



# Hindi Numeral Font Recognition Using Neural Network

Saima Rafiq<sup>1</sup>, Jayshree Boaddh<sup>2</sup>, Durgesh Wadbude<sup>3</sup>

PG Student, Dept. of Computer Science, Mittal College of Technology, Bhopal, India <sup>1</sup>

Assistant Professor, Dept. of Computer Science, Mittal College of Technology, Bhopal, India <sup>2</sup>

Head of Department, Dept. of Computer Science, Mittal College of Technology, Bhopal, India <sup>3</sup>

**ABSTRACT:** Handwriting has continued to persist as a means of communication and recording information in day-to-day life even with the introduction of new technologies. The constant development of computer tools lead to the requirement of easier interface between the man and the computer. Handwritten character recognition may for instance be applied to Zip-Code recognition, automatic printed form acquisition, or cheques reading. The importance to these applications has led to intense research for several years in the field of off-line handwritten character recognition. 'Hindi' the national language of India (written in Devanagri script) is world's third most popular language after Chinese and English. Hindi handwritten character recognition has got lot of application in different fields like postal address reading, cheques reading electronically. Recognition of handwritten Hindi characters by computer machine is complicated task as compared to typed characters, which can be easily recognized by the computer. This paper presents a scheme to recognize hindi number numeral with the help of neural network.

**KEYWORDS:** Hindi Numerals, Hindi Font recognition, Neural Network, Training, Testing, Images etc.

## I. INTRODUCTION

English Character Recognition (CR) has been extensively studied in the last half century and progressed to a level, sufficient to produce technology driven applications. But same is not the case for Indian languages which are complicated in terms of structure and computations. Digital document processing is gaining popularity for application to office and library automation, bank, publishing houses communication technology, postal services and many other areas. With ever increasing requirement for office automation, it is necessary to provide practical and effective solutions. Devanagri character recognition is becoming more and more important in the modern world. It helps human ease their jobs and solve more complex problems over the few past years, the numbers of companies involved in research on handwritten recognition are increasing continually. So Devanagri being the base of many Indian languages should be given special attention so that document retrieval and analysis of rich ancient and modern Indian literature can be effectively done. Development of a Character recognition system for Devanagari is difficult because (i) there are about 350 basic modified ("matra") and compound character shapes in the script and (ii) the characters in a word are topologically connected which is not in case of English characters. Here focus is on the recognition of offline handwritten Hindi characters that can be used in common applications like bank cheques, commercial forms, government records, bill processing systems, postcode recognition, signature verification, passport readers, offline document recognition generated by the expanding technological society.

Challenges in handwritten characters recognition lie in the variation and distortion of offline handwritten Hindi characters since different people may use different style of handwriting, and direction to draw the same shape of any Hindi character. This overview describes the nature of handwritten language, how it is translated into electronic data, and the basic concepts behind written language recognition algorithms.

Handwritten Hindi character are imprecise in nature as their corners are not always sharp, lines are not perfectly straight, and curves are not necessarily smooth, unlike the printed character. Furthermore, Hindi character can be drawn in different sizes and orientation in contrast to handwriting which is often assumed to be written on a baseline in an upright position. Handwritten characters also depend upon the mood of the person who is writing. Therefore, a robust offline Hindi handwritten recognition system has to account for all of these factors. The work that has been done in the



# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 9, September 2016

area of Devanagari text recognition is limited to only characters, no work has been reported for word, sentence or the entire document identification .

This paper recognized Devanagari numerals in a handwritten Devanagari curve script using ANN (Artificial Neural Network approach). Artificial Neural Network (ANN), often called as neural network (NN), is a mathematical model or structure or we can also say computational model that is inspired by the functional aspects and structure of biological neural networks. Neural networks have been implemented successfully in various fields like voice recognition, iris recognition, odor recognition and clustering. They are used to solve complicated problems. It is an effort in the field to make computers as intelligent as human beings i.e. it makes the computer act more like a human beings and answer “what if questions “to the users. ANNs are being used in a vast domain of pattern recognition, one of the areas of pattern recognition is handwritten text recognition. To recognize a handwritten text using ANN approach, we need to design a multilayer Perceptron i.e. a neural network with input layer, hidden layer and output layer. The size of these three layers may vary according to the problem. To recognize a handwritten Hindi curve script the first we that we required is the sample images of handwritten text, which was obtained by scanning handwritten text script through a scanner and saving the file in bitmap format. Before this image is feeded to the algorithm, there are certain steps that are required to be executed.

