# Study of Approaches to Remove Show-Through and Bleed-Through in Document Images

**Pritish Sahu**

**Department of Computer Science and Engineering**

**National Institute of Technology Rourkela**

**Rourkela-769 008, Orissa, India**

# Study of Approaches to Remove Show-Through and Bleed-Through in Document Images

Thesis submitted in partial fullfillment

of the requirements for the degree of

## BACHELORS OF TECHNOLOGY

IN

## COMPUTER SCIENCE AND ENGINEERING

### BY

### PRITISH SAHU (107CS044)

Under the Guidance of

**Prof. Pankaj Kumar Sa**



Department of Computer Science and Engineering

National Institute of Technology Rourkela

Rourkela-769 008, Orissa, India

Department of Computer Science and Engineering

**National Institute of Technology Rourkela**

Rourkela-769 008, India.   www.nitrkl.ac.in

**Pankaj Kumar Sa**

Assistant Professor

May 17, 2011

# Certificate

This is to certify that the project entitled, STUDY OF APPROACHES TO REMOVE SHOW-THROUGH AND BLEED-THROUGH IN DOCUMENT IMAGES submitted by Pritish Sahu, Roll No: 107CS044 in partial fulfillment of the requirements for the award of Bachelor of Technology Degree in Computer Science and Engineering at the National Institute of Technology, Rourkela is an authentic work carried out by him under my supervision and guidance. To the best of my knowledge, the matter embodied in the thesis has not been submitted to any other university / institute for the award of any Degree or Diploma.

**Pankaj Kumar Sa**

# ACKNOWLEDGEMENT

This thesis has benefited in various ways from several people. Whilst it would be simple to name them all, it would not be easy to thank them enough.

We avail this opportunity to extend our hearty indebtedness to our guide **Prof. Pankaj Kumar Sa**, Computer Science Engineering Department, for their valuable guidance, constant encouragement and kind help at different stages for the execution of this dissertation work.

We also express our sincere gratitude to **Prof. A.K. Turuk**, Head of the Department, Computer Science Engineering, for providing valuable departmental facilities.

We would like to thank all our friends for helping us and would also like to thank all those who have directly or indirectly contributed in the success of our work.

Last but not the least, big thanks to NIT Rourkela for providing us such a platform where learning has known no boundaries.

**Submitted by :**

**Pritish Sahu**

**Abstract**

The wok implemented describes a study of approaches to restore the nonlinear life mixture of images, which occurs when we scan or photograph and the back page shows through. We generally see this to occur mainly with old documents and low quality paper. With the presence of increased bleed-through, reading and deciphering the text becomes tedious. This project executes algorithms to reduce bleed-through distortion using techniques in digital image processing. We study the algorithm knowing the fact that in images the high frequency components are sparse and stronger on one side of the paper than on the other one. Bleed-through effect and show-through effect was removed in one time processing, with no iteration. Here the sources need not require to be independent or the mixture to be invariant.Hence it is suitable for separating mixtures such as those produced by bleed-through.

# Contents

# List of Figures

# Chapter 1

# Introduction

# 1   Introduction

Image processing is any form of signal processing for that the input is an image, such as a photograph or video frame; the output of image processing may be either an image or, a set of parameters related to the image.

## 1.1   Image Processing

Image processing totally refers to processing digital images by means of digital computer.An image may be described as a 2-D function I.

$$I = f(x, y) \tag{1.1}$$

where x and y are spatial coordinates. Amplitude of f at any pair coordinates (x, y) is called the intensity I or gray value of the image. When spatial coordinates and amplitude values are all finite, wecall the image a digital image. Digital image processing may be classified into various sub-categories based on methods whose:[1]

- both input and output are images.
- inputs may be images, where as outputs are attributes extracted from those images.

Following is the list of different image processing functions based on the above two classes.

- Image Acquisition

- Image Enhancement

- Image Restoration

- Color Image Processing

- Multi-resolution Processing

- Compression

- Morphological Processing

- Segmentation

- Representation and Description

- Object Recognition

For the first seven categories the inputs and outputs are images where as for the other three the outputs are attributes of the input images. With the exception of image acquisition and display, most image processing functions are implemented in software.In image processing the technique that works well in one area can be inadequate with another, as it is characterized by specific solutions. It requires significant research and development to find the actual solution of a specific problem.[2]

From the above ten sub-categories of digital image processing, this thesis deals with wavelet and multiresolution processing and also image enhancement. Here, different types of wavelets are used and various inputs are restored using these methods. Image enhancement techniques can be divided into two broad categories:

1. Spatial domain methods, it refers to image plane itself and approaches in this category are based on direct manipulation in an image.

