

A Kannada Document Image Retrieval System based on Correlation Method

Chandrakala H T
Acharya Institute of Technology
Bangalore, India,
chandrakl80@gmail.com

ABSTRACT

The present growth of digitization of books and manuscripts demands an immediate solution to access them electronically. This requires research in the area of document image understanding, specifically in the area of document image retrieval. There is an immense scope for such a retrieval system for a digital library of Kannada Document Images. This paper presents an efficient Content Based Document Image Retrieval system for a Kannada document image collection. A recognition-free approach is followed because recognition based approach is inefficient in terms of performance. The data is pre-processed and segmented for faster matching and retrieval. An efficient search technique - Correlation method is used to search in large collection of document images. Performance evaluation using different datasets of Kannada documents shows the effectiveness of the approach.

Keywords

Content based image retrieval, Correlation coefficient, Median Filtering, Kannada Document Images, Segmentation.

1. INTRODUCTION

Modern technology has made it possible to produce, process, store, and transmit document images efficiently. In an attempt to move towards the paperless office, large quantities of printed documents are scanned, digitized and stored as images in databases.

In recent years, there has been much interest in the research area of *Document Image Retrieval (DIR)*. The ever increasing amount of multimedia data creates a need for new sophisticated methods to retrieve the information one is looking for. *Content-Based Image Retrieval (CBIR)* uses image own information, usually by the similarity of image features (color, texture, shape and structure of the layout, and so on) and the comparable characteristics of each image for retrieval.

Document Image Retrieval systems are available for printed Roman, Japanese, Korean, Chinese, English and other oriental scripts. However, the availability of such products for Indian scripts is still a rarity. The present work addresses the issues involved in designing a font style and size independent Document Image Retrieval System for printed Kannada text. Kannada is the official language of the south Indian State-Karnataka. Retrieval from Kannada document images is more difficult than many other Indian scripts due to higher similarity in character shapes, a larger set of characters and higher variability across fonts in the characters belonging to the same class.

Aimed at document image, in this paper the image similarity is defined as having a similar correlation coefficient. The rest of the paper is organized as follows- First the related work in this field is described. Then characteristics of Kannada script are introduced. This is followed by the Present work. Experimental results and analysis are presented to show the effectiveness of the system.

2. RELATED WORK

2.1 Segmentation

Segmentation is the process of extracting objects of interest from an image. The Text Segmentation works reported till now basically involve three levels of Segmentation namely: Line Segmentation, Word Segmentation and Character Segmentation. A new word segmentation technique based on an efficient distinction of inter-word and intra-word distances was proposed by G. Louloudis et al 2009[11]. Siddhaling Urolagin et al 2010[23] found that a multi-channel Gabor filter decomposition represents an excellent tool for image segmentation and texture analysis and proposed a character segmentation method using Gabor filters.

Casey R.G. et al 1996[5] and Liang et al 1994 [16] have used the three elementary methods-- The Classical Approach, Recognition Based Segmentation, and Holistic methods of segmentation. Segmentation using smearing method is done by Wahl F.M et al 1996. G. Louloudis et al 2006[10] applied a block based Hough transform approach taking into account the gravity centers of parts of connected components, which are called blocks. Three different recognizers based on Hidden Markov Models are designed, and results of writer dependent as well as writer-independent experiments are reported by F. Luthy et al 2007[9].

2.2 Feature Extraction

Features are a set of numbers that capture the salient characteristics of the segmented image. B. Vijaykumar et al 2002[24] proposed a system which uses ANN classifier, Discrete Cosine Transform, Discrete Wavelet Transform for character recognition. Keerthi S S et al 2000 [14] proposed a fast iterative nearest point algorithm for support vector machine classifier design. O'Gorman L et al 1995[19] present a technique for design of vectoriser and feature extractor. Bansal V et al 1999[13] exploited certain features of the script in both reducing the search space and creating a reference with respect to which correspondence could be established, during the matching process. An OCR (Optical Character Recognition) system for Kannada script is described by Ashwin T V et al 2002[2] using Support Vector Machine (SVM) classifiers.

