# Offline Handwritten Sanskrit Character Recognition Using Hough Transform and Euclidean Distance

*Sujata S. Magare[1] and Ratnadeep R. Deshmukh[2]*

[1]Department of CS-IT,
Dr. Babasaheb Ambedkar Marathwada University,
Aurangabad, Maharashtra, India

[2]Department of CS-IT,
Dr. Babasaheb Ambedkar Marathwada University,
Aurangabad, Maharashtra, India

**ABSTRACT:** In this paper, we describe the methodology derived for offline handwritten Sanskrit character recognition. This paper will provide a way for researcher to develop a dataset and techniques for offline handwritten Sanskrit character recognition. This paper describes basics of dataset; challenges associated with character system and proposed techniques to recognize Sanskrit Compound Characters.

**KEYWORDS:** Offline Character Recognition, Hough Transform, Prewitt Edge Detection, Enclosed Region, Normalization, Euclidean Distance Classifier.

## 1    INTRODUCTION

Recognizing handwritten compound character is the most challenging task in recognition system. Many works has been done in character recognition, but still it needs much more improvements. We have developed a system that will recognize compound characters of Sanskrit language.

Optical character recognition is the process of recognizing optically scanned characters. Character recognition has two types: Offline and Online. One of the challenging problems in pattern recognition is Offline character recognition. Offline character recognition takes scanned image of required document paper. Offline character recognition can be done in two ways: Handwritten and Printed.

Handwritten character recognition is abbreviated as HCR; handwritten characters have number of variations as different people have different writing styles. HCR can recognize offline character and online characters. Offline HCR takes input from scanned image of paper document and Online HCR takes input from digital pen. There are many handwritten historical documents exist in electronic form, HCR is used to recognize such documents.

### 1.1    SANSKRIT

Sanskrit is an historical language, many works has been written in Sanskrit Language and are need to be recognized. In these documents recognizing compound characters are more challenging. Sanskrit language is written in Devanagari script. Each character in Sanskrit language consists of a Horizontal line at the top of character. This horizontal line is called as Shirorekha. Compound character is a combination of two characters. These compound characters are difficult to recognize due to variation in their shapes. Therefore it is required to develop such system that would recognize these compound characters.

---

## 2   LITERATURE REVIEW

In paper [1], they have proposed to recognize handwritten Sanskrit word using a Prewitt's operator for the edge detection. In this they have use Freeman chain code (FCC) as the representation technique of an image character. They have chosen Support vector machine (SVM) for the classification step.

Chain coding used to extract chain code features at the feature extraction stage and use Combined MLP and Minimum Edit Distance Classifier for classification [2].

In paper [3] Median and Wiener filters are used for denoising. They have used Structural segmentation algorithm for segmentation purpose and for feature extraction they have used Zone based approach.

A work on multi-stage character recognition system for an Indian script is reported by [4]. They have used MLPs for classification and have fuzzy features as inputs. The MLP outputs represent the belongingness of an input pattern to different fuzzy character pattern classes.

In paper [5], they have presented a combined DCT-DWT approach for tri-script identification at block level for the handwritten documents. KNN classifier is used in recognition phase that yielded better results for k=1.

In Paper [6] various Feature Extraction Method has been described, such as Template matching, Deformable templates, Graph description, Discrete features, Zoning and Fourier descriptor. They found that Real- Valued feature vectors are ideal for statistical Classifier.

## 3   HCR SYSTEM

Handwritten character recognition system consists of Preprocessing, Segmentation, feature extraction, classification and recognition.

### 3.1   PREPROCESSING

Preprocessing technique is used to do improvement of image data that enhances some image features required for processing and suppresses unwanted noise and distortion from image data and aims to correct degradation in an image.

- Binarization

Binarization is the process of converting grayscale image in to binary (Black and White) image, so that image data will only contain 0 and 1. Binarization technique is usually used for separating foreground from background using required level of thresholding.

