LAMP-TR-065 CAR-TR-962 CS-TR-4218 4400019848 N660010028910/IIS9987944 February 2001

# Morphological Degradation Models and their Use in Document Image Restoration

Qigong Zheng and Tapas Kanungo

Language and Media Processing Laboratory Center for Automation Research University of Maryland College Park, MD 20742-3275 {qzheng,kanungo}@cfar.umd.edu

## Abstract

Document images undergo various degradation processes. Numerous models of these degradation processes have been proposed in the literature. In this paper we propose a model-based restoration algorithm. The restoration algorithm first estimates the parameters of a degradation model and then uses the estimated parameters to construct a lookup table for restoring the degraded image. The estimated degradation model is used to estimate the probability of an ideal binary pattern, given the noisy observed pattern. This probability is estimated by degrading noise-free document images and then computing the frequency of corresponding noise-free and noisy pattern pairs. This conditional probability is then used to construct a lookup table to restore the noisy images. The impact of the restoration process is then quantified by computing the decrease in OCR word and character error rate.

We find that given the estimated degradation model parameter values, the restoration algorithm decreases the character error rate by 16.1% and the word error rate by 7.35%. In some categories of degradation (e.g. model parameters that give rise to broken characters) there is a 41.5% reduction in character error rate and a 20.4% reduction in word error rate.

This research was funded in part by the Science Applications International Corporation under Contract 4400019848, the Defense Advanced Research Projects Agency under Contract N660010028910, and the National Science Foundation under Grant IIS9987944.

Report Documentation Page			Form Approved OMB No. 0704-0188		
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE FEB 2001	DATE 2. REPORT TYPE		3. DATES COVERED 00-02-2001 to 00-02-2001		
4. TITLE AND SUBTITLE				5a. CONTRACT	NUMBER
Morphological Degradation Models and their Use in Document Image			ment Image	5b. GRANT NUMBER	
Restoration				5c. PROGRAM E	LEMENT NUMBER
6. AUTHOR(S)				5d. PROJECT NU	JMBER
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Language and Media Processing Laboratory,Institute for Advanced Computer Studies,University of Maryland,College Park,MD,20742-3275					GORGANIZATION ER
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/M	ONITOR'S ACRONYM(S)
			11. SPONSOR/MONITOR'S REPORT NUMBER(S)		
12. DISTRIBUTION/AVAII Approved for publ	LABILITY STATEMENT ic release; distribut	ion unlimited			
13. SUPPLEMENTARY NO	DTES				
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF: 17. LIMIT			17. LIMITATION OF	18. NUMBER	19a. NAME OF
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified	- ABSTRACT OF PAGES 20		KESPONSIBLE PERSON

Standard Form 298 (Rev. 8-98) Prescribed by ANSI Std Z39-18 LAMP-TR-065 CAR-TR-962 CS-TR-4218

4400019848 N660010028910/IIS9987944 February 2001

# Morphological Degradation Models and their Use in Document Image Restoration

Qigong Zheng and Tapas Kanungo

# Morphological Degradation Models and their Use in Document Image Restoration

Qigong Zheng and Tapas Kanungo

Language and Media Processing Laboratory Center for Automation Research University of Maryland College Park, MD 20742-3275 {qzheng,kanungo}@cfar.umd.edu

#### Abstract

Document images undergo various degradation processes. Numerous models of these degradation processes have been proposed in the literature. In this paper we propose a model-based restoration algorithm. The restoration algorithm first estimates the parameters of a degradation model and then uses the estimated parameters to construct a lookup table for restoring the degraded image. The estimated degradation model is used to estimate the probability of an ideal binary pattern, given the noisy observed pattern. This probability is estimated by degrading noise-free document images and then computing the frequency of corresponding noise-free and noisy pattern pairs. This conditional probability is then used to construct a lookup table to restore the noisy images. The impact of the restoration process is then quantified by computing the decrease in OCR word and character error rate.

