# Handwritten Malayalam Character Recognition System using Artificial Neural Networks

Vaisakh V K
*Dept. of ECE*
*NIT Calicut*
Kerala, India
Email: vaisakhvkofficial@gmail.com

Lyla B Das, member of IEEE
*Dept. of ECE*
*NIT Calicut*
Kerala, India
Email: lbd@nitc.ac.in

*Abstract*—This paper presents a system that recognizes Malayalam in handwritten form using artificial neural networks. A system for recognizing handwritten text (HCR) is a technique that is used for recognizing human handwritten text in any language. HCR is one of the research areas of recognition of patterns, which is still very challenging as perfect solutions have not yet been found. For certain foreign languages like English, Japanese, Chinese, etc, HCRs have been developed, which are reasonably good. But it is still premature for languages in India, especially for languages in south India. Because of the large character set, compound characters, presence of modifiers, and the curvature of characters in these languages, the task is quite complicated. This project aims to convert the photograph containing handwritten script into corresponding text. In this approach a trained ANN is used to identify the handwritten characters. The recognition system has been developed in python. The OpenCV library is used for performing different operations on the input image.

This paper pertains to the first part of a work where individual characters alone are recognized. The continuation of the work which is ongoing, is to recognize complete sentences.

*Index Terms*—Malayalam, handwritten recognition, Neural Networks, Deep Learning, Python, OpenCV

## I. Introduction

We all know that we are living in a world of automation, in that artificial intelligence (AI) is the fuel. The developments are happening in the field of AI, and deep learning is very rapid. There are many fields that are utilizing this development in AI.The capability of a computer to understand handwritten matter obtained from different sources is generally considered as a system for recognition of characters. In the areas of Image processing and associated ML systems, such an area of research is a very prominent one in which hundreds of research papers are being published as new methods come up. It provides an immense contribution to the development of an automation process and helps to close the man-machine interface gap. There are mainly two kind of handwritten systems are there. One is online HCR and second one is offline HCR. In online HCR system letters are recognized while the user is writing. It also contain the information like the order in which the user writes. But in an offline recognition system the only input is the image. It becomes even more difficult when the people has different style of writing.

## II. Literature survey

An English handwritten recognition system using deep neural networks by Surya Gunawan et al [1] has recognized words and digits using the MNIST dataset.This paper is by Surya Gunawan et al. Detection accuracies for digits are found to be 97.7%, and for characters is 88.8%. They used stacked encoders for handwritten recognition. The number of layers in the network is three, which includes two hidden layers and one softmax layer, and these layers are stacked together to create the network.

In another English handwritten recognition system by P. M. Pimpale,,S. Satra, D. Trivedi and R. Vaidya [2], the NIST dataset is used for recognition. The network model they used is a convolutional neural network. They used TensorFlow to build the neural network. They recognized the letter with an accuracy of 94%.

Aiquan Yuan proposed [3] a novel English handwritten recognition system. Which is a segmentation based system, and the segmentation is done using a modified online technique. Later a convolutional network is introduced for the recognition part. They used UNIPEN lower case dataset is used as the database and recognized with an accuracy of 92%.

K. Dutta et al [4] developed data-set named IIIT-HW-DB is used for Devanagari handwritten recognition. A mixed SCNN-RNN is used for the recognition part.

An Offline recognition of Malayalam handwritten Text [5] by Shajana Cand, Ajay james proposed a method to convert Malayalam handwritten text to editable text format. The line segmentation is done using a variation of horizontal projection method. Using vertical projection method words are converted to characters. Then for classifying the letters they used SVM classifier, and they were able to achieve accuracy rate upto 82%

Pranav P Nair proposed [6] a Malayalam character recognition of only 6 characters using LeNet. G Raju et al. [7] proposed a recognition system for malayalam using feature gradients and the count of Run lengths. The authors have proposed another character recognition scheme, for the recognition of isolated Malayalam characters, they fusing global and local features.

## III. THE PROPOSED SYSTEM

In this part, the character recognition model is explained. The Malayalam language consists large character set and different characters having similar features. The Malayalam language consists of 15 vowels, 36 consonants and five pure consonants and 12 dependent vowels. Hence differentiating Malayalam characters is a challenging task. The proposed system consists of following stages as shown in Fig 1.
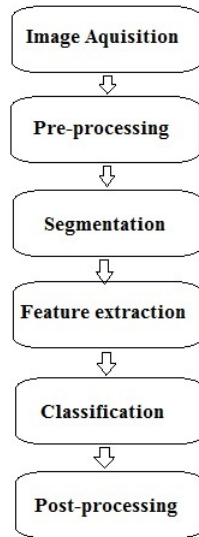


Fig. 1. System proposed – Schematic.

### A. Image acquisition

In this section the system gets the scanned image as input image. The input image must be acquired in specified format such as JPEG, JPG or PNG. The image is acquiesced in python using openCV as a JPG file.