- Noise Removal

Digital image consist of variety of noises. These noises are required to be removed from an image for better processing. Morphological operation, Median filter and Weiner filter is used to remove noise from an image. Median filter reduces blurring of edges. We have used Morphological Opening Operation. Morphological opening operation is Erosion followed by Dilation, using the same structuring element for both the operations. This operation removes small objects from an image while preserving the shape and size of larger objects in the image.

- Thinning and Filling

Smoothing implies both Filling and Thinning. Thinning reduces width of character while Filling eliminates gap, small breaks and holes in digitized character.

- Normalization

To obtain characters of uniform size, rotation and slant Normalization is applied on image. To improve the accuracy of character recognition Normalization reduces shape variation.

- Skew Detection and Correction

During the digitization of document page it is often that image is not align correctly or it may be happen by human while writing document. To make in correctly align Skew detection and correction technique is used.

Skew detection technique can be classified in to groups: Analysis of Projection profile, Hough transform, clustering, connected component and correlation between line techniques.

### 3.2    SEGMENTATION

Segmentation of an image is the process of subdividing image into number of parts. Segmentation takes the form as Paragraph Segmentation, Line Segmentation, Word Segmentation and Character Segmentation.

Paragraph wise segmentation divides the document into paragraph. Line wise segmentation divides paragraph into line. Line wise segmentation can use a horizontal projection profile based techniques Word wise segmentation divides line into word. Finally, Character wise segmentation divides words into characters. Chain code histogram can be used for each segment. Horizontal projection file method is used for segmentation.

### 3.3    FEATURE EXTRACTION

Feature extraction technique is aims to extract the essential and important features and characteristic of the given image. In Pattern recognition this is one of the difficult stages to implement. Selection of right feature extraction technique leads to achieving high performance for recognition.

Feature extraction technique is divided into three groups: Distribution of points, Transformation & series expansion and Structural analysis. Structural analysis extracts the feature which represents geometric and topological structure of character. Structural analysis gives feature with high tolerance of noise and style variation. Commonly used features are intersection between lines and loops.

Structural classification is used to detect presence of vertical line, its position in the character and presence of holes in character image [7].

### 3.4    CLASSIFICATION

After selection of the features, next step is to classify them according to its properties. Training and testing is done at the classification phase. Number of classifier can be used to train the character, such as Neural Network, K-NN, Support Vector Machine and Euclidean Distance Classifier.

- Neural Network:

    Neural network is one of the well known classifier used for character recognition system. Neural network are advantageous of their adaptive nature. Feed forward NN and Back propagation NN is used for character recognition mostly.

- SVM:
    Support vector machine construct the hyper-planes in high or infinite dimensional space. SVM is based upon statistical learning theory. The SVM was defined for the two class problem and it looked optimal hyper-plane, which maximized the distance, margin, between the nearest examples of both classes.

## 4    PROPOSED SYSTEM

Proposed system is shown in figure-1. It takes input image from user then preprocessing operations is performed on that image such as Binarization, Morphological opening, Character thicken and pass it to Structural classification. Structural classification stage will detect presence and position of vertical line and enclosed region within character image. Character is then normalized and Euclidean distance classifier applied on character image. At final stage, proposed system will give result for character recognition.

## 4.1    DATASET DEVELOPMENT

Researcher has to develop own character dataset collected from minimum 10-15 people, because there is no standard dataset available for handwritten characters. Different person has different writing styles; it includes variations in dataset which will be useful while training and testing phase for character recognition.

Offline Handwritten character recognition takes scanned image of required document paper. For this purpose we have taken a blank paper. Small blocks created on that blank paper. We have also added Name, Age and Sub-code field to the document. Although these fields are optional, but it is good to keep information about user which will help us to recognize handwritten styles with different age groups. After this block creation document paper is ready to take data from user. We have taken dataset from 26 people. For the reorganization process this image is converted to grayscale and then Binarization is applied on grayscale image. So that image can contain information only in 0 or 1.



