Gamma Enhanced Binarization - An Adaptive Nonlinear Enhancement of Degraded Word Images for Improved Recognition of Split Characters

H. R. Shiva Kumar, and A. G. Ramakrishnan

Dept. of Electrical Engineering, Indian Institute of Science, Bangalore, India shivahr@gmail.com, agr@iisc.ac.in

Abstract—Recognition performance of any OCR suffers because of the merged and split characters that occur in the scanned images of degraded printed documents. We propose an elegant method of non-linearly enhancing such degraded, gray-scale word images. This connects the broken strokes of the characters, so that binarization of the processed word images gives components with better connectivity for most characters or recognizable units. From an initial value of one, the value of gamma, the parameter determining the enhancement, is decreased in powers of 2 and the right value of gamma is chosen based on the recognition score of our character classifier. We have created a benchmark dataset of 1685 degraded word images obtained from scanned pages of several old Kannada books. The word images have been recognized before and after the proposed nonlinear enhancement. There is an absolute improvement of 14.8% in the Unicode level recognition accuracy of our SVM-based character classifier on the above dataset due to the proposed enhancement of the gray-scale word images. Even on the Google's Tesseract OCR for Kannada, our gamma enhanced binarization results in an improvement of 5.6% in the Unicode level accuracy.

Index Terms—Split characters, printed text, power law transformation, gamma enhancement, binarization, Kannada, OCR, Tesseract, word images, old books.

I. INTRODUCTION

In the analysis and classification of images, irrespective of their domain of origin (medical, document, etc.), the key issue is one of reliable segmentation of the region of interest (tumor, abscess, word, characters, etc.) [1], [2]. The major challenges in dealing with the recognition of historical or severely degraded documents are: robust segmentation of characters in the presence of merged symbols and broken characters [3] and noise of different types. Approaches exist in the literature for addressing the issue of merged characters [4], [5]. This work addresses the issue of character splits caused during binarization due to poor printing or degradation of the paper due to aging. Figure 1 illustrates this problem with example word images from several Indic scripts. In this paper, we mainly consider split characters in Kannada script and test our proposed gamma enhanced binarization algorithm on a benchmark dataset created by us, and made publicly available [6].



Fig. 1: Examples of degraded, binary, word images in various Indian languages, exhibiting splits (cuts) in characters.

A. Literature Survey

There have been different approaches in the literature to address the problem of recognition of broken characters. However, to the best of our knowledge, very few works have been reported on this problem for the Dravidian scripts, namely, Malayalam [7], Telugu [8], Kannada and Tamil. In fact, very few works have been reported in the literature on OCRs for Dravidian languages [9] such as Tamil [10]–[12] and Kannada [13], [14], [15]. Further, compound characters in Kannada have complex, twodimensional arrangement [16], [17] and hence, segmentation of characters even from a good document is not a trivial task. Sachan et al. [18] have proposed a number of heuristic approaches, such as following the end points of cut strokes in the gray image, degree of overlap between connected components and the shortest distance between pairs of connected components (CC) to treat a set of CCs as belonging to a single character. Another obvious approach is to explore different combinations of consecutive connected components, recognize and then select the set of combinations that maximizes the recognition probability of a word. Following such an approach, Sumetphong and Tangwongsan [19] [20] treated this as a set-partitioning problem and propose a partition-growing

algorithm to group the broken pieces. The partition with the best posterior probability is chosen and the recognized output is corrected using a dictionary. Another possible approach is to start from the binarized symbols and apply morphological dilation and/or closing in an attempt to merge the broken pieces. Peerawit et al. [21] employed heuristics and morphological closing to reconstruct broken Thai characters only in the middle zone of the text line, whereas the cuts in the upper and lower zones were not considered. Yu and Yan [22] used conditional dilation, morphological analysis and stroke extension to reconstruct handwritten, broken digits.

Sulem and Sigelle [23] proposed the use of two, coupled dynamic Bayesian networks to model the interactions between the two streams of image columns and rows. They demonstrate a performance improvement of 1.1%over support vector machines (SVM) in recognizing broken characters in digitized Renaissance festival books. Droettboom [24] used graph combinatorics to join broken components of Roman characters from the Statistical Accounts of Scotland and evaluated the connected subgraphs using k-nearest neighbour classifier. Drira et al. [25] proposed a partial differential equation based formalism, combining the coherence-enhancing Weickert tensor driven diffusion filter and the singularities-preserving Perona-Malik scalar diffusion filter. The technique was applied on images from Gazette de Leyde to enhance the quality of degraded documents in Roman script and thus to reduce their OCR errors.

