# Recognition of Gurmukhi character style through Neural Networks

## Reena Joshi[1] and Gaurav Agnihotri[2]

[1]Reena Joshi, Rayat-Bahra Institute of Engg. and Neno Technology, Hoshiarpur

[2]Gaurav Agnihotri, Assistant Prof GNA University, Phagwara

## Abstract

This paper defines recognition of font of various styles for Punjabi font recognition with neural network. The font recognition is depending upon extracted features so these features are known as core part of font recognition. Its major plan is to the identification of the font of various character image based on local features that using post prior approach. The technique used for recognition is nourishing forward neural network classifier and operates on a given set of known fonts. The presentation of the proposed method is evaluated by a set of tests made on a database of characters combining different type's fonts of Punjabi characters that are mostly used in Gurmukhi. The font recognition accuracy is about 87%. Benefit of this approach is that noise does not disturb it.

*Keywords: Artificial Neural Network (ANN), Optical Character Recognition (OCR), Optical Font Recognition (OFR), Arabic Font Recognition (AFR), Gurmukhi Script.*

## 1. Introduction

It is mostly a process which take unprocessed data to generate an action based on the category of the pattern is called Pattern Recognition. Optical Character recognition is one of the main applications of pattern recognition. OCR indicates mechanical or electronic conversion of scanned or photo images of typewritten or printed text into machine-encoded/computer-readable text file. The most important advantage to use OFR are as Logical Document arrangement Recognition, Document replica, Improve detection rate of OCR, Helpful in Retrieval System too. Only a few researches have been done on the identifications of the type face named optical font recognition (OFR) compared to the vast researches in the OCR field .In [3], main objective is identification of features such as typeface, weight, slope and size of the text image from an image block and uses multivariate Bayesian. Its rate of recognition was 95%. The authors of [2] use texture-analysis-based approach with global features and use multichannel Gabor filter and WED classifier toward font recognition. In paper [5] presented an algorithm for priori Arabic optical Font Recognition (AFR). They obtain overall font recognition rate 79.40 which was very low. The result was improved in paper [6] by extract 48 features value is used to learn the decision tree. A set of 33 fonts was investigated. The overall success rate is 89.8%. In paper [8] they use global texture analysis and Gabor filter in machine-printed document images on Persian font recognition. Experiment was taken out with two classifier and average accuracy are 85% with WED and 79% with SVM .Same approach is used for English language [9], with SVM classifier shows an average accuracy 91.43% .To our knowledge, there has been no study of the Punjabi Font Recognition (PFR) problem. Available studies deal with various foreign languages as discuss in previous.

## 2. Introduction To Punjabi Language And Fonts

Punjabi language is the official language of Punjab state in India. In Pakistan, even though Punjabi has no official status, this language is most spoken language and is the local language of Punjab (Pakistan) the second largest and the most populous province of Pakistan. In Punjab gurmukhi script is used. It contains 41 consonants and 12 vowels and 3 half characters,

*International Journal of Scientific Engineering and Applied Science (IJSEAS) - Volume-1, Issue-9, December 2015*
*ISSN: 2395-3470*
*www.ijseas.com*

which lie at the feet of consonants. There is no concept of upper or lowercase characters. A line of Gurmukhi script can be partitioned into three categories such as horizontal zones namely, upper zone, middle zone and lower zone. There are large number of fonts exist in Punjabi language. Here use different fonts that use mostly used in many applications [13].

## 3. Proposed Method

This paper deals with identify given text input belong to which font. The model for font recognition is as shown in fig.1
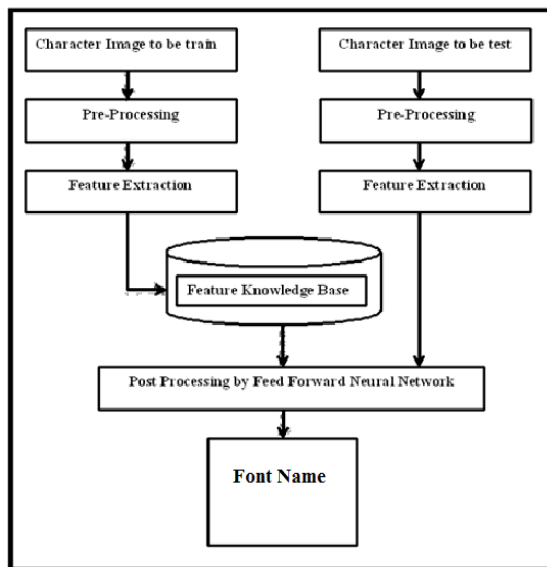


Fig 1: Recognition Model

2.1 Pre-processing

It is the First phase; its aim is to produce data that are easy for the systems to operate accurately. It takes a raw image and improves it by performed the various pre-processing steps. The pre-processing is as represent in fig. 2. The purpose of this phase is to make image ready for feature extraction. It includes the following steps:

1**.** *Image Acquisition*: It is a process to attain character image by image sensor such as image sensing device (Scanner, Digital Cameras).