- a) Preprocessing- In this phase the image is converted into grayscale and then binary image, then the image is made noise free i.e. removing any unwanted bit of pattern from the image, once the image is made noise free it is sent to a routine that skeletonizes (thinning) it. After skeletonizing the image the pixels required for the recognition are mapped into a fixed size matrix, in our project we have taken the size of the matrix as 10\*15.
- b) After completion of the preprocessing steps the image is segmented into individual characters. In the case of Hindi words the Shirrekha of the word has to be removed first and then the individual characters are extracted. So we developed an algorithm to remove the Shirrekha from each individual word in the document.
- c) Before starting the recognition process the neural network was to be trained with dataset (that we prepared manually for this project).Once the network was trained with the datasets, it was ready to identify the handwritten texts in any handwriting.
- d) The network was tested with about 40 sample images and of different individuals.

## II. LITERATURE SURVEY

OCR work on printed Devanagari script started in early 1970s. Among the earlier pieces of work, some of the efforts on Devanagari character recognition are done by Sinha and Mahabala (1979). A syntactic pattern analysis system and its application to Devanagari script recognition is discussed in his doctoral thesis. They also presented a syntactic pattern analysis system with an embedded picture language for the recognition of handwritten and machine printed Devanagari characters. The system stores structural description for each symbol of the Devanagari script in terms of primitives and their relationships. For recognition, an input character is labeled and compared it with stored description. To increase the accuracy of the system and reduce the computational costs, contextual information regarding the occurrences of certain primitives and their combinations and restrictions are used. They also demonstrated how the spatial relationship among the constituent symbols of Devanagari script plays an important role in the interpretation of Devanagari words. There are a number of constraints on these spatial relationships which characterize Devanagari script composition syntax. When the word composition is not found to be syntactically correct, the symbols are substituted with their resembling counterparts. The symbol substitution rules are mostly heuristic in nature.

Most Indian languages are very inflectional in nature. Because of this inflectional behavior, development of OCR error detection and correction technique is not an easy task. The complex character grapheme structure of some Indian scripts also creates difficulty in recognition error detection and correction. An OCR error correction scheme for the Devanagari text is proposed by Bansal and Sinha (1997). They used a partitioned word dictionary to reduce the search space besides preventing forced match to incorrect word. The envelope information of words consisting of number of top, lower, core modifiers along with the number of core characters form the second level partitioning feature for short words partition. The remaining words are further partitioned using a string of fixed length associated with each partition. A distance matrix for assigning penalty for a mismatch is incorporated in the search process.

The ability to identify machine printed characters in an automated or a semi-automated manner has obvious applications in numerous fields. Since creating an algorithm with a one hundred percent correct recognition rate is quite



## International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 9, September 2016

probably impossible in our world of noise and different font styles, it is important to design character recognition algorithms with these failures in mind so that when mistakes are inevitably made, they will at least be understandable and predictable to the person working with the program. Brown (1992) explores one such algorithm and tests it on two different fonts using a third font as a reference. The results are discussed and several improvements are suggested. He describes an algorithm that attempts to work with a subset of the features in a character that a human would typically see for the identification of machine-printed English characters. Its recognition rate is currently not as high as the recognition rates of the older, more developed character recognition algorithms, but it is expected that if it were expanded to work with a larger set of features this problem would be removed. If it were expanded to use more features, it would be made correspondingly slower; with the advent of faster microprocessors this fact is not viewed as a crippling problem. The procedure for extracting these feature points utilized by this algorithm is fairly straightforward. Since an eight by eight character consists of only sixty-four pixels, it is viable to simply loop through the entire character and examine each pixel in turn. If a pixel is on, its eight neighbors are checked. Since each neighbor can also only be on or off, there are merely 256 possible combinations of neighborhoods. Of these 256, fifty-eight were found to represent significant feature points in a fairly unambiguous way. Extracting feature points thus reduced to calculating a number between zero and 256 to describe a pixel's neighborhood and then comparing that number against a table of known feature points. While it is true that this method does not always catch every feature point (some can only be seen in a larger context) it catches the majority. Missing feature points is certainly not a limiting factor in the algorithm's accuracy. It also does not suffer from labeling too many uninteresting points as being feature points. It has virtually no false positives. The feature point extractor is thus fast and reliable.