We find that given the estimated degradation model parameter values, the restoration algorithm decreases the character error rate by 16.1% and the word error rate by 7.35%. In some categories of degradation (e.g. model parameters that give rise to broken characters) there is a 41.5% reduction in character error rate and a 20.4% reduction in word error rate.

This research was funded in part by the Science Applications International Corporation under Contract 4400019848, the Defense Advanced Research Projects Agency under Contract N660010028910, and the National Science Foundation under Grant IIS9987944.

#### 1 Introduction

Document images are usually corrupted by various types of noise during document generation and copying processes. We wish to design a filter to restore a class of document images that have similar structural features and degradation conditions. A traditional approach to image restoration is to use linear filters [Jai89]. Although linear filters are mathematically simple, their use usually results in distortion of many important image characteristics. In this paper we propose an algorithm to create a look-up-table that can be used for restoring degraded document images.

The issue of morphological filter design has been studied by many researchers. Dougherty [Dou92] proposed a method of characterizing the optimal binary morphological filter in terms of the Matheron representation. Using the Matheron representation, any binary morphological filter can be expressed as a union of binary erosions. The filter design procedure is thus essentially the problem of finding structuring elements that yield statistically optimal representations. To mitigate the computational burden of filter design, Loce [LD92] adds some constraints like the number of erosions, window size, and structuring element libraries to minimize search. As a result, his filter design is suboptimal. Schofeld and Goutsias [SG91] consider the set-difference distance as a measure of comparison between images, and by using this function, they prove that the class of alternating sequential filters is a set of parametric, smoothing morphological filters that best preserves the crucial structure of input images in the least mean difference sense. Liang and Haralick [LH96] present a method of restoring document images degraded by subtractive or additive noise, given a constraint on the size of the filters. The improvement of their algorithm is shown by the increased accuracy of an OCR system.

One of the common limitations of the above-mentioned algorithms lies in the lack of prior statistical information or an adequate image noise model, which makes them computationally complex. This suggests that greater improvement of restoration algorithms may be achievable by using an image noise model.

A survey of document image degradation models proposed in the literature can be found in [Bai99]. We use the model proposed by Kanungo et al. [KHP94, KHB<sup>+</sup>00] for our restoration algorithm.

## 2 Document Degradation Model

Our degradation model [KHP94] has six parameters:  $\Theta = (\eta, \alpha_0, \alpha, \beta_0, \beta, k)$ . We model the probability of a pixel flipping from foreground to background or vice versa as an exponential function of its distance from the nearest boundary pixel. The foreground and background 4-neighbor distances are computed using a standard distance transform algorithm. The flipping probabilities of foreground and background pixels are controlled by  $\alpha_0, \alpha$  and  $\beta_0, \beta$  respectively. The parameters  $\alpha_0, \beta_0$  are the initial values for the exponentials, and the decay speeds of the exponentials are controlled by the parameters  $\alpha, \beta$ . Parameter  $\eta$  is the constant probability of flipping for all pixels. Parameter k is the size of the disk used in the morphological closing operation. This operation normally simulates the correlation introduced by the point-spread function of the optical system. The procedure for degrading an ideal binary image is as follows:

I DEHEVIOL OF FILESE SYS	A PREVENT OF PHOSE BAD	I DETERATOR OF STEEDE BAD
ed by sets of coupled	ed by sets of coupled	ed by sets of coupled
is for formal neurons	is for formal neurons	is for formal neurons
;ation of essential featu	sation of essential featu	gation of essential featu
y and adaptation dep	y and adaptation dep	y and adaptation dep
mathematical theory	mathematical theory	mathematical theory
nd efficient analysis of	nd efficient analysis of	nd efficient analysis of
t theoretical research	t theoretical research	t theoretical research
adopted are frequent	adouted are framment	adapted are frequent
(a)	(b)	(c)

Figure 1: (a) A typical ideal image; (b) Degraded version of (a) with parameters (1.0, 0.7, 1.0, 3.0); (c) Degraded version of (a) with parameters (1.0, 3.0, 1.0, 0.7).

