

A Computational Framework for Ink Bleed Suppression in Handwritten Document Images

Thesis submitted in partial fulfilment
of the requirements for the degree of

*Master of Science
in Computer Science and Engineering
by Research*

by

Shrikant Baronia

200702046

shrikant.baronia@research.iiit.ac.in



International Institute of Information Technology, Hyderabad

Deemed to be University

Hyderabad - 500 032, INDIA

June 2025

Copyright © Shrikant Baronia, 2025
All Rights Reserved

International Institute of Information Technology
Hyderabad, India

CERTIFICATE

It is certified that the work contained in this thesis, titled “A Computational Framework for Ink Bleed Suppression in Handwritten Document Images” by Shrikant Baronia, has been carried out under my supervision and is not submitted elsewhere for a degree.

Date

Advisor: Prof. Anoop M Namboodiri

To my late father,

whose unwavering belief in aiming higher continues to guide me, even in his absence.

To my mother,

whose silent strength and love have been my foundation.

To my wife,

for her endless support, patience, and partnership through every challenge.

To my twin daughters and son,

whose presence fills my life with purpose and joy.

To my Guru, Sri Sri Ravi Shankar ji,

an eternal source of inspiration — the embodiment of dedication, commitment, and boundless grace.

This work is as much yours as it is mine.

Acknowledgments

I'm very grateful to my advisor Prof. Anoop M Namboodiri, who has been a beacon of inspiration and guidance over the years. His kindness, patience, grace and empathy will forever stay with me and push me to be a better person. I am lucky to have been associated with him.

I remain deeply grateful to my fellow batch mates, Sankalp Khare and Varun Kuchibotla, for their invaluable support and encouragement throughout the completion of this thesis.

So much love to my parents who kept faith in me through thick and thin. They did not agree with many of my decisions yet they let me do this at my own pace, and their unwavering support played a huge part.

There are many more people who played a part in this journey and I am thankful to all of them.

Special mention to Dr. P. Kumaraguru for his efforts towards completion of the course for long-pending dual degree students like me.

IIIT Hyderabad has given me so much. The people, the campus and the times spent here will forever be a part of me, and I hope in some small way I will also remain a part of this great institution.

Abstract

Ink-bleed is a common form of degradation in handwritten document images, where ink from the reverse side or adjacent lines seeps into the visible content, impairing readability and adversely affecting downstream tasks such as Optical Character Recognition (OCR). Traditional image processing techniques often fall short in preserving the fine structure of handwritten strokes while removing bleed-through noise.

This thesis presents a machine learning-based computational framework for ink bleed suppression using a layer separation approach. The proposed method models the document image as a combination of content and bleed-through layers and employs a Dual Layer Markov Random Field (DL-MRF) architecture to learn the separation in a supervised setting. A synthetic dataset with controlled bleed artifacts was constructed for training, along with augmentation strategies to simulate various bleed intensities and patterns.

The model was evaluated on both synthetic and real handwritten document images. The results demonstrate significant improvement over traditional filtering techniques and baseline learning models, preserving content integrity while effectively reducing ink bleed.

This work contributes towards robust document image restoration, with applications in digital archiving, historical manuscript preservation, and pre-processing pipelines for handwriting analysis systems.

Publications

List of Publications:

- Shrikant Baronia and Anoop Namboodiri. Ink-Bleed Reduction using Layer Separation, *12th International Conference on Document Analysis and Recognition 2013*.

Contents

Chapter	Page
1 Introduction	1
1.1 Background and Motivation	1
1.2 Problem Statement	2
1.3 Research Objectives	2
1.4 Scope of the Work	2
1.5 Contributions	3
2 Related Work	4
2.1 Classical Image Processing Techniques	4
2.2 Signal Separation and Probabilistic Models	4
2.3 Paired Image Methods and User-Guided Approaches	5
2.4 Machine Learning-Based Document Enhancement	5
2.5 Summary	5
3 Document Image Model	7
3.1 Document Generation Model	7
3.2 Document Degradation Model	9
4 Estimation of Degradation	10
4.1 Complexity of the problem	10
4.2 Details of Algorithm	10
4.2.1 Input and Pre-processing	10
4.2.2 Background vs Non-Background Classification	10
4.2.3 Background and Ink Estimation	11
4.2.4 MRF Formulation	11
4.3 Synthetic Dataset	14
4.4 Results and Analysis	14
4.5 Summary	15
5 Applications and Future Work	16
5.1 Applications	16
5.1.1 Optical Character Recognition (OCR)	16
5.1.2 Digital Archiving and Restoration	16
5.1.3 Preprocessing for Handwriting Recognition	16
5.1.4 Historical Manuscript Analysis	17

