An Optical Character Recognition System from Printed Text and Text Image using Adaptive Neuro Fuzzy Inference System

Mustain Billah

Department of Information and Communication Technology Mawlana Bhashani Science and Technology University Tangail,Bangladesh Sajjad Waheed

Department of Information and Communication Technology Mawlana Bhashani Science and Technology University Tangail,Bangladesh

Abu Hanifa

Department of Information and Communication Technology Mawlana Bhashani Science and Technology University Tangail,Bangladesh

ABSTRACT

This is the age of digital systems. Now a days, everything is being computerized. Peoples are using mobile phones, laptop, computer, camera, notebook, pdf reader etc digital systems too much than ever. Use of papers and pen, printed books are decreasing. Rather peoples are using digital means of communication, study, documentation. Optical character recognition is an application of these digital systems. There are many ways and systems to recognize a character from a printed document. Many research works have been done for OCR in different languages using neural network, support vector machine, markov model. In this paper, an Adaptive Neuro Fuzzy Inference System (ANFIS) based methodology has been proposed for recognizing characters from printed documents such as text image. Purpose of this paper is not making a software but proving the usability and effectivity of ANFIS for optical english character recognition and propose a anfis system with less error and more efficiency. Proposed system can detect every character perfectly wihout error.

General Terms

Computer science, Image processing

Keywords

Optical character recognition, OCR, adaptive neuro fuzzy inference system, ANFIS, Printed text, artificial intelligence, image processing

1. INTRODUCTION

Day by day technology is advancing. Newer mehods are being discovered to make our life easier. The most advancement has been done in the field of communication and digital systems. Especially digital systems of communication[7], documentation and computation are going ahead too much. Now people read books in their mobile phone or pdf reader or tablet pc irrespective of place and time. They have not to buy huge and larger books. Instead they load a pdf version of the book and read it. Again they can capture or scan an image of printed books or documents and convert as intended format using optical character recognition[8] systems such as pdf converter, image converter. Thus it saves time and cost and work load.

The process of converting an text image or printed text into computer redable and editable format[6] starts with scaning and loading the image and recognising each character independently and accurately. Human can easily detect a character at the first see, but a computer or digital system can not do it unless it is familiar or trained with previous records. There are many algorithms, ways and methods of doing such job. Artificial intelligence[2] has a field of pattern recognition for recognising patterns such as character,face, image,voice etc. Artificial neural network[12], support vector machine[11] and various data mining techniques are being used for character recognition.

Adaptive neuro fuzzy inference system (ANFIS)[5] is now a advance field of artificial intelligence. Various types of problems are being solved by ANFIS. Being a current important issue, there are a lot of research works on optical character recognition. For purpose of our research work, we have analyzed many papers. In the paper "Optimization of optical character recognition for printed devanagari text using anfis techniques" Ganash s. Sable & Sheetal Arun Irve proposed an anfis based method for recognising devanagari text with 92.66accuracy[10]. Agazzi, Oscar E, and Shyh-shiaw Kuo proposed a OCR model based on hidden markov model in the paper "Hidden Markov model based optical character recognition in the presence of deterministic transformations"[1]. Cash, Glenn L, and Mehdi Hatamian in the paper "Optical character recognition by the method of moments" used two dimensional moments as features for recognizing characters. They applied the method to six machineprinted fonts[3].Barve, Sameeksha used back propagation neural neywork for recognising typed english characters in the paper "Optical character recognition using artificial neural network" [4]. Patil, Mrs VV, Rajharsh Vishnu Sanap, and Rohini Babanrao Kharate examines the use of neural networks to accomplish optical character recognition in the paper "Optical Character Recognition Using Artificial Neural Network" [9]. Singh, Raghuraj used artificial neural network for recognising davanagari character in "Optical character recognition (OCR) for printed devnagari script using artificial neural network". In this paper, an ANFIS based system or model has been shown for character recognition. For testing proposed method's accuracy and validation, english characters of different fonts and faces from different text images have been used. Testing results show that, ANFIS can recognize characters properly and accurately from printed text and text image.