It is not an exaggeration, if we say that each of the above papers have used a different and custom database to report their results, which are not publicly available. Thus, in the absence of a common research database, there is no way one can compare the effectiveness, pros and cons of these different methods. In this context, it makes a difference to this field, if there are standard databases made publicly available for these kinds of degraded images, annotated with the corresponding ground truth.

B. Contributions of the paper

- Our proposed algorithm enhances the histogram of the gray level word image nonlinearly, in order to obtain the best binarized image for recognition.
- Though tested here only on Kannada word images, the proposed algorithm can be expected to work on any other script, provided the classifier in the OCR used provides a reliable recognition score for each recognized symbol.
- A standard, annotated database of 1685 degraded word images of Kannada script has been created, which has been made freely available [6] for benchmarking various algorithms against one another.



Fig. 2: Example word images that illustrate that strokeconnectivity information that was available in gray level image is sometimes lost during binarization.

II. Adaptive nonlinear enhancement of the image to minimize character splits

Stroke-connectivity information that is there in the gray level document image is sometimes lost during binarization as shown by Fig. 2. Hence, we explored gray-level image enhancement techniques that can improve binarization to minimize character splits. Deepak and Ramakrishnan [26] proposed a technique to nonlinearly stretch the histogram of an image using power-law transformation (PLT) and applied it to obtain the best recognition result in the literature on the ICDAR2011 Robust Reading Competition Challenge-1: Word Recognition Task on born digital dataset. Later, they extended it to the different colour planes of camera-captured scene word images [27] and once again, obtained the best recognition results on the ICDAR2011 robust reading competition challenge 2. However, PLT was actually proposed with values of gamma greater than one, for the reverse problem of splitting character merges in born-digital word images. We propose a method for merging/handling character splits in degraded word images using PLT, where appropriate fractional values of gamma are chosen automatically by our algorithm for each word image, to obtain the best recognition result.

Figure 3 shows the flow chart of how gamma enhanced adaptive binarization is applied on each degraded image to obtain the best recognized result. The intensity of every pixel is modified as,

$$y_o[r][c] = (y_i[r][c])^{\gamma}$$
 (1)

where $y_i[r][c]$ and $y_o[r][c]$ are the gray level intensities of the pixel at r^{th} row and c^{th} column, before and after gamma enhancement, respectively. Here, the intensities are expressed as a value in the interval [0,1). The value of γ is varied in steps as decreasing integer powers of two, starting from $2^0 = 1$. The resulting gray level images are binarized using the Otsu global thresholding algorithm [28]. The recognized result with the best recognition score returned by our classifier is chosen as the final output. Figure 4 illustrates, with an example, the mechanism of how this improves the quality of a degraded Kannada word image. The gray level images obtained for different γ values, their histograms and the corresponding binarized images are all shown in Fig. 4.

Depending upon the level of degradation of the word images, different values of γ work optimally for different



Fig. 3: Block schematic of Gamma-enhanced, locally adaptive binarization of gray level word images from scanned, old Kannada books, for improved OCR performance.

images. This is illustrated in Fig. 5, which shows seven word images with different levels of degradation, each of which is recognized correctly after enhancement with distinctly different values of γ .

III. NEW KANNADA DATASET FOR BENCHMARKING

For rigorously evaluating the performance of the proposed algorithm in handling character splits, we have created a benchmarking dataset of 1685 degraded word images obtained from old Kannada books. Figure 6 illustrates some sample images from this benchmarking dataset. The ground truth text for each of the test images is provided in a separate Unicode text file, as shown in Table I.

Performance of the proposed algorithm is evaluated using the Levenshtein distance between the recognized and the ground truth texts, across the entire benchmarking dataset. Let N denote the number of Unicodes in the ground-truth, and S, I and D denote the number of substitutions, insertions and deletions with respect to the ground truth, respectively. The Unicode recognition accuracy is determined as:

$$Accuracy = (N - S - I - D)/N \tag{2}$$

We are making this dataset publicly available at [6]. In Indian language processing, standardized datasets are a rarety. Hence, researchers report results on different datasets and it becomes difficult to assess the improvement

TABLE I: Filenames of some samples from the Kannada degraded word image dataset, together with the corresponding ground truth (GT) text for each test image. The GT has been provided in an accompanying, separate Unicode text file in the format shown below.