5.1.5	Low-Cost Document Imaging	17
5.2	Future Work	17
5.2.1	End-to-End Deep Learning Model	17
5.2.2	Multi-Page and Double-Sided Bleed Modeling	17
5.2.3	Real-Time Document Enhancement	17
5.2.4	Generative Modeling and Diffusion-Based Restoration	17
5.2.5	Color and Multi-Script Documents	18
5.2.6	Semi-Supervised or Unsupervised Learning	18
5.3	Summary	18
	Bibliography	19

List of Figures

Figure		Page
3.1	(a) Document generation model: The combination of paper color with ink is modeled $I'^{\alpha}.BG'$. The degradation process (b) creates bleed of ink from one face to another and also changes the color of ink and paper.	8
4.1	Examples from the synthetic dataset showing varying α and β values.	14
4.2	Comparison of DL-MRF-SVM of [7] and the proposed method(From left to right input image, DL-MRF-SVM and proposed method). The proposed method is able to preserve the background and the finer details of foreground.	15

Chapter 1

Introduction

1.1 Background and Motivation

Handwritten documents are rich sources of cultural, historical, and linguistic knowledge. However, many of these documents suffer from physical degradation over time due to environmental exposure, ink chemistry, and paper aging. One of the most common and visually disruptive forms of degradation is **ink bleed-through**, where ink from the reverse side or adjacent lines seeps into the visible content of the document.

Ink-bleed not only reduces the legibility of handwritten texts but also severely affects downstream processing tasks such as document binarization, layout analysis, and optical character recognition (OCR). Early efforts to address this problem focused on classical image processing techniques such as thresholding [1], wavelet filtering [14, 18], and morphological operations. While effective in low-noise settings, these techniques often fail in cases of strong bleed where foreground and background blend non-linearly.

More advanced methods include unsupervised classification techniques [5], blind signal separation using Independent Component Analysis (ICA) [15], and Markov Random Fields (MRFs) [19]. These methods aim to preserve foreground strokes while suppressing bleed-through artifacts, but often require assumptions like linear mixing or paired front-and-back images, which limits their real-world applicability.

In recent years, the use of machine learning, particularly supervised learning with probabilistic modeling, has shown promise in improving ink-bleed removal while maintaining fine structural details. User-assisted models such as dual-layer MRFs [7] and iterative classification techniques [9] have further advanced the robustness of bleed separation methods, enabling them to work even on complex historical manuscripts.

Despite this progress, most existing methods treat ink-bleed suppression as a pixel classification task. They label pixels as either foreground, background, or bleed-through and attempt to suppress the bleed using post-processing steps. This approach often discards subtle ink patterns or introduces new artifacts. Our work takes a different route by modeling the contribution of background, ink, and bleed-through at every pixel using a layer separation framework guided by machine learning.

1.2 Problem Statement

The goal of this research is to develop a machine learning-based framework that can suppress ink bleed artifacts from handwritten document images without compromising the integrity of the foreground content. Specifically, we aim to:

- Separate the ink bleed layer from the clean content layer using supervised learning.
- Preserve the structural fidelity of handwritten strokes.
- Ensure that the cleaned images are OCR-friendly and visually appealing.

1.3 Research Objectives

This thesis aims to advance the field of document image restoration through the following key contributions:

1. Develop a computational framework based on **layered image decomposition** to effectively separate foreground content from degradation such as ink bleed-through.
2. Incorporate **probabilistic modeling** techniques to capture the underlying structure of document degradation and guide the separation of overlapping visual components.
3. Propose a document generation model that **explicitly accounts for the interaction between ink, background, and bleed-through**, enabling more interpretable and controllable restoration outcomes.
4. Demonstrate **visual quality enhancement and fine detail preservation** through the proposed method, particularly in challenging scenarios involving densely written or highly degraded handwritten text.
5. Validate the effectiveness and robustness of the framework using both **synthetically generated** and **real-world handwritten document images** exhibiting varying levels of ink bleed.
6. Perform a comparative evaluation against existing classical and machine learning-based baseline approaches to highlight the strengths and limitations of the proposed method.

1.4 Scope of the Work

This work focuses primarily on:

- Monochromatic (grayscale) handwritten documents.

- Synthetic and scanned data of moderate to high ink bleed.
- Offline image processing (real-time performance is out of scope).

The methodology does not currently extend to colored documents, bleed originating from multi-page scanning artifacts, or real-time restoration pipelines.

1.5 Contributions

The key contributions of this thesis are:

- A novel machine learning framework for layer-based ink bleed suppression in handwritten documents.
- A synthetic data generation strategy for training the network with diverse bleed patterns.
- A performance evaluation framework covering both image quality and OCR effectiveness.
- A comparison with state-of-the-art document cleaning methods to establish the advantages of the proposed model.

Chapter 2

Related Work

Ink bleed-through in scanned handwritten documents has been a long-standing challenge in document image analysis. Researchers have explored a variety of techniques ranging from traditional image processing to probabilistic graphical models and machine learning. This chapter provides an overview of the literature across three major directions: classical image enhancement techniques, probabilistic and statistical models, and modern machine learning approaches.