2. Frequency domain methods, which operate on the various kind transform(fourier transform,short fourier transform,wavelet transform) of an image.

The rest of the chapter is organized as follows. Document Image Processing and Data Capture are discussed in Section 1.2 and Section 1.3 respectively. The problem definition is described in Section 1.4. Motivation behind carrying out the work is stated in Section 1.5. Organization of the thesis is outlined in Section 1.6.

## 1.2 Document

The objective of document image processing is to recognize text and graphics components in images of documents, and to extract the intended information as a human would. Two categories of document image processing can be defined:[3]

• **Textual Processing** deals with the text components of a document image. Some tasks here are: determining the skew (any tilt at which the document may have been scanned into the computer), finding columns, paragraphs, text lines, and words, and finally recognizing the text (and possibly its attributes such as size, font etc.) by optical character recognition (OCR).

• **Graphics Processing** deals with the non-textual line and symbol components that make up line diagrams, delimiting straight lines between text sections, company logos etc. Pictures are a third major component of documents, but except for recognizing their location on a page, further analysis of these is usually the task of other image processing and machine

vision techniques. After application of these text and graphics analysis techniques, the several megabytes of initial data are culled to yield a much more concise semantic description of the document.

## 1.3   Problem Definition

When we scan or photograph a document and the back page shows through. This effect is often due to partial transparency of the paper or paper of low quality (which we designate by show-through), another possible cause is bleeding of ink through the paper, a phenomenon that is more common in old documents, in which the ink has had more time to bleed.We call this phenomenon as bleed-through. Correction of both the effects remains one of the vitals part in Document Image Processing.

## 1.4   Motivation

Our work was primarily focused on enhancement of documents where bleed through or show through effects are shown.As it becomes difficult to use the document or interpret data from the document.

## 1.5   Thesis Organization

The rest of the thesis is organized as follows:

Chapter 2 Gives us a brief idea why fourier transform and short term fourier transform is not used. Chapter 3 introduces what is wavelet and why wavelet transform is used to remove bleed-through to a considerable level.

Chapter 4 outlines the structure, organization and working of the bleed-through removal algorithms described within.This work for removal which based on wavelet transform is explored. In addition, a method for efficient registration of documents through wavelet transform is studied.

Chapter 5 presents the concluding remarks, with scope for further research work.

## 1.6   Conclusion

This chapter gives us brief idea about the procedure to study a work.Firstly, it gives knowledge about the problem definition, then what are the reasons behind to go forward with this work, at last its organization.

# Chapter 2

# Fourier & Short-Fourier Transform

## 2 Fourier and Short-Fourier Transform

### 2.1 Introduction

Most of the signals in we use, are TIME-DOMAIN signals in their basic format.When we plot time-domain signals, we obtain a time-dependent variable (usually amplitude) representation of the signal.In many cases, the most distinguished information is hidden in the frequency content of the signal.The fourier transform tells us a great deal about a signal, it tells what frequency components exist in the signal.We measure frequency using *fourier transform*.

### 2.2 Fourier Transform

#### 2.2.1 Definition

The Fourier transform is a mathematical operation that decomposes a signal into its constituent frequencies.The original signal depends on time, and therefore is called the time domain representation of the signal, whereas the Fourier transform depends on frequency and is called the frequency domain representation of the signal. The term Fourier transform refers both to the frequency domain representation of the signal and the process that transforms the signal to its frequency domain representation. Let us see an example how fourier transform works:

$$x(t) = cos(2*pi*10*t) + cos(2*pi*25*t) + cos(2*pi*50*t) \tag{2.1}$$
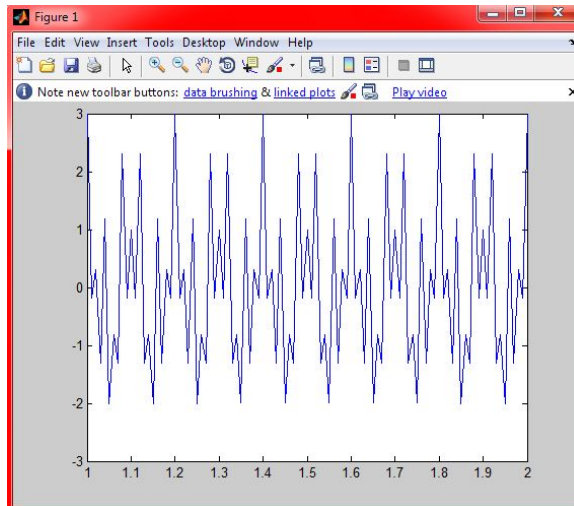


Figure 1: A stationary signal having 10Hz,25Hz,50Hz at any given instant.

By looking at the figure above, we don't get any idea about the frequency components.Now,lets find out the fourier transform of the signal.