A two-stage multi- network neural classifier and wavelets for feature extraction based OCR for Kannada is developed by R

Sanjeev Kunte et al 2007[22]. VijayaKumar B et al 2004[24] presented an OCR system for basic Kannada characters, selecting the scheme of feature extraction using the moments and RBF neural networks as classifiers to identify and classify the characters. Nagabhushan P, Pai et al 1999[18] proposed a modified region decomposition method and optimal depth decision tree for the recognition of printed Kannada characters. Khotanzad A 1998[15] considered Hu's invariant moments and Zernike moments for feature extraction. Ramachandra Manthalkar and P.K. Biswas 2002 [21] have proposed a method based on rotation- invariant texture features using multichannel Gabor filter for identifying six (Bengali, Kannada, Malayalam, Oriya, Telugu and Marathi) Indian languages.

2.3 Retrieval

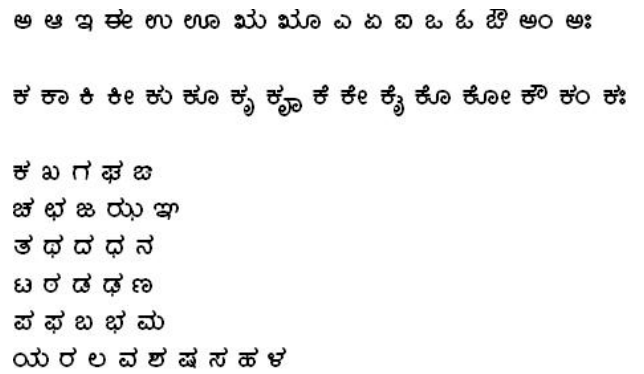
Two distinct techniques have been proposed in recent years for efficient search and retrieval from document image collections. They are namely Recognition Based Search and Recognition Free Search. Recognition Based Search approach is to convert the images into text and then apply a search engine. Recognition Free Search approach is to search directly in the images bypassing the recognition task. HOU Dewen et al 2010[13] describe a Content Based Document Image Retrieval System for document image database using hierarchical matching tree algorithm. E. Garcia [8] discusses Correlation Coefficient and its Confidence Intervals, and how to normalize its distribution through Fisher Transformation. An effective word image matching scheme that achieves high performance in the presence of script variability, printing variation, degradation and word-form variants is presented by Million Meshesha and C. V. Jawahar[17]

Most of the reported work is on the design of an OCR (Optical Character Recognition) system for Kannada script. A Content Based Document Image Retrieval System for Kannada Document Images of varying font styles and sizes is not yet available. Retrieval from Kannada document images is more difficult due to higher similarity in character shapes, a larger set of characters and higher variability across fonts in the characters belonging to the same class. Searching and retrieval from document image collections is challenging because of the segmentation errors, scalability issues and computational time. This paper presents an entirely new and fresh innovation in this field.

3. CHARACTERISTICS OF KANNADA SCRIPT

The Kannada alphabet is classified into two main categories: vowels and consonants. There are 16 vowels and 34 consonants as shown in Figure 1. Words in Kannada are composed of aksharas which are analogous to characters in an English word. While vowels and consonants are aksharas, the vast majority of aksharas are composed of combinations of these in a manner similar to most other Indian scripts. The aksharas of Kannada script may belong to any of these three categories: i) a stand alone vowel or a consonant ii) a consonant modified by an vowel iii) a consonant modified by one or more consonants and an vowel.

Figure 1. The aksharas of Kannada script



An vowel or a consonant can form a whole akshara. The aksharas of second category are formed by inflecting an vowel to a consonant. Here the akshara is formed by combining glyph corresponding to a consonant and vowel modifier's glyph which may occur independently or agglutinated with the consonant. The resulting aksharas may have several components. The vowel modifier glyphs are different from those of the vowels; all glyphs corresponding to vowels. The vowel modifier glyphs attach to the consonant glyphs up to three places corresponding to the top, right and bottom positions of the consonant.

In the third category, an akshara is formed by the combination of many consonants with an vowel. In practice at most three consonants are combined with an vowel to form an akshara. The first consonant is called as base consonant which is then combined with the modifiers of one or more consonants (these consonants are called as consonant conjunct or vatthu). Here an vowel modifier is attached to base consonant. The glyphs of many consonant conjuncts resemble those of the consonants. The consonant conjunct glyphs always appear below the base consonant and vowel combination.