- 1. Compute the distance d of each pixel from the nearest character boundary.
- 2. Flip each foreground pixel with probability

$$p(0|1, d, \alpha_0, \alpha) = \alpha_0 e^{-\alpha d^2} + \eta.$$

3. Flip each background pixel with probability

$$p(1|0, d, \beta_0, \beta) = \beta_0 e^{-\beta d^2} + \eta.$$

4. Perform a morphological closing operation with a disk structuring element of diameter k.

Figure 1 illustrates ideal and degraded images with different model parameters. Note that the two degraded images differ in the speed of decay of the exponential functions. If  $\alpha < \beta$ , more foreground pixels change to background so the images appear to be corrupted by subtractive noise. If  $\alpha > \beta$ , more background pixels change to foreground so the images appear to have additive noise.

## 3 The Estimation Algorithm

In this section, we briefly describe a parameter estimation algorithm [KZ01] for the degradation model described in the previous section. The basic assumption of this algorithm is that two document images with similar noise should have neighborhood pattern distributions that look similar. Thus we can estimate model parameters by degrading documents with various model parameter values and choose the one that gives rise to a neighborhood pattern distribution that is very close to that of the given degraded image.

Let P be a set of neighborhood bit patterns and p be an arbitrary element in the set P. If we choose a  $3 \times 3$  neighborhood, we will have a total of 512 different patterns. Let

 $H_R$  denote the pattern distribution of a degraded image R so that  $H_R(p)$ , where  $p \in P$ , is the number of times the pattern p occurs in the binary image R. Using mathematical morphology, we can define  $H_R(p)$  more precisely:

$$H_R(p) = \#\{R \ominus p\}. \tag{1}$$

We say that two images R and S are similar if the corresponding pattern distributions  $H_R$  and  $H_{S_{\theta}}$  are similar. To test the similarity of two pattern distributions, we use the Kolmogorov-Smirnov test [Mas51] of the two pattern distributions. Let  $KS(H_R, H_{S_{\theta}})$  denote the KS test *p*-value for the null hypothesis that the two distributions are the same. We will use this *p*-value as the objective function that the estimation process tries to maximize. That is,

$$\hat{\theta} = \arg\max_{\theta} KS(H_R, H_{S_{\theta}}) \tag{2}$$

Conventional optimization algorithms typically need the functional form of the objective function. However, in our case, since  $S_{\theta}$  is computed by simulation, it is impossible to use standard derivative approaches. We thus choose the simplex optimization algorithm [NM65] to minimize KS, which needs only function values to maximize or minimize functions. To prevent the problems of local minima, we select multiple random starting locations and pick the solution corresponding to the lowest *p*-value.

#### 4 The Restoration Algorithm

In this section we demonstrate that by using our degradation model, we can design filters in a more concise and efficient way, and the corresponding restoration procedure is thus simple and easily implemented.

Compared to other morphological restoration algorithms [LD92, LH96], our method is model-based. We always assume that the degraded image can be characterized by a set of parameters that can be estimated by using the algorithm described in the previous section. Our algorithm has two stages, a training stage and a restoration stage.

Suppose we have an ideal image I and a corresponding degraded image  $S_{\hat{\theta}}$  where  $\theta$  is the estimated parameter set used to generate  $S_{\hat{\theta}}$  from I. The training stage is responsible for computing the conditional distribution between the noise pattern pairs in the image pair  $(I, S_{\hat{\theta}})$ . During the training stage, we first scan  $S_{\hat{\theta}}$ . Next we obtain its noise pattern  $P_S(x, y)$  at location (x, y). We also obtain the point pattern at location (x, y) in the ideal image  $I: P_I(x, y)$ . From the pattern pairs  $(P_I(x, y), P_S(x, y))$ , we form the pattern distribution of an ideal image I conditioned on the degraded image  $S_{\hat{\theta}}: H_{\hat{\theta}}(P_I|P_S)$ .

