

Handwritten Tamil Character Recognition Using ResNet

R. Jayakanthan¹, A. Hiran Kumar², N. Sankarram³, B. S. Charulatha⁴, Ashwin Ramesh⁵

^{1,2,5}Student, Department of Computer Science & Engineering, Rajalakshmi Engineering College, Chennai, India

^{3,4}Professor, Department of Computer Science & Engineering, Rajalakshmi Engineering College, Chennai, India

Abstract: This paper describes the recognition of Tamil language characters by using Residual Neural Network (ResNet). ResNet is a kind of Artificial Neural Network (ANN) which is typically consists of double or triple layers skips that has nonlinearities and batch normalization. Tamil is one among most traditional languages especially found on the southern regions of India. Recognizing a non-digitalized character in Tamil is a very difficult process because of its large and compound character set. In this paper we describe an offline character recognition approach using ResNet Architecture which explains the ResNet layers and the difficulty in gaining the digitalized character from the ResNet Layer using deep learning. This involves the data set training and pre-processing of all the Tamil characters. There are totally 256 characters in Tamil language in which most of the letters are almost similar and only a slight change can be seen only at the end for most of the characters, so recognizing a particular character is quite difficult and require highly trained datasets for distinguishing each and every characters.

Keywords: Image processing, Residual Neural Network (ResNet), Deep learning, Machine learning, Character recognition.

1. Introduction

Image Processing is an important method for character classification and recognition but the character classification is a very difficult part especially in Tamil language because of its character similarity, structure and shape. Character classification is the process of identifying the Tamil characters from the data set which are trained initially and input is written by different users. One of the major obstacle in handwritten character recognition is unavailability of standard databases for training and testing purpose. So for overcoming this situation Hp Labs have designed a database samples for nearly 156 class of Tamil character it is called as HPL-Tamil-Iso-Char. By using this we can able to generate sample dataset which almost matches the handwritten characters. Some of the images in the dataset are created.

Manually and some of them are derived using the HPL datasets. These data sets are freely utilized for research purpose and are available in UNIPEN format which is the downloadable offline version of sample set.

HPL-Tamil-Iso-Char contains up to 156 classes of characters and in each and every class there will be 5 to 10 samples which are created by different Tamil scholars across different cities of

south India. Approximately there are 500 samples available for a particular class and in worst case scenario up to 250 classes will be created and making the system ready for the image acquisition and pre-processing.

2. Literature survey

U. Bhattacharya, S. K. Ghosh and S. K. Parui “A Two Stage Recognition Scheme for Handwritten Tamil Characters” 2007 [1] India is a multilingual multiscrypt country with more than 18 languages and 10 different major scripts. Not enough research work towards recognition of handwritten characters of these Indian scripts has been done. Tamil, an official as well as popular script of the southern part of India, Singapore, Malaysia, and Sri Lanka has a large character set which includes many compound characters. Only a few works towards handwriting recognition of this large character set has been reported in the literature. Recently, HP Labs India developed a database of handwritten Tamil characters. In the present paper, we describe an off-line recognition approach based on this database. The proposed method consists of two stages. In the first stage, we apply an unsupervised clustering method to create a smaller number of groups of handwritten Tamil character classes. In the second stage, we consider a supervised classification technique in each of these smaller groups for final recognition. The features considered in the two stages are different. The proposed two-stage recognition scheme provided acceptable classification accuracies on both the training and test sets of the present database.

M. A. Pragathi, K. Priyadarshini, S. Saveetha, A. Shavar Banu, K. O. Mohammed Aarif “Handwritten Tamil Character Recognition Using Deep Learning” 2019 [2]-Character recognition is developed for various patterns of handwritten or optical characters to be recognized digitally. There are many Tamil literatures in undigitized form. Using deep learning the undigitized Tamil literatures can be converted into readable format. Many researches were carried on character recognition using deep learning for languages like Arabic, Devanagari, Telugu, etc... Due to the larger category set and confusion in similarities between handwritten characters Tamil character recognition is a challenge. In this paper, we propose a character recognition system for handwritten Tamil characters using deep learning. Here, VGG 16 approaches is carried out. The

proposed work gives efficiency of 94.52% on our datasets.