Chandra et al (2006) proposed a system for the recognition of online handwritten characters for Indian writing systems. A handwritten character is represented as a sequence of strokes whose features are extracted and classified. Support Vector Machines (SVM) has been used for constructing the stroke recognition engine. The results have been presented after testing the system on Devanagari and Telugu scripts. An SVM machine is capable of learning to achieve good generalization performance, given a finite amount of training data, by striking a balance between the goodness of fit attained on a given training dataset and the ability of the machine to achieve error-free recognition on other datasets. With this concept as the basis, support vector machines have proved to achieve good generalization performance with no prior knowledge of the data. Burges (1998) presented a tutorial on SVM is to map the input data onto a higher dimensional feature space nonlinearly related to the input space and determine a separating hyper plane with maximum margin between the two classes in the feature space. This results in a nonlinear boundary in the input space. The optimal separating hyper plane can be determined without any computations in the higher dimensional feature space by using kernel functions in the input space.

An SVM in its elementary form can be used for binary classification. It may, however, be extended to multi class problems using the one-against-the-rest approach or by using the one-against-one approach. Arora and Bhattacharjee (2002) present a two stage classification approach for handwritten Devanagari characters. The first stage is using structural properties like shirorekha, spine in character and second stage exploits some intersection features of characters which are fed to a feed forward neural network. Simple histogram based method does not work for finding shirorekha, vertical bar (Spine) in handwritten Devanagari characters. So a new technique, differential distance based technique to find a near straight line for shirorekha and spine. This approach has been tested for 50000 samples and got 89.12% success.

Mishra and Rajput (2008) presented a system for recognizing hand written Indian Devanagari script. The system considers a handwritten image as an input, separates the lines, words and then characters step by step and then recognizes the character using artificial neural network approach, in which Creating a Character Matrix and a corresponding Suitable Network Structure is the most important step. In addition, knowledge of how one is deriving the Input from a Character Matrix must first be obtained before one may proceed. Afterwards, the Feed Forward Algorithm gives insight into the entire working of a neural network; followed by the Back Propagation Algorithm which comprises Training, Calculation of Error, and Modifying Weights. Once the characters are recognized they can be replaced by the standard fonts to integrate information from diverse sources.

Verma and Blumenstein (1997) presented a new intelligent segmentation technique is proposed that may be used in conjunction with a neural classifier and a simple lexicon for the recognition of difficult handwritten words. A heuristic



# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 9, September 2016

segmentation algorithm is initially used to over-segment each word. An Artificial Neural Network (ANN) trained with 32,034 segmentation points is then used to verify the validity of the segmentation points found. Following segmentation, character matrices from each word are extracted, normalized and then passed through a global feature extractor after which a second ANN trained with segmented characters is used for classification. These recognized characters are grouped into words and presented to a variable-length lexicon that utilizes a string processing algorithm to compare and retrieve words with highest confidences. This research provides promising results for segmentation, character and word recognition. In the proposed word recognition system, heuristic and intelligent methods are used for the segmentation of real world, handwritten words. Following segmentation, character matrices are extracted from the words and classified. Finally, to show how the segmentation technique may possibly be used in the context of an overall system, a lexicon is used to match each set of recognized characters (each set represents a single word) to potential correct words.

An algorithm for segmentation of touching Devanagari characters into its constituent symbols and characters proposed by Bansal and Sinha (1997) proposed algorithm extensively uses structural properties of the script. Statistical information about the height and width of character boxes, which are vertically separate from their neighbors, is used to hypothesize character boxes.

Sarkar (2006) showed that Artificial Neural Networks (ANNs) have been successfully applied to Optical Character Recognition (OCR) yielding excellent results. A method is used for segmentation of difficult handwriting with the use of conventional algorithms in conjunction with ANNs. The segmentation algorithm is heuristic in nature detecting important features which may represent a prospective segmentation point. An Artificial Neural Network is subsequently used to verify the authenticity of the segmentation points found by the algorithm. The C programming language, the SP2 supercomputer and a SUN workstation were used for the experiments. The algorithm has been tested on real-world handwriting obtained from the CEDAR database.

