

Review of Devnagari Handwritten Word Recognition

Mrs. Saniya M. Ansari
Research Scholar
Karpagam University, Coimbatore
Coimbatore

Dr. UdaysinghSutar
Professor and Head,
Department of Electronics Engg,
AISSMS COE,Pune, Maharashtra

ABSTRACT

Devnagari is the most popular script in India it is used by over 400 million people all over world. Recognition of Devnagari handwritten word has been a popular research area for many years because of its various applications. This paper describes different techniques for pre-processing, segmentation, feature extraction and classification which play an important role for recognition of word.

Keywords: Preprocessing, Feature extraction, Segmentation

1. Introduction

India is multilingual/multiscript country with various languages namely Gujarati, Marathi, Konkani, Bengali, Tamil, Telugu, Punjabi, Sanskrit, Urdu. Handwritten recognition is classified into two types as offline and online. In offline recognition document is scanned and complete writing is available as an image. Due to the availability of several computing devices such as Tablet PC, PDA and Smartphone in the market and affordable by common Indian, online handwritten word recognition gain enough attention. In online recognition input is given by Tablet PC, PDA and Smartphone which is equipped with pen based input technology. Input data to such as online handwriting recognition consist of (x, y) coordinates along with trajectory of the pen together with a few other possible information such as pen-up, pen-down etc.

1.1. Features of Devnagari Script

Devnagari is used in many languages like Marathi, Hindi, Konkani and Sanskrit which is used by approximately 400 million people in northern India and it is most widely used Indic script. Devnagari is written from left to right and it does not contain any lower and upper case letters. It consists of 11 vowels and 33 consonants.

अ आ इ ई उ ऊ ऋ ए ऐ ओ औ

Figure 1. Set of vowels

क ख ग घ ङ
च छ ज झ ञ
ट ठ ड ढ ण
त थ द ध न
प फ ब भ म
य र ल व श
ष स ह

Figure 2. Set of Consonants

Shirorekha or headline is the horizontal line at the upper part of the character or word. It does not contain any useful information so it should be detected and discarded. N. Joshi [14] describes a Shirorekha detection algorithm in the context of online Devnagari character recognition.

2. RECOGNITION OF DEVNAGARI HANDWRITTEN WORD

The schematic block diagram consists of various stages in Devnagari handwritten word recognition as shown in figure 3.

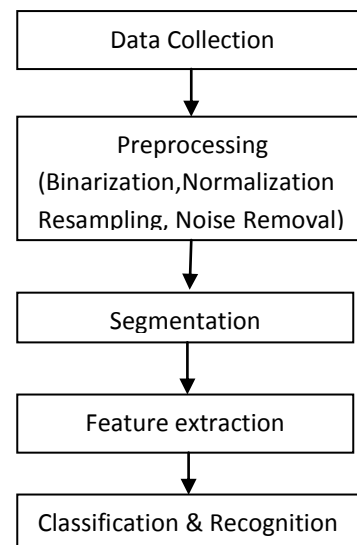


Figure 3. Stages in Handwritten word recognition

2.1. Pre-processing

For online recognition variations of handwriting occur due to various writers. For this preprocessing is required which involves elimination of noise, binarization of images, Size normalization, skew correction, thresholding and skeletonization techniques [5] [6]. While inputting data through a pen on the digitizer tablets, there may be certain noise and distortions present in the input text due to some limitations, which may make the recognition of input difficult. Irregular size, missing points due to fast movement of the pen, uneven distances of points from neighboring positions are various forms of noise and distortions. These noise and distortions present in the input text are removed in the second phase of online handwriting recognition i.e. pre-processing phase. The pre-processing phase includes five common steps [5] namely:

- Size normalization and centering
- Interpolating missing points
- Smoothing
- Slant correction
- Resampling of points

2.1.1. Binarization.

It is a method of transforming a gray scale image into a black and white image.

2.1.2. Size Normalization.

It is required so each segmented character is normalized to fit within a suitable matrix so that all characters have same data size.

2.1.3. Thresholding.

Thresholding is the process of reducing a gray scale image or colour image to a binary image.

2.1.4. Noise Removal.

It is necessary to eliminate imperfection like disconnected lines, gap of lines, etc. Median Filtering, Wiener filtering method and morphological operations can be performed to remove noise. Sobel technique is used to detect edges in binarized image [10].

