# A STUDY ON RESTORATION TECHNIQUES FOR OLD IMAGE DOCUMENTS

Roshan D Suvaris[1], Dr. S Sathyanarayana[2], Dhanush Shetty K[3] & Sandeep Raj[4]

**Abstract: Old documents are prone to being attacked by pests and insects. For this, most documents that are very old enough to just suddenly crumble at slight touch are sealed on tight containers to prevent any insect from being able to get to it at any probable way. Connecting past and present is essential in order for one to find the right path towards future.**
**Old document image restoration is the process of improving the appearance of the digital image of an old document. The aim of this paper is to introduce various techniques used in old document image restoration to the reader, who are just beginners in this field. There are various types of digital image restoration but this paper only discusses about wavelet based image restoration, image restoration using filter, using genetic algorithm and image restoration using texture inpainting and segmentation.**
**Keywords: Image processing, Inpainting, Wavelet, genetic algorithm, old document.**

## 1. INTRODUCTION

Historical documents are original documents that contain important historical information about a person, place, or event and can thus serve as primary sources as important ingredients of the historical methodology. Significant historical documents can be deeds, laws, accounts of battles (often given by the victors or persons sharing their viewpoint), or the exploits of the powerful.

The digital image database of historical documents is growing in the field of heritage studies. The work requires those images which are restored, enhanced and stored in a reasonable manner in order to simplify access and disseminate [5, 6]. In fact, the restoration and enhancement of degraded historical document images are considered a transformation process which concentrated to restore its original representation [7]. In addition, restoration and enhancement are desired to improve the results of subsequent segmentation and recognition [8].

The growing importance of digital image processing stems from two principal application areas and they are (i) Improvement of pictorial information for human interpretation and (ii) Processing of scene data for autonomous machine perception. Digital image processing techniques are now used to solve a variety ofproblems.

One such important problem in image processing is restoration. The goal of the restoration approach is to improve the given image, so that it is suitable for further processing. Restoration is a technique used to reconstruct or recover an image that has been degraded by using a priori knowledge of the degradationphenomenon.[6]

In this paper, we conduct a survey on the existing methods for image restoration of ancient document. Documents can be a valuable source of information but often they suffer degradation problems, especially in the case of historical documents, such as strains, background of big variations and uneven illumination, ink seepage, etc.[2] The purpose of the study is to preserve the information contained in old documents into a digital form, because the process to rescue the old images using physical approach is too slow. For ancient documents, major degradation factor is noise. In this paper, we have discussed the existing methods for image degradation for the ancient documents.

## 2. OLD DOCUMENT RESTORATION TECHNIQUES

*2.1 Using Segmentation and Texture Inpainting.*
According to the proposal made by [4] the system of restoration was divided into two modules.

- Superfluous content elimination module
- Missed content Inpainting module(uses the result of first module)

*2.2 Superfluous content elimination module*
Superfluous content is considered with annotations often stated in the border of the manuscript [4]

---

[1] Asst professor, Department of Information Technology, AIMIT (St. Aloysius), Mangaluru, Karnataka, India
[2] Asst professor, First Grade Womens College, Mysore.
[3] Department of Information Technology, AIMIT (St. Aloysius), Mangaluru, Karnataka, India
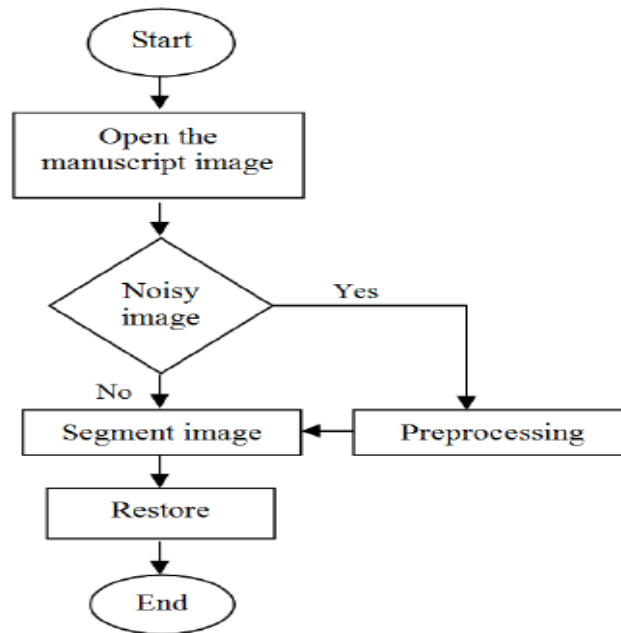[4] Department of Information Technology, AIMIT (St. Aloysius), Mangaluru, Karnataka, India

Figure 1. Flowchart of the superfluous content suppression module. [4]

*2.3 Preprocessing:*

Preprocessing is performed to enhance the image quality. Luminosity enhancement and thresholding are the two methods that comes under preprocessing.