4. PRESENT WORK

This work mainly aims at addressing some of the issues involved in effective and efficient retrieval of Kannada document images with variations in font-style and font-sizes. Figure 2 below shows the architecture of the system. Scanned Kannada document images are stored in the database. Each image is preprocessed to remove noise and extract co-ordinates. This system accepts a textual query in English from the user. The textual query is first converted to Kannada text and then an image by rendering. Segmentation is performed at word and character level. Features are extracted from these character images and then a search is carried out for retrieval of relevant documents. Results of the search are pages from Kannada document image collections containing the query word sorted based on their relevance to the query.

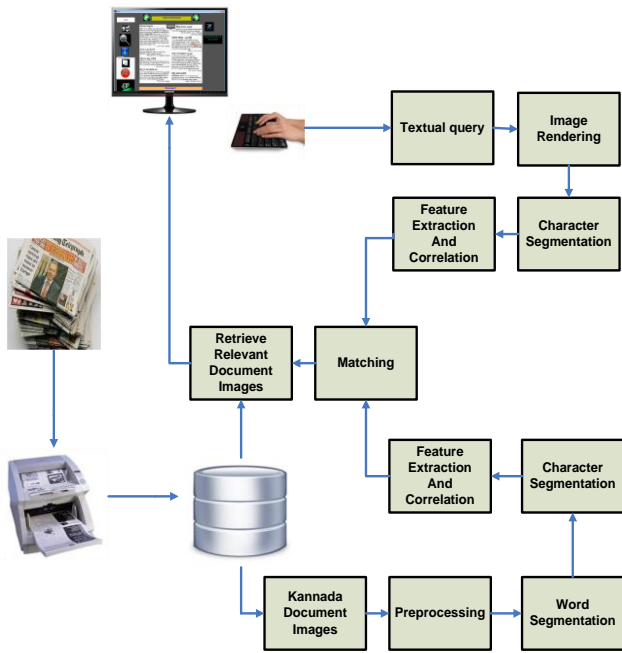


Figure 2. System Architecture: The textual query word is segmented upto character level. Characters of the query word are matched with those of the document images in the document image collection. The documents containing the words that match the query word are retrieved and displayed to the user.

4.1 Preprocessing

Preprocessing is required to prepare the source image for matching and retrieval. The scanned images are heavily infected by noise. To remove this noise, the images are passed through a median filter. The co-ordinate information is also extracted in order to speedup the matching process.

Median filtering is a nonlinear process useful in reducing impulsive, or salt-and-pepper noise. It is one of the better filters, as it also preserves edges in an image while reducing random noise. Salt-and pepper noise can occur due to distortion in the scanner. In a median filter, a window slides along the image, and the median intensity value of the pixels within the window becomes the output intensity of the pixel being processed. This is illustrated in the figure below

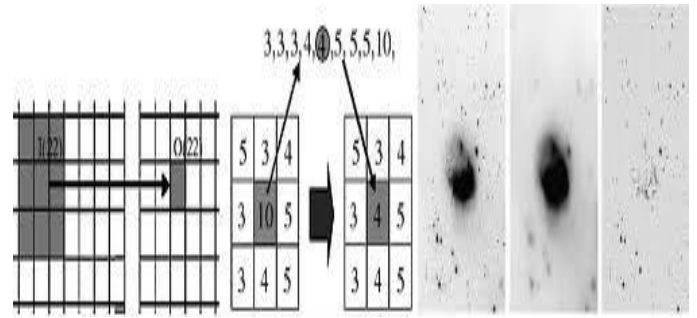
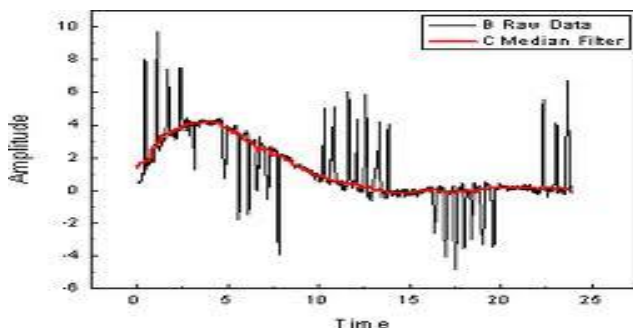


Figure 3 Median filtering smooths the image and is thus useful in reducing noise. Median filtering can preserve discontinuities in a step function and can smooth a few pixels whose values differ significantly from their surroundings without affecting the other pixels as shown in the images.