Manigandan T, Vidhya V, Dhanalakshmi V, Nirmala B “Tamil Character Recognition from Ancient Epigraphical Inscription using OCR and NLP” 2017 [3]- Recognition of ancient Tamil characters is one of the challenging task for Epigraphers as the language has evolved with different characters set. If the inscriptions are on stone walls, it adds even more complexity in identifying characters. This proposed work mainly focuses on recognition of various Tamil characters between 9th and 12th centuries using OCR and NLP techniques. In this work, the inscription images collected from Tamil Nadu, Archaeological Department are pre-processed and segmented. During the segmentation process the color images were converted to gray image and to binary image based on threshold value. From segmented, image features like number of lines, curves, loops and dots have been extracted using Scale Invariant Feature Transform (SIFT) algorithms for each letter to identify the exact character. Characters will be classified and constructed based on Vectors extracted, using Support Vector Machine (SVM) classifier and the patterns of the character will be matched with known characters and predicted using Trigram technique. Each identified character will be assigned with its corresponding Unicode value and it will be updated in the image corpus for further character identification, and to make the system in identifying the characters more effectively. Thus the proposed system can solve the major problems in reading the inscription images.

N. Prameela, P. Anjusha, R. Karthik “Off-line Telugu Handwritten Characters Recognition using optical character recognition” 2017 [4]- The Aim of the proposed paper is to recognize offline Hand written Telugu characters using Optical character recognition, OCR is one of the most popular and challenging topic of pattern recognition This paper proposes an OCR system for Telugu documents which comprises of three stages, namely pre-processing, feature extraction, and classification. In the preprocessing stage, we have employed median filtering on the input characters and applied normalization and skeletonization method over characters for extraction of boundary edge pixel points. In the feature extraction stage, initially each character is divided into 3x3 grids and the corresponding centroid for all the nine zones are evaluated. With this we can identify the characters of different styles. Thereafter, we have drawn the horizontal and vertical symmetric projection angel to the nearest pixel of the character which is dubbed as Binary External Symmetry Axis Constellation for unconstrained handwritten character. From which we have calculated the horizontal and vertical Euclidean distance for the same nearest pixel from centroid of each zone. Then we have calculated the mean Euclidean distance as well as the mean angular values of the zones. This is considered as the key feature values of our proposed system. Lastly, both support vector machine (SVM) and Quadratic discriminate Classifier (QDA) has been separately used as the classifier.

Chamila Liyanage, Thilini Nadungodage, Ruwan

Weerasinghe. “Developing a commercial grade Tamil OCR for recognizing font and size independent text” 2015[5]- Optical Character Recognition (OCR) of Indic scripts such as Tamil and Sinhala has lagged behind those for languages based on the Latin script. Several attempts to build commercial grade OCR for these languages have failed in the past owing to them not generalizing well. This paper describes a set of training regimes for Tamil using the Tesseract engine that have enabled us to develop a robust Tamil OCR system. We describe in detail our training regime, which results in a performance improvement of 12.5 % over the default Tamil module shipped with Tesseract on a set of ancient Tamil documents, which were part of an authentic project to digitize important Tamil manuscripts of Sri Lanka.

There are many different methodologies available for character recognition system. Some of the most commonly used methods are support vector machines(SVM), CNN using VGG-16, CNN using LeNet, CNN using GoogleNet and so on. This paper clearly explains the methodology of using Residual Neural Network(ResNet). Based on implement the different methodologies the efficiency and accuracy of system is analyzed. On overcoming all those disadvantages in the above methods, this working model can able to produce an accuracy of 96% which is very much higher when compared to implementation of other systems. The following table shows the accuracy parameters for different proposed systems.