2.1 Classical Image Processing Techniques

Initial efforts in ink bleed removal used global or adaptive thresholding to enhance readability of scanned documents [1]. Morphological operations, background normalization [12], and directional wavelet filtering [14, 18] were also explored. These techniques often assume a clear separation between foreground and background, making them ineffective in cases where bleed overlaps significantly with primary text.

Principal Component Analysis (PCA) and color normalization were applied to reduce noise dimensionality before thresholding [5]. However, their success largely depends on the document's ink density, color, and uniformity. Such approaches may discard useful stroke-level information in severely degraded or historical documents.

2.2 Signal Separation and Probabilistic Models

More advanced approaches treat bleed removal as a signal separation problem. Independent Component Analysis (ICA) was used to decompose an RGB image into estimated foreground, background, and bleed layers [15]. These methods are particularly effective for scanned color documents but suffer under the assumption of linear mixture models.

Markov Random Fields (MRFs) have been widely used to model spatial dependencies in degraded documents [19]. By formulating ink bleed as a labeling problem, MRF-based models can enforce local

consistency. Tonazzini et al. further extended MRFs to blind signal separation using Expectation-Maximization (EM) algorithms for parameter estimation [16, 17].

Sharma [11] introduced a reflectance-based model considering bleed spread functions and corrected show-through using adaptive linear filtering, but it relied heavily on image registration and prior knowledge of ink-paper interaction.

2.3 Paired Image Methods and User-Guided Approaches

Several methods assume availability of both front and back scans. For example, Huang et al. [7] proposed a dual-layer MRF using user-specified pixels for initial labeling. This was extended by Lu et al. [9], who introduced interactive guidance for improved classification of ambiguous pixels.

Moghaddam and Cheriet [10] used real-paper scans as background priors and applied a three-layer diffusion model, relying on manually tuned parameters and clean background references. Hanasusanto et al. [6] employed a functional minimization framework using the Chan-Vese active contour model [3] to handle strong bleed overlapping with foreground strokes.

These approaches, while effective, often require extensive user input or tightly controlled datasets, limiting their generalizability.

2.4 Machine Learning-Based Document Enhancement

Recent advancements in supervised learning and probabilistic classification have brought new capabilities to document restoration. Support Vector Machines (SVMs) have been used with intensity-ratio features to classify pixels into background, ink, and bleed categories [7]. The SVM implementation from LIBSVM [4] has been widely adopted due to its flexibility and efficiency.

Graph-based energy minimization methods such as graph cuts and α -expansion moves [2] have been employed to optimize MRF label assignments. Comparative studies on energy minimization methods for image modeling further validated these techniques [13].

Still, most ML methods treat ink bleed suppression as a classification task, ignoring the physical interactions between ink, bleed, and background. Our work aims to bridge this gap by modeling these interactions using a discrete MRF formulation with label sets corresponding to alpha-beta combinations, trained using SVMs on synthetic data.

2.5 Summary

Although significant progress has been made in ink bleed removal, existing methods either assume strong priors (e.g., paired scans, background templates) or rely on pixel classification that may ignore

layer-level interactions. The concept of **layer separation**—distinguishing ink, bleed, and background contributions at the pixel level—offers a more interpretable and fine-grained approach. This thesis extends previous work by proposing a computational framework rooted in supervised learning and MRF modeling to jointly estimate ink and bleed layers, preserving content fidelity across a wide range of degraded documents.

Chapter 3

Document Image Model

This chapter discusses the complete document image model, we begin with the discussion of document generation, here we discuss the process that goes when we write something over a paper, how the ink gets absorbed and how the background of a paper affects the intensity values. This state of the document that gets generated is something that we want to recover. In Document degradation model we discuss the several degradation processes that happens over a period of time and how the intensity values of each pixel undergoes changes.

3.1 Document Generation Model

If we go back to the time when the text was actually being written on the paper, then initially the paper was blank and for every pixel we only had the Bg value. This is the original Bg value without any degradation or any application of ink or ink-bleed.

Now, as the text is written over this paper by the application of ink, then on the basis of the ink properties various changes occur to the color values of these pixels. Some of these changes are instantaneous and some occur over a period of time. The immediate effect that is noticed is the visibility of the ink. The ink's opacity describes how opaque or transparent an ink is and to what degree the ink allows or prohibits the transmission of light through it and how well the background on which the ink has been printed can be seen. It is this property of the ink which determines the effect of the background on the color values of ink pixels.