4.2 Segmentation

The source image, which is in the RGB form, is binarized. Morphological dilation is performed on this binary image both vertically and horizontally. After dilation the words in the document image appear as connected components. These connected components of the binary image are extracted based on the average height width of the words. The Figure 4 given below shows an illustration of word segmentation.



Figure 4. A paragraph of text from a document image after the document is segmented up to word level. Each word is highlighted in a rectangular box.

Next using the spacing between the characters and pixel summation, character segmentation is done on the extracted words. This is illustrated in Figures 5 and 6 below:

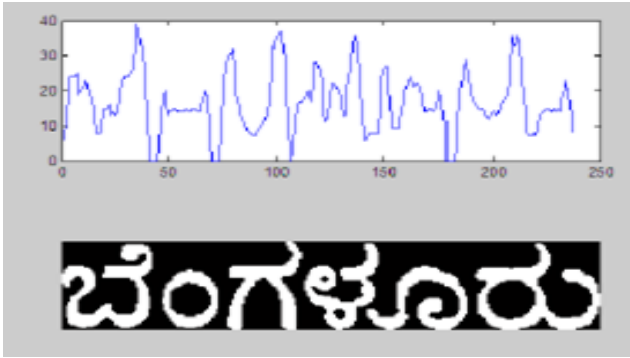


Figure 5. Valleys correspond to the spacing between the characters. Pixels between each two valleys are summed up to extract individual characters of a word.



Figure 6. Individual characters of a word “Bengaluru” after Character Segmentation

4.3 Matching and Retrieval

The individual characters of the query word and that of the words in the database documents extracted through character segmentation are matched using correlation method. Correlation measures the association between variables. Correlation is a measure of, the strength of relationship between random variables. The population correlation between two variables X and Y is defined as:

$$\rho = \frac{\text{Covariance}(X,Y)}{\sqrt{\{\text{Variance}(X) * \text{Variance}(Y)\}}}$$

ρ is called the Product Moment Correlation Coefficient or simply the Correlation Coefficient. It is a number that summarizes the direction and closeness of linear relations between two variables. The sample value is called r , and the population value is called ρ (rho).

The correlation coefficient, sometimes also called the cross-correlation coefficient, is a quantity that gives the quality of a least squares fitting to the original data. To define the correlation coefficient, first consider the sum of squared values SS_{xx} , SS_{xy} and SS_{yy} of a set of n data

points (x_i, y_i) about their respective means:

$$\begin{aligned} SS_{xx} &= \sum (x_i - \bar{x})^2 \\ &= \sum x^2 - 2\bar{x} \sum x + \sum \bar{x}^2 \end{aligned}$$

$$\begin{aligned} &= \sum x^2 - 2n\bar{x}^2 + n\bar{x}^2 \\ &= \sum x^2 - n\bar{x}^2 \\ SS_{yy} &= \sum (y_i - \bar{y})^2 \\ &= \sum y^2 - 2\bar{y} \sum y + \sum \bar{y}^2 \\ &= \sum y^2 - 2n\bar{y}^2 + n\bar{y}^2 \\ &= \sum y^2 - n\bar{y}^2 \\ SS_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) \\ &= \sum x_i y_i - \bar{x} \sum y_i - \sum x_i \bar{y} + \bar{x} \bar{y} \\ &= \sum xy - n\bar{x}\bar{y} - n\bar{x}\bar{y} + n\bar{x}\bar{y} \\ &= \sum xy - n\bar{x}\bar{y}. \end{aligned}$$

These quantities are simply unnormalized forms of the variances and covariance of X and Y given by

$$\begin{aligned} SS_{xx} &= N \text{ var}(X) \\ SS_{yy} &= N \text{ var}(Y) \\ SS_{xy} &= N \text{ cov}(X, Y) \end{aligned}$$