Image Name	Unicode Text (GT)	Image Name	Unicode Text (GT)
C0022.tif	ಇತಿಹಾಸದ	C0362.tif	ಅದಿಕ್ಯದಿಂದಾಗಿ
C0044.tif	ದೇಶಗಳ	C0365.tif	ಮುಂತಾದ
C0055.tif	ಗ್ರೀಕ್	C0386.tif	ಉಪಕಾರ
C0095.tif	ನುಗಳು	C0407.tif	ಮೂಲ
C0122.tif	ರಕ್ಷಣಾ	C0421.tif	ಮಟ್ಟ
C0250.tif	ಮೂರು	C0470.tif	ಸೇರಿದಂತೆ
C0274.tif	ಬೆಂಥಾಮ್	C0508.tif	ಮತ್ತು

of the new algorithms they propose. Making our benchmarking dataset open is a step towards addressing that problem.

IV. RESULTS AND DISCUSSION

The Kannada word images in the created dataset have been recognized by a custom built SVM classifier, using wavelet and autocorrelation features, using the LIBSVM package [29]. Table II lists the recognition accuracies before and after the application of the adaptive enhancement technique on the images. The proposed technique is able to improve the performance of the classifier from 74.5% on the original images to 89.3% after our processing. There is a significant increase of 14.8% in the overall accuracy at the Unicode level, which is encouraging.

We have tracked the actual value of γ that resulted in the best recognition of each of the test images. Figure 7 shows the histogram of the different values of γ automatically selected by our algorithm in enhancing the 1685 degraded images in the test database. It is interesting to see that the entire range of possible values has been used. Each step change in the value of γ decreases the dynamic range of the histogram and after eight steps, as shown in Fig. 4, there is no more information left to enhance.

Since Tesseract OCR supports Kannada script, we ran Tesseract v3 and v4 on each of the 1685 degraded word images in the benchmarking dataset, and computed the character recognition accuracy as per eqn. (2). Also, after we obtained the best enhanced version of each word image (based on our algorithm), we again ran it through both the versions of Tesseract Kannada OCR. Table III shows the recognition accuracies of Tesseract v3 on the test dataset, before and after applying our nonlinear enhancement on the word images. Our gamma enhanced binarization entails an increase of 5.6% in the recognition accuracy of Tesseract v3. Though Tesseract v4.0.0 gives 75.7% accuracy before enhancement, the accuracy after enhancing the images is unexplainably low (40.7%), and hence, we have not reported it in Table III.





Fig. 4: Illustration of the impact of the nonlinear enhancement due to a range of gamma values on the gray-level image of a degraded Kannada word with splits, its histogram and the corresponding binarized output image.

γ	Gray-level degraded image	Otsu binarized word image	Binarized image after GEB
2 ⁻¹ = 0.5	B. B B	ದ್. ಡಿಯೊ ಹೆ.ಡಿಯೊ	ದ್ರೌ.
2-2 = 0.25	ರಶ್ಷಣಾ	ರೆ ಕ್ಷ್ಮೆ ಣಾ ರಕ್ಷಣಠಿ	ರಸ್ಸೆ ಣಾ _{ರಕ್ಷಣಾ}
2-3 = 0.125	ದೇಶಗಳ	ದೇಶ ಗಳ ದೇಶ/ಸೊಳ	ದೇಶಗಳ ದೇಶಗಳ
2-4 =		*	a ma andad.
0.0625	0099440	್ ೧೯೧೪ ು ,ರಱಿಕೆ∕(ಳ,	ಗಿ ರಾಕಿಗಳು ಗಿರಾಕಿಗಳು
$2^{-6} = 0.015625$	ಭಾರೀ ಭಾರೀ	ಗೆ ೦ ಕ ೮ ಗಿ ೪೨ ,ರಱಿಕೆ∕(ಳ. ಧ್ ೪೮ೆ € -್ಇಫೀ	ಗಿರಾಕಿಗಳು ಗಿರಾಕಿಗಳು ಭಾರೀ ಭಾರೀ
$\begin{array}{c} 2^{-6} = \\ 0.0625 \\ \hline 2^{-6} = \\ 0.015625 \\ \hline 2^{-7} = \\ 0.0078125 \end{array}$	ಗಿರಾಕಿಗಳು ಭಾರೀ ದಿಕಾಸ	1:000190 ,ರಜಿಕೆ/(ಳ. -್ಇಫೀ ದಿ:5534 ವಿ''೬?ಳ	ಗರಾಕಿಗಳು ಗಿರಾಕಿಗಳು ಭಾರೀ ಭಾರೀ ನಿಕಾಸ ವಿಕಾಸ