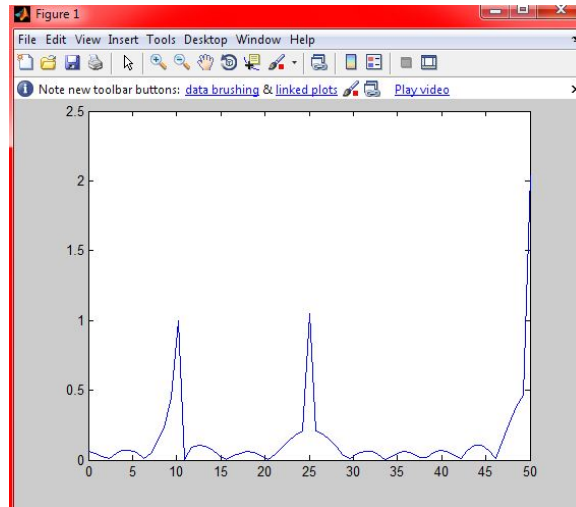The peaks in the figure above shows the frequency component i.e. 10Hz,25Hz,50Hz.



Figure 2: Fourier transform of the signal

### 2.2.2 Working Mechanism

FT decomposes a signal to complex exponential functions of different frequencies. The way it does this, is defined by the following two equations:

$$X(f) = \int_{-\infty}^{\infty} x(t) \bullet e^{-2j\pi ft} \mathrm{d}t \tag{2.2}$$

$$x(f) = \int_{-\infty}^{\infty} X(f) \bullet e^{2j\pi ft} \mathrm{d}f \tag{2.3}$$

In the above equations, t stands for time, f stands for frequency, and x denotes the signal at hand. Note that x denotes the signal in time domain and the X denotes the signal in frequency domain. This convention is used to distinguish the two representations of the signal. Equation (1) is called the Fourier transform of x(t), and equation (2) is called the inverse Fourier transform of X(f), which is x(t) [4].

### 2.2.3 Why fourier transform is not suitable?

We see that when a fourier transform of a stationary signal [1] is done we get the frequency components and we know they are present at any instant of time, but for non-stationary signal [2] we get the freqeuncy components but it fails to give time where it occured.We

---

[1]Stationary Signal: Signals whose frequency content do not change in time are called stationary signals.

[2]Non-Stationary Signal: Signals whose frequency content vary in time are called as non-stationary signals.

nextproceed to *short term fourier* transform.

## 2.3   Short Term Fourier Transfrom

### 2.3.1   Definition

Short Term Fourier Transform(STFT) is a revised version of the Fourier transform. In this the signal is divided into small enough segments, where these segments (portions) of the signal can be assumed to be stationary. For this purpose, a window function "w" is chosen. The width of this window must be equal to the segment of the signal where its stationarity is valid.

### 2.3.2   Mathematical Approach

STFT can be best described as:

$$STFTx(t) \equiv X(\tau, \omega) = \int_{-\infty}^{\infty} x(t)w(t-\tau)e^{-j\omega t}\mathrm{d}t \tag{2.4}$$

Where,x(t) is the signal to be transformed.

w(t) is the window function.

X($\tau$,$\omega$) is essentially the Fourier Transform of x(t)w(t-$\tau$).[4]

### 2.3.3   Resolution Issues

One of the downfalls of the STFT is that it has a fixed resolution. The width of the windowing function relates to how the signal is representedit determines whether there is good frequency resolution (frequency components close together can be separated) or good time resolution (the time at which frequencies change). A wide window gives better frequency resolution but poor time resolution. A narrower window gives good time resolution but poor frequency resolution.
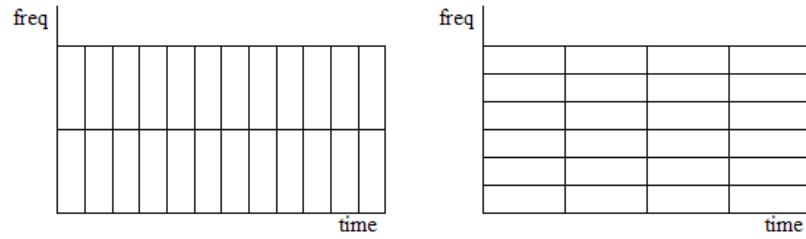
Figure 3: STFT resolution. Better time resolution in left, and Better frequency resolution in right.

## 2.4   Conclusion

We have seen in this chapter that *Fourier Transform* and *Short Term Fourier Transform* are not suitable to find good frequency resolution with good time resolution.For that we move on to wavelet transform, here we get good time and poor frequency resolution at high frequencies, and good frequency and poor time resolution at low frequencies.