2.2. Segmentation

Segmentation is the phase in which data is represented as character or stroke level so that nature of each character or stroke can be studied individually. The preprocessing stage yields a "clean" document in the sense that a sufficient amount of

shape information, high compression, and low noise on a normalized image is obtained.

The next stage is segmenting the document into its subcomponents. Segmentation is an important stage because the extent one can reach in separation of words, lines, or characters directly affects the recognition rate of the script. There are two types of segmentation: external segmentation, which is the isolation of various writing units, such as paragraphs, sentences, or words, and internal segmentation, which is the isolation of letters, especially in cursive written words.

In segmentation, pre-processed image is segmented into lines, words and characters. Segmentation process involves three steps namely line segmentation, word segmentation and character segmentation. Marathi word can be split into character by removing Shirorekha and then recognize [4].

2.3. FEATURE EXTRACTION

Feature extraction is a very crucial step as the success of a recognition system is often attributed to a good feature extraction method. The feature extractor determines which properties of the preprocessed data are most meaningful and should be used in further stages. For online recognition Vertical position of a point, curvature, PenUp/PenDown, writing direction, aspect, and slope are amongst the various features extracted in [2]. The most important aspect of handwriting recognition scheme is the selection of good feature set, which is reasonably invariant with respect to shape variations caused by various writing styles [4]. In [9] zoning method is used in which Diagonal feature extraction scheme is used to extract features from each zone. Features can be extracted using Freeman's direction code in which the change in directions while moving from one point to the next one of the sample are computed and quantized into one of the 8 possible values, via 1, 2... 8. Histograms of these direction codes are computed. In [15], global features are extracted from word which is useful to recognize the word.

The objective of feature extraction is to capture the essential information from data. This is an important

stage as its effective functioning improves the recognition rate and reduces the misclassification. In [10] Diagonal feature extraction scheme is used for recognizing offline handwritten character. For online recognition NPen++ recognition system [1] is used for feature extraction. Some feature extraction methods are Moments, Zoning and Projection Histogram.

2.4. CLASSIFICATION AND RECOGNITION

Recognition is the most important phase of the online recognition system and uses the features extracted in the previous stage to identify the input character according to preset rules. In many of the emerging new applications, it has become clear that no one method for recognition can be considered as optimal and usually a combination of multiple methods yield better recognition results. Combination of two or more recognition methods is a common practice now-a-days and is referred to as hybrid methods for recognition. The four best known approaches for pattern recognition are: 1) template matching, 2) syntactic or structural matching, 3) neural networks and 4) statistical classification.

The decision making stage of a recognition stage is classification stage and it uses the features extracted from previous stage. A number of classification methods were proposed by different researchers some of these are template matching, SVM classifiers and artificial neural network.

2.4.1. Template Matching.

This is the simplest approach of pattern recognition. Given pattern that is to be recognized is compared with stored patterns.

2.4.2. SVM Classifiers.

Support vector machines (SVM), when applied to text classification provide excellent precision, but poor recall. SVM have achieved excellent recognition results in various pattern recognition applications. Different types of kernel functions of SVM are: Linear kernel, Polynomial kernel, Gaussian Radial Basis Function and Sigmoid.

2.4.3. Artificial Neural Network.

Neural network is a computing architecture that consists of a massively parallel interconnection of adaptive "neural" processors. Because of its parallel nature, it can perform computations at a higher rate compared to the classical techniques. Neural network architectures can be classified as, feed forward and feedback (recurrent) networks. The most common neural networks used in the OCR systems are the multilayer perceptron (MLP) of the feed forward networks and the Kohonen's Self-Organizing Map (SOM) of the feedback networks.

3. COMPREHENSIVE STUDY OF PREVIOUS RECOGNITION TECHNIQUES:

Shailendra Kumar (2013) [1] proposed method Real Time Recognition of Handwritten Devnagari

Signatures without Segmentation Using Artificial Neural Network in this data collection is done by digital tablet. Various features of signature such as height, length, slant, Hu's moments etc are extracted and used for training of the NN. The highest accuracy rate achieved by for Devnagari handwritten signature recognition system was 96.12 %.

