

READABILITY ENHANCEMENT OF HISTORICAL DOCUMENT IMAGES



MSC IN COMPUTER SCIENCE

University of Moratuwa, Sri Lanka.

Electronic Theses & Dissertations

www.lib.mrt.ac.lk

D.D.G.A.C.W. SAPARAMADU

Department of Computer Science and Engineering

University of Moratuwa

12/2007

READABILITY ENHANCEMENT OF HISTORICAL DOCUMENT IMAGES

D.D.G.A.C.W. SAPARAMADU



University of Moratuwa, Sri Lanka.
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

**This Dissertation was submitted to the Department of Computer
Science and Engineering of the University of Moratuwa in Partial
Fulfillment of the requirements for the Degree of MSc in
Computer Science**

**Department of Computer Science and Engineering
University of Moratuwa**

12/2007

Declaration

I, D.D.G.A.C.W. Saparamadu hereby declare that the work included in this dissertation in part or whole has not been submitted for any other academic qualification at any institution.

Dr. Chatura De Silva

C.W. Saparamadu

Supervisor



University of Moratuwa, Sri Lanka.
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

Abstract

The Department of National Archives of Sri Lanka posses a huge collection of historical documents and those can be considered as a national heritage as it provides documental evidences for country's history. The documents anyway are subjected to many degradations due to natural aging, acidification, manual handling and lack of maintenance. The department is in the process of digitizing the documents as a preservation strategy and faced the problem of illegible document images after the digitization process. This research is conducted to find the solutions for this problem and particularly paid attention on ink-bleed-through problem, which occurs in aged, double sided documents.

We have evaluated many traditional image enhancement techniques like contrast stretching, histogram equalization and also paid attention on some other novice technologies like k -Means clustering and Principle Component Analysis. In order to conduct the experiments, we have used Image Processing Toolkit of Matlab and written an application to explore the performance of algorithms in depth. Further we have extended our experiments even to many different color spaces and results were analyzed in detail. This flexible approach has helped us to practically evaluate an existing algorithm with actual document images and to understand the pros and cons of each approach.

We found that some traditional image enhancement techniques are contributing positively on historical document image enhancement, while some others are not. The application of commonly used image registration based techniques are found not feasible in real-time environment and certain other techniques like color segmentation seems capable of providing a reasonable output within a very short time. k -Means clustering coupled with Principle Component Analysis is yet another combined solution and its performance and the limitations are described.

Based on the analysis, we have reached certain conclusions and some recommendations were made for a real-time historical document image enhancement strategy.

Acknowledgements

Throughout this research, the Department of National Archives of Sri Lanka has extended their fullest corporation and without it, this work not had been successful. The Director of National Archives, Dr. Saroja Wettasinghe and Planning Assistant, Mr. Chaminda Mihindukula Suriya were always accessible and their cordial support and guidance helped immensely to better understand the problem. They were able to supply us with necessary document images within a very short time and put lot effort in finding many ink-bleed-through related examples as that was the most challenging task that we had to address.

Project supervisor, Dr. Chatura De Silva has shown lot of flexibility in consulting him either personally, through emails or even via phone as it saved lot of time and effort and led us in the correct direction. Once we were in desperate need of certain research material, his intervention helped to contact authors of the material and get certain things clarified. I would really like to appreciate his support extended to date from the day the project has initiated.

Dr. Damith C. Rajapakse of National University of Singapore is appreciated warmly as he personally intervened in getting certain research material and Dr. Sanath Jayasena also helped immensely in this regard. Another big thank is deserved by Dr. Sanath Jayasena for his relentless efforts on scheduling progress meetings and guiding us on a methodical project progress.

Table of Contents

Abstract.....	i
Acknowledgements.....	ii
Table of Contents.....	iii
List of Figures.....	v
List of Tables.....	vi
List of Symbols, Notations, Abbreviations and Acronyms.....	vii
1 Introduction.....	1
1.1 Overview of the Research problem.....	1
1.2 Overview of the Method of Study.....	2
1.3 Previous Work.....	3
2 Literature Review.....	6
2.1 Thresholding Techniques.....	6
2.2 Non-Blind approaches.....	7
2.3 Blind Approaches.....	7
2.3.1 Directional Wavelet Approach.....	8
2.3.2 Edge Detection Approach.....	8
2.3.3 Non-supervised segmentation approach.....	8
2.3.4 Advantages of Blind approaches.....	10
2.4 Taxonomy of solutions for ink-bleed-through problem.....	11
2.5 Color Segmentation based Image Analysis.....	11
3 Research Methodology.....	13
3.1 Overview of method of study.....	13
3.2 Traditional Image Enhancement Techniques.....	14
3.2.1 Histogram Equalization.....	14
3.2.2 Adaptive Histogram Equalization.....	15
3.2.3 Contrast Stretching.....	16
3.2.4 Decorrelation Stretching.....	17
3.2.5 Spatial and Frequency domain filters.....	18
3.3 Thresholding Techniques.....	18
3.3.1 Global Thresholding.....	19
3.3.2 Adaptive Thresholding.....	19
3.3.3 Role of Illumination in Thresholding.....	20
3.4 Image Registration based ‘Non-Blind’ approaches.....	20
3.5 Single Image used ‘Blind’ approaches.....	21
3.5.1 Principle Component Analysis (PCA).....	22