# Chapter 3

## Wavelets

# 3 Wavelet

## 3.1 Introduction

Although Fourier transform has been there from 1950s for transform-based image processing,but a recent transformation *wavelet transform* is making it even easier to compress,transmit and analyze images. Wavelet transforms are based on small waves called *wavelets* of varying frequency and limited duration.It aslo helps in MRA(multi resolution analysis) i.e. it analyzes the signal at different frequencies with different resolutions.
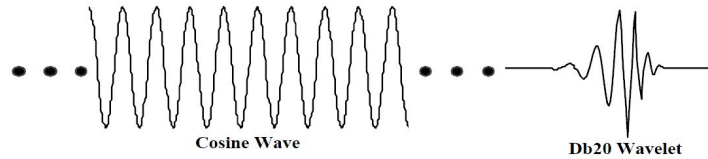


Figure 4: A portion of an infinitely long sinusoid (a cosine wave is shown here) and a finite length wavelet. Notice the sinusoid has an easily discernible frequency while the wavelet has a pseudo frequency in that the frequency varies slightly over the length of the wavelet[5]

.

## 3.2 Why Wavelet?

In our work, we focus on removing bleed-through and show-through.The mixtures images used in this work are nonlinear and noisy, and the letter and partiture mixtures are spatially variant. From looking at the pictures, we can say that high-frequency components of images are sparse (and that high-frequency wavelet coefficients are also sparse), and we used this fact to competition based on the observation that source image strongly represents itself in one of the mixture components than in the other one.

WT gives us a good time, and poor frequency resolution at high frequencies, and good frequency, and poor time resolution at low frequencies.

The figure below shows that lower scales [3] (higher frequencies) have better scale resolution that correspond to poorer frequency resolution . Similarly, higher scales have scale frequency resolution that correspond to better frequency resolution of lower frequencies.

---
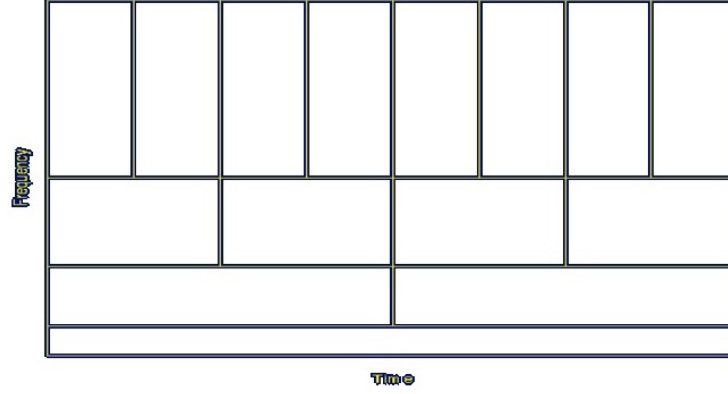
[3]scale is reciprocal of frequency

Figure 5: Wavelet resolution.

## 3.3   Continuous Wavelet Transform

The continuous wavelet transform was developed in a view to overcome the resolution problem which couldnot be addressed by short time Fourier transform. The wavelet analysis and STFT both work in a similar way when applied on signals, in the sense that we multiply the signal with a function, *the wavelet*, similar to the window function in the STFT, and the computation work is separated for different segments of the time-domain signal.

The continuous wavelet transform is defined as follows:

$$CWT_x^{\psi}(\tau, s) = \Psi_x^{\psi}(\tau, s) = \frac{1}{\sqrt{|s|}} \int x(t) \psi * (\frac{t - \tau}{s}) \qquad (3.1)$$

As seen in the above equation , tau and s are the two variables in the function to calculate transformed signal , the translation and scale parameters, respectively. The transforming function is psi(t), and we call it the mother wavelet .

## 3.4   Stationary Wavelet Transform

The Stationary wavelet transform (SWT) is a special kind wavelet transform algorithm designed to overcome the lack of translation-invariance of the discrete wavelet transform (DWT).This is used inherently as a redundant scheme as the output of each level of SWT contains the same number of samples as the input  so for a decomposition of N levels there is a redundancy of N in the wavelet coefficients.
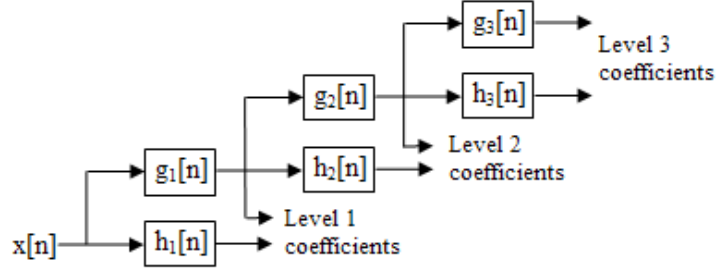
Figure 6: Wavelet decomposition.

x[n] is the original signal. $g_i[n]$ is the high pass filter. $h_i[n]$ is the low pass filter.