The restoration stage takes place after estimating the model parameters of the degraded image. Let Q represent the restored image version of  $S_{\hat{\theta}}$ . Given the pattern  $P_S(x,y)$  at location (x,y) of the degraded image  $S_{\hat{\theta}}$ , the restored pattern  $P_Q(x,y)$  in Q is computed as

$$P_Q(x,y) = \arg\max_{p \in P_I} H_{\hat{\theta}}(p|P_S(x,y))$$
(3)

Equation (4.1) is essentially the Maximum Likelihood (ML) estimate of the pattern based on the known parameter  $\theta$ . Figure 2 shows an ideal image and its degraded versions with two different parameter sets. Figure 3 shows four typical noise patterns in the degraded image in Figure 2(b) and its conditional pattern distribution based on the corresponding ideal image in Figure 2(a).

I NOTIONEDT OF PERDOC DAD	T NETTOATOT OF STEED BAD
ed by sets of coupled	ed by sets of coupled
is for formal neurons	is for formal neurons
ation of essential featu	ation of essential featu
y and adaptation dep	y and adaptation dep
mathematical theory	mathematical theory
nd efficient analysis of	nd efficient analysis of
t theoretical research	t theoretical research
adapted are frament	adapted are frequent
(b)	(c)
	ed by sets of coupled is for formal neurons gation of essential featury y and adaptation dep- mathematical theory id efficient analysis of t theoretical research adopted are frequent (b)

Figure 2: (a) A typical ideal image; (b) Degraded version of (a) with parameters (1.0, 0.7, 1.0, 3.0); (c) Degraded version of (a) with parameters (1.0, 3.0, 1.0, 0.7).



Figure 3: Four typical noise patterns are shown in the leftmost column. The pattern entries in the other columns show possible ideal patterns and the corresponding probabilities. The ideal image was degraded with parameter set (1.0, 0.7, 1.0, 3.0).

#### 5 Experimental Protocol and Results

The experiment is outlined illustrated in Figure 5. The basic idea is to compare the OCR result of the degraded image with that of the restored one. The evaluation software is provided by the University of Maryland. It compares the OCR outputs and the corresponding groundtruth information and generates statistical information such as character-level or word-level accuracy in a batch mode. We believe that the OCR accuracy rate is a good and objective indicator for showing how well our algorithm improves the overall image quality.

E NUMAVIOL OF SHORE BYD ed by sets of coupled s for formal neurons ation of essential feature y and adaptation dep mathematical theory nd efficient analysis of t theoretical research adapted are fragment (a) and of anode by a ed by sets of coupled s for formal neurons ation of essential featu y and adaptation depmathematical theory id efficient analysis of t theoretical research adapted are frequent (b)

Figure 4: (a) The restored version of the image shown in Figure 2(b). (b) The restored version of the image shown in Figure 2(c).

The test images were 100 one-column pages of English Bible that were typeset using **E**T<sub>F</sub>X. The image size is A4 with 12-point font size. One additional image was typeset to generate pattern distributions for the estimation process. While its text content was different from that of the 100 test images, its font and bigram symbol probabilities had the characteristics of the test images. The 100 test images were degraded and then categorized into ten groups with each group possessing a unique parameter set. The OCR product was FineReader4.0, manufactured by ABBYY. Tables 6–15 give the OCR accuracy before and after our restoration algorithm. Figures 6–15 show typical degraded images and restored images with the corresponding parameter sets. We also compute the image noise level (absolute mean error) for the purpose of comparison with morphological filter based algorithms. The foreground noise level (FNL) is an indicator that measures how many black (foreground) pixels in the original image change to white (background), and the background noise level (BNL) is used to detect how many white pixels change to black. They can be computed by doing logical operations between the ideal image (I)and degraded image (D) or restored image (R). The number of flipping pixels (EFP) basically summarizes both kinds of noise. Mathematically, the above three metrics can be represented in terms of set operations:



Figure 5: Illustration of the experimental setup to compare OCR accuracy on restored versus unrestored images.

$$FNL = \frac{\#\{(I \oplus D) \cap I\}}{\#\{I\}} \tag{4}$$

$$BNL = \frac{\#\{(I \oplus D) \cap I^{c}\}}{\#\{I\}}$$
(5)

$$EFP = \frac{\#\{I \oplus D\}}{\#\{I\}} \tag{6}$$

where  $\oplus$  denotes the XOR operation and # is the cardinality of the set (i.e. the number of foreground pixels in a binary image).