3.5.2	<i>k</i> -Means Clustering	22
3.5.3	Color spaces considered for <i>k</i> -Means Clustering and PCA	23
3.5.4	Limitations of some other ‘blind’ approaches	26
3.6	Color Image Segmentation	27
3.7	Color Reduction	28
3.8	The need of a flexible approach.....	29
4	Observations and Results	30
4.1	Thresholding results.....	30
4.2	Image Registration based image enhancement results.....	33
4.3	Principle Component Analysis (PCA)	36
4.4	<i>k</i> -Means Clustering	39
4.5	Combination of PCA with <i>k</i> -Means clustering.....	43
4.6	Color Image Segmentation	45
4.7	Color Reduction	46
5	Analysis and Discussion of Results	47
5.1	Thresholding based image enhancement	47
5.2	Image Registration based techniques.....	48
5.2.1	Drawbacks of Image Registration based Non-Blind approaches	48
5.3	Single image used ‘blind’ approaches.....	49
5.3.1	Principle Component Analysis (PCA)	49
5.3.2	<i>k</i> -Means Clustering	51
5.3.3	Combined approaches.....	52
5.3.4	Color Image Segmentation	53
5.3.5	Color Reduction.....	54
5.4	Impact of Pre-processing	54
5.5	Advantages of a flexible approach.....	55
6	Conclusions and Recommendations	58
6.1	Recommendations for method selection.....	60
	Appendix.....	61
	References.....	62

List of Figures

Figure 2.1 - Recursive Non-supervised classification technique.....	10
Figure 2.2 - Taxonomy of solutions for ink-bleed-through problem.....	11
Figure 3.1 - Histogram Equalization Example	15
Figure 3.2 - Contrast Stretching.....	16
Figure 3.3 - Contrast Stretching Example.....	17
Figure 3.4 - Decorrelation Stretching Example	17
Figure 3.5 – Histogram for Thresholding	18
Figure 3.6 - RGB Color space representation.....	24
Figure 3.7 - HSI color space representation (in color).....	25
Figure 3.8 - Lab color representation.....	26
Figure 3.9 - Document images with non-slanted writing style.....	27
Figure 3.10 - Example image with strong interfering background.....	27
Figure 3.11 - Spherical data region for RGB vector segmentation.....	28
Figure 4.1- Image used in PCA (in color).....	37
Figure 4.2 - Image used for Clustering (in color).....	39
Figure 5.1 - Outputs of PCA in Lab and CMYK color spaces	50
Figure 5.2 - Thresholding applied on PCA output.....	50
Figure 5.3 - Impact of pre-processing on PCA.....	51
Figure 5.4 - Impact of number of clusters (k value)	51
Figure 5.5 - Impact of Color space on k-Means clustering.....	52
Figure 5.6 - Impact on Tolerance value on RGB segmentation (in color)	53
Figure 5.7 – Impact of Pre-processing on Historical Document Images (in color).....	55
Figure 5.8 - Matlab Application ‘ <i>DocImage</i> ’	57
Figure 6.1 - Color picking for RGB segmentation (in color)	59

List of Tables

Table 4.1 - Thresholding Results with Image Pre-processing (in color)	31
Table 4.2 - Thresholding with different document images (in color).....	33
Table 4.3 - Point Pattern Matching Results	35
Table 4.4 - PCA in different color spaces (in color).....	38
Table 4.5 - PCA results of different images obtained using different colors (in color) ...	39
Table 4.6 - k -Means clustering in different color spaces (in color)	42
Table 4.7 - k -Means Clustering results (in RGB, $k=3$) (in color).....	43
Table 4.8 - PCA with k -Means clustering (in color)	44
Table 4.9 - RGB Color segmentation results (in color).....	45
Table 4.10 - Color Reduction results (in color)	46
Table 5.1 - Impact of PCA on clustering ($k=2$)	53



University of Moratuwa, Sri Lanka.
Electronic Theses & Dissertations
www.lib.mrt.ac.lk

List of Symbols, Notations, Abbreviations and Acronyms

PCA - Principle Component Analysis

HSI - Hue-Saturation-Intensity color space

RGB – Red-Green-Blue color space

CIE - International Commission on Illumination

Lab - CIE 1976 (L^* a^* b^*) color space or CIELAB

CMYK – Cyan-Magenta-Yellow-Black color space



University of Moratuwa, Sri Lanka.
Electronic Theses & Dissertations
www.lib.mrt.ac.lk