



Tamil Handwritten Text Recognition using Convolutional Neural Networks

Deepa. M¹, Deepa. R², Meena. R³, Nandhini. R⁴
Assistant Professor¹, Student^{2,3,4}

Meenakshi Sundararajan Engineering College, Chennai, India

Abstract:

Handwritten character recognition has been one of the active and challenging research areas in the field of image processing and pattern recognition. It has numerous applications which include, reading aid for blind, bank cheques and conversion of any handwritten document into structural text form. In this paper an attempt is made to recognize handwritten characters for Tamil alphabets without feature extraction using multilayer Convolutional Neural Networks. Each character data set contains Tamil alphabets. The trained network is used for classification and recognition. In the proposed system, each character is resized into required pixels, which is directly subjected to training. That is, each resized character has predetermined pixels and these pixels are taken as features for training the neural network. The results show that the proposed system yields good recognition rates which are comparable to that of feature extraction-based schemes for handwritten character recognition.

I. INTRODUCTION:

Optical Character Recognition systems have been effectively developed for recognizing the printed characters of many non-Indian languages like English, Chinese, French. Now various efforts are on the way for the development of efficient systems for recognizing the Indian languages, especially for Tamil, a south Indian language widely used in Tamilnadu, Pudhucherry, Singapore, Srilanka. The process of character recognition of any script can be broadly broken down into three stages; pre processing, text extraction, classification. Typical text extraction includes a collection of operations that apply successive transformations.

II. EXISTING SYSTEM:

The existing systems are based on character recognition using feature extraction techniques.

- Firstly, approaches that can recognize the segmented text by proposing their own features with classifiers training which works well for specific languages and specific data.
- Secondly approaches that can recognize and binarize the text without segmentation of text lines using multiple hypothesis frames work.
- Thirdly approaches that can improve recognition rate by enhancing the text through binarization.

DISADVANTAGES:

- To overcome the problems low light, reflection and broken text etc.
- To overcome the main challenges associated with the natural scene images like complex background, different font styles of the text, sizes of the text and orientation of the text.

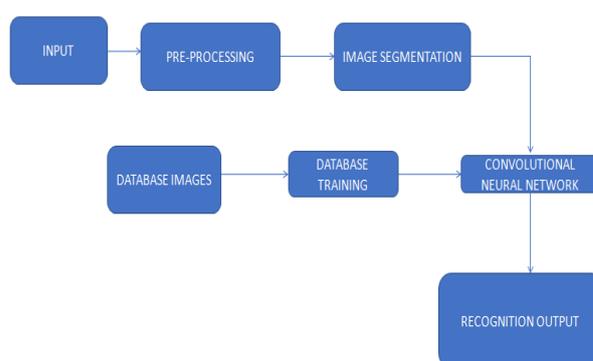
III. PROPOSED SYSTEM:

- ❖ The proposed approach is robust to different kinds of text appearances, including font size, font style, color, and background.
- ❖ Combining the respective strengths of different complementary techniques and overcoming their

shortcomings, the proposed method uses efficient character detection and classifier training based on neural networks.

- ❖ Neural Networks are used to achieve the highest accuracy in the lowest processing time possible.

IV. BLOCK DIAGRAM:



V. PRE-PROCESSING:

Preprocessing is the first step in the processing of scanned image. The scanned image is checked for noise, skew, slant etc. There are possibilities of image getting skewed with either left or right orientation or with noise such as Gaussian. Here the image is first convert into grayscale and then into binary. Hence we get image which is suitable for further processing. After pre-processing, the noise free image is passed to the segmentation phase, where the image is decomposed into individual characters. The binarized image is checked for inter line spaces. If inter line spaces are detected then the image is segmented into sets of paragraphs across the interline gap. The lines in the paragraphs are scanned for horizontal space intersection with respect to the background. Histogram of the image is used to detect the width of the horizontal lines. Then the lines are scanned vertically for vertical space intersection. Here histograms are used to detect the width of the words. Then the words are decomposed into characters using character width computation. Feature extraction follows the segmentation phase of OCR where the individual image glyph

is considered and extracted for features. First a character glyph is defined by the following attributes like height of the character, width of the character. Classification is done using the features extracted in the previous step, which corresponds to each character glyph. These features are analyzed using the set of rules and labeled as belonging to different classes. This classification is generalized such that it works for single font type. The height of the character and the width of the character, various distance metrics are chosen as the candidate for classification when conflict occurs. Similarly the classification rules are written for other characters. This method is a generic one since it extracts the shape of the characters and need not be trained. When a new glyph is given to this classifier block it extracts the features and compares the features as per the rules and then recognizes the character and labels it.