Gray-level Test Image	Otsu Binarized Image	Gray-level Image	Otsu Binarized
ಇತಿಹಾ ನದ	ಇತಿಹಾ ನದ	ವ:ತ್ತು	ವಃತ್ತ್ತು
ದೇಶಗಳ	ದೇ ಶಗಳ	వి.ల్లనూ	ವಿ ಲ್ಲವೂ
1) (2.	r_ (8*	30	್ಟ್
ನಂಗಳು	<u>ನಂಗಳು</u>	2002	೧ ೮೪%
00, 550	ರ ಶೈಣಾ	31260	an fir
5.60	ಸ್ಪಂದು	52206	ಭಾರೀ
ນິວແລະປີ.	బిందె కరి.	ස්ළුදා	ಬೆಲೆ %್)
ಆರ್ಥಿ ದಿಂದಾಗಿ	16000, Tab	NG38ಗಳು	1.0.814.
8	2	82089	ອ ະ ຕໍ ເອ
1000000	ມອອວ.ນ ໜາສາວ ກ	ಗೇಣಿಯಲ್ಲ.	ಗೇಶಯ ್ದು.
ສະເດຍ	ವ್ಯೂಲ	8, nabr	ส,กรุธส์
చున్న	ను స్ట	ుత్ర	ింతు
ಸೇರಿದಂತೆ,	ಸೇರಿದಂಕೆ:	88.19 B	ಹೆ.ಇ ಶ

Fig. 5: Illustration of the fact that different words with varying levels of degradation require different values of gamma for optimal binarization. The last column shows the best binarized image, as determined by the recognition score and the corresponding recognized word.

Fig. 6: Some samples from the benchmark dataset created by us, containing degraded (with cuts) word images of Kannada language. The corresponding binary images, obtained by Otsu global thresholding, are also shown, illustrating the level of splits in the different characters.

TABLE II: Comparison of the Unicode recognition accuracies (in %) before and after gamma enhanced binarization (GEB), on Kannada benchmarking dataset of 1685 degraded word images. N, M: # of Unicodes in ground-truth and OCR output. N = 15,486. S, I and D: # of Unicode substitutions, insertions and deletions.

	Accuracy	М	S	Ι	D
Before GEB	74.5	$17,\!058$	1,962	1,782	210
After GEB	89.3	$15,\!954$	927	599	131
Improvement	+14.8	-1,104	-1,035	-1,183	-79



Fig. 7: Histogram of the different gamma values used by the algorithm for the images in the test database. It illustrates the extent to which different gamma values are automatically selected as being optimal for the different word images from the benchmark dataset of 1685 images.

TABLE III: Recognition accuracy (in %) of Tesseract OCR (v3.04.00) on the benchmark dataset, without and with gamma enhanced binarization (GEB). All the notations are the same as in Table II.

	Accuracy	М	S	Ι	D
Before GEB	36.7	19,136	5,582	$3,\!937$	287
After GEB	42.3	18,772	5,069	3,579	293
Improvement	+5.6	-364	-513	-358	+6

Figures 8 and 9 show the recognition results on all the representative samples from the benchmark Kannada degraded word image dataset, illustrated in Fig. 6, before and after the application of GEB. Errors in the recognized text are highlighted in red color. Figure 8 shows the cases, where GEB has managed to correct all the recognition errors that existed before applying our nonlinear enhancement to the word images. Figure 9 shows the cases, where the degradation/loss of information is so severe that even GEB does not suffice in the correct recognition of degraded words. In this figure, there are example images, where there is no improvement, there is partial correction of errors and also one case, where the error has increased after the automated, adaptive image enhancement.

