

Writer Identification using Texture Features in Kannada Handwritten Documents

Praveen Bangarimath
M.Tech(IT),Dept of ISE
SDM College of Engineering
&Technology, Dharwad
Karnataka, India.

Deepa Bendigeri
Assistant Professor,ISE Dept
SDM College of Engineering
&Technology, Dharwad
Karnataka, India.

Jagadeesh Pujari
PhD, HOD, ISE Dept
SDM College of Engineering
&Technology, Dharwad
Karnataka, India

ABSTRACT

Writer Identification has great scope in emerging technology due to its usage in various types of applications in biometric and forensic science. Aim of this work is to identify the writer from script which is handwritten and obtained as scanned images. Features of textures will be elicited from wavelet decomposed images based on co-occurrence histograms. These will get the information about the relations among the sub bands of less frequency and that in sub bands of higher frequency at the particular level of the transformed image. If the co-relation between the sub bands has the same resolution then that indicates a stronger resolution. The relationship will indicate as information was essential considered to differentiating the textures. The proposed methodology will be executed with kannada handwritten document images by considering 14 different writers. Ability of features from texture in identifying writers is indicated through outcome achieved in experimentation.

Keywords

Wavelet, Texture Feature ,Document images ,Scanned images ,Co-occurrence histograms.

1. INTRODUCTION

Writer identification is a framework designed for recognizing the writer of a manually written record. An arrangement of reports from known authors must be known ahead of time to select another archive to one of this writer[1]. To begin with, components (features) are registered on the writer of a reference report and after that these feature parts will be appeared differently in relation to the ones which exist in the database set. An author having most noteworthy closeness to the existing is related to archive[3]. Writer identification distinguishes the handwriting methodology by taking into account of unknown handwriting image which continues by coordinating obscure handwritten image against a database of composed examples with known writer[1].In this manner, writer identification is imperative with numerous applications, for example, report examination, security, and monetary action, measurable and utilized as access control. The difficulties for writer identification and retrieving the writer incorporates from various pen features, that differs in the writer handwriting style of their composition, if certainty which the author (writer) has composed content in the rush or not, furthermore that single word is uncommon composed the very same way twice. Writer identification technique falls in two categories: Text dependent and Text independent.[1,2] Text-dependent approaches require handwritten tests taking into account on a particular content, or expect handwriting recognizer accessible for checking realness of writer. Writer identification utilizing signature is most prominent occurrence of these sorts of methodologies. Text-dependent approaches

have advantages that they utilize the information of the substance of the information to isolated style from substance. This will build the precision of text-dependent frameworks. The significant issue of text-dependent frameworks is that they are non-relevant to situations when the content is not accessible, for example, in criminal equity frameworks includes examination between content arch[1]ives with various substance. Also, message subordinate frameworks are much inclined to falsification when same information is displayed during testing. These sorts of frameworks can be actualized in the agreeable environment, where significant concern is exactness and writer might be asked to author particular substance to demonstrate their personality[2]. These strategies are fundamentally the same as signature check procedures which includes the correlation between individual characters or words which have known semantic substance. Subsequently these techniques require earlier restriction and division of the right data[3,5], which is generally performed by human client collaboration. The text-independent writer identification framework displays the style data, free of the content, which is utilized to recognize the writer in view of any given text content. This requires the insights of features which are computed from large extensive quantity of data to avoid anomalies due to specific content[7]. Identification of writer and verification techniques utilizes measurable elements which are extracted from entire image containing a text block. Base measure of manual written (e.g. few content lines in section) is extremely important to determine stable features components which will be hard to the text content of the examples.[4] Consequently the technique will fall in this later class. Writer identification includes input of two types: on-line and off-line[5]. Online technique includes catching of pen development of author, where the style is compelled. Online writer identification framework uses temporal succession code, which tracks pressure and velocity(speed) varies in handwriting, and pattern(shape codes) that relay on direction of trajectory in writing was developed for Chinese and English language [5]. It works better for little number of characters. Online text-independent writer recognition framework [4] for the language Thai depends on speed of pen pointer utilizes Fourier change technique.In offline writer identification framework, scanned image of the author composing is utilized which delineates his behaviour. Offline text independent writer identification using Hidden Markov Model [6] works on the basis of computing the score unknown author and comparing it scores of every individual author[2,3]. The score of every individual author is computed by recognizer in view of hand writing. The recognizer with the most astounding score is assigned as unknown writer. In offline writer identification framework, the hand written text of the writer is filtered and utilized for feature extraction. In that capacity offline writer identifications postures more