The steps involved in pre-processing are

- Read input image.
- Image Resize.
- Grey scale conversion.
- Noise removal.
- Binary image conversion.

However, extensive research, it is not easy to design series general-purpose systems. This is because there many possible sources of variation when extracting text. Shaded from the textured background or, from the low-contrast or complex images, or images with variations in font size, style, color, orientation, and alignment. This variation makes the problem very difficult to draw automatically. Generally text-detection methods can be classified into three categories. The first one consists of connected component-based methods, which assume that the text regions have uniform colors and satisfy certain size, shape, and spatial alignment constraints. However, these methods are not effective when the text have similar colors with background. The second one consists of the texture based methods, which assume that the text regions have special texture. Though these methods are comparatively less sensitive to background colors, they may not differentiate the texts from the text-like backgrounds. The third one consists of the edge-based methods. The text regions are detected under the assumption that the edge of the background and the object regions are sparser than those of the text regions. However, this kind of approaches is not very effective to detect texts with large font size compared the Support Vector Machines (SVM) based method with the multilayer perceptrons (MLP) based one for text verification over four independent features, namely, the distance map feature, the grayscale spatial derivative feature, the constant gradient variance feature and the DCT coefficients feature. They found that better detection results are obtained by SVM rather than by MLP. Multi-resolution-based text detection methods are often adopted to detect texts in different scales.

VI. SEGMENTATION:

- Image segmentation is the process of partitioning a digital image into multiple segments.
- Image segmentation is a vital part of image analysis process. It differentiates between the objects we want to inspect further and the other objects or their background.
- Connected Component analysis: Once region boundaries have been detected, it is often useful to extract regions which are not separated by a boundary. Any set of pixels which is not separated by a boundary is call connected. Each maximal region of connected pixels is called a

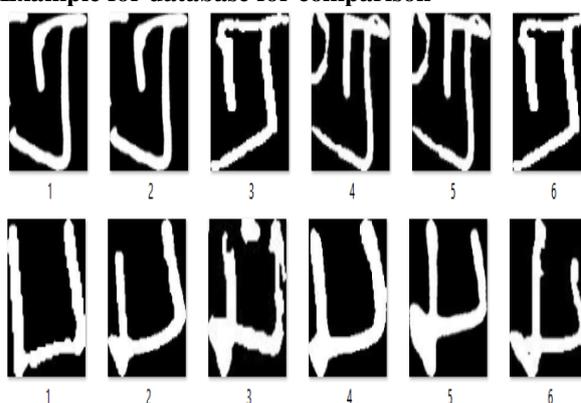
connected component. The set of connected components partition an image into segments.

VII. DATABASE:

- The database used here consist of 2 parts -

1. database for comparison
 2. output database
- The database used for comparison and has images of different letters in different styles of writing.
 - The output database consists of the images of letters to be displayed in the output.
 - The name of the images used should match with each other.
 - The images in the comparison database should be of same size.
- The images in output database can be of any varying size since it is only used for output display

Example for database for comparison



Example for ouput database



VII. CONVOLUTIONAL NEURAL NETWORK:

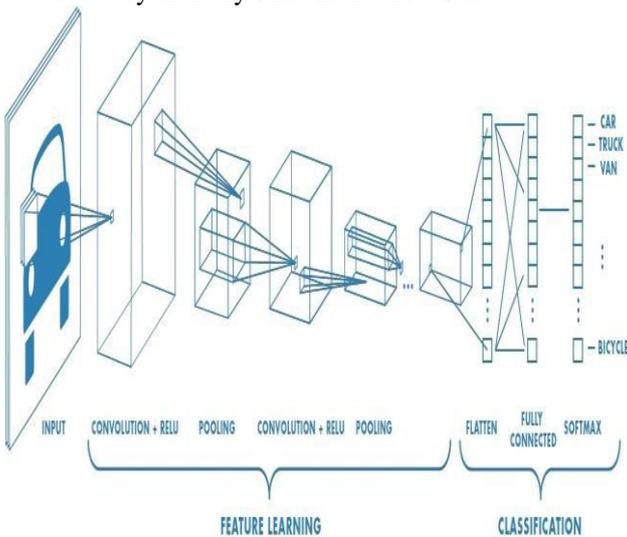
- Back propagation is a method used in artificial neural networks to calculate a gradient that is needed in the calculation of the weights to be used in the network.
- The disadvantage of Back Propagation Neural Network is, the features of the image should be extracted manually.
- A convolutional neural network (CNN) is a specific type of artificial neural network that uses perceptrons, a

machine learning unit algorithm, for supervised learning, to analyze data.

