

Digital Enhancement of Palm Leaf Manuscript Images using Normalization Techniques

Zhixin Shi, Srirangaraj Setlur and Venu Govindaraju

Center of Excellence for Document Analysis and Recognition (CEDAR)

State University of New York at Buffalo, Amherst, USA

Email: zshi@cedar.buffalo.edu

Abstract

Palm leaves were one of the earliest forms of writing media and their use as writing material in South and Southeast Asia has been recorded from as early as the fifth century B.C. until as recently as the late 19th century. Palm leaf manuscripts relating to art and architecture, mathematics, astronomy, astrology, and medicine dating back several hundreds of years are still available for reference today thanks to many ongoing efforts for preservation of ancient documents by libraries and universities around the world. Palm leaf manuscripts typically last a few centuries but with time the palm leaves degrade and the writing becomes illegible to be useful in any form. Image processing techniques can help enhance the images of these manuscripts so as to enable retrieval of the written text from these degraded documents. In this paper we propose a transform based method for enhancing digital images of palm leaf manuscripts. The method uses a dynamically selected pivoting background color in a linear transform to enhance the legibility of the foreground text. Then a combination of two other image processing algorithms viz., histogram normalization and background normalization developed in our earlier research, are applied to the transformed image. The algorithms can be mathematically combined into one or two transformations for computational efficiency. The method is tested on a set of palm leaf images from various sources and the preliminary results show significant improvement in readability. The techniques can also be used to enhance images of ancient, historical, degraded papyrus and paper documents.

Keywords:

Image enhancement, image processing, document pre-processing, document recognition, historical documents, color system transformation, OCR.

1 Introduction

Palm leaves have been a popular writing medium for over two thousand years in South and Southeast Asia. Use of palm leaves for recording literary and scientific texts have been reported from about the fifth century B.C., with the oldest existing documents dating from about the second century A.D. Palm leaf manuscripts are produced from two main types of palms: palmyra and talipot. The manuscripts are typically created by using a metallic stylus to etch letters into the dried leaf and enhancing the contrast and legibility of the script by applying lampblack or turmeric mixed with aromatic oils chosen for their insect repellent qualities.

A survey by the Institute of Asian Studies, Chennai, India indicates that there are still about a hundred thousand palm leaf manuscripts surviving in South Indian repositories alone with many more scattered across India, Nepal, Myanmar, Laos, Thailand, Cambodia and other South-east Asian Countries. These manuscripts contain religious texts and treatises on a host of subjects such as astronomy, astrology, architecture, law, medicine and music. In the past, Indian kings, temple authorities, and other concerned individuals ensured that the oldest manuscripts were ritually disposed only after they had been copied onto new palm-leaves. When this age-old cycle was broken in the 19th century, the remaining corpus of palm-leaf manuscripts and the knowledge contained in them began a long slide into obscurity and destruction. Most of these palm-leaves are nearing the end of their natural lifetime or are facing destruction from elements such as dampness, fungus, ants and cockroaches. This has spurred many new preservation projects to protect these valuable historical documents.

Efforts, funded by many foundations, universities and other institutions, are now underway for recovering and preserving these valuable palm leaves. Besides the many programs for preserving the manuscripts in their physical form, scanning and digital photograph imaging have been used to preserve their content and current appearance for future studies. Despite the availability of advanced photography and scanning equipment, natural aging and deterioration have rendered many palm leaf images unreadable. The original leaves are aged, leading to deterioration of the writing media, with seepage of ink and smearing along cracks, damage to the leaf due to the holes used for binding the manuscript leaves and dirt and other discoloration. The process of capturing a digital image of the leaves also presents some difficulties. In order to best preserve the fragile originals, the digital images are sometimes captured by using digital cameras instead of platen scanners. Leaf manuscripts cannot be forced flat and the light source for digital cameras is usually uneven. These factors lead to a very poor contrast between the foreground text and the background. Digital image processing techniques are necessary to improve the legibility of the manuscripts.

