

IMAGE CHARACTER RECOGNITION USING DEEP CONVOLUTIONAL NEURAL NETWORK LEARNED FROM DIFFERENT LANGUAGES

Jinfeng Bai, Zhineng Chen, Bailan Feng, Bo Xu

Interactive Digital Media Technology Research Center,
Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China
{jinfeng.bai, zhineng.chen, bailan.feng, xubo}@ia.ac.cn

ABSTRACT

This paper proposes a shared-hidden-layer deep convolutional neural network (SHL-CNN) for image character recognition. In SHL-CNN, the hidden layers are made common across characters from different languages, performing a universal feature extraction process that aims at learning common character traits existed in different languages such as strokes, while the final softmax layer is made language dependent, trained based on characters from the destination language only. This paper is the first attempt to introduce the SHL-CNN framework to image character recognition. Under the SHL-CNN framework, we discuss several issues including architecture of the network, training of the network, from which a suitable SHL-CNN model for image character recognition is empirically learned. The effectiveness of the learned SHL-CNN is verified on both English and Chinese image character recognition tasks, showing that the SHL-CNN can reduce recognition errors by 16-30% relatively compared with models trained by characters of only one language using conventional CNN, and by 35.7% relatively compared with state-of-the-art methods. In addition, the shared hidden layers learned are also useful for unseen image character recognition tasks.

Index Terms— deep convolutional neural network, image character recognition, multi-task learning

1. INTRODUCTION

Text information in images and videos, such as names of shops, road signs, billboard and superimposed text, usually offers reliable high level semantics which are useful in many applications. Although current commercial optical character recognition (OCR) systems have already reached high performances for document images, the recognition of text in images or videos (e.g. superimposed text) is still a challenging problem. This is due to image text often suffering several degradations from blur, uneven illumination, complex

background, perspective distortion and so on. Recently, deep convolutional neural network implemented on GPU unfold its potential in many fields to deal with challenging tasks. It has strong tolerance to shift, scale and distortion, therefore large improvement has been observed [1, 2, 3] on image character recognition compared to other methods. However, all these methods only focus on the single task-specific data and the feature extraction and transformation methods for different language are developed respectively. Common character traits existed in different language characters such as strokes, which portrays the distinguishing characteristic against other vision object, fail to utilize.

Motivated by this, this paper proposes a shared-hidden-layer convolutional neural network (SHL-CNN) for image character recognition. In SHL-CNN, the hidden layers are shared across characters of different tasks while the final softmax layer is tasks dependent. Specifically, the shared hidden layers (SHLs) are considered as a universal feature extraction process that trained based on characters of multiple languages. It thus emphasizes on learning common character traits existed in different languages such as strokes. The final softmax layer is trained solely based on characters of one destination language. It plays a role that differs the destination language with other languages. Under the SHL-CNN framework, image character data of different languages could be used together to learn more discriminative SHLs, and thus benefits the character recognition task for all languages.

Under the SHL-CNN framework, we discuss several issues including architecture of the network, training of the network. Based on these, a suitable SHL-CNN model for image character recognition is empirically learned. We evaluate the SHL-CNN with two tasks: English image character recognition and Chinese image character recognition, respectively on the famous ICDAR-03 dataset [4] and a Chinese video character dataset collected by ourselves, which contains 5036 text lines collected from 38 new videos. Experimental results show that the SHL-CNN reduce recognition errors by 16-30% relatively compared with models trained by characters of only one language using conventional CNN, and by 35.7% compared with state-of-the-art methods on the ICDAR-03 dataset.

The work was supported by National Nature Science Foundation of China(grant No.61202326, No.61303175)

In addition, the SHLs is also useful for unseen image character recognition tasks, in which the SHLs initialized CNN achieves better results than that whose weights of hidden layers are randomly initialized.