Swapnil A. Vaidya et.al (2013) [2] proposed method in which add all the sample character image matrices and divide the resultant matrix by total number of matrices added, called as average matrix. Then subtract it from each sample character image matrix, which results in unique features because of their positional properties of pixels present in that image. They used singular value decomposition technique to get projection vector matrix then used generalized regression neural network for resulting feature vectors and obtain classification performance in the character recognition task. The proposed recognition scheme provided 82.89 percent and 85.62 percent accuracies on Devnagari and Kannada character databases respectively.

Bharat et.al (2012) [3] proposed method based on HMM for lexicon driven and lexicon free word recognition for online handwritten for feature extraction they used NPen++ feature [29] for curliness, linearity and slope. The two different techniques for recognition of word written in Devanagari Text based on Hidden Markov Models (HMM): lexicon driven and lexicon free. The lexicon-driven technique models each word in the lexicon as a sequence of symbol HMMs according to a standard symbol writing order derived from the phonetic representation. The lexicon-free technique uses a novel Bag-of-Symbols representation of the handwritten word that is independent of symbol order and allows rapid pruning of the lexicon. In combination with the lexicon-driven and lexicon free approach, high recognition accuracy 93.38 percent can be achieved.

Naveen S. et.al (2012) [4] proposed a method of detecting Devanagari text of a printed document by mapping word directly to Unicode sequence. Printed text is different than handwritten text because it does not vary with users. Here they consider Unicode as the main recognition unit and use a sequence transcription module to map the words features to corresponding Unicode. A variant of RNN known as BLSTM (Bidirectional long Short term memory) was used for the task. The total time taken to test a word image is approximately 0.25 seconds.

Ved Agnihotri (2012) [5] proposed a new technique of Chromosomes function generation and fitness function for classification by extracting diagonal features from zones of an image. Handwritten Devanagari script recognition system using neural network is presented in this paper. For diagonal

extraction image is divided into some zones. Diagonal based feature extraction is used for extracting features of the handwritten Devanagari script. After that these feature of each character image is converted into chromosome bit string of length 378. In the recognition phase and classification Genetic Algorithm is used. The precision of offline Devanagari system is 85.78% match, 13.35% mismatch.

Gunjan Singh et.al (2012) [6] proposed an offline handwritten Hindi character recognition system using neural network. Neural networks are good at recognizing handwritten characters as these networks are insensitive to the missing data. In this paper proposed approach is to recognize Hindi characters this is done in four stages—1) Scanning, 2) Preprocessing, 3) Feature Extraction and, 4) Recognition. The feature vector comprises of pixels values of normalized character image. A Back propagation neural network is used for classification. Experimental results show that back-propagation network yields recognition accuracy of 93%.

Mitrakshi B. Patil et.al (2012) [7] proposed method for recognition of offline handwritten Devanagari characters using segmentation and Artificial neural networks. In this input image is scan image so it can easily preprocess. The whole process of recognition includes two phases- segmentation of Characters into line, word and characters and then recognition through feed-forward neural network.

Jayadevan R. et.al (2011)[8] did a survey of the comparative study of recognition of printed as well as handwritten word recognition by different classification techniques like Artificial Neural Network, Hidden Markov Model, Support Vector Machine, MQDF. In this survey paper he compares all classifiers with their accuracy.

A handwritten character recognition system using multilayer feed forward neural network is proposed in (2010) [9]. Three different orientations, namely, horizontal, vertical and diagonal directions are used for extracting 54 features from each character. Trained neural network is used for recognition of character.

In [10] Shanti N.(2010) proposed system for Handwritten Tamil character recognition which is based on support vector Machine. In this system global features are extracted from the character on which SVM is applied.

In [11] Suresh Kumar (2010) proposed a system in which the scanned image is segmented into paragraphs using spatial space detection technique, paragraph into lines using vertical histogram. The extracted features are given to SVM, self-organizing Map, RCS, Fuzzy Neural network. Structure analysis suggested that the proposed system of RCS with back

propagation network is given higher recognition rate. With the combination of RCS and back propagation network, a high accuracy recognition system is realized. The training set consists of the writing samples of 25 users selected at random from the 40, and the test set, of the remaining 15 users. A portion of the training data was also used to test the system. In the training set, a recognition rate of 100% was achieved and in the test set the recognized speed for each character is 0.1sec and accuracy is 97%.

In [12] SandhyaArora (2010) characters are collected in a systematic manner from printed pages scanned on a HP7670 Scanjet scanner. Documents were skew corrected and components were extracted. After dimensionality reduction (PCA), Telugu needs around 150 features for representation, while Hindi needs only 50. SVM and ANN are used for classification. SVM-based classifiers are found to perform better than KNN. Experimental results show that recognition accuracy of 93.31%.

