# Intelligent Character Recognition of Handwritten Characters

## Kailash Kangne[1]

[1]Eight Semester, Information Science and Engineering, The National Institute of Engineering, Mysore

---***---

**Abstract -** *The goal of this work is to identify printed text in image or handwritten characters using a camera or scanning equipment by converting them into machine-readable text with the help of a neural network. This research is different from traditional OCR in fact that it provides accurate results as it is using supervised machine learning which helps in training the system well. The system input is digital images which contains some patterns to be classified. The analysis and recognition of the patterns in images are becoming more complex, yet easy with advances in technological knowledge. It can handle multiple language scripts in its dataset and provide results accurately up to 95% for alphanumeric characters.*

**Key Words:** intelligent character recognition, supervised machine learning, dataset visualization, feature extraction.

## 1. INTRODUCTION

Optical character reader (OCR) is a device which electronically converts images of handwritten or printed text into machine-readable text, maybe it from a scanned document, a nameplate, a photo or from subtitle text superimposed on an image.

It is not possible to use handwritten and printed documents that are either scanned or photographed as they are in an image format that proves to be a lot of hindrances and involves hard work to convert them into data editable for computers. This conversion can be costly for organizations that handle large amounts of data and in the resulting documents can also generate huge amounts of human errors. When labour is scarce and efficiency needs to be increased, a computer-aided conversion of such optical documents can prove. In the past, people struggled with data entry and digitization. Due to changes in data type, existing techniques to overcome these issues such as OMR sheets and form data have been outdated. A neural network-based approach is necessary to convert data of varied types into digital information to overcome these problems. Many organizations have spent a lot of money, time and work to convert present information on a paper into computer data. Using a trained tool, this small task can be easily accomplished with little or no work, making the company focus its resources and valuable workforce on its challenges. If done on similar types of documents, this computerized conversion can be very productive, saving valuable time.

## 2. RELATED STUDIES

Researchers are very concerned about OCR as they aim to develop a robust and efficient recognition system. Study was conducted to develop associative pattern matching technique and match algorithms with accurate and fast intelligent character recognition system. There have been approaches related to microscopic image texture characters related to surface roughness. This proves that with the roughness of the surface, different parameters change regularly. Efforts have been made to capture characters from paper forms.

Research has also been done in license plate recognition. Nevertheless, developing a system where characters are automatically transferred to excel sheet after testing is challenging. Segmentation is an important phase in character recognition and at this stage efforts have been made to improve accuracy also discuss segmentation in which trained classifier is used to segment the connected characters into text-based CAPTCHAs. With regard to machine learning, SVM has been used in modified freeman chain code, which is used to extract features using SVM classifiers. As many performance parameters depend on in, the algorithm is very necessary. This gives a high-performance algorithm as well as low memory consumption. Matching templates and extracting features using back propagation algorithm largely focuses on the noise removal aspect. There, algorithms are applied after transferring data to excel sheet or creating a CSV file. As previously stated, printed text is relatively easier to recognize, but handwritten text is difficult due to changes in parameters of writing. Therefore, a lot of work has been done previously to address the issues facing the recognition of handwriting, but no such efforts have been made to apply supervised learning algorithms and transfer data simultaneously to excel sheet after python programming.

## 3. ARTIFICIAL NEURAL NETWORK

It is the most widely used techniques for classification purposes. However, other techniques such as Naive Bayes, Support vector machine, rough set are also popular tools for different classification context. The suitability of artificial neural network (ANN) as a smart technique lies in its attempt to imitate human effort. Because of its vast capacity to process information, ANN gives better results for complex pattern recognition tasks. Neural networks are therefore preferred to other techniques, particularly for pattern recognition tasks such as identification of offline text.

---

In principle, ANN is defined by three parameter types: 1.The pattern of interconnection between neuron layers 2. The learning process to update interconnection weights 3. Activation features that map a weighted input of a neuron to its activation of output.