Test Image – Otsu Binarized	-GEB Recognized Text	Best γ	GEB Image	+GEB Recognized Text
ದೇ ಶಗಳ	ದೇಶ <mark>/ಸ</mark> ೊಳ	2 ⁻³	ದೇ ಶಗಳ	ದೇಶಗಳ
ಗ್ರೀಕ್	ಳ್ರ್ಕ್	2-8	ಗ್ರೀಕ್	ಗ್ರೀಕ್
ನಂಗಳು	ನ <mark>ಸ</mark> ೊಗಳು	2-1	ನಂಗಳು	ನುಗಳು
ರಶ್ಷಣಾ	රಕ್ಷಣ <mark>ර</mark> ಿ	2-7	ರಕ್ಸೆ ಹಾ	ರಕ್ಷಣಾ
ಸೇರಿದಂತೆ.	ಸೇರಿದ <mark>ದೊರೆ</mark> :	2-4	ಸೇರಿದಂತೆ,	ಸೇರಿದಂತೆ,
୍ବଶ୍ଚ	ರಲ್ಲಿ	20	୍ବଶ୍ଚ	ರಲ್ಲಿ
	-್ಇರ್ಥಿ	2-6	ಭಾರೀ	ಭಾರೀ
ಬೆಲೆ ್ರ	ಬೆಲೆ <mark>ಲ</mark> '	2-7	బిలిం ు	ಬೆಲೆಯ
1.0.8145	,ರಱಿಕೆ/(ಳ,	2-7	ಗಿರಾಕಿಗಳು	ಗಿರಾಕಿಗಳು
အနှင့်အ	<mark> </mark>	2-4	ಕಾರಣ	ಕಾರಣ
ಗೇಣಿಯ ್ದು.	ರ್/ಫಿಯಭೌ,,?	2-4	ಗೇಣಿಯಲ್ಲ.	ಗೇಣಿಯಲ್ಲ.
ಗೈಗಾರಿಗೆ	ಕೈಲೌತರಿಕೆ	2-6	สำการส	ಕೈಗಾರಿಕೆ
ಹೆ.ಾಶ	ಹೆ,ಡಿಯೊ	2-7	ಹೊಸ	ಹೊಸ

Fig. 8: Illustration of some sample cases, where GEB has been successful in the correct recognition of the degraded words. - **GEB**: binarization of the original word image and the recognized word. + **GEB**: binarization after GEB and the word recognized. Errors in the recognized text in column 2 have been highlighted in red color.

Test Image – Otsu Binarized	-GEB Recognized Text	Best γ	GEB Image	+GEB Recognized Text
ಇತಿಹಾ ಸದ	ಇತಿಹಾ <mark>ವು</mark> ದ	2-4	ಇತಿಹಾ ನದ	ಇತಿಹಾದ
ವರಿಂದು	;១,ಗ <mark>[</mark> ರು	2 ⁻⁸	ನ್ನೂರು	;ಮೊರು
బింథాకార్,	ಬೆಂಥಾ <mark>ಯಿರ</mark> -,	2-1	బింథాకాన్,	ಬೆಂಥಾ <mark>ದ್ರ</mark> -,
೧೯ದಂದ _ೇ ದಲ	ಆದಿ ಘ್ ಯಾದಿಂದಾಗಿ	2 ⁻⁸	<i>ಗ೯ದಂದ್ಮಿ</i> ಗಿದಿಅ	ಆದಿ <mark>ಗ್</mark> ಯಿದಿಂದಾಗಿ
র .১০ ভ ় র	<mark>ವೆ.ಂತಾವೆ</mark>	2 ⁻²	ನ ಎಂತಾನ	<mark>ವೆ.ಂತಾವೆ</mark>
en 1830.	ಪುಸಕಾವೆ.	2 -2	ಉ ಸಕಾರ .	ಉಪಕಾ <mark>ವೆ</mark> .
ສະເຈຍ	ಮೂಜೊ	2-2	ສະເຈຍ	ಮೂಪ
వ ుస్ట	ಮ/,ನೈ	2 ⁻⁸	ను స్ప	ಷಾಟ್ತ
ವಃತ್ತು	ವಭೌತ್ತು	2-1	ನ:ತ್ತು	ವಯೆತ್ತು
ವಿ ಲ್ಲ ನೂ	ವಿಲ್ಲವು	2-4	ವಿ ಲ್ಲ ನೂ	<mark>ವಿ</mark> .ಲ್ಲನೂ
Des V	చి" ఙ ?ళ	2 -3	೧ ೯೪%	ವಿಕಾ ಳ
anr.	ಎತೊ	2-1	សក្រ	ವಲರ್/
8:2 3 2	ಱೆ,ತ್ತ್	2-5	ಿುತ್ತು	s1త్త

Fig. 9: Illustration of some sample cases, where the degradation/loss of information is so severe that even after GEB, the word could not be correctly recognized. - **GEB**: binarization of the original word image and the word recognized. + **GEB**: binarization after GEB and the word recognized. Errors in the recognized text have been highlighted in red color. There are example images, where there is no improvement, there is partial correction of errors and also one case, where the error has increased after the automated, adaptive image enhancement.