- CNNs apply to image processing, natural language processing and other kinds of cognitive tasks.
- CNNs use relatively little pre-processing compared to other image classification algorithms.
- Other learning algorithms or models can also be used for image classification. However CNN has emerged as the model of choice for multiple reasons.
- These include the multiple uses of the convolution operator in image processing,.
- The CNN architecture implicitly combines the benefits obtained by a standard neural network training with the convolution operation to efficiently classify images.
- Further, being a neural network, the CNN is also scalable for large datasets, which is often the case when images are to be classified.

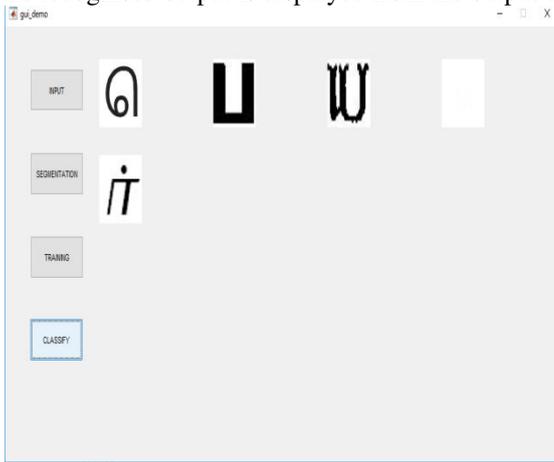
WORKING:

- CNN uses different kind of filters which is applied on the input 2-D image.
- It has a series of layers which performs the application of filters on the image.
- These results are convolved when they move from one layer to the other layer so that the results have high accuracy than any other neural network.



IX. RECOGNITION:

Thus the pre trained database is compared with the input image (after segmentation) by the convolutional neural network and the recognised output is displayed from the output database.



X. CONCLUSION:

A lot of research work exists in the survey for Tamil Handwritten recognition. However, there is standard solution to identify all Tamil characters with reasonable accuracy. Various methods have been used in each phase of the recognition process, Challenges still prevails in the recognition of normal as well as abnormal writing, slanting characters, similar shaped characters, joined characters, curves and so on during recognition process. In this paper, we have projected various aspects of each phase of the offline Tamil character recognition process. We have used minimal character set. Coverage is not given for different writing styles and font size issues. The following key challenges can be further explored in the future .A neural network based off line handwritten character recognition system without feature extraction has been introduced in this paper for classifying and recognizing the Tamil alphabets. The pixel values derived from the resized characters of the segmentation stage have been directly used for training the neural network. As a result, the proposed system will be less complex compared to the offline methods using feature extraction techniques. Of the several neural networks architectures used for classifying the characters, the one with two hidden layers each having 100 neurons has been found to yield the highest recognition accuracy of 90.19%. The handwritten recognition system described in this paper will find potential applications in handwritten name recognition, document reading, conversion of any handwritten document into structural text form and postal address recognition.

XI. REFERENCES:

[1]. Akshay Apte and Harshad Gado, “Tamil character recognition using structural features” , 2010

[2]. Banumathi P and Nasira G.M, “Handwritten Tamil Character Recognition using Artificial neural networks”, International Conference on Process Automation, Control and Computing (PACC), page(s): 1 – 5, 2011

[3]. Bhattacharya U, Ghosh S.K and Parui S.K, “A Two Stage Recognition Scheme for Handwritten Tamil Characters”, Ninth International Conference on Document Analysis and Recognition, Vol: 1 page(s): 511 – 515, 2007

[4]. Bremananth R and Prakash A, “Tamil Numerals Identification”, International Conference on Advances in Recent Technologies in Communication and Computing, page(s): 620 – 622, 2009

[5]. Hewavitharana S and Fernando H.C, “A Two Stage Classification Approach to Tamil Handwritten Recognition”, Tamil Internet, California, USA, 2002