Typically, two steps are involved in the process. The first step is the training phase in which the correct data class is identified. The extracted features would serve as the data for the neural network training. Once the NN is trained, by extracting features from new characters, new characters can be classified into known classes. For this purpose, which is driven by the principle of back propagation algorithm, multi-layer feed forward neural network is mostly used.

## 4. PROPOSED SYSTEM

This project approaches the problem of optical character recognition in two stages, the first being training stage in which the system is provided with handwritten characters of all alphabets in ascending order and then prediction stage in which an input image is given that is split into characters and each is individually predicted by the trained neural network. In both stages characters are processed by using a gradient change based feature extraction which involves slicing individual characters into multiple equally sized zones and extracting features such as starters, intersections and end points as well the lowest and highest points in every zone. The zones are then converted into feature vectors which gives the values of all the features in the zones. These feature vectors are used for training as well as during prediction during comparison.
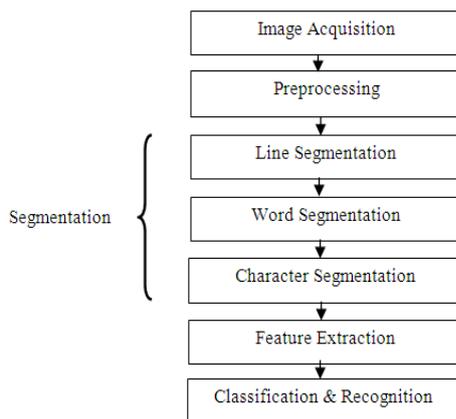


Fig 1: Block diagram of system.

## Document Accessing Module:

Input: Documents that contain handwritten or printed text.
Process: Using a scanner or digital camera, the submitted documents are read and converted into a PNG image. Also, the documents are labelled either as training or test data or as actual data. The naming conventions followed are' TR' or' TE' or' AC' respectively, followed by a number for training, testing and actual data.
Output: The acquired set of images are stored into the directory system.

## Pre-processing Module:

Input: Images stored in the directory
Process: Regardless of the image, it must be free from noise caused by degradation, physical damage and also from the scanning equipment boundaries. By padding or removing borders, the image is converted into a square image, then the image is binarized to facilitate processing. Averaging is also done by taking multiple scans of the same image to improve the quality of the image.
Output: Image ready for mathematical operations.

## Feature Improvement Module:

Input: Binarized Images.
Process: While binarization produces good results for printed text images, further processing of manuscript images is still required. There are still holes within an image's particles that must be filled and very small text dilated in order to detect them.
Output: Images with no holes and better contrast.

## Feature Extraction module:

Input: Detailed Binary Images
Process: The image characteristics here are recognized with spaces between words and each word is divided into characters. The characters are further segmented into several sub-functions, which then become a function vector. Elements can be form, slope, height, width, etc.
Output: Features of an image as vectors.

## Artificial Neural Network:

Input: Feature Vectors.
Process: The feature vectors are fed into the feed forward network's input layer, which uses the cluster log sigmoid weight function to classify the image features and determine the confidence level for each output node for a given feature. The trained data set is stored in the directory along with its probabilities and confidence levels.
Output: Trained Neural Network.
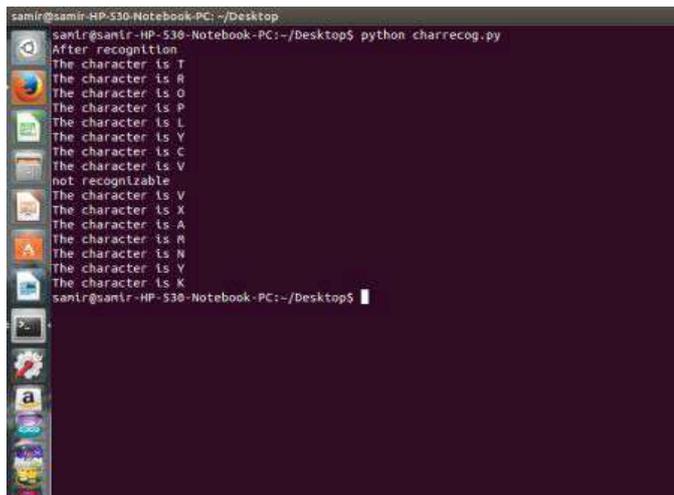
## User Input module:

Input: New Images
Process: Using this module, the user uses the trained neural network to identify the images he wants. Open a picture and send it for prediction to this module. The trained node neural network is not able to find the characters in the image.
Output: Recognized characters.