*2.4 Luminosity enhancement:*

The enhancing process is done as follows:

Each pixel of a image is represented by three components R (red component), G (green component) and B (blue component). For each pixel of the original image, and for each channel of the pixel, the pixel of the enhanced image is calculated as follows:

$$C`n = Cn + (S*rate)/256 \tag{1}$$

Where: $C`n$ is $n^{th}$ component of the enhanced image and n its channels: R, G and B. Cn: is the $n^{th}$ component of the mean of the three components of the original image:

$S = 1/n \ \Sigma \ (Cn)$ and rate = (TX*256)/100 (TX is a luminosity parameter).

*2.5 Thresholding*:

Thresholding is done following the original Otsu algorithm [6]. Since this method is applied to grey level images and uses histogram of images, the image is converted to grey level and then its histogram is computed.

- Segmentation with k nearest neighbors

[4] Defines a threshold of neighboring distance between gray level pixels of the image. If the distance between two pixels is less than this threshold then both pixels belong to the same class, otherwise one of them belongs to another class. The only parameter here is the neighboring distance which is chosen manually.

- Classification of segments

In this step, a label to each segment provided by the previous step. Essential content (the original one) and spurious content (stamps, annotations and ink spots) are the two classed taken by [4].

The rational of classification is that significant content is very often larger than the non-significant one and may appear more often. On this basis, for each segment, we compute the number of pixels which compose it which we call the weight of the segment. The segment with the largest weight is considered essential content and the remaining segments are classified into spurious content. This method has the advantage of being simple to implement and having a fast execution time.

- Restoration

In the restoration step, all pixels belonging to the spurious content class are eliminated. The image is then divided into background and foreground. Foreground is the essential content provided with the previous step, background is the rest. The mean of the background is then calculated and all pixels except the essential content are replaced with that mean. In the case of a printable version, background pixels are given the white value.

*2.6.Missed Content Inpainting Module*

Content of manuscripts may be lost during many stages and degradation may be physical and logical. To retrieve the lost content due to degradations, many approaches are proposed. We have used in our work the exemplar based approach to reconstruct our missing content in the manuscript image. Fig. 2 shows the flowchart of this module.
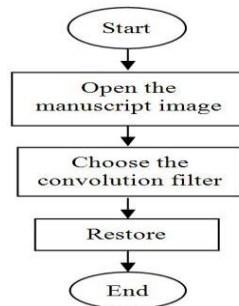


Figure 2.     Flowchart of the missed content inpainting module.

Restoration is performed in several steps. These steps are: Select the missing area, locate border pixels and inpaint by exemplar based.

- Select the missing area

Selection of the missing area can be done in two ways: automatically and manually.

In an automatic way, a zone is considered as degraded if its color is given a certain value. This value is chosen according to the type of degradation. Degradations are usually presented as spots and holes. Their color is usually very darker or lighter than the rest of the image and its amount is usually less than other colors. We have chosen in our experiments a blue color as a synthesized degradation and use Euclidian distance to decide whether the pixel is considered as a degradation or not.

- Locate border pixels

A convolution is applied on the missing area to get their border pixels. Convolution can be made using a Laplacian or a Sobel filter.

## 3. IMAGE FILTERING ALGORITHM

There are different image filtering algorithms used in image processing

In image processing, filters are mainly used to suppress either the high frequencies in the image, i.e. smoothing the image, or the low frequencies, i.e. enhancing or detecting edges in the image. Image restoration and enhancement techniques are described in both the spatial domain and frequency domain, i.e. Fourier transforms. Noise removal is easier in the spatial domain as compared to the frequency domain as the spatial domain noise removal requires very less processing time. Spatial processing is classified into point and mask processing. Point processing involves the transformation of individual pixels independently of other pixels in the image. These simple operations are typically used to correct for defects in image acquisition hardware, for example to compensate for under/over exposed images. On the other hand, in mask processing, the pixel with its neighborhood of pixels in a square or circle mask are involved in generating the pixel at (x, y) coordinates in the enhanced image. It is a more costly operation than simple point processing, but more powerful. The application of a mask to an input image produces an output image of the same size as the input. One of the most important requirements of noise removal algorithms is that they should provide satisfactory amount of noise removal and also help preserve the edges. For the stated conditions to be satisfied there are two types of filters with their significant advantages and disadvantages. The two types of filters are the linear and non-linear filters. The linear filters have the advantage of faster processing but the disadvantage of not preserving edges. Conversely the non-linear filters have the advantage of preserving edges and the disadvantage of slower processing.[8]