Previous image enhancement algorithms for historical documents have been designed primarily for segmentation of the textual content from the background of the images. An overview of the traditional thresholding algorithms for text segmentation are given in [Leedham *et al.*, 2002] which compares three popular methods, namely Otsu's thresholding technique [Otsu, 1979], entropy techniques proposed by Kapur *et al.* [Kapur and A.K.C.Wong, 1985] and the minimal error technique by Kittler and Illingworth [Kittler and Illingworth, 1986]. Another entropy-based method specifically designed for historical document segmentation [C.A.B.Mello

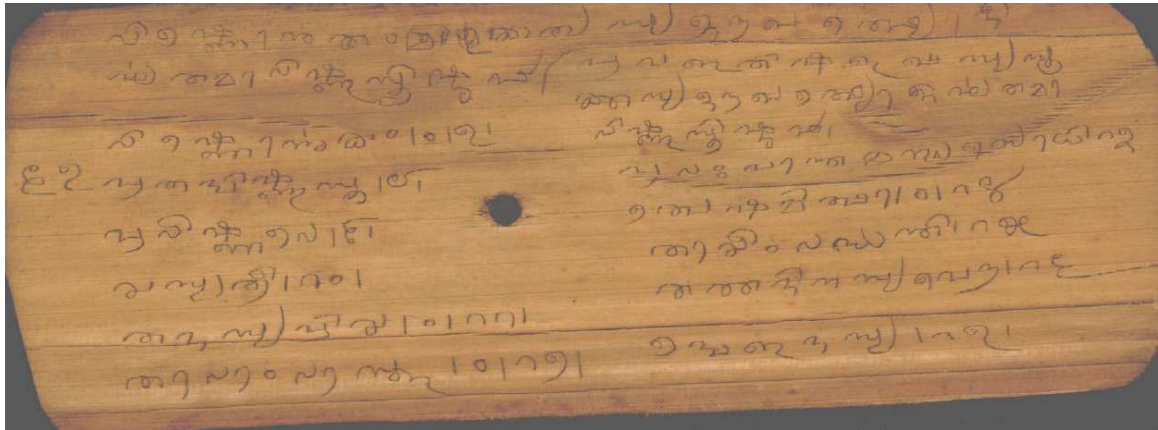


Figure 1: A sample palm leaf manuscript.

and R.D.Lins, 2000] deals with the noise inherent in the paper especially in documents written on both sides. Tan et al. presented methods to separate text from background noise and bleed-through text (from the backside of the paper) using direct image matching [Wang and Tan, 2001] and directional wavelets [Wang *et al.*, 2003]. These techniques are designed mainly as preparation stages for subsequent OCR processing. Other methods for historical image enhancement are driven by the goal of improving human readability while maintaining the original "look and feel" of the documents [C.A.B.Mello and R.D.Lins, 2002]. These methods do not produce satisfactory results in processing these palm leaf manuscripts since the contrast between the foreground and background is typically low and the color intensity of the background varies throughout the image.

In this paper we propose a transform based method for enhancing digital photograph images of palm leaf manuscripts. The method uses a dynamically selected pivoting background color in a linear transform to enhance the legibility of the foreground text. Then a combination of two other image processing techniques for histogram normalization and background normalization developed in our earlier research [Shi and Govindaraju, 2004] are applied to the transformed image. The algorithms can be mathematically combined into one or two transformations for computational efficiency. The method was tested using a set of palm leaf manuscript images from publicly available sources and the preliminary result shows significant improvement in readability.

In section 2 we present our enhancement algorithms. We present experimental results in section 3 and conclusions in section 4.

2 Proposed Techniques for Image Enhancement

The color of the treated palm leaves is light brown when the leaves are ready for *writing*. Damage to and deterioration of palm leaves are usually the result of staining, mechanical damage, splitting and cleavage, and insect and rodent activity. Palm leaf is susceptible to desiccation, losing its flexibility and becoming brittle. In many cases this dryness is treated by reapplying oil, which has a darkening effect if done too often. The lignified cells are particularly susceptible to degradation and discoloration. These processes often result in uneven background coloration across the image and the darkening of the background which reduces the contrast between the foreground text and the background color of the palm leaf.

We first target the problem of low contrast between the foreground and the background. We have designed a transform to bleach out the background colors of the leaf image to the extent possible. At this stage, the text color level is still close to the bleached background. To further enhance the contrast, we apply the histogram normalization algorithm on the transformed image to elevate the text color away from the background. Finally, to solve the uneven background problem, we apply a background normalization algorithm which smoothes out the background. The background normalization enhances the image, making it more legible to the eye as well as facilitating segmentation of the text from the non-text background.