V. CONCLUSION

Locally adaptive, nonlinear enhancement of gray-level word images has been proposed for improved binarization, and hence, better OCR recognition results on old Kannada documents, such as printed books. The proposed technique of gamma enhanced binarization has resulted in an improvement of 14.8% in the character recognition rate on the created benchmark dataset of 1685 degraded Kannada word images. Of course, there are severely degraded images, which could not be fully recognized correctly, even after the application of GEB. More work is needed to handle such difficult cases. The reported results do not make use of any dictionary or contextual information to correct the words, and the presence of the latter modules could improve the results further, when the proposed enhancement is embedded in a good OCR.

Even on the Google's Tesseract OCR for Kannada, our gamma enhanced binarization results in an improvement of 5.6% in the Unicode level accuracy.

Together with a good pre-processing that splits the merged characters [5], this enhancement is very promising to obtain better recognition performance from any Kannada OCR on old printed books of low quality. Further, the proposed technique is generic, and hence can be applied to degraded document images of any script, provided the character classifier used provides reliable recognition scores or posterior probabilities.

References

- N. Sinha and A. G. Ramakrishnan, "Blood cell segmentation using EM algorithm," in Proc. Indian Conf. on Comp. Vision, Graphics and Image Processing (ICVGIP-02). ACM, 2002.
- [2] D. Kumar, M. N. Prasad, and A. G. Ramakrishnan, "MAPS: Midline analysis and propagation of segmentation," in *Indian* Conf. Comp. Vision, Graphics and Image Proc. ACM, 2012.
- [3] M. K. Jindal, R. K. Sharma, and G. S. Lehal, "A study of different kinds of degradation in printed Gurmukhi script," in *Computing: Theory and Applications. ICCTA'07. International Conf. on.* IEEE, 2007, pp. 538–544.
- [4] A. Madhavaraj, A. G. Ramakrishnan, H. R. Shiva Kumar, and N. Bhat, "Improved recognition of aged Kannada documents by effective segmentation of merged characters," in *Signal Processing and Communications (SPCOM)*, International Conference on. IEEE, 2014, pp. 1–6.
- [5] H. R. Shiva Kumar, A. Madhavaraj, and A. G. Ramakrishnan, "Splitting merged characters of Kannada benchmark dataset using simplified paired-valleys and L-cut," in *Proc. 25th National Conference on Communication*, 2019.
- [6] MILE-IISc. (2018, Dec) Kannada benchmarking dataset of degraded word images, with character splits. [Online]. Available: https://github.com/MILE-IISc/DegradedWordsKannada
- [7] S. Dutta, N. Sankaran, K. P. Sankar, and C. V. Jawahar, "Robust recognition of degraded documents using character n-grams," in *Document Analysis Systems (DAS)*, 10th IAPR International Workshop on. IEEE, 2012, pp. 130–134.
- [8] P. P. Kumar, C. Bhagvati, A. Negi, A. Agarwal, and B. L. Deekshatulu, "Towards improving the accuracy of Telugu OCR systems," in *Document Analysis and Recognition (ICDAR)*, Int. Conf. on. IEEE, 2011, pp. 910–914.
- [9] D. Arya, C. V. Jawahar, C. Bhagvati, T. Patnaik, B. Chaudhuri, G. S. Lehal, S. Chaudhury, and A. G. Ramakrishnan, "Experiences of integration and performance testing of multilingual OCR for printed Indian scripts," in *Proc. Joint workshop on multilingual OCR and analytics for noisy unstructured text data*. ACM, 2011, p. 9.