*Fig. 1.    Proposed System*



*Fig. 2.    Datasheet for Sanskrit Compound Characters*

*Fig. 3.    Sample of Cropped Compound Character*

## 4.2    STRUCTURAL CLASSIFICATION

Structural Classification of a character implies detection of Vertical line, its Position in an image and Enclosed Region detection. For vertical line detection we have used Hough Transform and Prewitt Edge Detection algorithm. Enclosed region is detected with help of boundary tracing algorithm.

### 4.2.1    HOUGH TRANSFORM

Hough transform is applied on the binarized edge map to generate the Hough image of it. The Hough Transform was originally developed for detecting lines and it is still popular in this particular area. Therefore, a large amount of different line parameterizations and various algorithmic modifications exist [8].

In automated analysis of digital images, a frequently arising problem is detecting the simple shapes like straight line, circle or ellipse. In most of the cases an edge detector can be used as a pre-processing stage to obtain image points or image pixels that are on the desired curve in the image space. But due to imperfections in either the image data or the edge detector there may be missing or isolated or disjoint points or pixels on the desired curves as well as there may be spatial deviations between the ideal lines or circle or ellipse and the noisy edge points as obtained from the edge detector. For these reasons, it is often non-trivial to group the extracted edge features to an appropriate set of lines, circles or ellipses. The purpose of the Hough transform is to address this type of problem by making it possible to perform groupings of edge points into object candidates by performing an explicit voting procedure over a set of parameterized image objects.

### 4.2.1.1    PREWITT EDGE DETECTION

The Prewitt operator is used in image processing, particularly within edge detection algorithms. Technically, it is a discrete differentiation operator, computing an approximation of the gradient of the image intensity function. the operator calculates the gradient of the image intensity at each point, giving the direction of the largest possible increase from light to dark and the rate of change in that direction. The result therefore shows how abruptly or smoothly the image changes at that point and therefore how likely it is that part of the image represents an edge, as well as how that edge is likely to be oriented.



*Fig. 4.    Vertical Line Detection*

- *Position of Vertical Line:*

Once the vertical edge has been detected, next task is to find its position in character image such as Middle Line and End Line. For this reason, character image is divided in two parts vertically. If edge is exist then using its beginning point line position is find out. If point position is less than image half then line is called as Middle Line. If point position is greater than image half then line is called as End Line.

- *Position of End Point:*

Detection of End point is done with the analysis of position of line end points. For this reason image is divided horizontally, then lower part is again divided vertically and checked for the end points. If end point position is less than lower half then End point is in Fourth Quadrant. If point position is greater than lower half then End Point position is in Third Quadrant.

#### 4.2.2    DETECTION OF ENCLOSED REGION

In Sanskrit language, there are many characters formed with the enclosed region. Hence, we have detected enclosed region within the character too. We have traced boundaries within binary image with the help of connected component.



*Fig. 5.    Enclosed Region Detection*

### 4.3    CHARACTER NORMALIZATION

Character Normalization is the process of making images to be uniform in their size. After Structural classification, character normalization is performed. We have normalized character in to 32x32 size. Distance measurement can be efficiently applied on 32x32 image and these can minimize computation time for distance calculation.

### 4.4    EUCLIDEAN DISTANCE CLASSIFIER

Euclidean Distance classifier is used to classify features. We have calculated Euclidean distance between the input image and training images. After calculation, we have considered those classes whose Euclidean distance is minimum. Accordingly, testing has been done and Output is displayed.

## 5    RESULTS AND DISCUSSION

We have collected a database of 34 compound characters from individuals of different groups.  Resulting into 60 samples for each character. There are 2040 total characters into the database. These samples are preprocessed and applied for segmentation. At preprocessing stage Morphological opening operation fill the gap between characters and removes unwanted strokes from character image. And then it is thinned. At the segmentation stage Vertical lines are detected using Hough transform and Prewitt Edge detection.