2.1 Pivoting Color Bleaching Transform

The first step in the proposed process is to identify a base color for the background in order to wash out the background colors. Since the background colors on different leaf manuscripts are different due to differences in age or material, the background color has to be dynamically determined for each individual leaf. The simplest way of determining the background color is by calculating a color histogram. Our assumption is that the most dominant colors from a leaf are from the background. We first locate a range for the most frequently occurring colors on a leaf, then take the mean of the colors in the range as our background base color.

After calculating the base background color as (r_0, g_0, b_0) , we design a linear model in terms of the following transform:

$$L = \frac{R}{r_0} + \frac{G}{g_0} + \frac{B}{b_0} \quad (1)$$

This transforms the background colors to a range around number 3. Re-scaling by a factor of 88, the background colors are transposed to a range near 255. This moves all the other darker colors to lower levels and colors lighter than the background are transformed to have values greater than 255. A transformed grey-scale image is constructed by truncating the above values at corresponding pixel positions to the normal range of 0 to 255.

The transform in (1) is constructed by putting the contrast contributions from all the R, G and B color component channels together. In each channel, the component color of the background is mapped to 1 and other darker colors will have mapped values smaller than 1. The above

scaling and truncating operations are equivalent to a washing-out process for removing the predominant background color and lighter colors i.e. colors which we would want mapped to white in a binarization process. Although the transform works very well in the case of the palm leaf manuscript images, care should be taken in implementing the transform for any other application. In the extreme case when one of the RGB components of a background base color has a value 0, the transform will have undesirable effects on the image. However, in this extreme case, the component channel which creates the problem does not have any meaningful contribution to the contrast (background is so dark, that there would be no text visible in the foreground). In this extreme case, the corresponding term in (1) for that component channel can simply be omitted.

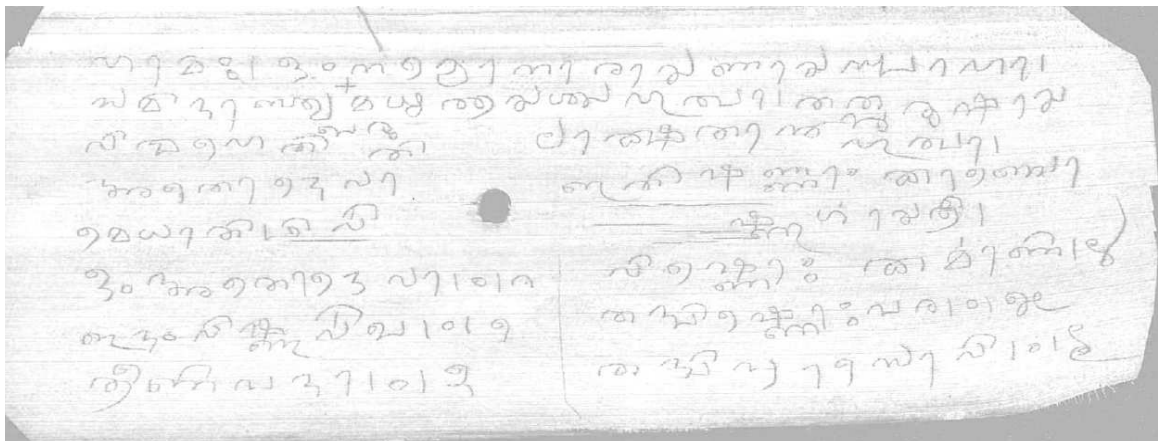


Figure 2: Result of the color-bleaching transform applied to the palm leaf image in Figure 1.

2.2 Histogram Normalization

A grey-scale bleached or brightened image is created from the original color image using the color bleaching transform described in 2.1. The original background colors are mapped into a color range very close to white and all other colors darker than the estimated original background colors map to darker grey levels. Since the original background color in many of the aged leaf-manuscripts are so dark that the foreground text colors are very close to the background, the transformed grey-scale images - while brighter than the original dark image - need further enhancement to increase the contrast between the text and the background. A histogram normalization algorithm is applied to the transformed grey-scale images to effect this contrast enhancement.