2. RELATED WORK

The motivation of this study is to use cross knowledge learned from multiple tasks to improve the performance of image character recognition. Since CNN has strong tolerance to shift, scale and distortion, it has shown good performance on image character recognition [1, 2, 3] and been applied to other recognition tasks [5]. Saidane [1] proposed a CNN based character recognition method that through recognizing character on color image directly, it outperformed the traditional binary character recognition method. On a same test set with [6, 7], Saidane's method got better result. Zhu [2] first find an optimum projection from color to grayscale conversion, and then using CNN to do scene character recognition. Jacobs [3] also proposed a CNN based character recognition method. Its CNN based character recognizer worked on grayscale character images and outperformed the binary character recognition based commercial OCR software on low-resolution documents captured by camera. Chinese image character recognition based on CNN has received relatively few concerns. We argue that this situation is mainly attributed to there is no public large labeled Chinese character database available. Although a few methods [8, 9] were proposed to deal with this tough problem, the thousands of categories, each corresponding to a Chinese character, as well as the tiny and varied differences among categories, prohibit the problem of Chinese character recognition far from addressed. However, the feature extractions and transformations of all these works are only learned from single and monolingual task. Common character traits existed in different language characters such as strokes, which portrays the distinguishing characteristic against other vision object, fail to utilize.

For cross tasks learning in the neural network, [10, 11] has been shown that a resource-limited tandem ASR system can benefit from the model trained with a resource-rich language under multilayer perceptions (MLPs) framework. Recently, Plahl et al [12] found out that resource-rich languages can also benefit from cross-lingual MLP features by changing the network topology to the bottleneck structure (BN). Huang [13] demonstrates that the benefit of out-of-language data is not limited to low-resource languages, and the degree of kinship between the source and the target language becomes unimportant if the neural network is powerful enough in the CD-DNN-HMM framework. Training procedure in [14] is performed in a sequential way. Firstly, the network is trained using one source language. Then, the top layer is substituted with the phoneme set for another source language for weight retraining. But Huang [13] shows that the simultaneous training strategy is superior to the sequential training strategy.

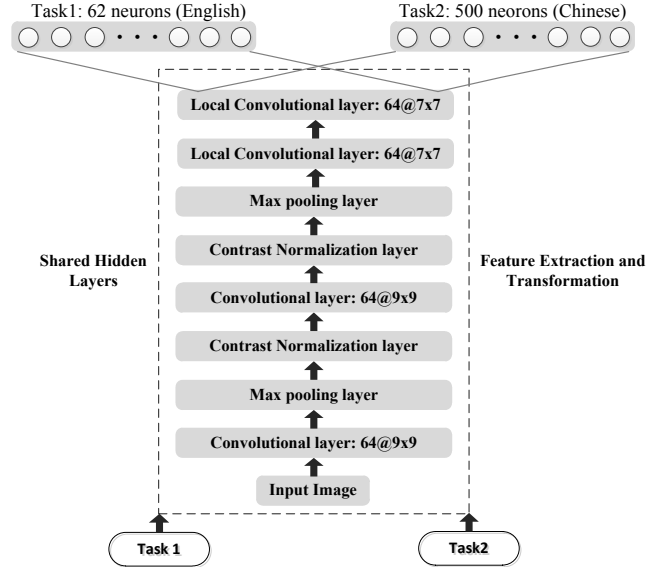


Fig. 1. Architecture of the shared-hidden-layer CNN

3. SHARED-HIDDEN-LAYER CNN

Fig. 1 depicts the architecture of the proposed SHL-CNN. In this architecture, the input and hidden layers are shared across all the tasks that the SHL-CNN can recognize, and can be considered as a common feature extractions and transformations module. The top softmax layer, however, is not shared. Instead, each task has its own softmax layer to estimate the posterior probabilities of the characters specific to that task.