For linear least squares fitting, the coefficient b in

$$y = a + bx$$

is given by

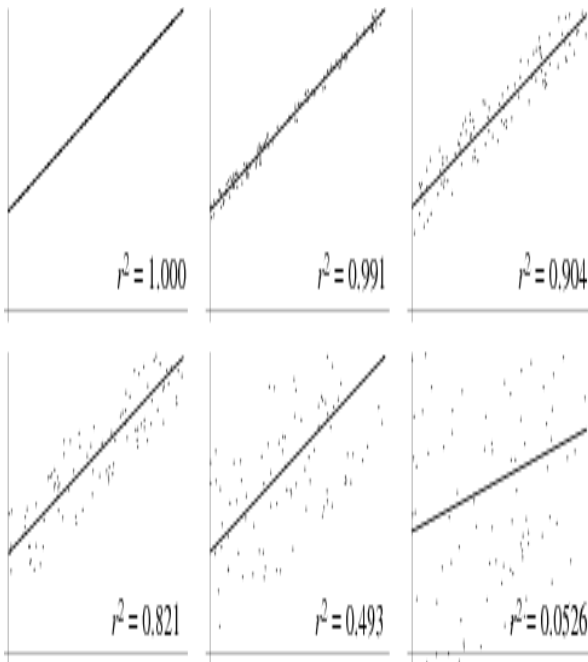
$$\begin{aligned} b &= \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} \\ &= \frac{SS_{xy}}{SS_{xx}} \end{aligned}$$

and the coefficient b' in

$$x = a' + b' y$$

is given by

$$b' = \frac{n \sum xy - \sum x \sum y}{n \sum y^2 - (\sum y)^2}$$



The correlation coefficient r (sometimes also denoted by R) is

then defined by

$$r^2 = bb'$$

$$= \frac{SS_{xy}^2}{SS_{xx}SS_{yy}}$$

The correlation coefficient is also known as the product-moment coefficient of correlation or Pearson's correlation. The correlation coefficients for linear fits to increasingly noisy data are shown above.

The correlation coefficient has an important physical interpretation. To see this, define

$$A = \left[\sum x^2 - n\bar{x}^2 \right]^{-1}$$

and denote the "expected" value for y_i as \hat{y}_i . Sums of \hat{y}_i are

then

$$\begin{aligned} \hat{y}_i &= a + bx_i \\ &= \bar{y} - b\bar{x} + bx_i \\ &= \bar{y} + b(x_i - \bar{x}) \\ &= A \left(\bar{y} \sum x^2 + (x_i - \bar{x}) \sum xy - n\bar{x}\bar{y}x_i \right) \\ &= A \left(\bar{y} \sum x^2 - \bar{x} \sum xy + x_i \sum xy - n\bar{x}\bar{y}x_i \right) \end{aligned}$$

$$\begin{aligned} &= A^2 [n\bar{y}^2 (\sum x^2)^2 - n^2 \bar{x}^2 \bar{y}^2 (\sum x^2) \\ &\quad - 2n\bar{x}\bar{y} (\sum xy) (\sum x^2) + 2n^2 \bar{x}^{-3} \bar{y} (\sum xy) \\ &\quad + (\sum x^2) (\sum xy)^2 + n\bar{x}^2 (\sum xy)] \end{aligned}$$

$$\begin{aligned} \sum y_i \hat{y}_i &= \\ &= A [n\bar{y}^2 \sum x^2 + (\sum xy)^2 + n\bar{x}\bar{y}] \\ &= A [n\bar{y}^2 \sum x^2 + (\sum xy)^2 - 2n\bar{x}\bar{y} \sum xy] \end{aligned}$$

The sum of squared errors is then

$$\begin{aligned} &= \sum (\hat{y}_i - \bar{y})^2 \\ &= \sum (\hat{y}_i^2 - 2\bar{y}\hat{y}_i + \bar{y}^2) \\ &= A^2 \left(\sum xy - n\bar{x}\bar{y} \right)^2 \left(\sum x^2 - n\bar{x}^2 \right) \\ &= \frac{(\sum xy - n\bar{x}\bar{y})^2}{\sum x^2 - n\bar{x}^2} \\ &= \frac{bSS_{xy}}{SS_{xx}} \\ &= \frac{SS_{xy}^2}{SS_{yy}} r^2 \\ &= b^2 SS_{xx} \end{aligned}$$

and the sum of squared residuals is

$$\begin{aligned} SSR &= \sum (y_i - \hat{y}_i)^2 \\ &= \sum (y_i - \bar{y} + b\bar{x} - bx_i)^2 \\ &= \sum (y_i - \bar{y})^2 + b^2 \sum (x_i - \bar{x})^2 - 2b \sum (x_i - \bar{x})(y_i - \bar{y}) \\ &= SS_{yy} + b^2 SS_{xx} - 2b SS_{xy}. \end{aligned}$$