The paper is organised as follows: Section II contains a brief description of theory of ANFIS. Section III describes the proposed model. It has some sub sections describing each parts of the model. In the section IV, ANFIS has been trained. Results are analyzed in section V.

2. THEORY OF ANFIS

ANFIS derives its name from adaptive neuro-fuzzy inference system.It works similarly to that of neural networks. Using a given input/output data set, ANFIS constructs a fuzzy inference system (FIS) whose membership function parameters are tuned (adjusted) using either a backpropagation algorithm alone or in combination with a least squares type of method. This adjustment allows the fuzzy systems to learn from the data they are modeling. [4] presents the brief theory of ANFIS about how it works:

Assume that the fuzzy inference system has two inputs x and y and one output z. A first-order Sugeno fuzzy model has rules as the following:

Rule1:

If x is A1 and y is B1, then f1 = p1x + q1y + r1

Rule2:

If x is A2 and y is B2, then f2 = p2x + q2y + r2

Here,

$$output = \frac{w1.f1 + w2.f2}{w1 + w2} \tag{1}$$



Fig. 1. Sugeno model



Fig. 2. Anfis architechture with two input variable

Layer 1 - I

 $O_{1,i}$ is the output of the ith node of the layer 1. Every node i in this layer is an adaptive node with a node function.

$$O_{1,i} = \mu Ai(x)$$
 for $i = 1, 2, or$ (2)

$$O_{1,i} = \mu B i - 2(x)$$
 for $i = 3, 4$ (3)

x (or y) is the input node i and Ai (or Bi-2) is a linguistic label associated with this node. Therefore $O_{1,i}$ is the membership grade of a fuzzy set (A1,A2, B1, B2).

Layer 2

Every node in this layer is a fixed node labeled "Prod". The output is the product of all the incoming signals.

$$O_{2,i} = wi = \mu Ai(x) \cdot \mu Bi(y)$$
 for $i = 1, 2, ..$ (4)

Each node represents the fire strength of the rule Any other T-norm operator that perform the AND operator can be used.

Layer 3

Every node in this layer is a fixed node labeled Norm. The ith node calculates the ratio of the ith rule's firing strength to the sum of all rule's firing strengths.

$$O_{3,i} = wi = \frac{wi}{w1 + w2}$$
 for $i = 1, 2, ...$ (5)

Outputs are called normalized firing strengths.

Layer 4

Every node i in this layer is an adaptive node with a node function:

$$O_{3,i} = \bar{wi} fi = \bar{wi} (px + qi.y + r)$$
 for $i = 1, 2, ...$ (6)

wi is the normalized firing strenght from layer 3. $\{pi, qi, ri\}$ is the parameter set of this node. These are referred to as consequent parameters.

Layer 5

The single node in this layer is a fixed node labeled sum, which computes the overall output as the summation of all incoming signals:

$$O_{5,i} = \sum \bar{w}i.fi = \frac{\sum \bar{w}i.fi}{\sum \bar{w}i}fi \text{ for } i = 1, 2, ...$$
 (7)

Infact, ANFIS is the combination of both ANN fuzzy logic. ANN algorithms are also used for anfis training, learning.

3. PROPOSED MODEL

Actually all the research works related with OCR have about the same mehodology in image processing section. Variation comes in the feature extraction sections and used tools. However, our system also comprises of such methodology except that we have used image resize algorithm. A flow chart of our proposed model is given below:



Fig. 3. Our proposed model for OCR using ANFIS

3.1 Image acquisition

In image acquisition process, digitized images are obtained from real world sources. Scanner, PDA, digital camera, camcorder etc are used for this purpose. Pdf books are also sources of optical characters. However, pdf documents have been used as the sources of character images as we are going to build a system to recognise optical characters. In the following figure 4, a sample image is represented.