The algorithm applied is as follows. The distribution of the grey levels is computed and a small percent of values at both ends of the grey spectrum (black and white) are cut off. The cut off pixel values are folded to the nearest cut-off levels. The pixel values are re-scaled to stretch the grey levels to the range of 0 to 255. The resultant image shows appreciable contrast enhancement

(See Figure3).

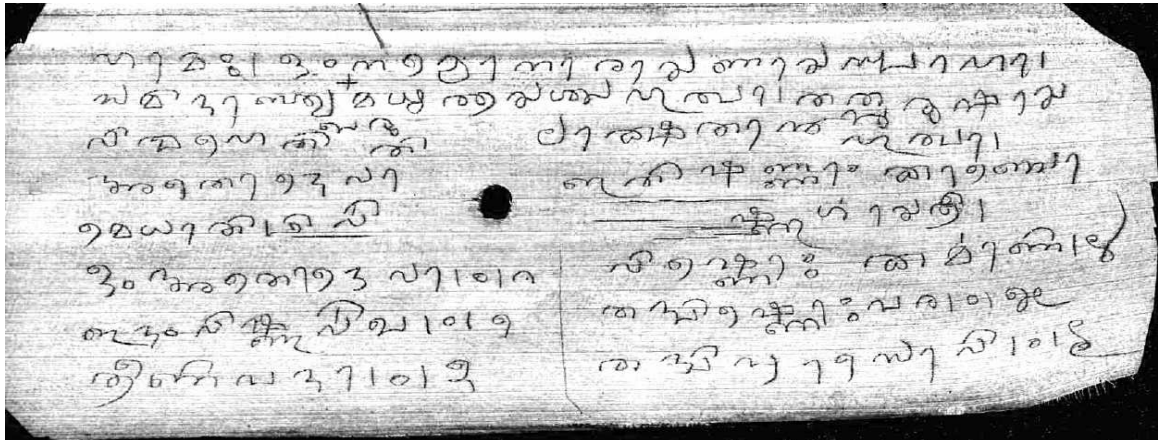


Figure 3: Histogram normalization further enhances the contrast.

2.3 Background Normalization

The problem of uneven background color intensity across the image is often seen in historical document images. We propose a new background normalization method building on our earlier work [Shi and Govindaraju, 2004] using a nonlinear approximation. We approximate the uneven background by a nonlinear curve that best fits the background color values. For efficiency, we compute a non-linear approximation for the image background color along each scanline, as shown in Figure 4.

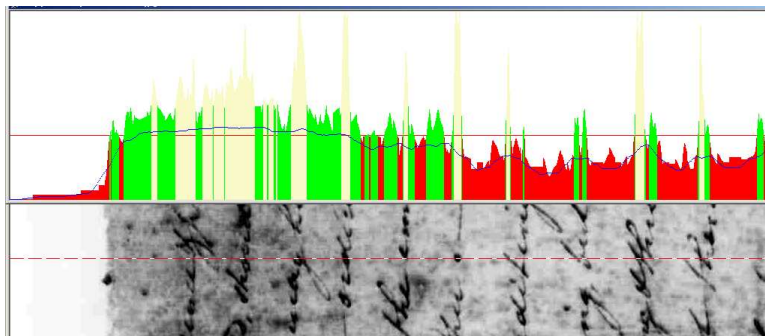


Figure 4: Scanline histogram and background approximation. The histogram of black pixel intensity along a selected scanline is shown above the grey-scale document. The horizontal line is the calculated average level. The curve is the approximation of the background.

Consider the histogram of foreground pixel color intensity. The histogram exhibits taller peaks

with higher variations at text locations. The non-text locations in the histogram, on the other hand, appear as a lower and less variant distribution. Another fact to be noted is that the number of background pixels in the document image is significantly larger than the number of foreground pixels for text. Based on the above observations, we first compute the mean or average level of the histogram. Then we use the mean level as a reference guideline to set a background level at each pixel position along the scanline. We scan the scanline from left to right. If the pixel level at the current position is less than the mean, then we take the value of the level for the next computation of our approximation and update a variable *previousLow* with the value of the current level. If the current level is higher than the mean, we retain the value in *previousLow* as the background level at the current location for the following computation of our approximation.