2. *Cropping the image:* Cropping is done as to eliminate all the white pixels and have the character fully stretched inside the box.

3. *Convert the image into Gray Scale Image*: In this step convert the image into gray scale image. Gray scale image represent the gray level intensity of the character image i.e. glow at particular pixel.

4. *Convert the image into Binary image:* This step converts the Character image into binary image which represent the image in the form of zero's and ones.

5. *Resize the Image*: This step change the size of image into given dimensions.

2.2 Feature Extraction

It is second phase that deals with extraction of features from text image which is act as an input for next stage i.e. Post processing. It is the process of extracting the characteristics or attributes from a character image. The whole system performance is depends upon this stage. Generally, the feature extraction technique consists of two categories as shown in fig 2.
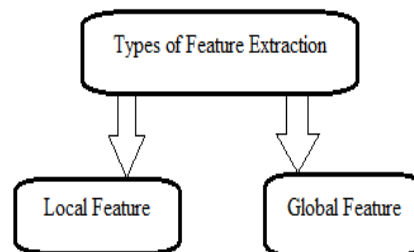


Fig 2: Types of features

a) *Local Features*: It is extraction from individual printing of image and deal with posteriori font recognition.
b) *Global Features*: These features are extracted from text image and can be termed as word, line, and paragraph.

Here deals with selected features according to which we can recognize from character images like area, length, width etc. These features are local features and the features that we use are such as given below:

*International Journal of Scientific Engineering and Applied Science (IJSEAS) - Volume-1, Issue-9, December 2015*
*ISSN: 2395-3470*
*www.ijseas.com*

1. *Area*: Actual number of pixels in the region.
2. *Length:* It specifies the length of the bounding box along each dimension.
3. *Breadth:* it specifies the width of the bounding box along each dimension.
4. *Corr2 (2-D correlation coefficient):* It computes the correlation coefficient between two images.

## 2.3 Post processing

The previous feature extraction phase should extract the features or attributes from the reference set and after than generate a template. In this stage the characters are classified by artificial neural networks to recognize**.** An Artificial Neural Network (ANN) is an information processing concept that is encouraged by the way natural nervous systems, such as the brain, process information. The ANN is collected as large number of highly interconnected processing elements (neurons) working in unity to solve specific problems. It is most of similar to people, such as they learn by example. Which include two steps: -

1. In the Training phase, System trained according to template character images which was created in feature Extraction phase and used for recognize when the system reads a new input image. Recognition is done according to the trained template. The training of neural network is as shown in figure 3.The features are extracted from template images and store in dataset. After this the neural network are trained according to extracted features which is stored in dataset. After than a threshold value 90% is taken for recognition.
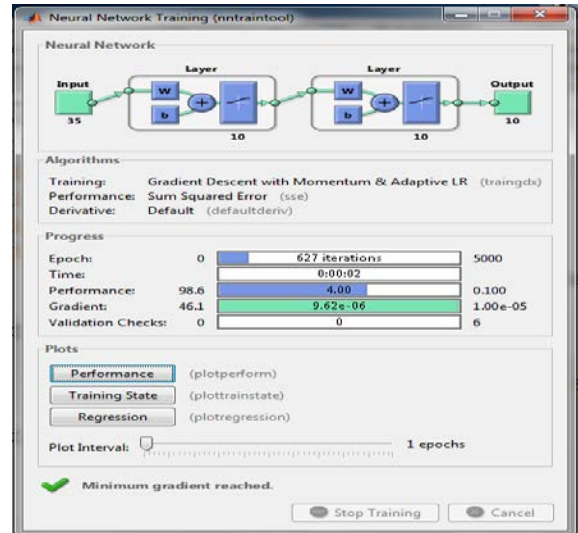


Fig 3: Train network for font recognition system

2. In testing phase the image to be tested depend on after pre processing and extracted feature after tat matching and make decision.

## 4. Result

In this paper we evaluate an approach by construct a dataset which consists of Punjabi character images. These documents are noiseless and without skew and are written in various different font faces .The method has been implemented and evaluated using MATLAB tool, with a dual core, 2.60 GHz Intel Core i5 with 4GB RAM memory. The Common criteria for performance evaluation are based on the Accuracy (A).The accuracy of a measurement system is the degree of closeness of measurements of a quantity to that quantity's actual (true) value.

$$A = \frac{\text{No.of words correctly recognised}}{\text{Total No of Words}} \quad \text{......... Eq. (1)}$$

Table1: Punjabi Typeface Confusion

Matrix (Percentage)

| **Font** | Amrit | Anmol | joy | Prabhki | Bulara | Anmol |
|---|---|---|---|---|---|---|

| Name | | lipi | | | | neon |
|---|---|---|---|---|---|---|
| Amrit | 85 | 4 | 2 | 1 | 1 | 2 |
| Anmo lipi | 0 | 88 | 2 | 0 | 2 | 0 |
| joy | 1 | 0 | 89 | 1 | 0 | 1 |
| Prabhki | 1 | 1 | 2 | 86 | 2 | 4 |
| Bulara | 2 | 2 | 1 | 0 | 88 | 1 |
| Anmol Neon | 1 | 0 | 2 | 2 | 1 | 86 |

In the Table 1, the row represents sample's actual font, and the column represents the result of the font recognition. For example, the data 2 % in the first column of the fifth row represents that there are 2% characters of Bulara font have been recognized as characters of amrit font in the test. The data in the table show that the rates of fonts are all above 87 %, which indicates the validity of this method.