Table 1
 Comparison of accuracies of different methodologies

S. No.	Methodology Implemented	Accuracy
1	Support Vector Machines(SVM)	82.04%
2	CNN-using VGG-16	94.52%
3	CNN using Googlenet	95.01%
4	Gabor Filters and SVM	92.5%
5	CNN using ResNet(Proposed System)	96.014%

3. Proposed system

A. Data Set

The initial image of the dataset will be irregular and the pixel quality of the picture is very low. The following sample shows the data batch of input images.

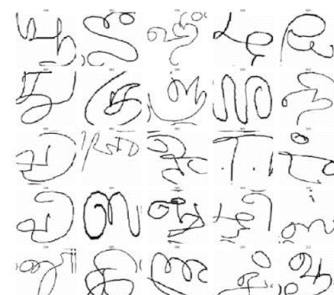


Fig. 1. Data batch of input images

B. Data Augmentation

Data augmentation is the important regularization process used when training a model or dataset in computer vision.

Generally, the system used to feed with similar images instead of feeding the same image a slight change in the pixel values of the image is done. In our working model after the pre-processing is done the images are fed into data augmentation process and thus producing a trained dataset of pixel value (128x128). The irregular pixel image is converted into 128x128 definite pixel images. The get-transform function is used in the operation of data augmentation.

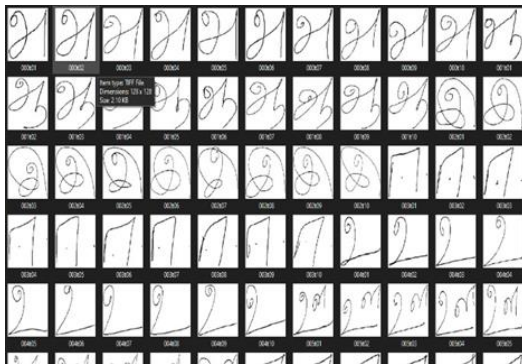


Fig. 2. Data set containing 128x128 images

4. Methodology

A single stage character recognition scheme is used in the case of small dataset problems but this system fails when there are many number of sample datasets. So, on overcoming this situation multi stage recognition scheme is utilized, which explain the different stages of processing an acquired image. The first stage of the system is image acquisition.

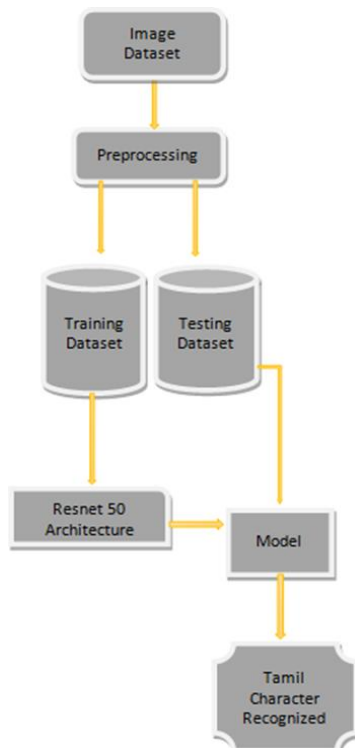


Fig. 3. Architecture of the proposed model

5. Image Acquisition

For the above proposed system, we have considered a standard database of Tamil character named as HPL-Tamil-Iso-Char database. In this database they had collected various 15000 numbers of images of character of different shapes and styles of Tamil characters. But in our case only 5000 images alone are considered for 156 Tamil characters.

A. Image Pre-processing

Image Pre-processing is the process of removing or reducing the unwanted noise present in the particular character image. These noises which are present in the input image affects the quality and efficiency of the character recognition system so we have to remove those. At first, the obtained input image with unwanted and noisy areas are cropped and neglected for creating a better framing. Then they are checked for skewing. Skewing is the process of aligning the image with proper coordination. Skewing is done by certain skew functions and detection checks. They check for an angle orientation of +/- 15 degrees in the actual image with the true horizontal lines. If there occurs any deviation certain skew angle is noted and for that particular skew angle a simple rotation of the Image is carried out till the lines meet the horizontal lines.