Thus far, we have set an approximate background level for each pixel position on the scanline. This rough background is not very accurate for two reasons. First, at the foreground pixel location, the foreground level is set using a previously remembered background level which may be used multiple times for a consecutive run of foreground pixels. Second, due to the low image quality even the real background pixels may be locally very distant from the desired globally dominant document background level. We therefore propose to use this roughly selected and estimated background (SEB) to obtain a better approximation of the normalized background level.

Using the selected and estimated background (SEB) pixel levels on a scanline, the approximation of the normalized document background level can be achieved in two ways. One approach is to use a *moving window* paradigm. At each pixel position, the approximated background level is computed from an average of the SEB values in the local neighborhood of the pixel position.

A better approximation is computed using a best fitting straight line in each of the above neighborhoods. At each position, we use all the SEB values in its neighborhood to find a best fitting line using least squares. Then the approximation value for the pixel position is calculated from the straight line corresponding to the position. If the final approximation of the background is a curve, the line segments going through each point on the curve at the corresponding pixel position form an envelope of the approximation curve.

Continuing the processing of the palm leaf manuscript image after the histogram normalization, the gray scale image can be further normalized using the nonlinear approximation described above. For any pixel at location (x, y) with pixel value z_{orig} , the normalized pixel value is then computed as

$$z_{new} = z_{orig} - z + c \quad (2)$$

where z is the corresponding pixel value on the approximated background; c is a constant fixed at some number close to the white color value 255. The resultant normalized image is shown in Figure 5.

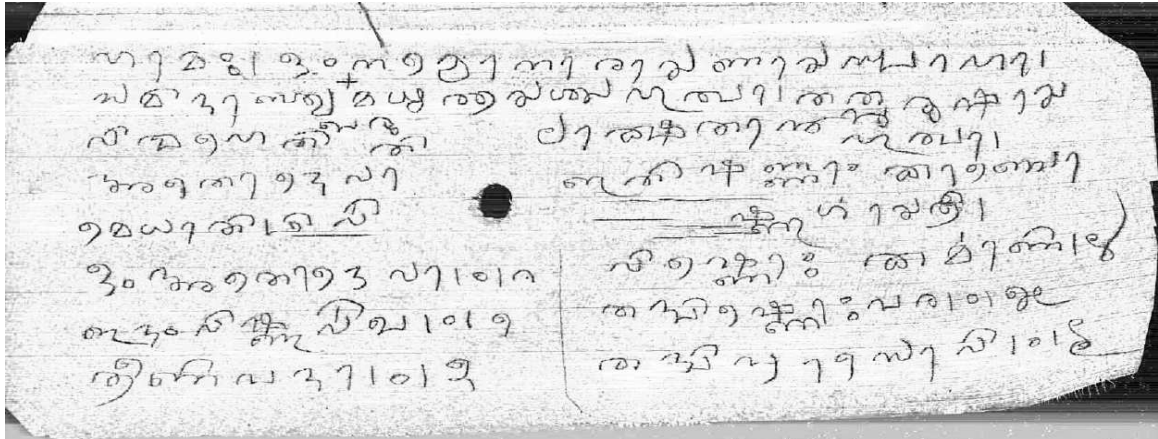


Figure 5: Background normalization evens out the media background for easier binarization.

3 Experiments

Over 100 historical palm leaf manuscript images were downloaded from many online repositories. A large number of the images have obvious uneven background problems and low contrast. Visual inspection of the enhanced images produced by the proposed techniques show a marked improvement in image quality for human reading.

Further, 20 images from the set were selected that could only be segmented (binarized) very poorly using a global threshold value. The proposed method successfully finds a better binarized image in all these cases. The binarized image produced from our example palm leaf manuscript is shown in Figure 6.

The techniques described in this paper were also used to process images of other historical documents such as papyrus manuscripts and aged, stained or otherwise discolored paper documents and were found to generate binarized images of very high quality with very little text degradation.

One of the goals of automatic document segmentation is to facilitate document OCR, and we propose to test the performance of in-house OCR systems on papyrus and aged paper documents containing Roman character text processed using the segmentation techniques proposed in this work. A longer term objective is to also process palm leaf manuscripts in Indic scripts using Indic OCR systems currently under development.