### B. Preprocessing

The first step performed on the image is preprocessing. In this stage, the image is processed using standard techniques. For removing the noise in the image, the denoising operator in the Opencv is applied on the image. Then the image is converted to a binary image to get the region of interest. For that, an adaptive thresholding technique is used.

### C. Segmentation

For any image in this application, segmentation is an important step. Segmentation is used to locate the object and boundaries. In this stage, the image is segregated into individual characters. Contours are drawn around letters and cropped out for recognition.

### D. Feature extraction and Classification

**Dataset** The major backbone of any handwritten character recognition is a good database. There is no public domain Malayalam character database available.

Recently a team of research scholars developed an open-source handwritten Malayalam character image database named Amrita_MalCharDb. This data was accumulated by using handwritten material from 77 people in the ages between 20 and 55, whose native language is Malayalam. In this data set, the training data is of 59 writers and testing is done using the remaining images. We have taken this set as our data set.

**Amrita_MalCharDb:** This dataset [8] contains 85 Malayalam character classes which has the basic characters plus the modifiers associated with the language, See Fig 3. The dataset is saved along with the labels in CSV files. In each row, the values of the image pixels and the associated class label is available. In the first row, we see the label of the character class, and the rest of the columns represent the image in the form of vectors (a $32*32$ image has been converted to a vector of length 1024). There are 17236 samples for training and 6630 images for testing. Fig 2 illustrates the characters of this dataset.



Fig. 2. Malayalam character classes included in the current database

### Artificial Neural networks

For feature extraction and classification of the letters, an artificial neural network is designed and implemented. A Convolutional Neural Network (CNN) is implemented in python. This model includes convolutional layers, maxpooling layers and a flattening layer, as shown in the Fig 4. The CNN is built as a sequential model, which works sequentially and will focus on one layer at a time.
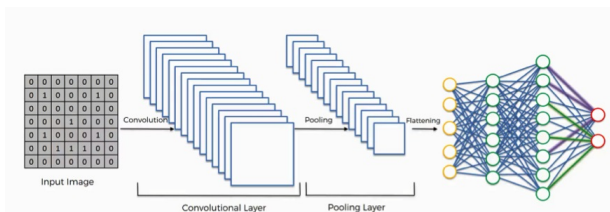


Fig. 3. Model of the artificial neural network

CNN's are special kinds of neural networks. It has mainly two functional parts. One is feature extraction, and the next one is classification. Each section consists of multiple layers. For feature extraction the system uses combination of Convolutional and Maxpooling layers as shown in Fig 5. Using these learned features, the letter is classified into the corresponding class by the dense layer. The dense layer is also known as a fully connected layer. In addition to this, there is an input layer

that accepts the input image and output layer, which contains the classes of characters that we are trying to classify.
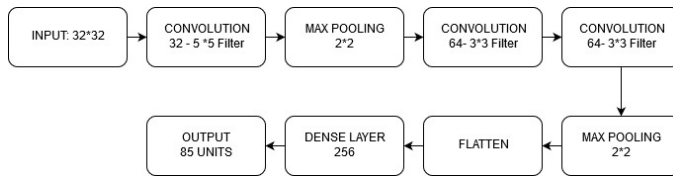


Fig. 4. CNN structure

For feature extraction, the system contains three convolutional layers and two max-pooling layers. The last max-pooling layer output is fed to the flatten layer and then to the dense layer for classification of the letter, as shown in Fig 4.

The shape of input layer is exactly same as the input image. In this case, it is 32*32=1024. In the hidden layers, the Relu activation function is used to find the compact representation of input images. At the end, the output units contain a vector of size 85, which includes the probability matrix of values between 0 and 1 and the unit that has a maximum probability value will determine the class label. Then the corresponding recognized letter will be written to the text file.

## IV. TRAINING AND TESTING

The model is trained using the training dataset containing 17236 samples and tested using the testing dataset containing 6230 samples. After training, the final model is saved.. The model's accuracy is increased with the number of epochs, as shown in Fig 5.
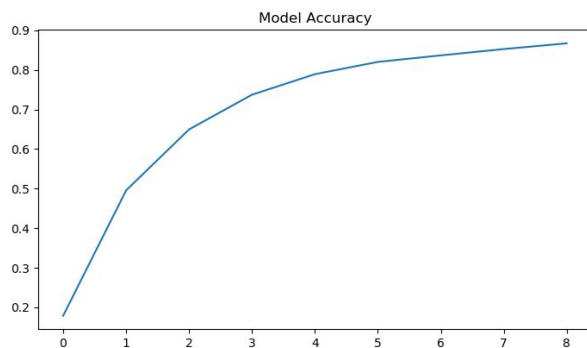


Fig. 5. Model accuracy with no of epochs

## V. RESULTS

The individual precision, recall, and f1-score are varying for each letter from 0.8 to 1, as shown in the Table 1 (confusion matrix) . The total accuracy of the classifier is 91%.