3.1. Architecture of the SHL-CNN

The description of SHL-CNN architecture used for the experiments is given in the following way: 3x48x48-64C9-MP2-CN-64C9-CN-MP2-64L7-32L7-62N (500N) represents a net with 3 channel input images (i.e. color images) of size 48x48, a shared weight convolutional layer with 64 maps and 9x9 filters, a max-pooling layer over overlapping regions of size 3x3 at stride in 2 pixels, a local contrast normalization layer with size 9x9, a shared weight convolutional layer with 64 maps and 9x9 filters, a local contrast normalization layer with size 9x9, a max-pooling layer over overlapping regions of size 3x3 at stride in 2 pixels, an unshared weight convolutional layer with 64 maps and 7x7 filters, an unshared weight convolutional layer with 32 maps and 7x filters, and two fully connected output layers with 62 neurons (one per class) or / and 500 neurons. We use a logistic activation function for convolutional and fully connected layers, a linear activation function for max pooling layers and a softmax activation function for the output layer. All CNNs are trained using on-line gradient descent with an annealed learning rate. During training, images are continually translated, scaled, stretched and ro-

tated, whereas only the original images are used for validation. Training ends once the validation error is zero or when the learning rate reaches its predetermined minimum. Initial weights are drawn from a normal distribution with mean zero and standard deviation 0.01.

3.2. Training of the SHL-CNN

The key to the success of the SHL-CNN learning is to train the model for all the tasks simultaneously. The conventional training method of deep convolutional neural network is the classical stochastic gradient descent algorithm described in [15]. However, a slight modification has to be made for SHL-CNN, since an independent softmax layer is adopted for each different task. During the training phase, every training sample is only used to update gradient of the SHLs and the language-specific softmax layer (i.e. the softmax layer associated with target language of the training sample). Other irrelevant softmax layers are kept invariant. The SHLs serve as a universal feature extractor and a structural regularization to the model. The entire SHL-CNN and its training procedure can be regarded as an example of multi-task learning.

After being trained, the SHL-CNN can be used to recognize character of any task used in the training process. By sharing the hidden layers in the SHL-CNN and by using the joint training strategy, we can improve the recognition accuracy of all the tasks over the conventional CNNs trained using only the task specific data.

4. EXPERIMENTS AND RESULTS

To assess the performance of the proposed scheme, we train a SHL-CNN model for both English and Chinese image character recognition and conduct extensive experiments elaborated as follows.

4.1. Databases

We choose the ICDAR 2003 character database for the English recognition task, which a benchmark dataset used for the robust reading task of the ICDAR 2003 competition. The objective of the competition is to addressing text images that are rather challenge to current commercial OCR packages. The ICDAR 2003 character database is divided into three subsets: a train subset (containing 6185 images), a test subset (containing 5430 images), and a sample subset (containing 854 images). These character images are of different sizes (e.g., 5x12, 36x47, 206x223), fonts, colors, and present with varying kinds of distortion. For simplicity, Yokobayashi[6] merge the train and test set into a new train set, adopt the sample set as the test set and select only number and alphabet for use. Therefore, the new training set has 11492 samples including 6113 from original training set and 5379 from original test set. The new test set includes 698 images (bad samples are

discarded) and is further classified to seven groups according to the type of image degradations, as shown in the Fig. 2. In following experiments, we also adopt the classification scheme.

For Chinese image character recognition task, a news video character database is collected from 38 different Chinese TV news programs respectively as no public database is available. Each of them is about 30-minute long and the total time is about 50 hours. We have manually labeled text appearing in these videos including its position and content, eventually forming total 5036 text lines with 69340 characters. These characters include superimposed text, scrolling caption and a small quantity of scene text. These characters are ranked according to their frequency, and based on the ranking list, two character subsets, one formed by the top 500 categories and the other one formed by the next top 500 categories, are obtained. The two subsets are named as TV-500-1 and TV-500-2, which contain 49480 and 9997 occurrences of characters, respectively. Both subsets are divided into two equal parts, one for training and the other for testing.

4.2. SHL-CNN V.S. Con-CNN

Experiments are conducted on the ICDAR-2003 and TV-500-1 depicted before. The inputs to the conventional CNN (Con-CNN) and SHL-CNN are both 48x48 RGB color image. For ICDAR-2003 task, the output has 62 neurons including 10 numerals and 52 alphabets. For TV-500-1 task, 500 output neurons are included. Since the sizes of local receptive fields are sensitive for different tasks, three different models (listed in Table 1 below) are used in the experiment in order to testify the proposed scheme more comprehensively.