### 3.4.1 What swt does?

The procedure starts with passing this signal (sequence) through a half band digital low-pass filter with impulse response h[n]. Filtering a signal corresponds to the mathematical operation of convolution of the signal with the impulse response of the filter. The convolution operation in discrete time is defined as follows:

$$x[n] * h[n] = \sum_{k=-\infty}^{\infty} x[k] \bullet h[n-k] \qquad (3.2)$$

A half band low-pass filter removes all frequencies that are above half of the highest frequency in the signal. Low-pass filtering halves the resolution, but leaves the scale unchanged. The signal is then sub-sampled by 2 since half of the number of samples are redundant. This doubles the scale.

$$y[n] = \sum_{k=-\infty}^{\infty} h[k] \bullet x[2n-k] \qquad (3.3)$$

Similarly for High-pass & Low-pass filter

$$y_{high}[n] = \sum_{n} x[n] \bullet g[2n-k] \qquad (3.4)$$

$$y_{low}[n] = \sum_{n} x[n] \bullet h[2n-k] \qquad (3.5)$$

The above procedure, which is also known as the subband coding, can be repeated for further decomposition. At every level, the filtering and subsampling will result in half the number of samples (and hence half the time resolution) and half the frequency band spanned (and hence double the frequency resolution)[4] .

## 3.5    Conclusion

There are a lot of wavelet transform available and to know which will work best is difficult, so the analysis on the wavelets done earlier was necessary. Analysis showed that swt is good approach for removing nonlinear mixture images ,when compared to cwt or dwt.

# Chapter 4

# Algorithms

# 4 Algorithms

## 4.1 Overview

This is non-iterative method to eliminate the bleed-effect and show-through effect. The Stationary wavelet transform (SWT) is a wavelet transform algorithm designed to overcome the lack of translation-invariance of the discrete wavelet transform (DWT).Stationary wavelet transform(swt) is used mainly for denoising.Here we have studied methods to separate non linear image mixtures.The study involves digitizing both sides of the document and attempting to perfectly align both sides manually.

## 4.2 Initial Process

For the application of algorithm, we need the digitized image. So, first scanning of both side of the image is done, with manual alignment techniques.We need the components mixture to be precisely registered with one another to apply the separation method .At first we need to scan one side of the image then horizontally flip it.After that, an alignment procedure was applied, to correct misalignments due to the different positions of the paper during the two scanning acquisitions. In order for it to be precise, the alignment had to be performed locally. This local alignment was needed even for documents that were not wrinkled, probably due to some geometrical imperfections of the scanner[6]. In the case of the air-mail letter, which was significantly wrinkled, the local alignment was even more important.

Let us see some examples of show-through:



Figure 7: First pair of tracing paper mixtures

## 4.3  Separation Methods

Instead of assuming independence of the source images, the method that we study uses a property of common images and a property of the mixture process to perform the separation[6]. These properties are:

• High-frequency components of common images are sparse. This translates into the fact that high-frequency wavelet coefficients have sparse distributions. Consequently, the high-frequency wavelet coefficients from two different source images will seldom both have significant values in the same image location.

• In the mixture processes considered here, each source is represented more strongly in one of the mixture components (the one acquired from the side where that source is printed or drawn) than in the other component.
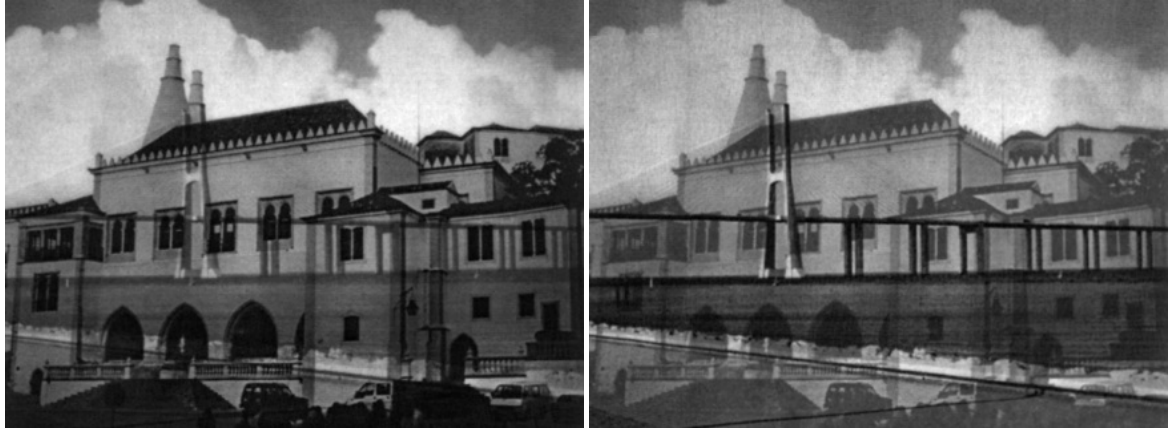
Figure 8: Second pair of tracing paper mixtures