Very good recognition results have been obtained as shown in Fig 6.

|    | precision | recall | f1-score | support |
|----|-----------|--------|----------|---------|
| 0  | 0.00      | 0.00   | 0.00     | 0       |
| 1  | 0.95      | 0.88   | 0.91     | 60      |
| 2  | 0.87      | 1.00   | 0.93     | 102     |
| 3  | 0.93      | 0.93   | 0.93     | 58      |
| 4  | 0.93      | 0.91   | 0.92     | 69      |
| 5  | 0.89      | 0.95   | 0.92     | 62      |
| 6  | 0.98      | 0.89   | 0.93     | 56      |
| 7  | 0.92      | 0.80   | 0.86     | 56      |
| 8  | 0.90      | 1.00   | 0.95     | 64      |
| 9  | 0.92      | 0.90   | 0.91     | 63      |
| 10 | 0.88      | 0.85   | 0.86     | 105     |
| accuracy |     |        | 0.91     | 6360    |

Fig. 6. Confusion matrix

### A. Real Time Testing

To test in real time, an input image is given to the system for recognition, as shown in Fig8. The image was read and a bounding box was drawn around the letter, as shown in Fig 8. After getting a region of interest, the letter is cropped, as shown the Fig 9. Then the image is binarized and converted to a 32*32 image, as shown the Fig 10. The resulted 32*32 image is applied to the classifier to predict the character, which outputs a probability vector of size 85. The probability vector finds the maximum probability and identified the corresponding letter from the stored file and prints that letter to the output text file, as shown in Fig 11.
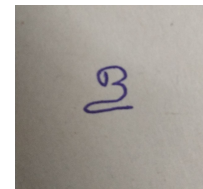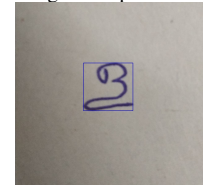


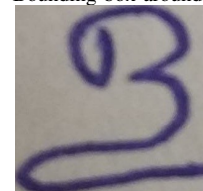Fig. 7. Input letter



Fig. 8. Bounding box around the letter
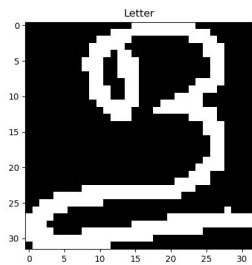


Fig. 9. Cropped image

Fig. 10. The input image for the classifier

## B. Output text file

The corresponding letter is printed in the text file as shown in the Fig 12.
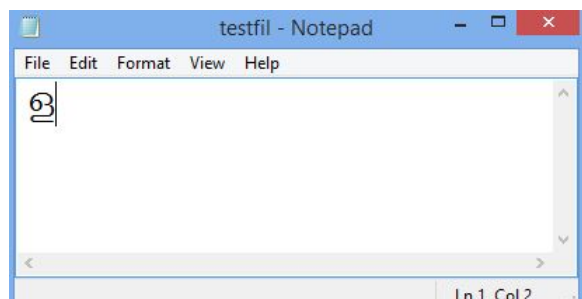


Fig. 11. The generated output letter

## VI. CONCLUSION AND FUTURE WORK

There is a wide variety applications for HCR. Office automation is one of them. In classification and recognition of images, CNN gives the best results. As of now the system is recognizes Malayalam characters only. It is proposed to continue this work to identify words, sentences and paragraphs.

## ACKNOWLEDGEMENT

## REFERENCES

[1] T. Gunawan, A. Noor, and M. Kartiwi, "Development of english handwritten recognition using deep neural network," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 10, pp. 562–568, 05 2018.

[2] R. Vaidya, D. Trivedi, S. Satra, and P. M. Pimpale, "Handwritten character recognition using deep-learning," in *2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT)*, April 2018, pp. 772–775.

[3] A. Yuan, G. Bai, P. Yang, Y. Guo, and X. Zhao, "Handwritten english word recognition based on convolutional neural networks," in *2012 International Conference on Frontiers in Handwriting Recognition*, Sep. 2012, pp. 207–212.

[4] K. Dutta, P. Krishnan, M. Mathew, and C. V. Jawahar, "Offline handwriting recognition on devanagari using a new benchmark dataset," in *2018 13th IAPR International Workshop on Document Analysis Systems (DAS)*, April 2018, pp. 25–30.

[5] C. Shanjana and A. James, "Offline recognition of malayalam handwritten text," *Procedia Technology*, vol. 19, pp. 772–779, 12 2015.

[6] P. P. Nair, A. James, and C. Saravanan, "Malayalam handwritten character recognition using convolutional neural network," in *2017 International Conference on Inventive Communication and Computational Technologies (ICICCT)*, March 2017, pp. 278–281.

[7] E. E. Swartzlander, "Truncated multiplication with approximate rounding," in *Proc. 33rd Asilomar Conf. Signals Syst. Comput.*, vol. 2, Oct 1999, pp. 1480–1483 vol.2.

[8] K. M. Manjusha K and S. Kp, "On developing handwritten character image database for malayalam language script," *Engineering Science and Technology, an International Journal*, vol. 22, 02 2019.