## 5. Conclusion and Future Work

The main objective of this paper is to study and design Punjabi Font Recognition system in MATLAB. In this paper we describe approach for recognition of Punjabi language font style and type of font. We use neural network tool for testing and training so that system's performance can be evaluated in a better way. It is based on local features. Extensive experiments have shown that the algorithm performs very well. The common recognition accuracy of different fonts is 87%. In future, the work may be extended to get more accuracy by using different set of features in recognizing and reorganization done by using word.

**REFERENCES:**

[1] AbdelwahabZramdini and Rolf Ingold, "Optical Font Recognition Using Typographical Features", Transactions on pattern analysis and machine intelligence, VOL. 20, No. 8,IEEE AUGUST 1998.

[2] Min-Chul Jung, Yong-Chul Shin and Sargur N. Srihari, "Multifont Classification using Typographical Attributes", Center of Excellence for Document Analysis and Recognition, 1999.

[3] Chi-Fang, Lina,Yu-Fan, Fang,Yau-TarngJuang, "Chinese text distinction and font identification by recognizing most frequently used characters", Image and Vision Computing 19, 329-338,2001.

[4] Yong Zhu, Tieniu Tan, and Yunhong Wang, "Font Recognition Based on Global Texture Analysis", Transactions on pattern analysis and machine intelligence, VOL. 23, No. 10, IEEE, October 2001.

[5] Ibrahim Abuhaiba, "Arabic Font Recognition Based on Templates", The International Arab Journal of Information Technology, Vol. 1, No. 0, July 2003.

[6] Ibrahim S. I. Abuhaiba, "Arabic Font Recognition using Decision Trees Built from Common Words", Journal of Computing and Information Technology - CIT 13, 3, 211–223, 2005.

[7] Ming-Hu Ha, Xue-Dong Tian, Zi-Ru Zhang , "Optical Font Recognition based on Gabor filter", Proceedings of the Fourth International Conference on Machine Learning and Cybernetics, Guangzhou, August 2005.

[8] A.Borji, and M. Hamidi, "Support Vector Machine for Persian Font Recognition", World Academy of Science, Engineering and Technology, 2007.

[9] R.Ramanathan, L.Thaneshwaran, V.Viknesh,, Dr. K.P.Soman, "A Novel Technique for English Font Recognition Using Support Vector Machines", International Conference on Advances in Recent Technologies in Communication and Computing, 2009.

[10] Usha Rani, Er. Balwinder Singh, Er. Ravinder Singh, "Machine Printed Punjabi Character Recognition Using Morphological Operators on Binary Images", International Journal of Engineering Research & Technology (IJERT), ISSN: 2278-0181, Vol. 1 Issue 3, May, 2012.

[11] FouadSlimane, Slim Kanoun , Jean Hennebert , Adel M. Alimi , Rolf Ingold ," A study on font-family and font-size recognition applied to Arabic word imagesat ultra-low resolution, Pattern Recognition Letters 34, 209–218, 2013.

[12] Harjit Singh, "Detection of Bold and Italic Character in Gurmukhi Script", IOSR Journal of Computer Engineering (IOSRJCE) ISSN: 2278-0661 Volume 1, Issue 6, PP 28-31, July-Aug 2012.

[13] www.sikhnet.com/Gurmukhi-Fonts.

[14] YaghoubPourasad ,HoushangHassibi , AzamGhorbani ," Farsi Font Face Recognition in Letter Level", Procedia Technology 1, 378 – 384, 2012 .

[15] R Sanjeev Kunte and R D Sudhaker Samuel, A simple and efficient optical character recognition system for basic symbols in printed Kannada text, Vol. 32, Part 5, pp. 521–533, October 2007.

**Ms Reena Joshi** , Masters of Technology from DAVIET Jalandhar, Punjab, India. She did her Bachelors of Technology (B-Tech) from College of Engineering & Management Kapurthala, Punjab. She is currently working as an assistant professor in Computer Science department at Rayat and Bhara institute of Neno Technology, Hoshiarpur Punjab, India. She is more than 5 years teaching experience. Her Research area of interest includes Natural Language Processing, Machine Learning, programming language and database management system.

**Gaurav Agnihotri** is a Astt. Professor in department of information technology at GNA University, Phagwara, Punjab, India. He has done B.Tech in Computer Science & Engineering, M.Tech in Information Technology and currently pursuing PhD degree in computer engineering from the Punjabi University, Patiala. Mr. Gaurav has more than 6 years teaching and research experience. He has supervised more than 20 M.Tech. students in Data mining , Data Structure and cloud computing