Now our paper will have some *ink* pixels and the rest are the original Bg pixels. Lets call the intensity values at these original *ink* and Bg pixels as ink' and Bg' . The color values of these *ink* pixels will lie in a range depending on the amount of ink applied, how much the ink penetrates into the paper and the paper's thickness. Here we are assuming that the ink will not penetrate enough to become visible on the opposite side and this process will happen at some later time. Lets define I' as the intensity referring to the Ink applied. We can define the intensity of *ink* at each pixel p as:

$$Ink'_p = I'^{\alpha'_p} \cdot Bg'_p$$

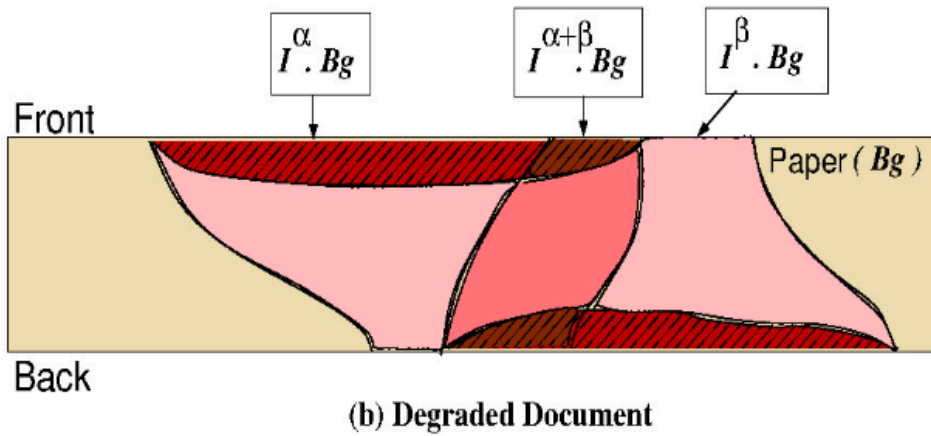
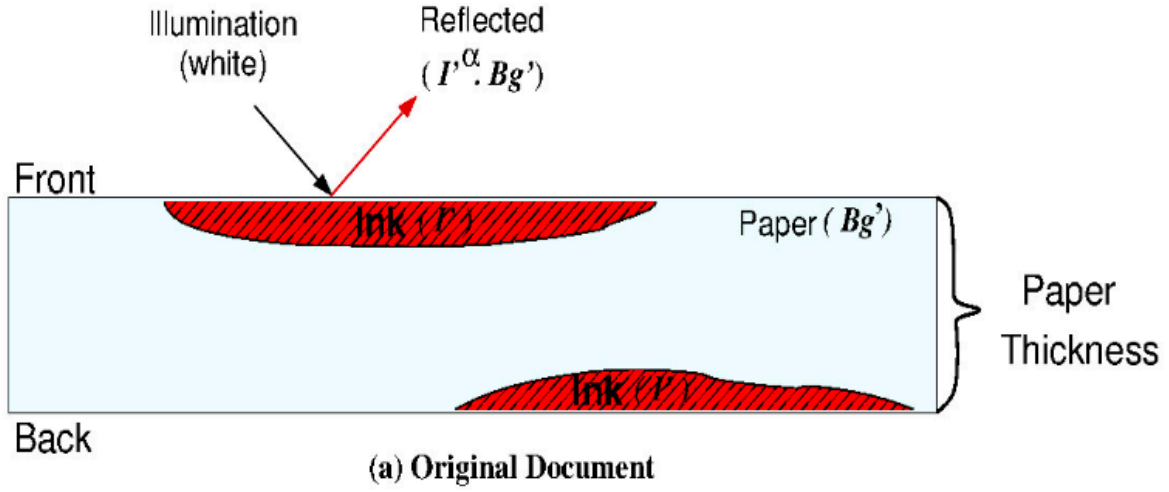


Figure 3.1: (a) Document generation model: The combination of paper color with ink is modeled $I'^{\alpha} . Bg'$. The degradation process (b) creates bleed of ink from one face to another and also changes the color of ink and paper.

Here, α'_p is defined for each pixel p and it refers to α to the amount of *ink* present at each pixel. Most of the ancient documents are written by the use of a quill, the ink intensity is highest at the first point of contact and then it keeps on decreasing until the quill is dipped again in the ink. Here we have defined α at each pixel to capture this phenomena. Note that I' remains constant for all the pixels, based on the assumption that the ink is not changed while writing a complete document. For Bg pixels $\alpha' = 0 = 0$ and for Ink pixels α' can lie in a range $(0,1]$.

3.2 Document Degradation Model

To achieve a degraded document image restoration of acceptable quality, a better understanding of the different types of image degradation processes is required.

The permanence of an ink pigment is the property which determines the extent to which an ink will retain its color strength and brightness with time or upon the exposure to light. This property of ink will change the value of I' over time. With the passage of time and due to variety of other factors including the ink's chemical makeup, the paper's physical and chemical construction and the amount of humidity in the environment, the value of α' will decrease over time. Due to the ageing of the document, the paper also undergoes changes, it tends to become brittle and the brightness of the paper decreases. Various chemical reactions and the oxidation of the cellulose results in the yellowing of the paper. Due to all these reasons the Bg' value also changes over time. Lets remove superscripts for the new values of α , Bg , Ink and I . Now for every pixel p , we have

$$Ink_p = I^{\alpha_p} \cdot Bg_p$$

One more prominent phenomena that will happen over time is that the ink will seep through the paper to interfere with the opposite side of the paper. In our generation model we only had background pixels and ink pixels. Now after the process of degradation we will have background pixels, ink pixels and some ink-bleed pixels.