Position of a vertical bar and presence of enclosed region is detected at the Feature extraction. After extracting features, image is then normalized to 32x32 images. This image is then proceeding for a Classification. Minimum Distance classifier used to classify features. We have calculated Euclidean distance between input image and training images. Class, whose distance is Minimum, is qualified for the recognition. Based upon these testing has been made and Corresponding Recognition result is displayed.

Proposed system is designed using MATLAB R2013a. Proposed system provides overall 94% recognition rate.

*Table 1.  Recognition Rate of Individual Character*

| Characters | Correctly Classified | Misclassified | Percentage of Recognition |
|---|---|---|---|
| Bba | 58 | 2 | 96 |
| Bda | 58 | 2 | 96 |
| Bja | 54 | 6 | 90 |
| Chya | 56 | 4 | 93 |
| Dha | 54 | 6 | 90 |
| Dhya | 55 | 5 | 91 |
| Dya | 56 | 4 | 93 |
| Gwa | 57 | 3 | 95 |
| Jva | 56 | 4 | 93 |
| Jya | 57 | 3 | 95 |
| Kta | 58 | 2 | 96 |
| Kva | 60 | 0 | 100 |
| Kya | 60 | 0 | 100 |
| Mha | 55 | 5 | 91 |
| Nha | 55 | 5 | 91 |
| Nma | 56 | 4 | 93 |
| Nta | 56 | 4 | 93 |
| Nya | 54 | 6 | 90 |
| Pta | 57 | 3 | 95 |
| Pya | 58 | 2 | 96 |
| Rya | 58 | 2 | 96 |
| Shka | 59 | 1 | 98 |
| Shta | 58 | 2 | 96 |
| Ska | 58 | 2 | 96 |
| Sma | 56 | 4 | 93 |
| Spa | 55 | 5 | 91 |
| Sta | 54 | 6 | 90 |
| Swa | 56 | 4 | 93 |
| Sya | 56 | 4 | 93 |
| Tma | 57 | 3 | 95 |
| Tpa | 58 | 2 | 96 |
| Tta | 60 | 0 | 100 |
| Tva | 59 | 1 | 98 |
| Tya | 59 | 1 | 98 |

## REFERENCES

[1] Namita Dwivedi, Kamal Srivastava and Neelam Arya, "Sanskrit Word Recognition Using Prewitt's Operator And Support Vector Classification", *IEEE International Conference on Emerging Trends in Computing, Communication and Nanotechnology (ICECCN)*, pp. 265-269, 2013

[2] S. Arora, D.Bhattacharjee, M. Nasipuri, D. K. Basu & M. Kundu, "Recognition of Non-Compound Handwritten Devanagari Characters using a Combination of MLP and Minimum Edit Distance", *International Journal of Computer Science and Security (IJCSS)*, vol. 4, issue 1

[3] Veena Bansal and R.M.K. Sinha, "Segmentation of touching and fused Devanagari character", *Pattern Recognition*, vol. 35, issue 4, pp. 875-893, 2002

[4] Shamik Surul, "Recognition of an Indian Script using Multilayer Perceptron and Fuzzy Features", *Proc. of 6th International Conference on Document*, 2003

[5] U. Garain, B.B. Chaudhari, " Touching Characters in Printed Devanagari and Bangla Scripts Using Fuzzy Multifactorial Analysis", *IEEE Transaction* on vol. 32, issue 4, pp. 449-459, 2002

[6] Olivind Due Tier, Anil K Jain, Torfin Tax, "Feature Extraction Method For Character Recognition: A Survey", *Pattern Recognition*, vol. 29, no. 4, pp. 641-662, 1996

[7] Sushama Shelke and Shaila Apte, "A Multistage Handwritten Marathi Compound Character Recognition Scheme using Neural Networks and Wavelet Features", International Journal of Signal Processing, Image Processing and Pattern RecognitioInternational Recognition, vol. 4, no. 1, 2011.

[8] Princen, J., Illingowrth, J., Kittler, J., "Hypothesis testing: a framework for analyzing and optimizing Hough transform performance", *IEEE Transaction Pattern Analysis and Machine Intelligence*, p.p. 329–341, 1994.