Fig. 4. Scanned image

3.2 Image preprocessing

3.2.1 segmenting individual character. There are thousands of characters in an text image. To recognise them individually they must be segmented for analysing them. In this part, each characters is separated.



Fig. 5. Separated individual character

3.2.2 Converting to grayscale. Next step is to convert image from rgb to grayscale image. Matlab function has been used for converting.

3.2.3 Image filtering. Scanned image may contain unwanted noise. So used median filter is used to remove noise.

3.2.4 converting to binary image. In this stage, each image is converted from grayscale to it's binary form which contain 1 for white portion and 0 for black portion.

3.2.5 *Image resize.* Now each separated charcter is in it's binary form. But to find the actual shape of any character it should be resized in it's actual smallest shape. It will remove all unnecessarry background and keep only the smallest rectangle around the actual character's rectangular boundary. These part is needed when image separation is not done perfectly. Following figure will clear the concept:



Fig. 6. Image resizing

For accomplishing above task a simple algorithm is designed and named as "resize algorithm". Followung algorithm 1 is the resize algorithm.

3.3 Feature Extraction

Now the main and critical point comes. It is feature extraction, the main part on which proposed system?s accuracy depends. After getting the resized image, it is time to find out important features from the image. we can gain many information from an image but all of them do not influence the outcome of a system. Many features are found but choosen only seven of them, as large number of input to ANFIS makes the computer system complex to train the ANFIS creating large sets of rules.

Area. Total areas of images are different for each characters as we have brought the character image into it's actual shape. For example, character A and I have different area.

Algorithm 1 Algorithm for resizing Image
Input the Image
[row column] = size(Image)
% delete the row with all white background%
for $i = 1$ to row do
count = 0
for $j = 1$ to $column$ do
if $Image(i, j) = 1$ then
count = count + 1
end if
end for
if $count = column$ then
delete the ith row
end if
end for
% delete the column with all white background%
[row column] = size(Image)
for $i = 1$ to $column$ do
count = 0
for $j = 1$ to row do
if $Image(i, j) = 1$ then
count = count + 1
end if
end for
if $count = row$ then
delete the ith column
end if
end for
return resized Image

Total black pixel. As size of characters are different, so the number of black pixels(number of 0?s) are also different.

$$Totalblackpixel = row * column - nnz(Image)$$
(8)

nnz is the number of nonzero elements.

Entropy. Entropy is a statistical measure of randomness that can be used to characterize the texture of the input image.

$$entropy = entropy(Image)$$
 (9)

GLCM. A statistical method of examining texture that considers the spatial relationship of pixels is the gray-level co-occurrence matrix (GLCM). The GLCM functions characterize the texture of an image by calculating how often pairs of pixel with specific values and in a specified spatial relationship occur in an image, creating a GLCM, and then extracting statistical measures from this matrix. GLCM contains four measures of an image.

Contrast. Returns a measure of the intensity contrast between a pixel and its neighbor over the whole image.

Correlation. Returns a measure of how correlated a pixel is to its neighbor over the whole image.

Energy. Returns the sum of squared elements in the GLCM.

Homogeneity. Returns a value that measures the closeness of the distribution of elements in the GLCM to the GLCM diagonal. For example seven feature values of discussed input is shown below:

Table 1. Features sample of B and C	mple of B and C
-------------------------------------	-----------------

	_	
Features	В	С
Area	795	1248
Total Black pixel	1270	1800
Entropy	0.9959	0.8893
Contrast	0.0912	0.0615
Correlation	0.8157	0.8553
Energy	0.4220	0.5171
Homogeneity	0.9544	0.9692

4. TRAINING ANFIS

Input dataset is ready which contains seven features value of each individual characters. A specific output value is assigned to each characters input dataset. For example, a is assigned 1, b is 2, c is 3, A is 27, B is 28. For experiment convenience, font times new roman, 48, bold, A-Z and a-z characters are choosen.