Various filtering techniques:

The purpose of smoothing is to reduce noise and improve the visual quality of the image. Often, smoothing is referred to as filtering. There are two types of filters that are fund useful:

*A. Spatial filter*

*B. Temporal filter*

Spatial channels are connected to both static and dynamic pictures, though transient channels are connected just to dynamic images. The easiest smoothing system is the nine-point smooth. The nine-point smooth will take a 3-x-3 square of pixels (aggregate of nine) and decide the quantity of tallies in every pixel. The tallies per pixel are then arrived at the midpoint of, and that esteem is relegated to the focal pixel (Figure 5). This same operation can be rehashed for the whole PC screen or limited to an assigned territory. Comparable operations can be performed with

| 5 | 7 | 3 |
|---|---|---|
| 4 | 2 | 5 |
| 6 | 2 | 2 |

| 5 | 7 | 3 |
|---|---|---|
| 4 | 4 | 5 |
| 6 | 2 | 2 |

5-x-5 or 7-x-7 squares.3-x-3 square values after smoothing
3-x-3 square values before smoothingSimple Nine-Point Smooth Schematic
5+7+3+4+2+5+6+2+2=4(average value)


*3.1 Spatial channels:-*
A wide cluster of techniques, and a few devoted „spatial" econometric methods for the factual investigation of geo referenced information is accessible in the writing. These strategies are valuable while breaking down territorial joblessness information, as for our situation ponders and especially when the last point is to create estimating models for some local scale. 10 Among customary spatial econometric techniques, spatial auto relapse is a capable strategy normally utilized. Spatial autoregressive systems consider spatial impacts by methods for geographic weights lattices that give measures of the spatial linkages (reliance) between estimations of geo referenced factors.

*3.2 Transient Filtering:-*

Transient sifting permits lessening signals that are not associated from casing to outline. It can viably diminish commotion when joined with movement pay, as movement remuneration corresponds the picture content from casing to outline. This makes this preparing reasonable to enhance the productivity of consequent encoders [7]. It is executed utilizing a recursive channel since it gives a superior selectivity at bring down expenses. The general objective of transient separating is to build the flag to-clamor proportion. Because of the generally poor worldly determination off MRI (Functional attractive reverberation imaging), time arrangement information contain minimal high-recurrence commotion. They do, in any case, frequently contain moderate recurrence variances that might be random to the flag of intrigue. Moderate changes in attractive field quality might be in charge of part of the low-recurrence flag saw in MRI time arrangement.

## 4. GENETIC ALGORITHM.
Hereditary Algorithms (GAs) are versatile heuristic inquiry calculation in view of the developmental thoughts of regular choice and hereditary qualities. All things considered they speak to a clever misuse of an irregular inquiry used to tackle enhancement issues. Albeit randomized, GAs are in no way, shape or form arbitrary, rather they abuse chronicled data to coordinate the hunt into the locale of better execution inside the pursuit space. The essential methods of the GAs are intended to reproduce forms in normal frameworks fundamental for advancement, uniquely those take after the standards initially set around Charles Darwin of "survival of the fittest.". Since in nature, rivalry among people for sparse assets brings about the fittest people ruling over the weaker ones.

*4.1 Chromosome Structure*
This usage of hereditary technique utilizes chromosome structure proposed in [2]. The structure utilizes arranged cluster of irregular whole number numbers from 0 to 255 that speaks to dark level range in a 8-bit grayscale picture. Size of every chromosome is equivalent to n, which speaks to the quantity of dark levels in input picture. This implies each picture will have distinctive length of chromosome portrayal as indicated by its current dim levels
$T(G(k))=Ci(k)$                                                                                         (1)
*Where k=1,2,....n*
The remapping of chromosome structure into digital image form follows the transformation in (1) where *T* is thefunction that will change the original image gray levels, *G* is the array of input gray levels in ascending order, *k* is number of gray levels in the input image and *G(k)* is *kth* gray level of them, *Ci* is the *ith* chromosome in the population, and *Ci(k)* is the value of *kth* cell. This transformation simply means that the first gray level in original image is replaced by the value of the first cell of the enhanced chromosome and so on. Before creating initial population, firstly, the original image chromosome representation's length should be obtained, that is the number of input image gray levels (*n*).
After that, each chromosome is generated using the following steps:
1) Create *m*, where *m* represents population count, empty arrays of length *n*.
2) Set each element of each empty array with random integers ranging from 0 to 255. Although the length of chromosomes may vary according to number of existing gray values in each image, values to be assigned to new chromosomes don't depend on gray level values on input image. It can be any value, from 0 to 255.
3) Sort each array or chromosome in ascending order.