After skewing, for the skewed image the smoothing process is done. Smoothing is the process of removing the unwanted noise over the image. After smoothing, the image is then resized to the size of the actual image by resizing process. The processed image then passed into the residual neural network (ResNet) for recognizing the particular.

B. Training Dataset

ResNet is the key feature of the Tamil character recognition because it makes the system possible to train hundreds and thousands of layers and still acquires with the most efficient performance and results. Generally, most of the neural network train their datasets by means of back propagation. When more number of layers keep on adding or the depth of the layers increasing then their multiplicative derivative make the gradient very small, which implies that when more number of layers have been added in the ResNet there wont be any improvement in the overall performance.

For overcoming this situation ResNet uses “identity short connection” layer which does nothing but makes the system to skip the identical layer and reuses the activation layer from the previous layer. This main advantage in ResNet systems help the system to have only fewer layers and no layer can be found very deeper.

Multiple deep learning models and pre trained images are available in “torchvision.models” PyTorch. We can able to build the basic building blocks of ResNet using the PyTorch. PyTorch is based on the Torch library, it is an open source library which comes under the Machine Learning concept which is more specifically used in the fields of computer vision and natural language processing. It is a free and open source platform developed by facebook’s research lab.

epoch	train_loss	valid_loss	accuracy	time
0	1.242802	0.703643	0.786976	04.04
1	0.745869	0.384966	0.880555	04.09
2	0.502490	0.259989	0.919867	04.11
3	0.376871	0.197480	0.936027	04.12
4	0.302228	0.180675	0.941574	04.17

Fig. 4. Accuracy of datasets before

C. Transfer Learning

Transfer learning is the process of gaining knowledge over solving a problem then that particular gained knowledge is utilized to solve similar but different problem. This can be used in the convolutional network, because it is very rare to obtain all the datasets in a system, so the pre-trained network weights help in solving the problem through transfer learning. Freezing and finetuning are the important feature helps in transfer learning.

6. Fine Tuning

Fine Tuning is the process of giving a new data set in the pre-trained convolutional neural network to get a fine tuned data set. If a new dataset is almost similar to the old dataset for pre-training then same weights can be used to extract the features of new dataset. In this case if the new dataset is very much larger, then the whole network is retrained with initial basic weights from pretrained model.

epoch	train_loss	valid_loss	accuracy	time
0	0.539130	0.317683	0.897196	05.40
1	0.405956	0.225525	0.926017	05.55
2	0.262287	0.165592	0.947965	05.55
3	0.197185	0.133194	0.958276	05.56
4	0.158299	0.123576	0.960145	05.56

Fig. 5. Accuracy of datasets after unfreezing, fine tuning

7. Results

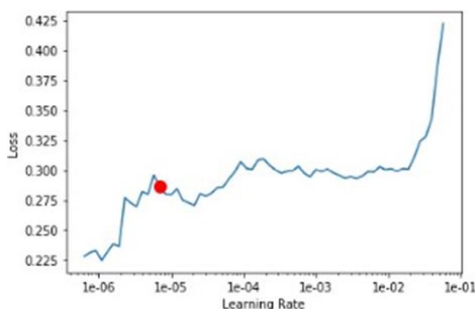


Fig. 6. Initial learning rate graph

In simple words learning rates are hyper-parameters used to determine the flow of the weights in the particular input image. Initially the input image which is being fed has many dynamic variations of weight. If the learning rate is too high derivative may miss the 0 slope point or if the learning rate is too low, then

it may take forever to reach that point.