Table 1. Configuration of the models

Name	Configuration
Model-1	3x48x48-64C5-MP2-CN-64C5-CN-MP2-64L3-32L3-62N (500N)
Model-2	3x48x48-64C7-MP2-CN-64C7-CN-MP2-64L5-32L5-62N (500N)
Model-3	3x48x48-64C9-MP2-CN-64C9-CN-MP2-64L7-32L7-62N (500N)

Table 2 compares the recognition error rate (RER) obtained on the task specific test sets using the Con-CNN (trained using only the data from that task) and the SHL-CNN (whose hidden layers are trained using data from both two tasks). From the table we can observe that the SHL-CNN outperforms the Con-CNN with a 16-30% relative RER reduction (RR) across both the tasks for all models. Since the model structure and training procedure are same between the Con-CNN and SHL-CNN, we ascribe the gain of SHL-CNN to cross-task knowledge. Additionally, since most of samples come from superimposed text we can see that Chinese character recognition task has a low RER.

Group	Clear	Background Design	Multi-color	Uneven Light	Little Contrast	Blurring	Distortion
Number	199	130	54	40	37	210	28
Example							

Fig. 2. Classification and examples of ICDAR-2003 OCR sample database

Table 2. RER Comparison of Con-CNN and SHL-CNN

RER	ICDAR-2003			TV-500-1		
	Con-CNN	SHL-CNN	RR	Con-CNN	SHL-CNN	RR
Model-1	0.156	0.116	25.6%	0.053	0.037	30.1%
Model-2	0.140	0.098	30.0%	0.036	0.030	16.7%
Model-3	0.122	0.090	26.2%	0.029	0.023	23.1%

4.3. SHLs for Unseen Characters

The SHLs extracted from the SHL-CNN can be considered as an intelligent feature extraction module jointly trained with data from multiple character recognition tasks. As such they carry rich information to distinguish character classes in multiple tasks and can be carried over to distinguish characters in new tasks for image character recognition.

In this experiment, only the model behaving best performance before (i.e. 3x48x48-64C9-MP2-CN-64C9-CN-MP2-64L7-32L7-500N) is adopted. The procedure of training a model for new task TV-500-2 is simple. All SHLs from the SHL-CNN are extracted and an untrained softmax layer with 500 nodes is added on top of them. During training process, the weights of the SHLs keep intact and only the softmax layer are updated using training data from training set of TV-500-2. Using this method we achieve an RER of 3.6% on the TV-500-2 test set. For comparison, the experiment of the model randomly initialized weights without making use of SHLs from SHL-CNN are also performed. The training process is nothing different except that all weights are updated using training data. We only achieve an RER of 6.3% on the TV-500-2 test set.

According to the results, it indicates that the model leveraging the SHLs extracted from SHL-CNN are more effective than randomly initialized weights. An additional absolute 2.7% RER reduction has been observed by doing so.

4.4. Comparison with existing methods

In order to compare our system to the recent works in detail, we test the seven subset of new test set of ICDAR-2003 OCR database mentioned before respectively. Fig. 3 shows the results of our method compared with those of [6, 7, 1, 2]. The overall RER of our system reaches 9.03% compared to 14.04% of the best method [2] and it ranges from 32.14% for seriously distorted images to 5.03% for clear images. Com-

pared with the other methods, the performance of our system behaves better especially in cases such as multi-color, uneven light, little contrast and blurring. Since the large improvement were made by cross-task training, we attribute it to cross-task knowledge learned by proposed SHL-CNN, which may includes more variations knowledge of character images.

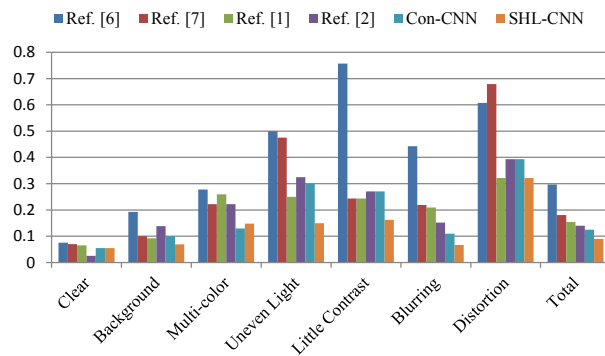


Fig. 3. Comparison of RER for each group of the test set

5. CONCLUSION

In this paper we firstly propose the SHL-CNN for image character recognition. With the idea that sharing hidden layers common across different tasks while keeping the final softmax layer task dependent, the learning of CNN benefits greatly from characters of multiple languages. Experiments on English and Chinese image character recognition tasks demonstrate that the SHL-CNN reduce errors on both the tasks by 16-30% relatively, over the CNNs trained individually. In addition, we also show that the learned shared hidden layers are also useful for unseen image character recognition tasks.

