 (Historical Document Image Binarization Contest 2016) datasets.

The DIBCO dataset consists of images of degraded documents, with challenges such as noise, low contrast, and uneven lighting, commonly encountered in real-world scanning conditions. This dataset is widely used in the research community for assessing binarization algorithms.

The HDIBCO 2016 dataset, on the other hand, focuses on historical document images, which present additional difficulties due to faded text, intricate background patterns, and distortions typical of older documents. By utilizing these datasets from previous research, this study aims to thoroughly evaluate the proposed binarization techniques across a range of document types, from modern to historically degraded. This approach allows for meaningful comparisons with existing methods and demonstrates the effectiveness of the proposed technique.

Data analysis:

The analysis was carried out on the grayscale images obtained after converting the Original colour documents into shades of grey. Following this, two binarization techniques—global and local—were applied to assess their effectiveness in document image binarization. For global binarization, a single threshold was used across the entire image to separate foreground text from the background. Local binarization applied different thresholds to smaller regions of the image, allowing for more flexibility in handling documents with uneven illumination.

Overall, the analysis indicated that while global binarization was faster and simpler, local binarization provided superior results in terms of accuracy and visual quality, especially in cases with significant document degradation or inconsistent lighting. The trade-off between speed and accuracy was a key consideration in the performance evaluation of both methods.

Implementation

1. Preprocessing

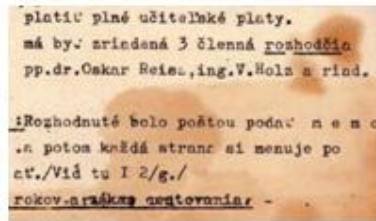
Pre-processing is an essential step in document image binarization that enhances the quality of the input image before applying the binarization algorithm. It involves operations like resizing, noise removal, and grayscale conversion. These steps help reduce distortions such as blur, uneven lighting, and smudges, making the text regions more distinct and improving the accuracy of binarization.

1.1 Grayscale Conversion:

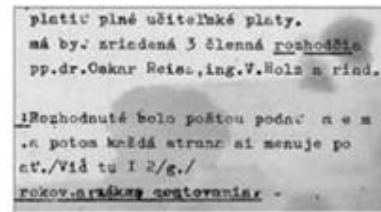
Grayscale conversion is essential before binarization because it reduces image complexity by eliminating colour, allowing the algorithm to focus solely on intensity differences. Binarization works by classifying pixels as either black or white based on their intensity, so having a single-channel grayscale image ensures more accurate thresholding and reduces processing time. It also helps in handling variations in lighting and background noise more effectively.

Most binarization algorithms are designed to work on grayscale images, so RGB images are typically converted to grayscale beforehand. Although simple conversion techniques like the luminosity formula ($g = 0.21r + 0.72g + 0.07b$) are commonly used, they can sometimes reduce the contrast of colored foreground text. This issue becomes more noticeable when the image contains multiple ink colours in the foreground [6].





(a) Degraded document



(b) Grayscale document

IV. TRADITIONAL BINARIZATION TECHNIQUES

4.1. Global threshold method

The Otsu algorithm, developed in 1979, is a prominent method of a global thresholding technique. The algorithm aims to determine an optimal threshold value, denoted as T , by analysing the grayscale properties of an image. This process involves partitioning the image into foreground and background segments. The objective is to minimize the gap between the two segments while maximizing the difference between them. The difference in grayscale distribution serves as a measure of the contrast between foreground and background, with a larger difference indicating an easier segmentation. The Otsu algorithm is also commonly known as the maximized difference between classes method [4]. The optimal threshold for the desired image is the value that maximizes the gap between categories, and it can be expressed as follows:

$$T' = \arg \max_{0 \leq T \leq L} \omega_0(T)\omega_1(T)(\mu_0(T) - \mu_1(T))^2 \quad (1)$$

we represent the image pixel in the grey level of the image, the image has L -order grey level, $\omega_0(T)$ and $\omega_1(T)$ are the probability distribution of the target and background when the threshold value is T , $\mu_0(T)$ and $\mu_1(T)$ represent the average grey value of the pixel of the target and background, respectively, if the pixel value of the input image is greater than T' . The pixel value is set to white, or otherwise it is black [1].

The Otsu algorithm partitions the entire image based on a single threshold, allowing for the determination of the optimal threshold for the image at once.