4) Set the first element of each chromosome to 0, and the last element of each chromosome to 255. The number of individuals in the population is constant. This implementation use *K*-elitist scheme with *K* = 10 which means 10 individuals are considered having better fitness values determined by the selection algorithm, and are forwarded to the next generation. These processes are performed until it meets the terminating condition. The terminating criteria used here is a determined number of generations.

*4.2 Fitness Function*

Reference [2] utilizes number of edges and their general power as a wellness work, a dark picture with great visual differentiation incorporates numerous serious edges.[1] have connected the wellness work from [2] to this issue, yet it didn't portray the normal for the needed arrangement. The space issue of this examination is old archive data should have been protected is spoken to by long edges. Yet, it doesn't really say that every single short edge are the commotions. Along these lines, [1] needed to evacuate those which are considered as clamors without losing the data. In light of past contemplations, number of edges identified by Sobel administrator that contain just 1 pixel is utilized as a wellness work. The best wellness esteem is the one with the most similitude with the first picture chromosome portrayal's wellness esteem, since we would prefer not to lose excessively data. For instance, in the Figure below we have a double picture with 2 edges. The principal edge is of length 16, while the other is of length 1..
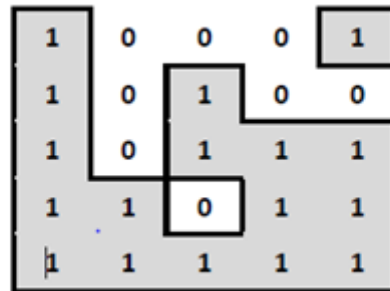
Figure: A binary image with fitness value 1

*4.3. Selection*

Roulette wheel rule is used as the selection operator. As mentioned in the previous sub-section, the selection operator will return 10 best individuals based on the fitness criteria. To perform selection operation, the fitness value of all individuals in the population should be calculated first. Furthermore, the selection mechanism is as follows:

1) Calculate standard deviation ($\sigma$) of population's fitness value, including fitness value of the original imagechromosome representation (*b*).

2) Remove all individuals with fitness value *a*, where *a<(b-$\sigma$)*or *a>(b+$\sigma$)*. This operation removes outliers,which are images with too high or too low fitness value, because based on the experimental observation, outliers always give bad results.

3) Give score for each remaining individuals. Individuals with more similar fitness value to fitness value of inputimage will have higher score. This score represents the probability a chromosome will be selected and survive to thenext generation.

4) Generate random number from 0 to the sum of given score for all individuals in step 3. This random number willdetermine which individual is selected. This is where roulette wheel selection is implemented.

5) Repeat step 4 if the number of selected individuals is less than 10. For example, suppose we have 6 individuals survive from step 2, and each will be given score in step 3. Let's say the score of the 1st to the 6th individual are 0.5, 0.7, 1.1, 1.2, 0.8, and 0.7, then the sum of given score is 5. From step 4, we will have a random generated number with value ranging from 0 to 5. If the random number value falls between:

- [0, 0.5], the 1st individual is selected.
- (0.5, 1.2], the 2nd individual is selected.
- (1.2, 2.3], the 3rd individual is selected.
- (2.3, 3.5], the 4th individual is selected.
- (3.5, 4.3], the 5th individual is selected.
- (4.3, 5], the 6th individual is selected

This selection process returns 10 best individuals thatwill be used as initial population for the next generation.