From the test statistics, we see that our restoration algorithm decreases both the OCR error rate and image noise level. For instance, the decreases in OCR accuracy error rate at the character and word levels range from 3.4% to 41.5% and from 1.0% to 20.4% respectively, depending on what model parameters are associated with the degraded images. In particular, we find that our algorithm performs better in restoring images suffering from broken characters (Figures 8 and 9) than those that have blurred characters (Figures 12 and 13). This gives us the impression that the OCR product seems to be more vulnerable to broken characters which have more subtractive noise. In addition to the OCR error rate, our algorithm significantly decreases the image noise level by amounts, ranging from 13.1% to 52.7%.

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	veloped for research and
nical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
ns as the operations on t	ns as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so:	To achieve our goal, so:
(a)	(b)

Figure 6: (a) A sample degraded image with parameters (0.6, 0.8, 1.0, 3.0); (b) Restored image of (a).

Table 1: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (0.6, 0.8, 1.0, 3.0).$ 

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	24660	24580	
Num. of Correct Chars	23885	23910	
Num. of Char Errors	775	670	13.5%
Num. of Words	4855	4855	
Num. of Correct Words	3762	3816	
Num. of Word Errors	1093	1039	4.9%
Foreground Noise Level	16.1%	11.8%	
Background Noise Level	0.19%	0.19%	
Num. of Error Flipping Pixels	502659	409992	18.4%

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	veloped for research and
nical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
ns as the operations on t	ns as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so:	To achieve our goal, so:
(a)	(b)

Figure 7: (a) A sample degraded image with parameters (0.8, 0.8, 1.0, 3.0); (b) Restored image of (a).

Table 2: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (0.8, 0.8, 1.0, 3.0).$ 

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	24391	24806	
Num. of Correct Chars	23935	23999	
Num. of Char Errors	996	807	19.0%
Num. of Words	4953	4953	
Num. of Correct Words	3737	3846	
Num. of Word Errors	1216	1107	9.0%
Foreground Noise Level	22.2%	14.7%	
Background Noise Level	0.18%	0.24%	
Num. of Error Flipping Pixels	625228	516481	17.4%

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	veloped for research and
nical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
ns as the operations on $t$	ns as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so:	To achieve our goal, so:
(a)	(b)

Figure 8: (a) A sample degraded image with parameters (1.0, 0.8, 1.0, 3.0); (b) Restored image of (a).

Table 3: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (1.0, 0.8, 1.0, 3.0)$ .

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	25651	25262	
Num. of Correct Chars	23973	24280	
Num. of Char Errors	1678	982	41.5%
Num. of Words	4973	4958	
Num. of Correct Words	3397	3703	
Num. of Word Errors	1576	1255	20.36%
Foreground Noise Level	28.8%	15.7%	
Background Noise Level	0.17%	0.30%	
Num. of Error Flipping Pixels	768872	584919	23.9~%

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	veloped for research and
nical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
as the operations on $t$	ns as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so:	To achieve our goal, so:
(a)	(b)

Figure 9: (a) A sample degraded image with parameters (1.0, 0.6, 1.0, 2.0); (b) Restored image of (a).