### III. ARTIFICIAL NEURAL NETWORK

An Artificial Neural Network (ANN) is an information processing structure that is adapted from biological nervous systems, such as the nervous system, brain. The basic element of this structure is the new structure of the information processing system. It consists of many highly interconnected information processing elements (neurons) working together to solve specific problems. Just like people, ANNs learn by example. An ANN is trained for a specific application, such as pattern recognition or data classification, by learning process. In a biological system learning means adjusting the synaptic connections between the neurons. The same is done in ANN. A biological neural network is made up of a group of chemically connected components or functionally associated neurons. A single neuron is connected to many other neurons and there may be a large number of neurons or connections. Connections between the neurons, called synapses, are formed from axons to dendrites. The structure and functioning of neural networks are extremely complex. Artificial intelligence and algorithms associated with cognitive system try to simulate some properties of biological neural networks. Although both are similar in techniques, but biological neuron aims at solving a particular tasks, while an artificial neuron aims to build mathematical models from biological neural systems. Artificial neural network have been applied successfully in artificial intelligence field, they have been applied successfully to voice recognition, image processing and others, with a purpose to construct software agents (in computer and video games) or autonomous robots. All the areas like Artificial intelligence, neural networks cognitive modeling, and are information processing structure inspired by the working of biological neural systems.

Neural networks have ability to derive meaning from complicated or imprecise data, which is why it can be used to extract patterns and detect trends that are very complex to be noticed by either humans or other computer technology. A trained neural network can assumed to as an "expert" in the given area for which it has been trained. Other advantages include:

**Adaptive learning:** It is an ability to learn how to do tasks that is based on the data given for training or initial experience.

**Self-organization:** It is a property of ANN that it can create its own organization or representation of the information which receives during learning time.

# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 9, September 2016

**Real time functions:** All the ANN calculations may be carried simultaneously, and special hardware devices are being designed and manufactured which take up advantage of this capability of ANN.

**Fault tolerance by redundant information coding:** If there is partial destruction in the neural network, the entire functioning does not stop but instead it continues to work with a bit low performance.

Component of a neuron is shown in Fig. 1. and its synapse is shown in Fig. 2.

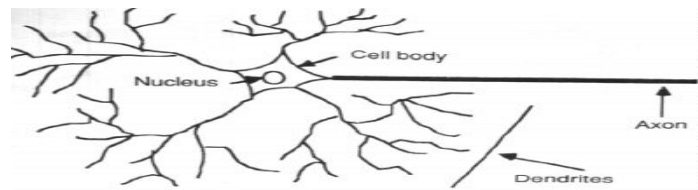


Fig. 1. Components of a neuron

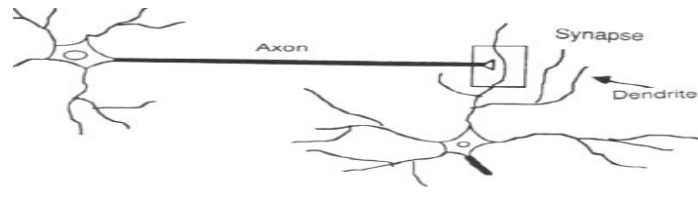


Fig. 2 The synapse

## Structure of an artificial neurons

In ANN we first try to take out the essential features neurons for recognizing and their interconnections. We then program a computer or write algorithm to simulate these features. But since our knowledge of neurons is incomplete and our computing power is also limited, our models are only close to the model of real networks of neurons. A typical neuron model is shown in Fig. 3

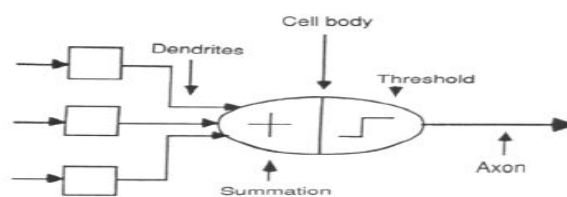


Fig. 3 A Neuron Model

## An engineering approach for neural network

A simple ANN is a structure with many inputs and one output. The neuron has two modes of operation; the training mode and the testing mode. In the training mode, the neuron is trained to fire, for particular input patterns. In the testing mode, when a taught input pattern is detected at the input, its corresponding output becomes the current output. If the input pattern does not belong in the taught list of input patterns, the firing rule is used to determine whether to fire or not. An example of simple neuron is shown in Fig. 4.