Fig. 7. Our trained ANFIS with 7 input

However, membership functions play a great rules in ANFIS. Different MFs with different numbers makes the result to vary from each to others. We have used a variety of MFs for gaining the best result. Following table shows some of the experimental results:

Table 2.	Results	of	ANFIS	training

MFs number	MFs name	Error
2 2 2 2 2 2 2 2	triangular	0.0294
$2\ 2\ 2\ 2\ 2\ 2\ 2\ 2$	gbellmf	0.04
$2\ 2\ 2\ 2\ 2\ 2\ 2\ 2$	gaussian	0.03
$2\ 2\ 2\ 4\ 2\ 2\ 2$	triangle	0.0059

From the table it is seen that, ANFIS system with 2 2 2 4 2 2 2 configuration of triangular membership functions shows least error.

5. RESULT AND DISCUSSION

For testing purpose, a test data set is prepared of different character's seven features. Each character has it's own output label which means a specific character. For example, 1 means a, 26 means z, 27 means A etc. After testing with test dataset, fraction output value is gained in some cases such as 5.88 or 9.31 etc.

But target is to recognize which character it is. So round values are needed. Floor and ceiling function is used to gain round figure



Fig. 8. Averaging testing error after training

which corresponds to a specific character. Table 1 shows ANFIS output and actual desired output.



Fig. 9. ANFIS Error

 Table 3. Anfis output vs actual output

ANFIS output	Floor/ceiling value	corresponding character
11.99961945	12	1
13.00118961	13	m
13.99914059	14	n
14.99964811	15	0
15.99966288	16	р

Thus proposed ANFIS system could recognize all the test character perfectly. But recognising capital letter 'I' and small letter 'I' is somewhat problamatic if input image contains much noise. However satisfactory results are gained which is usable for optical character recognition.

6. CONCLUSION

Optical character recognition is a complex task. It may not be easy all the time as image contains noise and complete noise removal is not possible. In this paper, a system is showed for recognising optical character using adaptive neuro fuzzy model (ANFIS). There are some works on OCR using ANFIS for different language. Proposed method works on english character and shows more than 98% accuracy. In future, we will work on the comparison of proposed ANFIS and neural network to find out the best soft computing technique.

7. REFERENCES

- Oscar E. Agazzi and Shyh-shiaw Kuo. Hidden markov model based optical character recognition in the presence of deterministic transformations. 26(12):1813–1826.
- [2] Rodney A. Brooks. Intelligence without representation. 47(1):139–159.
- [3] Glenn L. Cash and Mehdi Hatamian. Optical character recognition by the method of moments. 39(3):291–310.
- [4] Adriano Cruz and NCE Mestrado. Anfis: Adaptive neurofuzzy inference systems. IM, UFRJ, Mestrado NCE, 2009.
- [5] Jyh-Shing Roger Jang. ANFIS: adaptive-network-based fuzzy inference system. 23(3):665–685.
- [6] R. Jean-Marie. *Method and apparatus for converting bitmap image documents to editable coded data using a standard no-tation to record document recognition ambiguities.* Google Patents. US Patent 5,359,673.
- [7] Bhagwandas Pannalal Lathi. Modern Digital and Analog Communication Systems 3e Osece. Oxford university press.
- [8] Shunji Mori, Hirobumi Nishida, and Hiromitsu Yamada. Optical character recognition. John Wiley & Sons, Inc.
- [9] Mrs VV Patil, Rajharsh Vishnu Sanap, and Rohini Babanrao Kharate. Optical character recognition using artificial neural network.
- [10] GANESH S. SABLE and SHEETAL ARUN NIRVE. OP-TIMIZATION OF OPTICAL CHARACTER RECOGNI-TION FOR PRINTED DEVANAGARI TEXT USING AN-FIS TECHNIQUES.
- [11] Simon Tong and Daphne Koller. Support vector machine active learning with applications to text classification. 2:45–66.
- [12] Sun-Chong Wang. Artificial neural network. In *Interdisciplinary Computing in Java Programming*, pages 81–100. Springer.