## 4.4   Algorithm

Let us denote the recto of the document as f(x,y) and flipped side as g(x,y).Inorder to apply the technique both should have same dimension.  Now we need to apply PCA(Principal Component Analysis) to f(x,y) & g(x,y).

### 4.4.1   PCA

PCA stands for *principal component analysis*.PCA is a useful statistical technique that has found application in fields such as face recognition and image compression, and is a common technique for finding patterns in data of high dimension.The other main advantage of PCA is that once you have found these patterns in the data, and you compress the data, ie. by reducing the number of dimensions, without much loss of information.

### 4.4.2   Caluclate PCA[7]

1. Get the data set:

    Let us take data set of observations of M variables, and the reduced data for each observation can be described with only L variables, L < M. Let there are N vector set from $\mathbf{x}_1 \ldots \mathbf{x}_N$ with each $\mathbf{x}_n$ with each representing grouped observation of the M variables.

    - Write $\mathbf{x}_1 \ldots \mathbf{x}_N$ as column vectors, having M rows.
    - Construct a single matrix X of dimension MxN by putting the column vectors.

2. Calculate mean :

    - Calculate the empirical mean along each dimension m = 1, ..., M.

• Put the calculated mean values into an empirical mean vector **u** of dimensions M x 1.

$$u[m] = \frac{1}{N} \sum_{n=1}^{N} X[m,n] \qquad\qquad (4.1)$$

3. Find the deviation:

   Calculate the difference between the empirical mean vector **u** from each column of the data matrix X and store the data in the M N matrix B.

   $B = X - uh$

   where, h[n]=1   for all n=1 ... N

4. Calculate the covariance matrix:

   Find the M M empirical covariance matrix C from the outer product of matrix B with itself:

   $$C = \frac{1}{N} \sum B \bullet B^{*} \qquad\qquad (4.2)$$

5. Calculate the eigenvectors and eigenvalues of the covariance matrix:

   $V^{-1}$ CV=D                          (4.3)

   where D is the diagonal matrix of eigenvalues of C.

6. Rearrange the eigen vectors and eigen values, and Calculate the cumulative energy content for each eigenvector:

   Now the eigenvector matrix V and eigenvalue matrix D are sorted in order of decreasing eigenvalue. Make sure to maintain the correct pairings between the columns in each matrix.

7. Find out cumulative energy for each eigen vector:

   The sum of the energy content across all of the eigenvalues from 1 through m gives the cumulative energy content g for the mth eigenvector :

   $$g[m] = \sum_{q=1}^{m} D[q,q] \quad \text{for m=1 ... M} \qquad\qquad (4.4)$$

8. Determine a subset of the eigenvectors as basis vectors:

   Store the first L columns of V as matrix W of dimension MxL W[p,q]=V[p,q]    for p=1 ... M   q=1 ... L

9. Convert the source data to z-scores:

   Take the square root of each element along the main diagonal of the covariance matrix

C to find empirical standard deviation vector s

$$s=s[m]=\sqrt{C[p,q]} \quad \text{for p=q=m=1} \ldots M \tag{4.5}$$

Calculate the z-score matrix:

$$Z = \frac{B}{s \bullet h} \tag{4.6}$$

10. Project the z-scores of the data onto the new basis:

$$Y = W * \bullet Z \tag{4.7}$$

### 4.4.3 Resotarion Algorithm

1. After computing PCA, we need to calculate the value upto which level it will be decomposed using $SWT$.

2. Using Haar wavelet decompose the mixture images upto a certain level.

3. Now the corresponding high-frequency wavelet coefficients of both mixture images were subject to the following competition process:

$$m_i = \frac{1}{1+e^{\left(-a\frac{x_i^2-x_{3-i}^2}{x_i^2+x_{3-i}^2}\right)}} \tag{4.8}$$

$$y_i = x_i \, m_i \tag{4.9}$$

where i $\in$ {1,2} indexes the two sides of the paper, $x_i$ are the wavelet coefficients of a given type (for example, vertical coefficients at the first decomposition level) from the ith mixture image, $x_{3-i}$ are the corresponding coefficients from the other image of the same mixture, and $y_i$ are the coefficients that are used for synthesizing the ith separated image; a is a parameter that controls the strength of the competition[6].