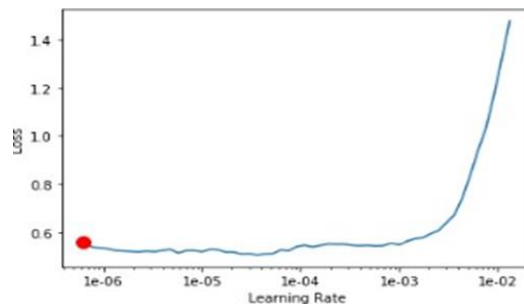


Fig. 7. Final learning rate graph

This dynamic variation in the above graph is trained to form a gradual and uniform increase in the weights of the learning rates and it also satisfies the stopping point.

Our proposed system for Tamil character recognition gives an accuracy of 96% and which is the current most efficient system utilized.

8. Conclusion and future work

We have presented a solution for Tamil character recognition system using ResNet-50 that includes a database, algorithm and an application. The dataset consists of more than 15,000 images each of 128x128 in dimensions. The dataset has many characters that are visually similar or written in a similar way by most of the people. Our proposed system for Tamil character recognition gives accuracy of 96%. This system could be implemented to design a complete handwritten documenting and digitizing system. The work can be extended to recognize the words in Tamil, and in future it will be merged with opensource Tamil database for recognizing all the 256 characters of the Tamil language thus becoming the repository for Tamil character recognition.

References

- [1] U. Bhattacharya, S. K. Ghosh and S. K. Parui "A Two Stage Recognition Scheme for Handwritten Tamil Characters" IEEE Transactions on Ninth International Conference on Document Analysis and Recognition, 2007.
- [2] M. A. Pragathi, K. Priyadarshini, S. Saveetha, A. Shavar Banu, K. O. Mohammed Aarif "Handwritten Tamil Character Recognition Using Deep Learning" IEEE Transactions on International Conference on Vision Towards Emerging Trends in Communication and Networking (ViTECoN), 2019.
- [3] Manigandan T, Vidhya V, Dhanalakshmi V, Nirmala B "Tamil Character Recognition from Ancient Epigraphical Inscription using OCR and NLP" IEEE Transactions on International Conference on Energy, Communication, Data Analytics and Soft Computing (ICECDS-2017).
- [4] N. Prameela, P. Anjusha, and R. Karthik "Off-line Telugu Handwritten Characters Recognition using optical character recognition" IEEE Transactions on International Conference on Electronics, Communication and Aerospace Technology, ICECA 2017.
- [5] C. S. Sundaresan, S. S. Keerthi, "A Study of Representations for Pen based Handwriting Recognition of Tamil Characters," IEEE Transactions on 1999.
- [6] Chamila Liyanage, Thilini Nadungodage, and Ruvan Weerasinghe. "Developing a commercial grade Tamil OCR for recognizing font and size independent text," IEEE Transactions on International Conference on Advances in ICT for Emerging Regions (ICTer), 2015.

- [7] K. B. M. R. Batuwita, G. E. M. D. C. Bandara “Fuzzy Recognition of Offline Handwritten Numeric Characters” 2006.
- [8] R. Ramanathan, S. Ponmathavan, N. Valliappan L. Thaneshwaran, Arun. S. Nair, K. P. Soman “Optical Character Recognition for English and Tamil Using Support Vector Machines” IEEE Transactions on International Conference on Advances in Computing, Control, and Telecommunication Technologies, 2009.
- [9] B. S. Charulatha, Paul Rodrigues, T. Chitrallekha, Arun Rajaraman “Mining Ambiguities using Pixel-based Content Extraction”, Advances in Intelligent Systems and Computing, pp. 537 – 544, 2016.
- [10] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, and K. Lang. Phoneme recognition using time delay neural networks. IEEE Transactions on Acoustics, Speech and Signal Processing, March 1989.
- [11] Shaohan Xu, Qi Wu, and Siyuan Zhang “Application of Neural Network in Handwriting Recognition” IEEE Transactions on International Conference of Stanford University.