*4.4. Crossover and Mutation Operator*

After comparing the 2 crossover operations presented in [3], Uniform R/C Crossover is selected over Random R/CCrossover, as the crossover operator, with the crossover rate*Pc*= 0.8. In this crossover operation, parent chromosomes will not be guaranteed as the best chromosomes because it is selected at random. After getting2 parent chromosomes, generate 2 random numbers ranging from 0 to chromosome's length. These numbers willdetermine the position or index in which elements of both chromosomes are to be substituted. 2 new offspring will be produced after substitution, and each individual will be sorted
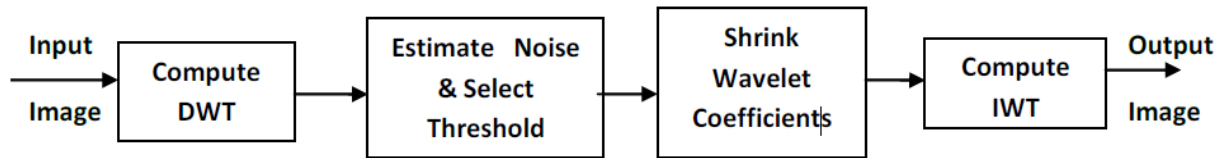
in ascending order to preserve structure [2]. The mutation operator that is selected is the one presented in [2], with a mutation rate $Pm = 0.1$. Individualthat will be mutated is selected at random. 5% of the individual chromosome elements are selected randomly formutation. Those elements will be replaced by randomly generated integer that should be less than or equal to the next element value and more than or equal to the previous element.

## 5. WAVELET BASED IMAGE RESTORATION

The principle of wavelet change is part up the flag into a bundle of signs, speaking to a similar flag, yet all relates to recurrence groups. The possibility of wavelet de-noising in view of the presumption that the abundance, as opposed to the area of the spectra of the flag to be as various as workable for that of clamor. This permits cutting, thresholding and contracting of the abundance of the co-proficient to isolate flags or evacuate noise.[4]

Steps in Wavelet De-noising [7]
1. DWT of the image is calculated
2. Resultant co efficient are passed through threshold testing
3. The coefficients less than threshold are removed, others shrinked
4. Resultant coefficients are used for image reconstruction with IWT

Input Image → Compute DWT → Estimate Noise & Select Threshold → Shrink Wavelet Coefficients → Compute IWT → Output Image

## 6. CONCLUSION

The main goal of this survey is to understand the usage of textureinpainting, genetic algorithm, wavelet and filtering in image restoration. These studies also explain about some MFs (Membership Functions) and also list out various inpainting techniques used for image Restoration. In this paper, we have presented the existing methods for digital image enhancement. There is the need of some more efficientapproach to enhance the document quality.

## 7. REFERENCES

[1]   Hilda Deborah and AniatiMurniArymurthy, "Image Enhancement and Image Restoration for Old Document Image using Genetic Algorithm", in Second International Conference on Advances in Computing, Control, and Telecommunication Technologies,2-3,2010

[2]   S. Hashemi, S. Kiani, N. Noorozi, M. E. Moghaddam, "An image enhancement method based on genetic algorithm," in Proc. of 2009 IEEE International Conference on Digital Image Processing, March 2009, pp. 167–171.

[3]   Y. W. Chen, Z. Nakao, K. Arakaki, "Genetic algorithms applied to neutron penumbral imaging," in Optical Review Journal, vol. 4, no. 1B, 1997, pp. 209–215.

[4]   InsafSetitra and AbdelkrimMeziane, "Old Manuscripts Restoration Using Segmentation and Texture Inpainting", in *InternationalJournal of Electrical Energy, 2-4, June 2014*

[5]   R. Hedjam and M. Cheriet, "Historical document image restoration using multispectral imaging system," *Pattern Recognition*, vol. 46, no. 8, pp. 2297-2312, 2013.

[6]   R. Hedjam, R. Moghaddam, and M. Cheriet, "A spatially adaptive statistical method for the binarization of historical manuscripts and degraded document images," *Pattern Recognition,* vol. 44, no. 9, pp. 2184-2196, 2011.

[7]   E. Dubois and S. Sabri, "*Noise reduction in images sequences using motion-compensated temporal filtering*," IEEE Trans. on Communications, vol. COM-32, no. 7, pp. 826–831, July 1984

[8]   RuchikaChandel and Gaurav Gupta, " Image Filtering Algorithms and Techniques: A Review", in International Journal of Advanced Research in Computer Science and Software Engineering 3(10), October - 2013, pp. 198-202

[9]   International Journal of Computer Applications (0975 – 888) Volume 48– No.12, June 2012

[10]  Prof. D.N. Satange, Ms. Swati S. Bobde and Ms. Snehal D. Chikate, "Historical Document Preservation using Image Processing Technique"