4.2. Local threshold method

The Niblack algorithm was developed to address the limitations of a fixed threshold by introducing a local binarization method. This approach involves utilizing a local window to calculate the mean and standard deviation within a small neighboring domain of each pixel. These values are then used to adjust the threshold for binarizing the image. The threshold calculation formula is expressed as follows:

$$T = m + k * s \quad (2)$$

Where, m represents the average grey value of pixels in the local area, s represents the standard deviation, and k is a constant, a correction factor, which can be adjusted according to the foreground and background conditions of the image [4].

Niblack performs better than other local binarization methods in grey images with low contrast, noise, and uneven background intensity [1].



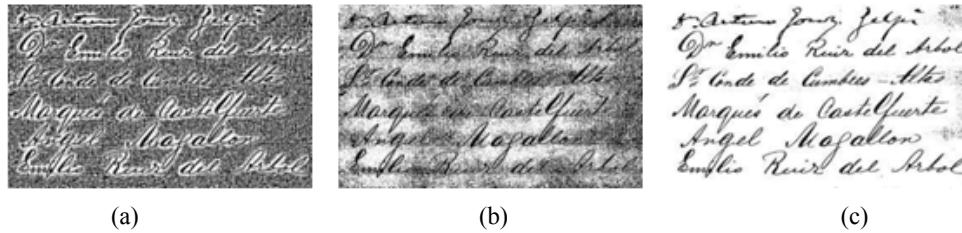


Figure 3. Niblack's binarization results

Classification	Algorithm	Description	Performance
Global Threshold	Otsu [31]	The grey level corresponding to the maximum inter-class variance is selected as the global threshold.	Low complexity and fast operation. However, it cannot handle complex degraded images and is suitable for processing high-quality document images.
Local Threshold	Niblack [32]	It calculates the mean and standard deviation of the pixel within a local window of anobvious noise can be seen in the output image, laying the foundation for local binarization methods.	The processing time has increased, and the foreground region. binarization image, which greatly increases the foreground region.

Table 1. Traditional binarization techniques (1).

V. CONCLUSION

This paper focused on the degraded document binarization technique. Document binarization is an important application of vision processing. The main objective of this paper is to evaluating the short comings of algorithms for degraded image binarization. It has been found that each technique has its own benefits and limitations; no technique is best for every case. And there is no universal binarization method that works for all types of document images.

This paper presents an overview of two binarization techniques. The evaluation shows that traditional algorithms work well for images with clean or simple backgrounds but are less effective on documents with complex backgrounds, especially those with mixed noise or heavy degradation. Considering the differences in characters and numerals across various languages, future work should aim to improve the cross-language adaptability of binarization methods so they can better handle different writing styles and stroke patterns. In near future we will propose a new algorithm which will use more reliable methodology to enhance the work.

REFERENCES

- [1] Yang, Z.; Zuo, S.; Zhou, Y.; He, J.; Shi, J. "A Review of Document Binarization: Main Techniques, New Challenges, and Trends." *Electronics* 2024, 13, 1394. <https://doi.org/10.3390/electronics13071394>
- [2] Bolan Su12*, Shijian Lu2+ and Chew Lim Tan. "Combination of Document Image Binarization Techniques". National University of Singapore, September 2011
- [3] Rohithkumar A, Rajath T ,Bipin Nair B.J, Dr.N.Shobha Rani. "Document Image Binarization" December 2020, 1533-9211
- [4] Zhengxian Yang , Rui Zhang , Yanxi Zhou , Jinlong He , Jianwen Shi , Shikai Zuo , "Summary of Document Image Binarization" Xiamen University of Technology, Fujian, China Date: 22 January 2024
- [5] Tarnjot Kaur Gill . "Document Image Binarization Techniques- A Review", *International Journal of Computer Applications* (0975 – 8887) Volume 98– No.12, July 2014
- [6] Chris Tensmeyer1 • Tony Martinez. "Historical Document Image Binarization: A Review" © Springer Nature Singapore Pte Ltd.16 May 2020. <https://doi.org/10.1007/s42979-020-00176-1>



[7] Wei Xiong, Lei Zhou, Ling Yue, Lirong Li and Song Wang. "An enhanced binarization framework for degraded historical document images". Xiong et al. EURASIP Journal on Image and Video Processing (2021) 2021:13 :xw@mail.hbut.edu.cn; songwang@cec.sc.edu