difficulties contrasted with [7] on-line writer identification method on account of the absence of extra features, which are accessible to online frameworks, is absent for offline system [3,4,5]. Statistical based Writer identification method for non consistently skewed handwriting images has been discussed in [6]. Different strategies for identification of writer depends on Contour based features, Hierarchical Shape Primitive elements, has been talked about in [5]. Offline text-independent writer identification is very imperative for measurable examination, archives approval, and calligraphic relics ID [6].

2. FEATURE EXTRACTION

2.1 Discrete wavelet transform

The continuous 1-D transform wavelet of (1D) signal $F(m)$ described as

$$(WaF)(b) = \int F(x)\psi_{a,b}^*(x)dx \quad (1)$$

ψ represents wavelet calculated according to Base ψ wavelet interpretation, expansion

$$\psi_{x,y}(M) = \frac{1}{\sqrt{|a|}} \psi\left(\frac{x-a}{b}\right) \quad (2)$$

Due to certain conditions, base wavelet ψ fulfills limitation which includes zero mean. It is differentiated through control of x,y to distinct cross section ($x=2^y$, $b \in \mathbb{I}$). Normally, forced it changes to non excess, integrated and that it should include multi-resolution primitive signal view. Expansion of Two-D is generally operates utilizing the result of One-D wavelet channels. The Haar wavelet is defined as

$$\psi(t) = \begin{cases} 1, & 0 \leq t \leq 1/2 \\ -1, & 1/2 \leq t < 1 \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

Above equation is much easier, where it checks

$$\left\{ \psi_{n,a} = \left(2^{-\frac{n}{2}} \psi \left(2^{-\frac{n}{2}} s - a \right) \right) \right\}_{n, a \in \mathbb{Z}}$$

is orthogonal and unit vectors premise

for $L^2(\mathbb{R})$. In this, discretization $x=2^a$ and $y=s2^a$ is utilized, which will be sought after all through this area. This wavelet is, verifiably, primarily identified wavelet.

The simplest way to compute a 2D discrete wavelet transform (DWT) of an image is to apply one-dimensional transform over image rows and columns separately and then to carry out down sampling. This transform decomposes an image with the overall scale factor of four, providing at each level one low resolution subimage and three wavelet coefficient subimages. [4]

$$\begin{aligned} A &= |S_x * |S_y * I|_{2,1}|_{1,2} \\ H &= |T_x * |S_y * I|_{2,1}|_{1,2} \\ V &= |S_x * |T_y * I|_{2,1}|_{1,2} \\ D &= |T_x * |T_y * I|_{2,1}|_{1,2} \end{aligned} \quad (4)$$

Here I is the input image. S_x, S_y and T_x, T_y represent low and high pass filters respectively, $*$ denotes the convolution operator and 2 denotes down sampling operation. The sub bands labeled H, V, D correspond to the detail images, while A corresponds to the approximation image.

2.2 Computational scheme

Script or style types usually vary from each other by the way they are assembled or grouped into words, and also the state of individual attribute, and so forth. This gives diverse scripts particularly distinctive visual appearance. Texture can be characterized in simple definitive form as “same pattern occurring repetitively” or something comprising of commonly related elements. This identification of script or writer style from the handwritten images consist features which are based on texture, extracted from handwritten images provided by writer in Kannada Language. The feature extraction method is described below. The computation scheme is extraction of features influenced by perception of individuals is talented for recognition among new writings simply in view of simple visualization analysis [3]. Classification of texture is processed by considering the Identification of the script. Hence, this is complicated visible texture made out derived by sub-pattern.