M. C. Padma (2009) [13] proposes a method which uses distinct features extracted from the top and bottom profiles of the printed text lines. Using learning algorithm propose system learnt training data set. From the experimentations on the test data set, the overall accuracy of the system has turned out to be 99.67%. This algorithm is also tested on another test data set constructed from the scanned document images. The overall accuracy of the system reduces to 98.5%.

PrachiMukherji et.al (2009) [14] proposed method in which thinned character is segmented into segments (strokes), using basic structural features like endpoint, cross point, junction points and adaptive thinning algorithm. The segments of characters are coded using our Average Compressed Direction Code (ACDC) algorithm. The knowledge of script grammar is applied to identify the character using shapes of strokes, mean row and column co-ordinates, relative strength, straightness and circularity. Their location in the image frame is based on fuzzy classification. Characters are pre-classified using a tree classifier. Subsequently unordered stroke classification based on mean stroke features is used for final classification and recognition of characters. Recognition accuracy for characters with top modifier (matra) is 71.68 % and without top modifier (matra) is 88.33 %. Somaya Alma (2006) [15] presented system for offline Arabic handwritten word recognition using neural network. In this system preprocessing is done by edge detection and thinning of slant and slope is done. After preprocessing feature extraction is done by extracting global features from whole word using that neural network recognizes word. The achieved accuracy is 63%.

5. COMPREHENSIVE STUDY

Below table shows the comprehensive study of different techniques used for handwritten character, word and script recognition.

Reference Paper	Preprocessing	Segmentation	Feature Extraction	Classification /Recognition
[1]	Preprocessing is done to normalize the position and size of the sample.	–	NPen++ features are used for curliness, linearity and slope.	Hidden Markov Model based lexicon free technique used.
[2]	Image Binarization Thinning of binarized image windowing	Character recognition by neural network	Replacing the recognized characters by Standard fonts.	Assembling all the Separated characters in the same order as they appeared in the input image to give final output.
[3]	Thresholding method used for Binarization	Lines are segmented by noting the valleys of projection profile	Vertical Feature Bar, Horizontal Zero, Crossing, Moments	Tree Classifiers
[4]	Morphological operation are used to noise removal Thinning algorithm is used to remove the distortions Bicubic interpolation are used for standard sized image	Differential distance based technique used for identifying the Shirorekha and spine	Top, bottom, left, right or on a Combination technique. A single or double vertical line called a Danda (Spine) was traditionally used to indicate the end of phrase or sentence	Preliminary classification is performed for better results.
[5]	Gaussian filter used to make input data strokes smoother And reduce noise.	–	Sequential floating search method used for Indic script	K-nearest neighbor and Support Vector Machine (SVM) used for recognition.
[6]	Edge Detection is done and thinning for slant and slope	–	features are extracted from whole word.	Artificial neural network
[7]	Noise removal	–	Five different features from a vertical strip width, using a sliding window.	Neural network classifier known as Bidirectional Long Short Term Memory
[8]	Smoothing, Resampling and computation the length of input stroke if it less than	Cursive stroke segmentation for line and word segmentation.	Histogram of the direction codes calculated for each sub-stroke. Obtain co-ordinates of	Modified Quadratic discriminate function (MQDF) classifier is used. It improves
	set a priori ignore it for next phases		centre of gravity and normalize	efficiency over QDF.

	this approach is for noise removal.		these value by width and height of Stroke	
[9]	Preprocessing is done to Normalize the position and size of the sample and to Remove local noise so that the extracted features from the sample become robust.	Horizontal projection file method is used for segmentation	Images scaled into height and width using bilinear interpolation technique	feed forward algorithm
[10]	Detection of edges in the binarized image using sobel technique,	Preprocessed input image is segmented into isolated characters by assigning a number to each character using a labelling process.	Diagonal feature extraction scheme is used to extract features from each zone.	A feed forward back propagation neural network used for classification
[11]	Gabor Thresholding and Otsu Thresholding methods(global) are used for Binarization	Horizontal and vertical profile method is used for segmentation	Zone based approach is used for Feature Extraction.	Support vector machine (SVM) method is used for classification.
[12]	Detection of edges in binarized image is done by canny technique.	Preprocessed input image is segmented into isolated characters by assigning a number to each character using a labelling process.	Diagonal feature extraction scheme is used to extract features from each zone.	Chromosome function generation and Chromosome fitness function are used for classification.
[13]	Thresholding method used For Binarization. Thinning algorithm used to thin the characters	Histogram method used to convert the image to glyph	Character height, width, no. of horizontal and vertical lines.	Support Vector Machine(SVM) used for classification
[14]	Threshold technique used for preprocessing.	–	Encoding binary variation method used for extract the features. Then text and tested recognize the characters.	Support Vector Machine(SVM) used for classification
[15]	Global thresholding approach was used to binarized the scanned gray scale image	–	Top and bottom profile based features are used for feature extraction.	Learning Algorithm is used for classification.

