

Offline Handwritten Gurmukhi Character and Numeral Recognition using Different Feature Sets and Classifiers - A Survey

Anoop Rekha*

*(Department of Computer Science, Punjabi University, Patiala, INDIA)

ABSTRACT

A great work has been done for printed Gurmukhi text but in case of handwritten Gurmukhi text very less work has been done. In recent years research towards Indian handwritten character is getting increasing attention. Many approaches have been proposed by the researchers towards handwritten Indian character recognition and many recognition systems for isolated handwritten characters and numerals are available in the literature. This paper presents an overview of various O.C.R systems for Gurmukhi script which are developed for handwritten isolated Gurmukhi characters and numerals.

Keywords -Gurmukhi script, handwritten isolated text, offline handwritten character and numeral recognition, OCR

I. INTRODUCTION

Today, many researches have been done to recognize Gurmukhi characters and numerals. But the problem of interchanging data between human beings and computing machines is a challenging one. Even today, many algorithms have been proposed by many researchers so-that these Gurumukhi characters and numerals can be easily recognize. But the efficiency of these algorithms is not satisfactory. Many researches have been done to solve handwritten character recognition problem in related areas such as Image Processing, Pattern Recognition, Artificial Intelligence, and cognitive science etc. Further researches are being done to improve accuracy and efficiency. Many techniques have been applied for recognition of handwritten Gurumukhi characters and numerals but still it is the case of less efficiency and accuracy of recognition. Optical Character Recognition has been one of the most challenging research area in the field of image processing in the recent years. Optical Character Recognition (OCR) is the process of converting scanned images of machine printed or handwritten text into a computer processable format. The process of optical character recognition has following stages:

1. Digitization: Digitization is the process whereby a document is scanned and an electronic representation of the original, in the form of a bitmap image, is produced. Digitization produces the digital image, which is fed to the pre-processing phase.

2. Preprocessing: Preprocessing is used for skew detection/correction, skeletonization, and noise reduction/removal. Skewness refers to the tilt in the bit mapped image of the scanned paper for OCR. Skeletonization is used for decreasing the line width of text from many pixels to single pixel. Noise removal is used to remove unwanted bit pattern which does not play any significant role in document.

3. Segmentation : Segmentation is used to break the script into lines, words and characters.

4. Feature Extraction: In the feature extraction phase, one can extract the features according to levels of text, e.g., character level, word level, line level and paragraph level.

5. Classification: The classification phase is the decision making phase of an OCR engine, which uses the features extracted in the previous stage for making the class memberships in pattern recognition system.

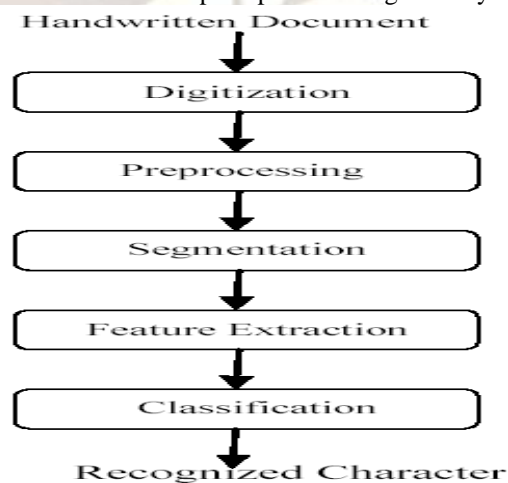


Fig. 1 Block diagram of OCR System

1.1 Introduction to Gurmukhi Script

Gurmukhi Script is used primarily for Punjabi language, which is the world's 14th most widely spoken language.

Following are the properties of Gurmukhi Script are:

- i. Writing style is from left to right.
- ii. No concept of upper and lower case characters.
- iii. Gurmukhi script is cursive.