Despite of fact that, the sub-patterns can have scarcity of a better mathematical standard, it is well entrenched that a texture is considered as completely only if all the sub-patterns are correctly defined [5]. It utilizes a multi-resolution method as elicitation of texture features according to DWT and by using minimum distance classifier, the classification of the textures is obtained. This extraction of features is illustrated as shown in Fig 1.

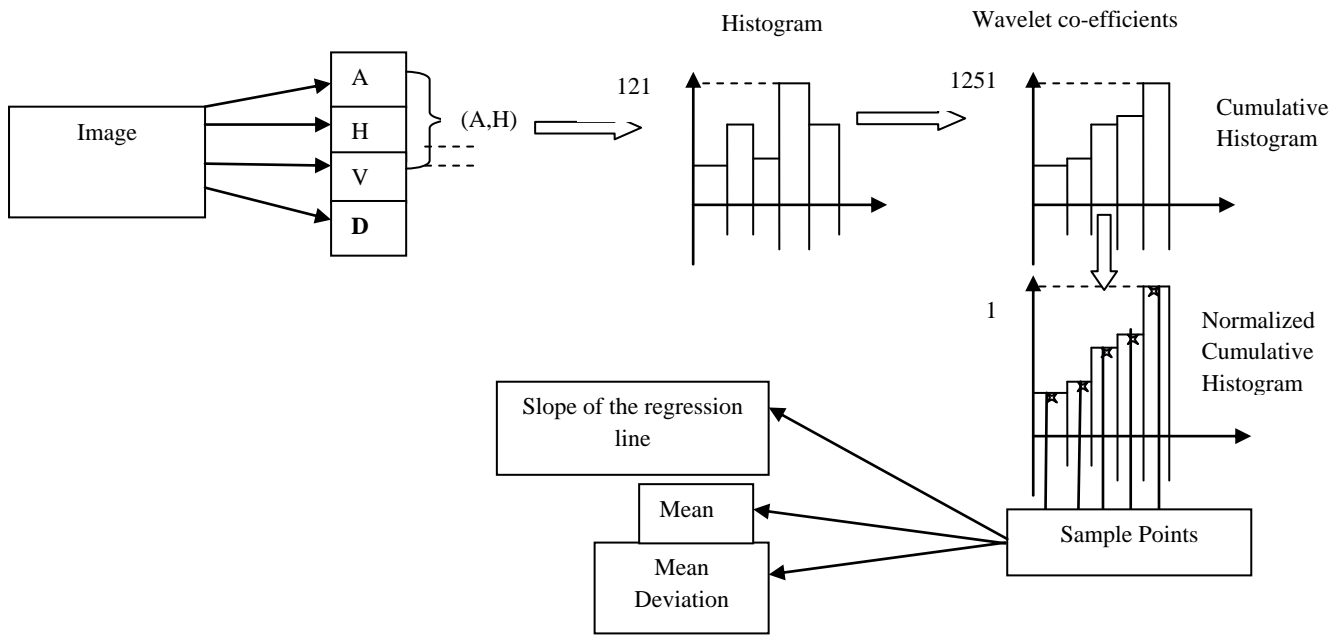


Fig 1: Schematic diagram of the feature extraction algorithm[1]

Handwritten image is considered as input image X and using Haar wavelet we apply 2D discrete wavelet transform (DWT), which provides us approximation sub-band image (I), and also detail sub-band images (H,V,D) (Fig 1). We take set of two images like (I, V), calculate the co occurrence histograms H1, H2 in provided direction. For every histogram, normalized cumulation is composed for all co- occurrence histogram and then enumerate the texture features, especially, mean, regression-line slope and mean-deviation described as above. The overall procedure will be looped in eight directions which yields two-histograms * three-features * eight-directions = 48 components features from each set (I, V). Correspondingly, for all sets such as, (I,H), (I,D), (I,abs(V-H-D)) by elicitation of features, we obtain 192 components to given handwritten image (I).

$$[\bar{I}] = [255 - I]$$

The detailed process will be repeated for the image which is complement image of I, represented as where i is gray-value pixel for image I. From I and features extracted for combining and obtaining a feature space containing dimension of 384. These are used for training of features and then classification. Fig 3.3 illustrates the feature extraction schema in detail. The schematic diagram of the feature extraction method is shown in the Fig. 1.