For the rest of this paper lets define p and p' to denote front image pixel and the corresponding back image pixel, q and q' will be used to denote the pixels corresponding to 4-N neighborhood of p and p' respectively. After taking ink-bleed into consideration the above equation will become

$$Ink_p = I^{\alpha_p} \cdot Bg_p \cdot I^{\beta_p}$$

Here, β refers to the amount of ink that is coming from the other side of the paper as a result of ink-bleed. Theoretically, the value of β can lie in a range of [0,1] and its magnitude tells about the severity of ink-bleed in the document. Note that the value of β_p depends on the value of $\alpha_{p'}$ and similarly the α_p influences the value of $\beta_{p'}$. It's this Ink_p that we observe from the input images and from this we want to get an estimate of Ink'_p . Figure 3.1 explains the complete process of document degradation in detail.

Chapter 4

Estimation of Degradation

In this Chapter we present a detailed overview about the complexity of the problem, and the assumptions taken in order to do the estimation.

4.1 Complexity of the problem

From previous chapter we know that $\alpha, \beta \in [0, 1]$. Estimating the values of α and β for each pixel is essentially a continuous class labelling problem with infinite number of classes where each α and β can take any value from zero to one. All the existing techniques that we have seen in Section I, more or less tries to label each pixel as either a foreground, background or an ink-bleed pixel. Our problem can be solved by defining an MRF over the image pixels and posing the problem as a piecewise continuous restoration task, wherein the single-site and pair-site clique potential functions are defined. The details of which can be found at [8]. We simplify this problem and define a discrete set of values for α and β . The details of which is given later.

4.2 Details of Algorithm

4.2.1 Input and Pre-processing

The input to our system is two high resolution images(2KX2K) corresponding to the front and the back of a page. These images need not be aligned and in addition may also suffer from 3D surface variations. Global and Local alignment of images are done as is described in [7]

4.2.2 Background vs Non-Background Classification

Once the images are aligned, we classify all the pixels as either a background pixel or a non-background pixel. For this we take assistance from the user and few background and non-background(background+ ink-bleed) pixels are marked by the users. Two features C_p , and ρ_p are used

to train an SVM classifier. Here C_p refers to grayscale pixel intensity of the front image pixels and ρ_p is defined as the ratio of intensities $\frac{C_p}{C'_p}$. Huang et al. [7] used the ratio feature alone for the initial classification of pixels as foreground, inkbleed and background based on the assumption that the ratio value for foreground data is mainly distributed between zero and one, inkbleed larger than one and background is mainly around one. This approach classifies foreground data with ratio values close to one as background. Whereas the use of ratio features along with the intensity values classifies the background and non- background(foreground + inkbleed) data quite robustly. We used libsvm [4] with the Radial Basis Kernel.

4.2.3 Background and Ink Estimation

For all the pixels labelled as non-background, we try to estimate the expected background value at these points. This process can be called as background generation, which in itself is an ill-posed problem as the exact RGB composition for the background is completely lost. We used the patch based approach to estimate the background pixels. A small patch is taken from the background region with similarities to a patch around the target pixel and thereby approximating the background values. The ink that we observe is the Ink_p of equation 3, from this we want to get an estimate of I . For this we took the help of user-assistance and marked few strongest Ink pixels such that the corresponding back image pixels does not belong to foreground pixels. We took the average of all these user marked pixels as an estimate of I .

4.2.4 MRF Formulation

After the initial classification of each pixel as either a background pixel or non-background pixel, we need to estimate the value of α and β for each non-background pixel p , corresponding to the intensity of ink and bleed at that pixel.

Instead of dealing with the continuous values of α and β , we release the assumption of continuity and define discrete values so that $\alpha \in \{0, 0.25, 0.50, 0.75, 1\}$ and $\beta \in \{0, 0.25, 0.50, 0.75\}$.

Here we are assuming that β can never take the value of one. Note that it doesn't mean that bleed is always lighter than the foreground strokes, it only suggests that bleed can never be equal to the strongest foreground strokes, there can be cases where $\beta > \alpha$.

Now we have 20 combinations of $\langle \alpha, \beta \rangle$ that can be associated to each image pixel. Lets define our label set as $L = \{l_1, l_2, \dots, l_{20}\}$, corresponding to 20 different combinations of $\langle \alpha, \beta \rangle$ and $l_p \in L$ is the label assigned to pixel p .

With these discrete definitions of α and β , the problem can now be formulated as a discrete labeling MRF where each pixel, p is assigned some combination of $\langle \alpha, \beta \rangle, l_p$.