Vowel Carriers:					
ੳ	ਅ	ੲ			
Consonants:					
ਸ	ਹ				
ਕ	ਖ	ਗ	ਘ	ਙ	
ਚ	ਛ	ਜ	ਝ	ਞ	
ਟ	ਠ	ਡ	ਢ	ਣ	
ਤ	ਥ	ਦ	ਧ	ਨ	
ਪ	ਫ	ਬ	ਭ	ਮ	
ਯ	ਰ	ਲ	ਵ	ੜ	
ਸ਼	ਖ਼	ਗ਼	ਜ਼	ਫ਼	ਲ਼

Table 1. Gurmukhi Characters

੦	੧	੨	੩	੪	੫	੬	੭	੮	੯
---	---	---	---	---	---	---	---	---	---

Table 2. Gurmukhi Numerals

1.2 Challenges of Handwritten Gurmukhi Script Recognition

Gurmukhi Script has following challenges[2]:

- i. Variability of writing style, both between different writers and between separate examples from the same writer overtime.
- ii. Similarity of some characters.
- iii. Low quality of text images
- iv. Unavoidable presence of background noise and various kinds of distortions.

II. FEATURE EXTRACTION METHODS

Feature extraction method analyzes a handwritten character image and selects a set of features that can be used for uniquely classifying the character.

2.1 Diagonal Feature Extraction:

Diagonal features are very important features in order to achieve higher recognition accuracy and reducing misclassification. These features are extracted from the pixels of each zone by moving along its diagonals.[8]

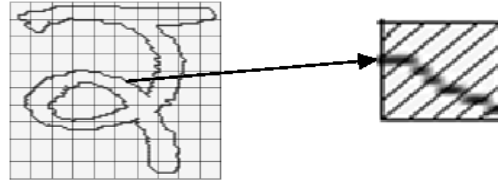


Fig. 2 Diagonal Feature Extraction

2.2 Transition Feature Extraction:

Transition feature extraction method was based on calculation and location of transition features from background to foreground pixels in the vertical and horizontal directions. To calculate transition information, image is scanned from left to right and top to bottom.[8].

2.3 Intersection and Open End Points Feature Extraction:

An intersection point is the pixel that has more than one pixel in its neighborhood and an open end point is the pixel that has only one pixel in its neighborhood.[6]

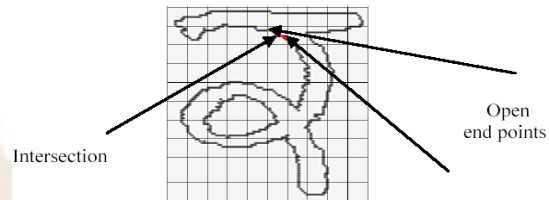


Fig. 3 Intersection and Open End points Feature Extraction

2.4 Zoning Density Features:

In zoning, the character image is divided into N×M zones. From each zone features are extracted to form the feature vector. The goal of zoning is to obtain the local characteristics instead of global characteristics. By dividing the number of foreground pixels in each zone by total number of pixels in each zone, we obtained the density of each zone.[1]

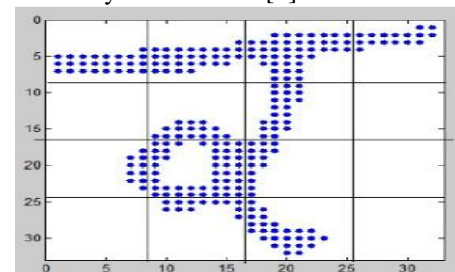


Fig.4 Zoning Density Features

2.5 Projection Histogram Features:

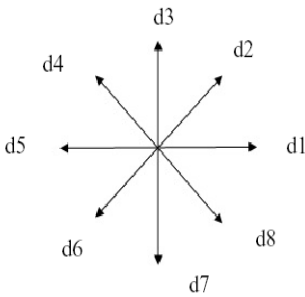
Projection histograms count the number of pixels in specified direction. There are three types of projection histograms i.e. horizontal , vertical , left diagonal and right diagonal.[1][7]

2.6 Distance Profile Features:

Profile counts the number of pixels (distance) from bounding box of character image to outer edge of character. In this approach, profiles of four sides left, right, top and bottom were used[1][7].

2.7 Background Directional Distribution (BDD) Features:

To calculate directional distribution values of background pixels for each foreground pixel, we have used the masks for each direction shown in figure . The pixel at centre ‘X’ is foreground pixel under consideration to calculate directional distribution values of background. The weight for each direction is computed by using specific mask in particular direction depicting cumulative fractions of background pixels in particular direction[1][7].