4. The main motive behind the procedure is ,it preserves the coefficients that are stronger in the mixture component than in the opposite component, and weakens the coefficients that are weaker than in the opposite component.

5. The computation process described in (4.8)-(4.9) was applied to all horizontal, vertical and diagonal wavelet coefficients at all decomposition levels where $H_j$-horizontal coefficients at level j,$V_j$-vertical coefficients at leve j,$D_j$-diagonal coefficients at level j.
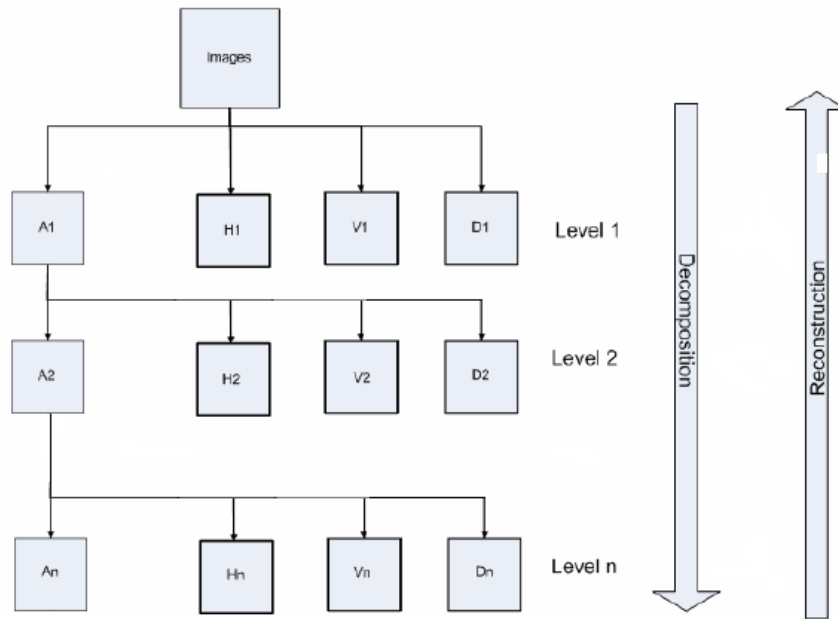
Figure 9: Schematic representation of the wavelet-based separation method

6. The reconstruction of the separated images were done by wavelet method using the high-frequency coefficients (Hj, Vj, Dj) after competition. For the low-frequency coefficients (An in Fig. 9) The coefficients which were obtained from the decomposition of the corresponding mixture image were used, with no change.

## 4.5 Result

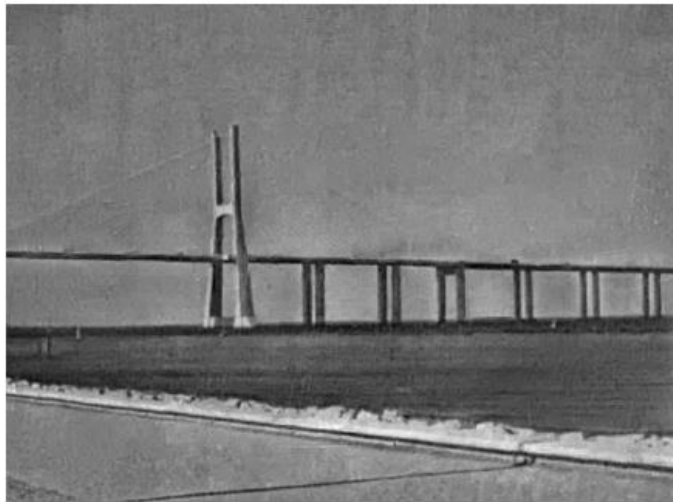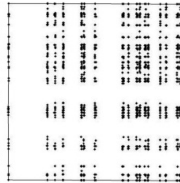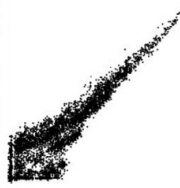Separated (page 1)



Separated (page 2)



Figure 10: Separated image

**Separation of nonlinear image mixtures**

When acquiring an image of a printed document, the image printed on the opposite page often shows through, due to partial transparency of the paper. Here we are dealing with quite a strong case of that effect, because we're using onion skin paper which is quite transparent.