Table 4: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (1.0, 0.6, 1.0, 2.0).$ 

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	27426	26370	
Num. of Correct Chars	22584	23455	
Num. of Char Errors	4842	2915	40.0%
Num. of Words	5040	5031	
Num. of Correct Words	2637	3089	
Num. of Word Errors	2403	1942	19.2%
Foreground Noise Level	31.7%	24.5%	
Background Noise Level	0.41%	0.43%	
Num. of Error Flipping Pixels	1026668	892519	13.1%

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	eveloped for research and
inical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
ns as the operations on t	ns as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so:	To achieve our goal, so:
(a)	(b)

Figure 10: (a) A sample degraded image with parameters (1.0, 0.8, 1.0, 2.0); (b) Restored image of (a).

Table 5: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (1.0, 0.8, 1.0, 2.0).$ 

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	25918	25771	
Num. of Correct Chars	24324	24408	
Num. of Char Errors	1594	1363	14.5%
Num. of Words	5037	5038	
Num. of Correct Words	3465	3581	
Num. of Word Errors	1572	1457	7.3%
Foreground Noise Level	24.3%	21.0%	
Background Noise Level	0.42%	0.37%	
Num. of Error Flipping Pixels	843493	758692	13.1%

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	eveloped for research and
inical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
as as the operations on $t$	ns as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so	To achieve our goal, so
(a)	(b)

Figure 11: (a) A sample degraded image with parameters (1.0, 1.0, 1.0, 2.0); (b) Restored image of (a).

Table 6: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (1.0, 1.0, 1.0, 2.0)$ .

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	25001	24950	
Num. of Correct Chars	23952	24003	
Num. of Char Errors	1049	947	9.7%
Num. of Words	4887	4889	
Num. of Correct Words	3614	3682	
Num. of Word Errors	1273	1207	5.2%
Foreground Noise Level	19.1%	18.6%	
Background Noise Level	0.42%	0.29%	
Num. of Error Flipping Pixels	750851	629294	16.2%

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	veloped for research and
inical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
ns as the operations on t	as as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so	To achieve our goal, so:
(a)	(b)

Figure 12: (a) A sample degraded image with parameters (1.0, 1.5, 1.0, 0.6); (b) Restored image of (a).

Table 7: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (1.0, 1.5, 1.0, 0.6).$ 

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	23612	23709	
Num. of Correct Chars	23065	23193	
Num. of Char Errors	548	516	5.8%
Num. of Words	4582	4586	
Num. of Correct Words	3639	3659	
Num. of Word Errors	943	927	1.7%
Foreground Noise Level	2.4%	17.4%	
Background Noise Level	1.96%	0.52%	
Num. of Error Flipping Pixels	1656108	783032	52.7%

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	veloped for research and
nical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
as the operations on t	as as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so	To achieve our goal, so:
(a)	(b)

Figure 13: (a) A sample degraded image with parameters (1.0, 1.5, 1.0, 0.8); (b) Restored image of (a).

Table 8: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (1.0, 1.5, 1.0, 0.8)$ .

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	24401	24558	
Num. of Correct Chars	23827	24037	
Num. of Char Errors	574	521	9.2%
Num. of Words	4737	4742	
Num. of Correct Words	3748	3787	
Num. of Word Errors	989	955	3.4%
Foreground Noise Level	3.8%	20.1%	
Background Noise Level	1.53%	0.4%	
Num. of Error Flipping Pixels	1337753	752212	43.7%

functions, formulae and	functions, formulae and
ystems for performing sj	ystems for performing sy
veloped for research and	veloped for research and
inical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
as as the operations on t	as as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so	To achieve our goal, so:
(a)	(b)

Figure 14: (a) A sample degraded image with parameters (1.0, 1.5, 1.0, 1.0); (b) Restored image of (a).

Table 9: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (1.0, 1.5, 1.0, 1.0)$ .