(a) 8 directions

0 0 1	0 1 2	1 2 1	2 1 0
0 X 2	0 X 1	0 X 0	1 X 0
0 0 1	0 0 0	0 0 0	0 0 0

d1 d2 d3 d4

1 0 0	0 0 0	0 0 0	0 0 0
2 X 0	1 X 0	0 X 0	0 X 1
1 0 0	2 1 0	1 2 1	0 1 2

d5 d6 d7 d8

(b) Masks used to compute different directional distributions

	d4	d3	d2
	d5	X	d1
	d6	d7	d8

(c) An example of 3*3 sample

Fig. 5 Computation of Background Directional Features

2.8 Combination of various features:

Each feature is used to form a feature vector hence if we use a combination of features then it will help us to derive the feature vectors with more elements which are helpful to increase the efficiency of recognition.

III. CLASSIFICATION METHODS

Classification method uses features extracted in the feature extraction stage to identify the unknown character.

3.1 K-Nearest Neighbour (KNN):

The k-nearest neighbor (k-nn) approach attempts to compute a classification function by examining the labeled training points as nodes or anchor points in the n-dimensional space, where n is feature size. We calculate the Euclidean distance between the test point and all the reference points in order to find K nearest neighbors, and then rank the obtained distances in ascending order and take the reference points corresponding to the k smallest Euclidean distances. A test sample is then attributed the same class label as the label of the majority of its K nearest (reference) neighbors. Euclidean distance is the straight line distance between two points in n-dimensional space[1][2][8].

3.2 Support Vector Machine (SVM):

The Support Vector Machine (SVM) is learning machine with very good generalization ability, which has been applied widely in pattern recognition, regression estimation, isolated handwritten character recognition, object recognition speaker identification, face detection in images and text categorization. SVM implements the Structural Risk Minimization Principal which seeks to minimize an upper bound of the generalization error.[] The standard SVM classifier takes the set of input data and predicts to classify them in one of the only two distinct classes. SVM classifier is

trained by a given set of training data and a model is prepared to classify test data based upon this model. For multiclass classification problem, we decompose multiclass problem into multiple binary class problems, and we design suitable combined multiple binary SVM classifiers. Different types of kernel functions of SVM: Linear kernel, Polynomial kernel, Gaussian Radial Basis Function (RBF) kernel and Sigmoid kernel[1][2][6].

3.3 Probabilistic Neural Network (PNN):

PNN is a multilayered feed-forward neural network classifier in which known probability density function (pdf) of the population is used to classify unknown patterns. PNN is closely related to Parzen window pdf estimator[1]. If the probability density function (pdf) of each of the populations is known, then an unknown, X, belongs to class “i” if:

$$f_i(X) > f_j(X), \quad \text{all } j \neq i, \quad f_k \text{ is the pdf for class } k.$$

3.4 Artificial Neural Network:

An Artificial Neural Network (ANN) is an information-processing paradigm that is inspired by the way biological nervous systems, such as the brain, process information. The key element of this paradigm is the novel structure of the information processing system. It is composed of a large number of highly interconnected processing elements (neurons) working in unison to solve specific problems. ANNs, like people, learn by example.[10] An ANN is configured for a specific application, such as pattern recognition or data classification through a learning process.

3.5 Neocognitron Neural Network:

The Neocognitron is a hierarchical multilayered neural network. Researchers found two types of cells in visual primary cortex called *simple cell* and *complex cell*, and also proposed a cascading model of these two types of cells. The local features are extracted by S-cells, and deformation of these features, such as local shifts, are tolerated by C-cells. Local features in the input are integrated gradually and classifying in the higher layers. The first version of the neocognitron was based on the learning without a teacher. This version is often called self-organized neocognitron. The main advantage of neocognitron is its ability to recognize correctly not only learned patterns but also patterns which are produced from them by using of partial shift, rotation or another type of distortion. The system is a neocognitron which recognizes handwritten characters of Gurumukhi script[5].