But

$$\begin{aligned} b &= \frac{SS_{xy}}{SS_{xx}} \\ r^2 &= \frac{SS_{xy}^2}{SS_{xx}SS_{yy}} \end{aligned}$$

so

$$\begin{aligned} SSR &= SS_{yy} + \frac{SS_{xy}^2}{SS_{xx}} SS_{xx} - 2 \frac{SS_{xy}}{SS_{xx}} SS_{xy} \\ &= SS_{yy} - \frac{SS_{xy}^2}{SS_{xx}} \end{aligned}$$

$$= SS_{yy} \left(1 - \frac{SS_{xy}^2}{SS_{xx}SS_{yy}} \right)$$

$$= SS_{yy}(1 - r^2)$$

and

$$SSE + SSR = SS_{yy}(1 - r^2) + SS_{yy} r^2 = SS_{yy}$$

The square of the correlation coefficient r^2 is therefore given

by

$$r^2 = \frac{SSR}{SS_{yy}}$$

$$= \frac{SS_{xy}^2}{SS_{xx}SS_{yy}}$$

$$= \frac{(\sum xy - n\bar{x}\bar{y})^2}{(\sum x^2 - n\bar{x}^2)(\sum y^2 - n\bar{y}^2)}$$

In other words, r^2 is the proportion of SS_{yy} which is

accounted for by the regression. If there is complete

correlation, then the lines obtained by solving for best fit (a, b) and (a', b') coincide since all data points lie on

them.

$$y = \frac{a'}{b'} + \frac{x}{b'} = a + bx.$$

Therefore, $a = -a'/b'$ and $b = 1/b'$ giving

$$r^2 = b b' = 1.$$

The correlation coefficient is independent of both origin and scale, so

$$r(u, v) = r(x, y)$$

where

$$u = \frac{x - x_0}{h}$$

$$v = \frac{y - y_0}{h}$$

The correlation coefficient can take values between -1 through 0 to +1. The sign (+ or -) of the correlation defines the direction of the relationship. When the correlation is positive ($r > 0$), it means that as the value of one variable increases, so does the other. If a correlation is negative ($r < 0$), it indicates that when one variable increases, the other variable decreases. This means there is an inverse relationship between the two variables. The formula for Correlation Coefficient is:

$$r = \frac{\sum_m \sum_n (A_{mn} - \bar{A})(B_{mn} - \bar{B})}{\sqrt{(\sum_m \sum_n (A_{mn} - \bar{A})^2)(\sum_m \sum_n (B_{mn} - \bar{B})^2)}}$$

where \bar{A} = mean(A), and \bar{B} = mean(B)

Correlation method gives the correlation coefficient between the individual characters of the query word and the individual characters of all the words in all the documents of the database. This correlation coefficient is the matching score between the words which are under comparison as shown in Figure 7 below. The words with a matching score higher than the threshold are highlighted and the containing documents are retrieved and displayed to the user based on their relevance to the query.

ಗಂಭೀರವಾಗಿದ್ದು, ಈಗಾಗಲೇ 224 ಗ್ರಾಮಗಳು ಮತ್ತು
ನಗರ ಸ್ಥಳೀಯ ಸಂಸ್ಥೆಗಳ ವ್ಯಾಪ್ತಿಯ 106 ವಾರ್ಡ್‌ಗಳಿಗೆ
ಟ್ರಾಂಕರ್ ಮೂಲಕ ನೀರು ಪೂರೈಸಲಾಗುತ್ತಿದೆ.
ಬರುವ ದಿನಗಳಲ್ಲಿ 4,853 ಗ್ರಾಮಗಳಲ್ಲಿ ಕುಡಿಯುವ
ನೀರಿನ ಸಮಸ್ಯೆ ಉಂಟಾಗಬಹುದು ಎಂದು ಅಂದಾಜಿಸಿ
ಲಾಗಿದ್ದು, ನೀರಿನ ಸಮಸ್ಯೆ ಮತ್ತಷ್ಟು ಉಲ್ಬಣಗೊಳ್ಳುವ
ಲಕ್ಷಣಗಳು ಗೋಚರಿಸುತ್ತಿವೆ. ಬರದ ಛಾಯೆ ಆವರಿಸಿ
ರುವ ಬೆನ್ನಲ್ಲೇ ಕುಡಿಯುವ ನೀರು ಮತ್ತು ಜಾನುವಾರು
ಗಳ ಮೇವಿಗೆ ಹಾಹಾಕಾರ ಉಂಟಾಗಿದೆ.
ಇದುವರೆಗೆ ಕೇವಲ ಎರಡು ಮೇವಿನ ಬ್ಯಾಂಕ್ ಮತ್ತು
34 ಗೋಶಾಲೆಗಳನ್ನು ತೆರೆಯಲಾಗಿದೆ. ಇನ್ನಷ್ಟು ಗೋಶಾ

Figure 7. The words with a Correlation Coefficient value above 0.6 are considered to be matching. Such words are highlighted in red color.