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	24717	24834	
Num. of Correct Chars	24095	24233	
Num. of Char Errors	622	601	3.4%
Num. of Words	4798	4804	
Num. of Correct Words	3757	3774	
Num. of Word Errors	1041	1030	1.0%
Foreground Noise Level	5.3%	21.0%	
Background Noise Level	1.2%	0.4%	
Num. of Error Flipping Pixels	1098220	700650	36.2%

functions, formulae and	functions, formulae and
ystems for performing sy	ystems for performing sy
veloped for research and	veloped for research and
nical sciences. However,	inical sciences. However,
rectly used for the analy	rectly used for the analy
as as the operations on t	as as the operations on t
hose involving an unspec	hose involving an unspec
definite summations, hav	definite summations, hav
To achieve our goal, so	To achieve our goal, so:
(a)	(b)

Figure 15: (a) A sample degraded image with parameters (1.0, 2.0, 1.0, 1.0); (b) Restored image of (a).

Table 10: OCR error improvement with parameters  $\alpha_0, \alpha, \beta_0, \beta = (1.0, 2.0, 1.0, 1.0)$ .

OCR Result	Degraded Image	Restored Image	Improvement
Num. of Chars	23604	23663	
Num. of Correct Chars	23049	23131	
Num. of Char Errors	555	532	4.1%
Num. of Words	4569	4572	
Num. of Correct Words	3614	3636	
Num. of Word Errors	955	936	2.0%
Foreground Noise Level	3.0%	18.2%	
Background Noise Level	1.17%	0.28%	
Num. of Error Flipping Pixels	1018179	604504	40.6%

## 6 Summary

A model-based document image restoration algorithm has been proposed based on the estimated parameters of the degradation model. We first use the degradation model to estimate the probability of an ideal binary pattern, given the noisy observed pattern. This probability is estimated by degrading noise-free document images and then computing the frequency of corresponding noise-free and noisy pattern pairs. This conditional probability is then used to construct a lookup table to restore the noisy images. The impact of the restoration process is then quantified by computing the decrease in OCR word and character error rates.

# Acknowledgments

We would like to thank Azriel Rosenfeld for his comments. This research was funded in part by the Science Applications International Corporation under Contract 4400019848, the Defense Advanced Research Projects Agency under Contract N660010028910, and the National Science Foundation under Grant IIS9987944.

# References

- [Bai99] H. S. Baird. Document image quality: Making fine discriminations. In Proceedings of the International Conference on Document Analysis and Recognition, pages 459-462, 1999.
- [Dou92] E. R. Dougherty. Optimal mean-square n-observation digital morphological filters. Part I: Optimal binary filters. Computer Vision, Graphics, and Image Processing, 55:36-54, 1992.
- [Jai89] A. K. Jain. Fundamentals of Digital Image Processing. Prentice Hall, Englewood Cliffs, NJ, 1989.
- [KHB<sup>+</sup>00] T. Kanungo, R. Haralick, H. Baird, W. Stuetzle, and D. Madigan. A statistical, nonparameteric methodology for document degradation model validation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22:1209– 1223, 2000.
- [KHP94] T. Kanungo, R. M. Haralick, and I. Phillips. Nonlinear global and local document degradation models. International Journal of Imaging Systems and Technology, 5:220-230, 1994.
- [KZ01] T. Kanungo and Q. Zheng. Estimation of morphological degradation model parameters. In Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing, 2001. To appear.
- [LD92] R. P. Loce and E. R. Dougherty. Facilitation of optimal binary morphological filter design via structuring element libraries and design constraints. Optical Engineering, 31:1008-1025, 1992.
- [LH96] J. Liang and R. M. Haralick. Document image restoration using binary morphological filters. In *Proceedings of SPIE*, volume 2660, pages 274–285, 1996.
- [Mas51] F. J. Massey. The Kolmogorov-Smirnov test for goodness. J. Amer. Stat. Assoc., 46:68-78, 1951.
- [NM65] J. Nelder and R. Mead. A simplex method for function minimization. Computer Journal, 7:308-313, 1965.
- [SG91] D. Schonfeld and J. Goutsias. Optimal morphological pattern restoration from noisy binary images. IEEE Transactions on Pattern Analysis and Machine Intelligence, 13:14–29, 1991.