The mixture that is obtained is rather nonlinear, as can be observed from the top figure on the right, which shows a scatter plot of the intensities of corresponding pairs of points from the two pages of a printed document. The scatter plot of the original images, shown in the bottom figure, filled a square, and had only a relatively small number of discrete intensity levels for each image. The fact that the shape of the scatter plot of Fig. 1 is very different from a parallelogram shows that the mixture was strongly nonlinear. The fact that this scatter plot becomes quite narrow in the upper-right corner (which corresponds to the lighter intensities in both images) indicates that, for those intensities, the mixture is close to singular. Finally, the fact that the discrete levels of Fig. 2 became largely blurred in Fig. 1 is due to noise in the process. The process leading from the sources to the observations involved printing the images, on both sides of a sheet of onion skin paper, at 1200 dpi, with a black and white laser printer (with the inherent halftoning of gray levels), and then scanning both sides of the printed sheet at 100 dpi. The noise is due, at least, to the printing process (including the halftoning), to the scanning process and to the non-uniformity in the onion skin paper, especially in its transparency.

The purpose of separation is to recover, from the mixed images that are obtained by scanning both faces of the printed document, the images that had been printed in each of its faces, with as little interference from the other image as possible.

In this example we are creating mixtures that involve natural images, printed text and graphs. The special characteristic of printed text and graphs is that they normally involve just two intensity levels (black and white) although, due to the above mentioned noise, these will appear, in the scanned images, as two clusters of intensity levels.

The separation of mixtures of two-level images, such as printed text, may be much easier than the separation of grayscale images. In fact, at least in the case of mixtures that are not too strong, a simple thresholding procedure may yield the desired results. Such a procedure can be easily performed by hand with most image processing programs, and should not be hard to automate. In such a case the use of more general blind source separation methods might be an overkill, both because it would involve a much larger amount of processing and because it might actually yield worse results. This is an extreme case in which prior knowledge about the sources can strongly simplify the separation process.

In the case of grayscale mixtures, the use of a separation method based on a good model of the physical mixing process should yield much better results that the use of a generic nonlinear separation method. A physical model could have a small number of parameters to be estimated, and would thus allow a much more precise estimation. Furthermore, it might avoid the inherent ill-posedness of nonlinear blind separation, which is currently addressed through regularization. The parameters of such a model could be estimated by an independent component analysis criterion.

Another issue of interest is the definition of separation criteria that are more suited for images or for printed documents than statistical independence. In fact, images and/or text from the opposite pages of a printed document can easily happen not to be independent from one other. For examples, images of landscapes tend to be lighter on the top than on the bottom, inducing a correlation between intensities of both. Also, in printed text with regularly spaced lines, the lines from both sides of the paper may happen to fall on top of each other, or the lines from one side may fall on the intervals of the lines from the other side, also inducing a significant correlation between intensities from both sides of the document. It would be interesting to use criteria based on a notion of image complexity, but these may not be easy to define, and may be even harder to use as criteria for optimizing a source separation system.



Figure 11: Separated image

# Chapter 5

# Conclusions

# 5 Preview

The work in this thesis, primarily focuses on bleed-through and show-through removal in scanned images. The work reported in this thesis is summarized in this chapter. Section 5.1 lists the achievements of the work. 5.2 lists the limations and 5.3 provides some scope for further development.

## 5.1 Acheivement

Using SWT(Stationary Wavelet Transform) and Haar wavelet, we removed bleed-through and show-through from scanned images.This is a faster method with no iteration required for separating real life nonlinear mixture of images.Here we use wavelet decompostion method to a deep level, the results are strong non-point wise separation. In contrast with previous solutions, this method does not assume the mixture to be invariant, and is therefore suitable for mixtures with varying local characteristics, such as those that result from bleed-through or from wrinkled documents.

## 5.2 Limitation of Work

The following limitations were encountered during the course of this project: The output obtained this method is applied to colour images are not good, better results could have been obtained by separately processing the three colour channels of the original images.

## 5.3 Further Development

Contrast compensation calculation becomes complex, it can be evaluated in future work.Color images can also be taken by separately processing the three color channels of the original images.

# References

[1] Rafael C. Gonzalez and Richard E.Woods. *Digital Image Processing*. Prentice Hall of India.

[2] Anil K. Jain. *Fundamentals of Digital Image Processing*. Prentice Hall of India.

[3] Lawrence OGorman Rangachar Kasturi and Venu Govindraju. Document Image Analysis: A primer. *Sadhana*, 27(1):3 – 22, February 2002.

[4] Robi Polikar. The Wavelet Tutorial. `http://users.rowan.edu/~polikar/WAVELETS/ WTtutorial.html`, january 2001.

[5] D.Lee Fugal. *Conceptual Wavelets In Digital Signal Processing*. Space and Signal Technical Publishing.

[6] Mariana S. C. Almeida and Luis B. Almeida. Wavelet-based Separation of Nonlinear Show-through and Bleed-through Image Mixtures. *Neurocomputing*, 72(1):57 – 70, October 2008.

[7] Lindsay I Smith. A Tutorial on Principal Component Analysis. February 2002.