- [10] A. G. Ramakrishnan and K. Mahata, "A complete OCR for printed Tamil text," in *Proc. Tamil Internet 2000*. INFITT, 2002, pp. 53–57.
- [11] K. G. Aparna and A. G. Ramakrishnan, "A complete Tamil optical character recognition system," in *International Workshop on Document Analysis Systems*. Springer, 2002, pp. 53–57.
 [12] A. Kokku and S. Chakravarthy, "A complete OCR system for
- [12] A. Kokku and S. Chakravarthy, "A complete OCR system for Tamil magazine documents," in *Guide to OCR for Indic Scripts*. Springer, 2009, pp. 147–162.
- [13] B. Vijay Kumar and A. G. Ramakrishnan, "Machine recognition of printed Kannada text," in *International Workshop on Document Analysis Systems*. Springer, 2002, pp. 37–48.
- [14] —, "Radial basis function and subspace approach for printed Kannada text recognition," in *International Conf. on Acoustics*, *Speech, and Signal Processing.* IEEE, 2004, pp. V-321.
- [15] T. V. Ashwin and P. S. Sastry, "A font and size-independent OCR system for printed Kannada documents using support vector machines," *Sadhana*, vol. 27, no. 1, pp. 35–58, 2002.
- [16] M. M. Prasad, M. Sukumar, and A. G. Ramakrishnan, "Divide and conquer technique in online handwritten Kannada character recognition," in *Proc. International Workshop on Multilingual OCR.* ACM, 2009, p. 11.
- [17] B. Nethravathi, C. Archana, K. Shashikiran, A. G. Ramakrishnan, and V. Kumar, "Creation of a huge annotated database for Tamil and Kannada OHR," in *Proc. Inter. Conf. on Frontiers in Handwriting Recognition (ICFHR)*. IEEE, 2010, pp. 415–420.
- [18] D. Sachan, S. Dutta, T. Naveen, and C. Jawahar, "Segmentation of degraded Malayalam words: methods and evaluation," in *Computer Vision, Pattern Recog., Image Proc. and Graphics* (NCVPRIPG), 3rd National Conf. on. IEEE, 2011, pp. 70–73.
- [19] C. Sumetphong and S. Tangwongsan, "Recognizing broken characters in Thai historical documents," in Advanced Computer Theory and Engineering (ICACTE), 3rd International Conf. on, vol. 1. IEEE, 2010, pp. V1–99.
- [20] —, "Effectively recognizing broken characters in historical documents," in Computer Science and Automation Engineering (CSAE), IEEE Int. Conf. on, vol. 3, 2012, pp. 104–108.
- [21] W. Peerawit, W. Yingsaeree, and A. Kawtrakul, "The utilization of closing algorithm and heuristic information for broken character segmentation," in *Cybernetics and Intelligent Systems, IEEE Conf. on*, vol. 2, 2004, pp. 775–779.
- [22] D. Yu and H. Yan, "Reconstruction of broken handwritten digits based on structural morphological features," *Pattern Recognition*, vol. 34, no. 2, pp. 235–254, 2001.
- [23] L. Likforman-Sulem and M. Sigelle, "Recognition of degraded characters using dynamic Bayesian networks," *Pattern Recognition*, vol. 41, pp. 3092–3103, 2008.
- [24] M. Droettboom, "Correcting broken characters in the recognition of historical printed documents," in *Digital Libraries, Joint Conf. on.* IEEE, 2003, pp. 364–366.
- [25] F. Drira, F. LeBourgeois, and H. Emptoz, "Document images restoration by a new tensor based diffusion process: Application to the recognition of old printed documents," in *Document Analysis and Recognition, ICDAR. 10th International Conf. on.* IEEE, 2009, pp. 321–325.
- [26] D. Kumar and A. G. Ramakrishnan, "Power-law transformation for enhanced recognition of born-digital word images," in Signal Processing and Communications (SPCOM), International Conf. on. IEEE, 2012, pp. 1–5.
- [27] D. Kumar, M. N. Anil Prasad, and A. G. Ramakrishnan, "NESP: Nonlinear enhancement and selection of plane for optimal segmentation and recognition of scene word images," in *Document Recognition and Retrieval XX*, vol. 8658. International Society for Optics and Photonics, 2013, pp. 865–806.
- [28] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Transactions on Systems, Man and Cybernetics*, vol. 9, no. 1, pp. 62–66, 1979.
- [29] C.-C. Chang and C.-J. Lin, "LIBSVM: a library for support vector machines, software available at http://www.csie.ntu. edu. tw/~ cjlin/libsvm," 2001.