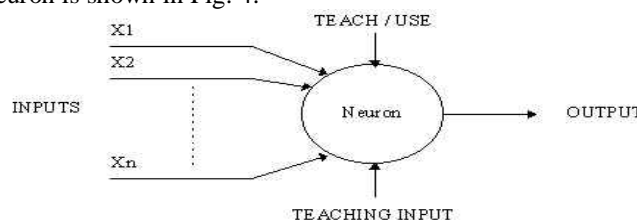


Fig. 4 A simple neuron

# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 9, September 2016

The previous neuron doesn't do anything that conventional computers don't do already. A more sophisticated neuron (Fig. 5) is the McCulloch and Pitts model (MCP). The difference from the previous model is that the inputs are 'weighted'; the effect that each input has at decision making is dependent on the weight of the particular input. The weight of an input is a number which when multiplied with the input gives the weighted input. These weighted inputs are then added together and if they exceed a pre-set threshold value, the neuron fires. In any other case the neuron does not fire.

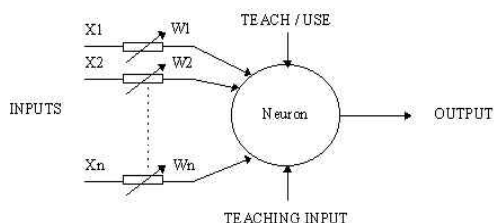


Fig. 5 An MCP neuron

In mathematical terms, the neuron fires if and only if;

$$X_1W_1 + X_2W_2 + X_3W_3 + \dots > T$$

The addition of input weights and of the threshold makes this neuron a very flexible and powerful one. The MCP neuron has the ability to adapt to a particular situation by changing its weights and/or threshold. Various algorithms exist that cause the neuron to 'adapt'; the most used ones are the Delta rule and the back error propagation. The former is used in feed-forward networks and the latter in feedback networks.

## How Neural Network Functions

According to Jain, Mohiuddin and Mao (1996) a neural network is a massively parallel distributed processor that has a natural propensity for storing experiential knowledge and making it available for use. It resembles the brain in two respects:

1. They adapt by learning process.
2. Knowledge is stored in interconnections between neurons known as synaptic weights.

Basically, learning is a process by which the free parameters (i.e. synaptic weights and bias levels) of a neural network are adapted through a continuing process of stimulation by the environment in which the network is embedded. The type of learning is determined by the manner in which the parameter changes take place.

Broadly learning can be classified into two categories:

- 1. Supervised Learning:** This form of learning assumes the availability of a labeled (i.e., ground-trusted) set of training data made up of N input-output.
- 2. Unsupervised Learning:** This form of learning does not assume the availability of a set of Training data made up of N input-output. They learn to classify input vectors according to how they are grouped spatially and try to tune its network by considering a neighbourhood.

## IV. PROPOSED WORKFLOW

The proposed work is carried out in following stages:

1. Image is taken as training database, and training is done using feed forward neural network algorithm.
2. Trained image is stored in database as a template which can be used for matching the similar images during testing.
3. Test image is taken.
4. Selection of the particular number is made.
5. The cropping of particular numbers are taken place.
6. Preprocessing is done to remove the noise in the cropped part.
7. The recognition shows the recognised digit from the trained database.

## V. RESULTS

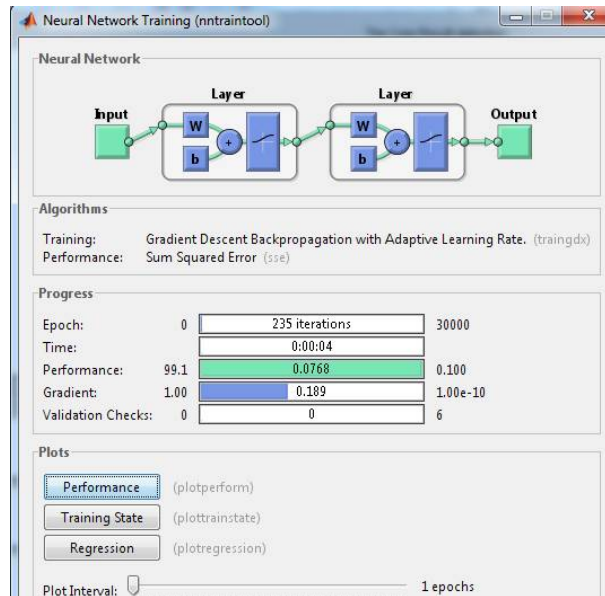


Fig 6. Neural Network training of the input or training set

Fig.6 shows the Neural N/W training set when an image is loaded or trained.