4 Conclusion

In this paper we present image enhancement techniques for historical palm leaf manuscript document images. The algorithm first converts the color image into a grey-scale image using

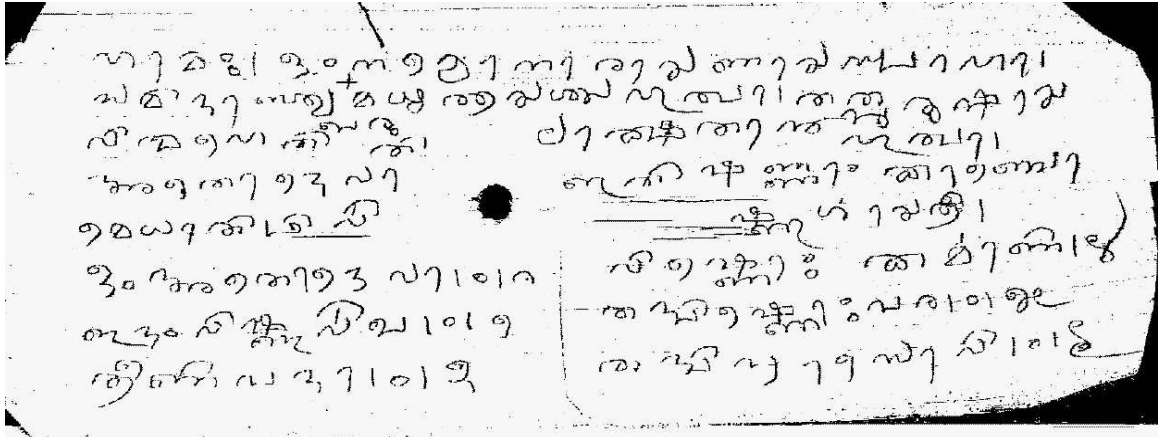


Figure 6: Binary image obtained using a simple global thresholding of the image in Figure 5.

a linear transform to brighten the text foreground by removing most of the background colors for better contrast. Then a nonlinear model to approximate the "flatness" of the background is applied after histogram normalization. The transformed leaf manuscript image is normalized by adjusting the pixel values relative to the background approximation. From our experiments and visual evaluation, the algorithm has been found to work successfully in improving readability of document images and produce high quality binarized images suitable for OCR, on not only palm leaf manuscripts but also on other aged and degraded documents such as papyrus and historical paper documents.

References

- [C.A.B.Mello and R.D.Lins, 2000] C.A.B.Mello and R.D.Lins. Image segmentation of historical documents. In *Visual2000, Mexico City, Mexico*, September 2000.
- [C.A.B.Mello and R.D.Lins, 2002] C.A.B.Mello and R.D.Lins. Generation of images of historical documents by composition. In *ACM Symposium on Document Engineering, McLean, VA, USA*, 2002.
- [J.Kittler and J.Illingworth, 1986] J.Kittler and J.Illingworth. Minimum error thresholding. *Pattern Recognition*, 19(1):41–47, 1986.
- [J.N.Kapur and A.K.C.Wong, 1985] P.K.Sahoo J.N.Kapur and A.K.C.Wong. A new method for gray-level picture thresholding using the entropy of the histogram. *Computer Vision, Graphics, and Image Processing*, 29:273–285, 1985.
- [Leedham et al., 2002] G. Leedham, S. Varma, A. Patankar, and V. Govindaraju. Separating text and background in degraded document images - a comparison of global thresholding techniques for multi-stage thresholding. In *Proceedings Eighth International Workshop on Frontiers of Handwriting Recognition*, September 2002.

- [N.Otsu, 1979] N.Otsu. A threshold selection method from gray level histogram. *IEEE Transactions in Systems, Man, and Cybernetics*, 9:62–66, 1979.
- [Shi and Govindaraju, 2004] Z. Shi and V. Govindaraju. Historical document image enhancement using background light intensity normalization. *17th International Conference on Pattern Recognition, Cambridge, UnitedKingdon*, 23-26 August 2004.
- [Wang and Tan, 2001] Q. Wang and C.L. Tan. Matching of double-sided document images to remove interference. In *IEEE Conference on Computer Vision and Pattern Recognition, Hawaii, USA*, 2001.
- [Wang *et al.*, 2003] W. Wang, T. Xia, L. Li, and C.L. Tan. Document image enhancement using directional wavelet. In *Proceedings of the 2003 IEEE Conference on Computer Vision and Pattern Recognition, Madison, Wisconsin, USA*, June 2003.