## 5. RESULTS

The captured image is scanned to detect the color and background of the characters. To reduce all irrelevant information, pre - processing of the image is necessary. Our research is based on alpha numerals and special characters being recognized.



Fig. 2. Results showing almost 95% accuracy

The accuracy was about 75 percent after testing as many as 100 samples. But the accuracy increased significantly after moving to nearly 400 samples and moved to nearly 95 percent.

The figure (Figure 2) shows 16 sample test results, out of which 15 samples show the correct character and one of the 16 characters cannot be recognized. Once again the recognizable characters are transferred to the Excel sheet.

## 6. CONCLUSION

Optical characters can be recognized in multiple ways, but using something that mimics the nature and character of the human brain is bound to succeed better than any other system. The proposed method uses neural network approach to the optical character recognition task, but using features rather than a wholesome comparison, which is inappropriate for handwritten documents, is not an effective way. In the future, the method of extraction of features may change to suit the needs of multiple languages. Such a system can also be used in robot vision as well as in understanding human written information. Such a work could be extended in the future so that it could be used to read hieroglyphics and other ancient languages.

## REFERENCES

1) Mrs. R. Rani, Prof. R. Dhir , G. Sinha, Recognition Using Zonal Based Feature Extraction Method and SVM Classifier , International Journal of Advanced Research in Computer Science and Software Engineering , ISSN: 2277 128X, Volume 2, Issue 6, June 2012.

2) V. Narawade, M. Patil, Recognition of Handwritten Devanagari Characters through Segmentation and Artificial neural networks , International Journal of Engineering Research and Technology, ISSN:2278-0181,Vol. 1 Issue 6, August 2012.

3) D. Bhattacharjee, M. Nasipuri, S. Arora, L. Malik and B. Portier, A Two Stage Classification Approach for Handwritten Devanagari Characters.

4) Fan Xiaoping, Jian, Zhang, and Huang Cailun. "Research on characters segmentation and characters recognition in intelligent license plate recognition system." Control Conference, 2006, Chinese. IEEE, 2006.

5) Ogden, Phil. "Applying intelligent character recognition in the 'real world'." Document Image Processing and Multimedia , IEE Colloquium on. IET, 1999.

6) Hussain, Rafaqat, "Recognition based segmentation of connected characters in text based CAPTCHAs." Communication Software and Networks , 2016 8th IEEE International Conference on. EEE, 2016.

7) Verma, Vivek Kumar, and Pradeep Kumar Tiwari. "Removal of Obstacles in Devanagari Script for Efficient Optical Character Recognition." Computational Intelligence and Communication Networks (CICN), 2015 International Conference on. IEEE, 2015.

8) Khaustov, P. A., V. G. Spitsyn, and E. I. Maksimova. "Algorithm for optical handwritten characters recognition based on structural components extraction." Strategic Technology (IFOST), 2016 11th International Forum on. IEEE, 2016.

9) Meher, Sukadev, and Debasish Basa. "An intelligent scanner with handwritten odia character recognition capability." Sensing Technology, 2011 Fifth International Conference on. IEEE, 2011.

10) T.M.Rath and R. Manmatha(2007), "Word spotting for historical documents" , International Journal on Document Analysis and Recognition , Vol.9, No 2 – 4, pp. 139- 152.

11) Ragot N., Salah, A. B., and Paquet, T. (2013). "Adaptive detection of missed text areas in OCR outputs: Application to automatic assessment of OCR quality in mass digitization projects" Proc. SPIE8658, Document Recognition and Retrieval XX, 865816.

12) T. K. Das (2016). Intelligent Techniques in Decision Making: A Survey, Indian Journal of Science and Technology, Vol. 9, No.12,pp.1-6.