Our implementation of MRF is motivated by the work of Huang et al. [7], where instead of defining MRF for a trivial task of labelling foreground, background and ink-bleed, a more complex set of labels

L is modeled. The optimal label assignment of an MRF is found by minimization of below energy function:

$$E = E_d + \phi E_s$$

where E_d is the energy associated with the data terms, which is the likelihood of assigning a label $l_p \in L$ to each pixel p , and E_s is the energy associated with the smoothness term, which is the cost of assigning different label values to the neighboring pixels. The scalar weight ϕ is used to balance the two terms. The readers are directed to [8] for the details about the MRF formulation. We define the Dual Layer MRF similar to [7], and the details about intra-layer and inter-layer edges can be seen in [7]. The Dual Layer MRF is chosen due to its effectiveness in preserving the foreground strokes. The details about the data term and smoothness term are given below:

Data Term E_d : The definition for the data cost E_d is the same for both the front and back image. Using the estimation of background and ink from above section and using Equation defined in Chapter 3, the value of $\alpha + \beta$ can be calculated as

$$\alpha_p + \beta_p = \log_I \left(\frac{Ink_p}{Bg_p} \right)$$

Let us define $\lambda_p = \alpha_p + \beta_p$. A 3-dimensional feature vector with values $\langle \lambda_p, \lambda_{p'}, |\lambda_p - \lambda_{p'}| \rangle$ is defined, where $|a - b|$ denotes the absolute difference between two values a and b . This 3-dimensional feature vector is capable of capturing the amount of ink present on the front and back sides of the page and their differences.

Using these features, a set of twenty SVM classifiers are trained corresponding to twenty different classes based on the one-against-rest method. To obtain training samples for these SVMs, we carefully generated synthetic data corresponding to these twenty classes.

For the generation of synthetic data, we downloaded two images of handwritten text and two background images. The text is first rendered over the background images, and then different sets of synthetic images are generated by varying the values of α and β .

Each pixel p is then classified using these SVMs, and its similarities (or probabilities) S_i , for i ranging from one to twenty, are calculated similar to [7]. The data cost, E_d , for each label is defined as:

$$E_d(l^p = l_i) = \frac{\sum_{j=1}^{20} (S_j - S_i)}{19 \cdot \sum_{j=1}^{20} S_j}$$

Note that the data cost, E_d , from the above equation lies in the range of zero to one, and

$$\sum_{i=1}^{20} E_d(l^p = l_i) = 1.$$

Smoothness Term E_s : In the dual-layer MRF setup, each pixel p has edges with its 4-neighborhood (intra-layer edges), and there is also an edge with the corresponding back page pixel p' (inter-layer

edge). The same setup is defined for both the front and back images. We can write the smoothness term as:

$$E_s = \sum_{(p,q) \in \mathcal{N}} V_1(l^p, l^q) + \sum_{(p,p') \in \mathcal{M}} V_2(l^p, l^{p'})$$

Here, V_1 corresponds to the intra-layer edge costs and V_2 corresponds to the inter-layer edge costs. Both of these terms are weighted equally.

Intra-Layer Edge Costs: The intra-layer edge costs are defined as:

$$V_1(l^p = \langle \alpha_p, \beta_p \rangle, l^q = \langle \alpha_q, \beta_q \rangle) = |\alpha_p - \alpha_q|^2 + |\beta_p - \beta_q|^2$$

These costs enforce smoothness constraints on neighboring pixels by assigning higher costs to neighboring pixels with large variations in their corresponding values of α and β .

Inter-Layer Edge Costs: Inter-layer edge costs, $V_2(l^p, l^{p'})$, are defined as:

$$V_2(l^p = \langle \alpha_p, \beta_p \rangle, l^{p'} = \langle \alpha_{p'}, \beta_{p'} \rangle) = \delta\alpha_p\beta_{p'} + \delta\alpha_{p'}\beta_p$$

where $\delta\alpha\beta$ is defined in the following table:

$\beta \backslash \alpha$	0	0.25	0.5	0.75	1
0	0	$\frac{\omega}{4}$	$\frac{\omega}{3}$	$\frac{\omega}{2}$	ω
0.25	∞	∞	$\frac{\omega}{4}$	$\frac{\omega}{3}$	$\frac{\omega}{2}$
0.5	∞	∞	∞	$\frac{\omega}{4}$	$\frac{\omega}{3}$
0.75	∞	∞	∞	∞	$\frac{\omega}{4}$

In the above table, we enforce the constraint that the intensity of bleed at any pixel p is always less than the ink intensity at pixel p' , thereby assigning infinite cost to all these cases. As the value of α for a pixel p increases, this in turn increases the likelihood for a pixel p' to have a higher value of β . We have captured this notion by assigning an increasing cost for the cases when the value of β at pixel p' is relatively less compared to the value of α at pixel p .

The only two parameters in our MRF formulation are ϕ and ω . For our current work, we have empirically fixed the values of ϕ and ω to one.