3. TEXTURE TRAINING AND CLASSIFICATION

3.1 Training

Training step includes extraction of features from different handwritten image samples, which are chosen randomly that belongs to every script using above feature extraction methodology. These extracted features will be saved in the feature library. Then these will be utilized for writer identification.

3.2 Classification

Classification step involves comparison of features values with that saved in feature library. The extraction of features for image I is done using feature extraction method, as explained above.[6] Then using distance vector code these will be compared with subsequent component values which will be saved in library,

$$D(N) = \sqrt{\sum_{j=1}^M [f_j(X) - f_j(M)]^2} \quad (5)$$

Herein, N represents overall components in component space f, f_j (Y) expresses jth component of texture of given example X, and f_j (N) reproduces lth texture value for Mth class in library. After this process, the written record is determined by minimum distance classifier.[7]

4. EXPERIMENTAL RESULTS AND DISCUSSION

Writer identification have effective approach in image document analysis, this observation is considered for grouping of composition. This approach is discussed in the project; efficiency of feature elicitation method for texture is articulated. We already discussed about different writers, who have different sense of writing in different state of mind, which also includes different styles. Hence we consider, a text block as distinct pattern. This study helps us to motivate the use of texture classification for identification of writer. The approach will not involve connected component method. Since, this is a global approach used for texture classification. We perform experiments by allowing the authors to write different scripts in Kannada language. These scripts have 512*512 dimensions. Execution of the method is done on 14 different writers. From each writer we have taken 50 images, resulting a total of 700 images. from which, we considered 50% of images for training, other 50% images for testing.

Average classification (%)
87.23%

Table 1: Average classification accuracies

5. CONCLUSION

In this paper, we propose a method based on the texture features for writer identification in a kannada handwritten document image. The co-occurrence histogram based texture features are extracted using the correlation between subbands at the same resolution of wavelet decomposed image, indicating that this information is significant in characterizing the writers based on the handwritten document image. The average classification accuracy is 87.23% for a single writer document with full text coverage for kannada language which decreases slightly with the increase in angle of orientation and decrease significantly with increase in number of authors. The experimental results demonstrate the efficacy of the proposed method and the potential of such a global approach for the writer identification in the document image analysis, which has significance in biometrics and forensic science.

6. REFERENCES

- [1] P.S Hiremath, Shivshankar S, Jagadeesh D Pujari, Mouneswara, "Script Identification in a handwritten document image using texture features", In Proc.IEEE 2nd International Advance Computing Conference, Patiala ,pp.110-114,2010
- [2] B.V.Dhandra, Vijayalaxmi M B 2014 , "Text and script independent and writer identification", International Conference on Contemporary Computing and Informatics, 2014.
- [3] Marius Bulacu "Text Independent Writer identification and verification using textural and allographic features", IEEE Transactions on Pattern analysis and machine intelligence, Vol 29 No.4 April 2007.
- [4] Marus Bulacu, Lambert Schomaker, Axel Brink, "Text-Independent Writer Identification and Verification on Offline Arabic hand writing", 2007.
- [5] C Naveena, V.N.Manjunatha, Aradhya, "Handwritten character segmentation for kannada scripts", 978-1-4673-4805-8/12/\$31_C 2012 IEEE
- [6] J Pradeep, E.Srinivasan, S.Himavathi, "Diagonal based feature extraction for handwritten character recognition system using neural network", 978-1-4244-8679-3/11/\$26 2011 IEEE.
- [7] Deepa Bendigeri, Rashmi Mundas, Dr.Jagadeesh Pujari, "Texture Features from Handwritten images for writer identification" International Journal On Recent and Innovation Trends in computing and communication Vol.4, Issue 7, ISSN 2321-8169, 76-79