Sl No	Type of query word	Precision Rate	Recall Rate
1	Single letter words	90%	93%
2	Two letter words	86%	84%
3	Three letter words	80%	84%
4	Above three letter words	80%	82%
5	Words with Touching characters	65%	70%
6	Words with Confusion characters	65%	68%

5. EXPERIMENTAL RESULTS AND ANALYSIS

The experimental data are the document images of articles of different Kannada newspapers which use different font styles and font sizes. The newspaper articles are scanned at 600 dpi, 256 gray-level. The image size is approximately [2904, 2808] and the quantity of data is 3.5MB. The coordinate information and width and height of the words in the source image are calculated and stored in the database as part of preprocessing. This facilitates faster matching and retrieval. Using Content Based Image Retrieval method, the system retrieves the documents which have words matching with the query word irrespective of their font style and size. Figure 8 shows is a sample output of the system. The Performance graph of the system is given in Figure 9. Table 1 gives the detailed analysis of the system performance.



Figure 8. A document from the document image collection, which contains the words that match query word -“bengaluru”.

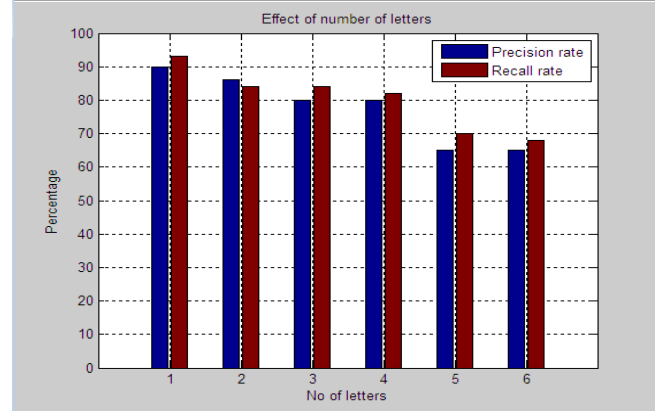


Figure 9: A graph showing the system performance

Table 1: Comparing the performance of the system for different kinds of queries. The performance is better for most of the query types except for words with touching and confusion characters

6. CONCLUSION

We have explored different approaches for addressing the document analysis and understanding problems with emphasis on scripts of Kannada language. After extensive script analysis, we have designed a new font style and size independent Content Based Image Retrieval algorithm for Kannada document image collection.

The similarity matching algorithm has high precision and less algorithm complexity. The performance of the method is presented with the help of experiments on real data sets from Kannada language documents. It can be concluded from the experiments that, the system performed efficient retrieval from large document image collections in few milliseconds.

Segmentation of Kannada Document Images is quite a challenge currently. Segmentation of the touching lines and characters may require some heuristic approaches. The option of search in multi-lingual documents could be another possible direction for research.

7. REFERENCES

- [1] A. Balasubrahmanian, Million Meshesha, C. V. Jawahar, "Retrieval from Document Image Collections", In Proc. of the 7th IAPR Workshop on Document Analysis Systems (DAS) (LNCS), pp 1-12, 2006.
- [2] Ashwin T V, Sastry P S 2002 "A font and size-independent OCR system for printed Kannada documents using support vector machines." *Sadhana* 27: 35–58.
- [3] Bansal V, Sinha R M K 1999 "On how to describe shapes of Devanagari characters and use them for recognition." In *Proc. Fifth Int. Conf. on Document Analysis and Recognition*, Bangalore (IEEE Computer Society Press) pp
- [4] B. Vijaykumar and A. G. Ramakrishnan, "Machine Recognition of Printed Kannada Text," in *Proceedings of the Fifth International Workshop on Document Analysis Systems*. 2002, pp. 37–48, Springer, Berlin.
- [5] Casey, R.G. and Lecolinet, E., "A Survey of Methods and Strategies in Character Segmentation". IEEE