Minimizing the Energy Function: For minimizing the energy function, we used the α -expansion move of graph cuts [?]. The dual layer MRF model is implemented by modifying Middlebury's MRF code provided by [?]. In all of our experiments, the energy minima is reached within four or five iterations.

4.3 Synthetic Dataset

We created a synthetic dataset by overlaying scanned handwriting samples on clean paper backgrounds and simulating varying degrees of ink bleed. The simulation uses controlled combinations of α (ink opacity) and β (bleed-through level) to model realistic degradation.

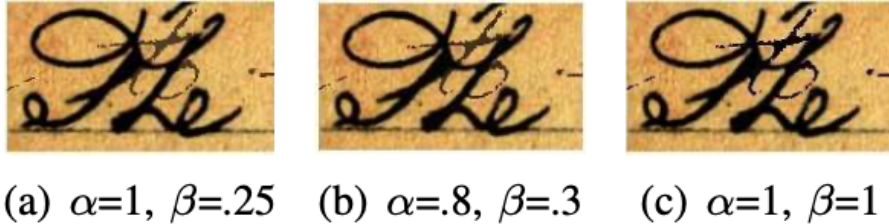


Figure 4.1: Examples from the synthetic dataset showing varying α and β values.

4.4 Results and Analysis

We ran our experiments over a dataset of Malayalam handwritten manuscripts. These documents are mostly 150 years old and contain conversations between the king and local groups. We compared our approach with the DL-MRF-SVM method of [7], which has shown better results than previous approaches.

One of the major contributions of this paper is the attempt to preserve background information and to extract the bleed even from the foreground strokes. As shown in Figure 3, our approach was able to recover the document in a better state than the corresponding DL-MRF-SVM of [7]. Figure 4 shows the evaluation on synthetic data and the failure cases of our algorithm. In fact, our algorithm fails every time the value of β exceeds 0.85.

As there is no ground truth available for the images, we manually evaluated the results. The error is calculated as the number of foreground words with undetected or missing strokes. We evaluated our results over a set of twenty images. Our method achieved a precision accuracy of 93.4%, whereas the DL-MRF-SVM gave an accuracy of 89.4%.

Precision is defined as:

$$\text{Precision} = \frac{W - W^F}{W + W^B}$$

where W is the total number of foreground words, W^F is the number of incorrectly classified foreground words, and W^B is the number of stroke-sized background or ink-bleed regions that were classified as foreground.



Figure 4.2: Comparison of DL-MRF-SVM of [7] and the proposed method(From left to right input image, DL-MRF-SVM and proposed method). The proposed method is able to preserve the background and the finer details of foreground.

4.5 Summary

Our results demonstrate that our model can generate outputs while preserving finer details. We were able to preserve most of the foreground strokes, and the background estimation is also of acceptable quality. Although better methods for background estimation may exist, we leave this as future work.

Our approach requires user markup for the initial classification of background versus non-background pixels and again for the estimation of ink. Consistent with most supervised learning approaches, we assume that the user-labelled data is correct.

The next chapter discusses broader applications and future directions.

Chapter 5

Applications and Future Work

This chapter discusses the practical applications of the proposed ink-bleed suppression framework and outlines possible future directions for extending this work. The ability to clean degraded handwritten documents has wide-ranging impact, especially in areas involving historical preservation, digital archives, and downstream machine learning tasks like OCR and handwriting recognition.

5.1 Applications

5.1.1 Optical Character Recognition (OCR)

One of the immediate benefits of ink-bleed suppression is improved accuracy in OCR. Bleed-through artifacts often confuse character segmentation and feature extraction stages, reducing recognition rates. Our method, by separating ink and bleed layers, produces cleaner inputs for OCR engines, as demonstrated in Chapter 4 with a substantial improvement in character-level accuracy.

5.1.2 Digital Archiving and Restoration

Many libraries and institutions hold valuable handwritten records that are centuries old. These documents often suffer from significant ink bleed and paper deterioration. The proposed framework enables automatic enhancement of such documents before digital storage or public access, improving readability and longevity without manual editing.

5.1.3 Preprocessing for Handwriting Recognition

Handwriting recognition models (especially deep learning-based ones) are sensitive to background noise and degradation. Preprocessing handwritten input using our bleed removal technique improves the quality of training data and inference, leading to more robust handwriting recognition pipelines.

5.1.4 Historical Manuscript Analysis

Philologists and historians analyzing ancient texts often require a faithful digital representation of original content. Layer separation allows scholars to view foreground text isolated from bleed artifacts, making textual transcription and interpretation more accurate.

5.1.5 Low-Cost Document Imaging

In low-resource settings where high-resolution scanners are unavailable, smartphone-captured document images often suffer from bleed artifacts. Integrating our framework into mobile document enhancement apps can significantly improve image clarity for education and record-keeping.

5.2 Future Work

5.2.1 End-to-End Deep Learning Model

The current framework uses an SVM classifier followed by MRF optimization. A promising future direction is to build an end-to-end deep neural network (e.g., U-Net or ResNet) that learns to perform bleed suppression directly from input images. Such models can be trained using synthetic data and validated on real scans for better scalability.