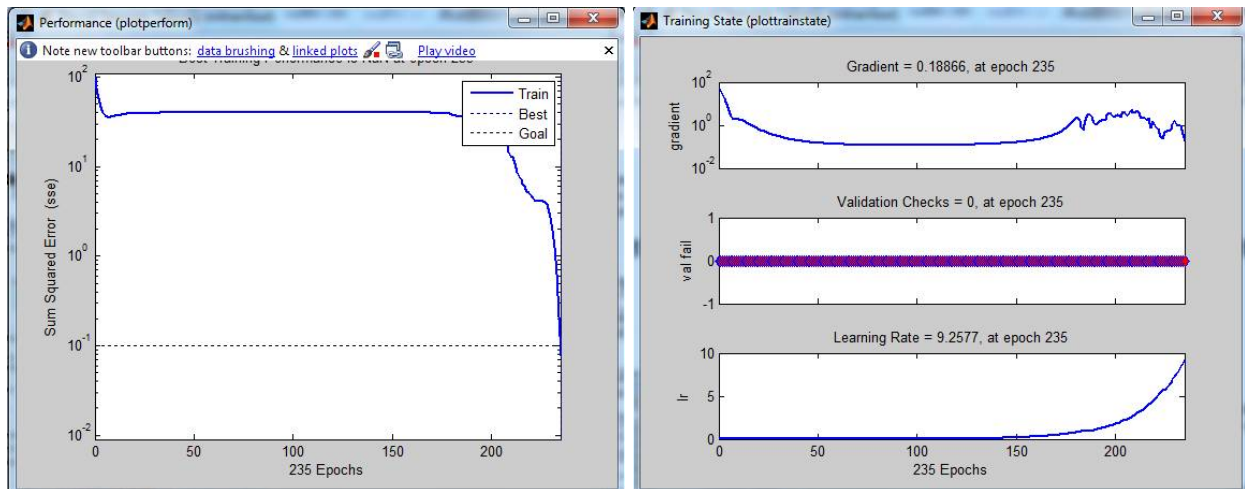


Fig 7. Performance Plot and Training state of the neural network training

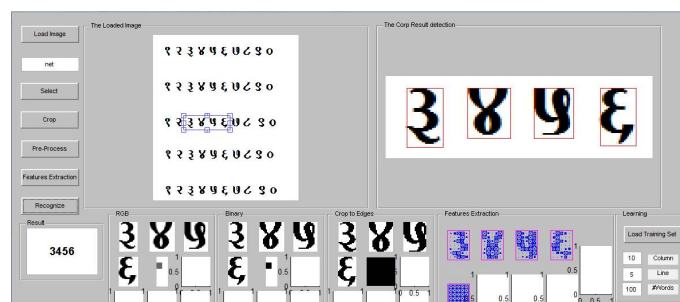


Fig 8. Selected, cropped, preprocessed and Recognised values



# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 5, Issue 9, September 2016

## VI. CONCLUSION & FUTURE SCOPE

Offline handwritten Hindi character recognition is a complex as well difficult problem, not only because of the the variations in human handwriting, but also, because of the overlapped and joined characters as in Hindi. Recognition approaches heavily depend on the nature of the data to be recognized. Since handwritten Hindi characters could be of various shapes and size, the recognition process needs to be much efficient and accurate to recognize the characters written by different users. This paper proposes a technique of applying Radial Basis Function for handwritten Devnagri numeral recognition. Since the database is not globally available, firstly we created the database, and then by the use of Principal Component Analysis we extracted the features of each image. At the hidden layer, centres are determined and the weights between the hidden layer and the output layer of each neuron are determined to calculate the output, where output is the summing value of each neuron.

It can be extended for the recognition of sentence and documents. Another research interest will be on the character images degraded or blurred by various reasons and skew detection and correction. This approach can be used in multilingual character recognition as well. It can be extended to recognize all the matras and ardhakshars in the Hindi script with better efficiency.