4. CONCLUSION

In this paper we have represented a survey of preprocessing, segmentation, feature extraction, classification and recognition techniques for handwritten Devnagari word recognition. This survey research paper helps researches and developers to understand various techniques which were implemented for recognition .

There is lot of issues in recognition such as variation of writer's style. It becomes interesting and challenging field for researcher in image processing and pattern recognition.

REFERENCES

- [1] ShailendraKumar, "Real Time Recognition of Handwritten Devnagari Signatures without Segmentation Using Artificial Neural Network", *IJ. Image, Graphics and Signal Processing*, 2013, 4, 30-37.
- [2] SwapnilVaidya and Balaji R. Bombade, "A Novel Approach of Handwritten Character Recognition using Positional Feature Extraction", *IJCSMC*, Vol. 2, Issue. 6, June 2013, pg.179 – 186.
- [3] Bharat, SriganeshMadhvanath, "HMM-Based Lexicon-Driven and Lexicon-Free Word Recognition for Online Handwritten Indic Scripts", *IEEE*, April 2012.
- [4] Naveen Sankaram, C. V. Jawahar, "Recognition of Printed Devnagari Text Using BLSTM Neural Netwrk", *ICPR*, Nov. 11-15, 2012.
- [5] VedAgnihotri, "Offline Handwritten Devnagari Script Recognition", *IJCSI International Journal of Computer Science Issues*, 2012.
- [6] Gunjan Singh, SushmaLehri, "Recognition of Handwritten Hindi Characters using Backpropagation Neural Network, (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 3 (4), 2012, 4892-4895.
- [7] Mitakshi B. Patil, VaibhavNarawade, "Recognition of Handwritten Devnagari Characters through Segmentation and Artificial neural networks, International Journal of Engineering Research & Technology (IJERT) ISSN: 2278-0181 Vol. 1 Issue 6, August – 2012.
- [8] R. Jayadevan, Satish R. Kolhe, Pradeep M. Patil and Umapada Pal, "Offline Recognition of Devanagari Script: A Survey in IEEE Transaction on Systems, Man and Cybernetics –Part C. Applications and Review, vol.41, pp-782-796, 2011.
- [9] J. Pradeep, E. Srinivasan, S. Himavathi, "Diagonal Feature Extraction Based Handwritten Character System Using Neural Network", *International Journal of Computer Applications*, Vol. 8– No.9, October 2010.
- [10] Shanthi N and Duraiswami K, "A Novel SVM -based Handwritten Tamil character recognition system", *Springer Pattern Analysis & Applications*, Vol-13, No. 2, 173-180, 2010.
- [11] Suresh Kumar C and Ravichandran T, "Handwritten Tamil Character Recognition using RCS algorithms", *Int. J. of Computer Applications*, (0975 – 8887) Volume 8– No.8, October 2010.
- [12] SandhyaArora et al., "Performance Comparison of SVM and ANN for Handwritten Devnagari Character Recognition", *IJCSI International Journal of Computer Science Issues*, Vol. 7, Issue 3, May 2010.
- [13] M.C. Padma, P.A. Vijaya, "Identification of Telagu, Devnagari and English Scripts using Discriminating features", *IJCSIT*, Vol.1, No 2, November 2009.
- [14] PrachiMukherji, Priti P. Rege, "Shape Feature and Fuzzy Logic Based Offline Devnagari Handwritten Optical Character Recognition, Journal of Pattern Recognition Research (2009), pp 52-68.
- [15] Somaya Alma, "Recognition of Off-Line Handwritten Arabic Words Using Neural Networks", *IEEE* 2006.