6. REFERENCES

- [1] Zohra Saidane and Christophe Garcia, "Automatic scene text recognition using a convolutional neural network," in *Proceedings of the Second International Workshop on Camera-Based Document Analysis and Recognition (CBDAR)*, 2007.

- [2] Yuanping Zhu, Jun Sun, and Satoshi Naoi, "Recognizing natural scene characters by convolutional neural network and bimodal image enhancement," in *Camera-Based Document Analysis and Recognition*, pp. 69–82. Springer, 2012.
- [3] Charles Jacobs, Patrice Y Simard, Paul Viola, and James Rinker, "Text recognition of low-resolution document images," in *Document Analysis and Recognition, 2005. Proceedings. Eighth International Conference on*. IEEE, 2005, pp. 695–699.
- [4] Robust word recognition dataset, "<http://algoval.essex.ac.uk/icdar/RobustWord.html>," .
- [5] Huiqun Deng, George Stathopoulos, and Ching Y Suen, "Error-correcting output coding for the convolutional neural network for optical character recognition," in *Document Analysis and Recognition, 2009. ICDAR'09. 10th International Conference on*. IEEE, 2009, pp. 581–585.
- [6] Minoru Yokobayashi and Toru Wakahara, "Segmentation and recognition of characters in scene images using selective binarization in color space and gat correlation," in *Document Analysis and Recognition, 2005. Proceedings. Eighth International Conference on*. IEEE, 2005, pp. 167–171.
- [7] Minoru Yokobayashi and Toru Wakahara, "Binarization and recognition of degraded characters using a maximum separability axis in color space and gat correlation," in *Pattern Recognition, 2006. ICPR 2006. 18th International Conference on*. IEEE, 2006, vol. 2, pp. 885–888.
- [8] Zhen Jin etc., "Ssift: An improved sift descriptor for chinese character recognition in complex images," in *CNMT.2009*, 2009, pp. 1–5.
- [9] Nongliang Sun etc., "Recognition of chinese character in gray image based on improved radiant-projection-transform," in *MASS*, 2004, pp. 168–171.
- [10] Samuel Thomas, Sriram Ganapathy, and Hynek Hermansky, "Cross-lingual and multi-stream posterior features for low resource lvcsr systems.," in *INTER-SPEECH*, 2010, pp. 877–880.
- [11] A. Stolcke, F. Grezl, Mei-Yuh Hwang, Xin Lei, N. Morgan, and D. Vergyri, "Cross-domain and cross-language portability of acoustic features estimated by multilayer perceptrons," in *Acoustics, Speech and Signal Processing, 2006. ICASSP 2006 Proceedings. 2006 IEEE International Conference on*, 2006, vol. 1, pp. I–I.
- [12] Christian Plahl, R Schluter, and Hermann Ney, "Cross-lingual portability of chinese and english neural network features for french and german lvcsr," in *Automatic Speech Recognition and Understanding (ASRU), 2011 IEEE Workshop on*. IEEE, 2011, pp. 371–376.
- [13] Jui-Ting Huang, Jinyu Li, Dong Yu, Li Deng, and Yifan Gong, "Cross-language knowledge transfer using multilingual deep neural network with shared hidden layers," in *Proc. ICASSP*, 2013.
- [14] Samuel Thomas, Sriram Ganapathy, and Hynek Hermansky, "Multilingual mlp features for low-resource lvcsr systems," in *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on*. IEEE, 2012, pp. 4269–4272.
- [15] Alex Krizhevsky, Ilya Sutskever, and Geoff Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*, 2012, pp. 1106–1114.