- Transactions on Pattern Analysis and Machine Intelligence, 1996, Vol.18, No.8, pp.690-706.
- [6] Choudhury B B, Pal U 1997 "An OCR system to read two Indian language scripts: Bangla and Devanagari." In *Proc. Fourth Int. Conf. on Document Analysis and Recognition* (IEEE Computer Society Press) pp 1011–1015
- [7] David Shen, Zaizai Lu "Computation of Correlation Coefficient and Its Confidence Interval in SAS"
- [8] E. Garcia "A Tutorial on Correlation Coefficients"
- [9] F. Luthy, T. Varga, H. Bunke, "Using Hidden Markov Models as a Tool for Handwritten Text Line Segmentation", *Ninth International Conference on Document Analysis and Recognition*, Curitiba, Brazil, 2007, pp. 8-12.
- [10] G. Louloudis, K. Halatsis, B. Gatos, I. Pratikakis, "A Block-Based Hough Transform Mapping for Text Line Detection in Handwritten Documents", *10th International Workshop on Frontiers in Handwriting Recognition (IWFHR 2006)*, La Baule, France, October 2006, pp. 515- 520.
- [11] G. Louloudis, B. Gatos, I. Pratikakis, C. Halatsis" Line and Word Segmentation of Handwritten Documents"031-3203/©2009 Elsevier Ltd.
- [12] Gonzalez R C, Woods R E 1993 *Digital image processing* (Boston, MA, USA: Addison Wesley Longman Publishing Co. Inc.)
- [13] HOU Dewen, WANG Xichang, LIU Jiang "A Content-Based Retrieval Algorithm for Document Image Database" 978-1-4244-7874-3/10/©IEEE 2010 Crown
- [14] Keerthi S S, Shevade S K, Bhattacharyya C, MurthyKRK2000 "A fast iterative nearest point algorithm for support vector machine classifier design." *IEEE Trans. Neural Networks* 11: 124–136
- [15] Khotanzad A 1998 "Rotation invariant pattern recognition using Zernike moments." *Proc. Int. Conf. on Pattern Rec.* 326–328
- [16] Liang, S., Shridhar, M. and Ahmadi, M., "Segmentation of Touching Characters in Printed Document Recognition". *Pattern Recognition*, 1994, Vol.27, No.6, pp.825-840.
- [17] Million Meshesha , C. V. Jawahar "Matching word images for content-based retrieval from printed document images" DOI 10.1007/s10032-008-0067-3
- [18] Nagabhushan P, Pai Radhika M 1999" Modified region decomposition method and optimal depth decision tree in the recognition of non-uniform sized characters—An experimentation with Kannada characters." *Pattern Rec. Lett.* 20: 1467–1475.
- [19] O’Gorman L, Kasturi R 1995 *Document image analysis* (IEEE Computer Society Press)
- [20] Pavlidis T 1986 "A vectorizer and feature extractor for document recognition." *Computer. Vision Graphics Image Processing.* 35: 111
- [21] Ramachandra Manthalkar and P.K. Biswas, "An Automatic Script Identification Scheme for Indian Languages", NCC, 2002.
- [22] R Sanjeev Kunte, R D Sudhaker Samuel., 2007. An OCR system for printed Kannada text using Two-stage Multi-network classification approach employing Wavelet Features. *Proc. International Conference on Computational Intelligence and Multimedia Applications* (IEEE Computer Society Press.), 349-353.
- [23] Siddhaling Urolagin, Prema K. V, V.Subba Reddy "A Gabor Filters Based Method for Segmenting Inflected Characters of Kannada Script" 978-1-4244-6653-5/10/©2010 IEEE
- [24] VijayaKumar B, Ramakrishnan A G 2004 "Radial basis function and sub-space approach for printed Kannada Text recognition." *Proc. IEEE ICASSP 2004* 5: 321–324.
- [25] Wahl F.M., Wong, K.Y., Casey R.G.: "Block Segmentation and Text Extraction in Mixed Text/Image Documents" *Computer Graphics and Image Processing*, 20 (19 82) 375- 390.

AUTHOR

Chandrakala H T

Assistant Professor
Department of Information Science and Engg. Acharya
Institute of Technology
Bangalore, India.