## REFERENCES

- [1] Arora S., Bhattacharjee D., Malik L. and Nasipuri M., 2007, A Two Stage Classification Approach for Handwritten Devanagari characters, International Conference on Computational Intelligence and Multimedia Applications, December 13-15, 2007, vol. 2, pp 399-403.
- [2] Bahlmann , Burkhardt , H. and Haasdonk, C.B., 2004, Online Handwriting Recognition With Support Vector Machine- A Kernel Approach, *IEEE Transaction on Pattern Analysis Machine Intelligence*, Vol. 26, Issue 3, pp 299-310.
- [3] Bajaj, R., Chaudhary, S. and Dey, L., 2002, Devanagari numeral recognition by combining decision of multiple connectionist classifiers, *Sadhna* Vol.27, Part 1, pp 59-72.
- [4] Bansal, V. and Sinha, R.M.K., 1999, On how to describe shapes of Devanagari characters and use them for recognition, Proceedings of the 5th Int. Conference on Document Analysis and Recognition, Bangalore, India, pp 410-413.
- [5] Bansal, V. and Sinha, R.M.K., 1995, On Devanagari Document Processing, Int. Conference on Systems, Man and Cybernetics, Vancouver, Canada, Oct 22-25, 1995, pp 1621 - 1626.
- [6] Bansal, V. and Sinha, R.M.K., 1997(a), Integrating Knowledge Sources in Devanagari Text Recognition System", Technical Report, I.I.T. Kanpur, India, pp 97-248.
- [7] Bansal, V. and Sinha, R.M.K., 1997(b), On Automating Trainer For Construction of Prototypes for Devanagari Text Recognition, Technical Report, I.I.T. Kanpur, India, pp 95-232.
- [8] Bansal, V. and Sinha, R.M.K., 1997(c), Partitioning and Searching Dictionary for Correction of Optically-Read Devanagari Character Strings, Technical Report, I.I.T. Kanpur, India, pp 97-246.
- [9] Bansal V. and Sinha, R.M.K., 1997(d), Segmentation of touching and fused Devanagari characters, Technical Report, TRCS, I.I.T. Kanpur, India, pp 97-247.
- [10] Bansal V and Sinha, R. M. K., 1996, Designing a Front End OCR System for Indian Scripts for Machine Translation - A Case Study for Devanagari, Symposium on Machine Aids for Translation and Communication (SMATAC-96), New Delhi, India.
- [11] Bin, Yong, Z.L and Shao-Wei, X., 2000, Support Vector Machine and Its Application In Handwritten Numeral Recognition, Proceedings of the 15<sup>th</sup> Int. conf. on Pattern Recognition, Barcelona, Spain, Sept 3-8, 2000, pp 720-723.
- [12] Blumenstein, M. and Verma B., 1998, "An Artificial Neural Network Based Segmentation Algorithm for Off-line Handwriting Recognition", International Conference on Computational Intelligence and Multimedia Applications (ICCAL4 '98), Melbourne, Australia.
- [13] Blumenstein, M. and B. Verma, 1999, A New Segmentation Algorithm for Handwritten Word Recognition", IEEE conference of IJCNN'99, Washington, U.S.A, Vol. 4, pp 2893-2898.
- [14] Brown, Eric W., 1992, Character Recognition by Feature Point Extraction, Northeastern University internal paper.
- [15] Burges, C. J. C., 1998, A tutorial on support vector machines for pattern recognition, *Data Mining and Knowledge Discovery*, Data Mining and Knowledge Discovery, Vol. 2, Issue 2, pp 121-167.
- [16] Casey, R. G. and Lecolinet, E., 1996, A survey of Methods and Strategies in Character Segmentation , *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 18, Issue 7, pp 690-706.
- [17] Chandra Sekhar, C., Jayaraman Anitha, Srinivasa Chakravarthy , Swethalakshmi V. H., 2006, Online Handwritten Character Recognition of Devanagari and Telugu Characters using Support Vector Machines, Tenth International workshop on Frontiers in handwriting recognition, 6 October 2006.
- [18] Chatterjee, B. and Sethi, I.K., 1976, Machine recognition of hand printed Devanagari Numerals, *Journal of Institution of Electronics and Telecommunication Engineers*, vol. 22 Issue 1, pp 532- 535.
- [19] Chaudhuri, B.B ., 1998, A complete printed Bangla OCR system, *Elsevier Journal of Pattern Recognition*, Vol. 31, Issue 5, pp 531-549.
- [20] Chitnis, S.D. and Khanale P.B., 2011, Handwritten Devanagari Character Recognition using Artificial Neural Network, *Journal of Artificial Intelligence*, Vol.4, Issue 1, pp 55-62.