5.2.2 Multi-Page and Double-Sided Bleed Modeling

While this work models ink bleed from the reverse side of a single page, future work can extend to multi-page scans, where bleed comes from both the reverse side and adjacent pages. Cross-page alignment and multi-channel modeling may be used to learn such interactions.

5.2.3 Real-Time Document Enhancement

Deploying the framework in real-time systems such as mobile apps or scanning kiosks would require computational optimization. Future work could involve converting the current method to lightweight models suitable for edge devices.

5.2.4 Generative Modeling and Diffusion-Based Restoration

Recent advances in generative modeling, such as denoising diffusion models (DDPMs), could be explored to reconstruct clean document images from heavily degraded inputs. These models can be conditioned on estimated background or text priors and trained on paired synthetic data.

5.2.5 Color and Multi-Script Documents

This work focuses on grayscale handwritten documents. Extending the framework to color manuscripts, multi-ink documents, and multilingual scripts (e.g., Devanagari, Arabic, Chinese) would require additional adaptations in the feature space and data generation pipeline.

5.2.6 Semi-Supervised or Unsupervised Learning

Training with fully supervised labels can be limiting, especially for rare historical data. Future directions could include using unsupervised clustering, self-supervised pretraining, or weak labels to improve generalizability with minimal annotation.

5.3 Summary

This chapter discussed practical domains where the proposed bleed suppression framework can be applied, and identified several avenues for future improvement. The model's integration into OCR systems, archival platforms, and mobile scanners demonstrates its versatility. Future enhancements using deep learning and generative models hold strong promise for further performance gains and real-world deployment.

Bibliography

- [1] J. Bescos. Image processing algorithms for readability enhancement of old manuscripts. In *Intl Electronic Imaging Exposition and Conference*, pages 392–397, 1989.
- [2] Y. Boykov and V. Kolmogorov. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(9):1124–1137, 2004.
- [3] T. Chan and L. Vese. Active contours without edges. *IEEE Transactions on Image Processing*, 10(2):266–277, 2001.
- [4] C.-C. Chang and C.-J. Lin. Libsvm: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):27, 2011.
- [5] F. Drira, F. L. Bourgeois, and H. Emptoz. Restoring ink bleed-through degraded document images using a recursive unsupervised classification technique. In *Document Analysis Systems VII*, pages 38–49, 2006.
- [6] G. Hanasusanto, Z. Wu, and M. Brown. Ink-bleed reduction using functional minimization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 825–832. IEEE, 2010.
- [7] Y. Huang, M. S. Brown, and D. Xu. A framework for reducing ink-bleed in old documents. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–7. IEEE, 2008.
- [8] S. Z. Li. *Markov Random Field Modeling in Image Analysis*. Springer, 2009.
- [9] Z. Lu, Z. Wu, and M. S. Brown. Directed assistance for ink-bleed reduction in old documents. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 88–95. IEEE, 2009.
- [10] R. Moghaddam and M. Cheriet. Low quality document image modeling and enhancement. *International Journal on Document Analysis and Recognition*, 11(4):183–201, 2009.
- [11] G. Sharma. Show-through cancellation in scans of duplex printed documents. *IEEE Transactions on Image Processing*, 10(5):736–754, 2001.
- [12] Z. Shi and V. Govindaraju. Historical document image enhancement using background light intensity normalization. In *17th International Conference on Pattern Recognition (ICPR)*, volume 1, pages 473–476. IEEE, 2004.
- [13] R. Szeliski, R. Zabih, D. Scharstein, O. Veksler, V. Kolmogorov, A. Agarwala, M. Tappen, and C. Rother. A comparative study of energy minimization methods for markov random fields. In *European Conference on Computer Vision (ECCV)*, pages 16–29. Springer, 2006.

- [14] C. Tan, R. Cao, and P. Shen. Restoration of archival documents using a wavelet technique. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(10):1399–1404, 2002.
- [15] A. Tonazzini, L. Bedini, and E. Salerno. Independent component analysis for document restoration. *International Journal on Document Analysis and Recognition*, 7(1):17–27, 2004.
- [16] A. Tonazzini, L. Bedini, and E. Salerno. A markov model for blind image separation by a mean-field em algorithm. *IEEE Transactions on Image Processing*, 15(2):473–482, 2006.
- [17] A. Tonazzini, E. Salerno, and L. Bedini. Fast correction of bleed-through distortion in grayscale documents by a blind source separation technique. *International Journal on Document Analysis and Recognition*, 10(1):17–25, 2007.
- [18] Q. Wang, T. Xia, L. Li, and C. Tan. Document image enhancement using directional wavelet. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 534–537, 2003.
- [19] C. Wolf. Document ink bleed-through removal with two hidden markov random fields and a single observation field. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(3):431–447, 2010.