IV. COMPARISON OF RESULTS

Many researchers have proposed various techniques for offline handwritten Gurumukhi character. The results are given below:

S.No	Feature Extraction	Classifier	Recognition Accuracy
1.	Zoning	KNN	72.54%
2.	Zoning	SVM(Poly.Kernel)	73.02%
3.	Structural Features	Neural Network	83.32%
4.	Transition Features	KNN	86.57%
5.	Diagonal Features	SVM(Linear Kernel)	90.29%
6.	Profiles, width, height, aspect ratio, neocognitron	Neocognitron Neural Network	92.78%
7.	Diagonal Features	KNN	94.12%
8.	Intersection and Open End Point Features	SVM(Poly.Kernel)	94.29%
9.	Diagonal and Intersection and Open End Point Features	SVM(Poly.Kernel)	94.29%
10.	Zoning and BDD Features	SVM(RBF Kernel)	95.04%

Table 3 Results of Gurumukhi character recognition

Many researchers have proposed various techniques for offline handwritten Gurumukhi numerals. The results are given below:

S.No.	Feature Extraction	Classifier	Recognition Accuracy
1.	Distance Profiles	SVM(RBF Kernel)	98%
2.	Zoning and BDD Features	SVM(RBF Kernel)	99.13%
3.	Projection Histograms	SVM(RBF Kernel)	99.2%

Table 4 Results of Gurumukhi numeral recognition

V. CONCLUSION AND FUTURE SCOPE

SVM(RBF Kernel) is used as a classifier in offline handwritten Gurmukhi character and numeral recognition. Accuracy of 95.04% is achieved using zoning and BDD features which is the highest accuracy achieved during Gurmukhi character recognition. Accuracy of 99.2% is obtained using Projection Histograms, which is the highest accuracy achieved during Gurmukhi numeral recognition. The work presented in this paper can be further extended for on connected characters or words, by first segmenting the words and then recognizing the so obtained characters. In the present work only consonants, while lie in the middle zone and digits have been considered, in future vowels, lying in upper and lower zone and other half characters can also be used. There are many feature extraction methods and classifiers which are not implemented in case of handwritten Gurmukhi script recognition e.g Zernike moments, Fourier Descriptors, HMM etc. So a lot of work can be done in the field of handwritten Gurmukhi character and numeral recognition.

REFERENCES

- [1] Kartar Singh Siddharth, Mahesh Jangid, Renu Dhir, Rajneesh Rani, "Handwritten Gurmukhi Character Recognition Using Statistical and Background Directional Distribution Features", *International Journal on Computer Science and Engineering* (0975-3397), Vol. 3 No. 6 June 2011.
- [2] Puneet Jhajj, D. Sharma, "Recognition of Isolated Handwritten Characters in Gurmukhi Script", *International Journal of Computer Applications* (0975-8887), Vol. 4, No. 8, 2010.
- [3] O. D. Trier, A. K. Jain and T. Text, "Feature Extraction Methods For Character Recognition- A Survey", *Pattern Recognition*, Vol. 29, No. 4, pp. 641-662, 1996.
- [4] G.S. Lehal and Chandan Singh, "A Gurmukhi Script Recognition System", *Proceedings of the International Conference on Pattern Recognition (ICPR'00)*, 1051-4651/00, 2000.
- [5] Ubeeka Jain, D. Sharma, "Recognition of Isolated Handwritten Characters of Gurmukhi Script using Neocognitron", *International Journal of Computer Applications* (0975-8887), Vol. 4, No. 8, 2010.
- [6] Munish Kumar, M.K.Jindal and R.K.Sharma, "SVM Based Offline Handwritten Gurmukhi Character Recognition" *International Conference on Image Information Processing (ICIIP 2011)*
- [7] Kartar Singh Siddharth, Renu Dhir, Rajneesh Rani, "Handwritten Gurmukhi Numeral Recognition Using Different Feature Sets", *International Journal of Computer Applications* (0975-8887), Vol. 28, No. 2, August 2011.
- [8] Munish Kumar, M.K.Jindal and R.K.Sharma, "k-Nearest Neighbor Based Offline Handwritten Gurmukhi Character Recognition" *International Conference on Image Information Processing (ICIIP 2011)*
- [9] U.Pal, T.Wakabayashi, F.Kimura, "Comparative Study of Devnagri Handwritten Character Recognition using Different Feature and Classifiers", *International Conference on Document Analysis and Recognition (ICDAR 2009)*
- [10] Naveen Garg, Karun Verma, "Handwritten Gurmukhi Character Recognition Using Neural Network", M.Tech Thesis, Thapar University, 2009[online]. Available: <http://dspace.thapar.edu:8080/dspace/bitstream/10266/788/1/thesis+report+final.pdf>