

ATTENTION BASED NETWORK FOR REMOTE SENSING SCENE CLASSIFICATION

Shaoteng Liu¹, Qi Wang^{1,2}, Xuelong Li³*

¹ School of Computer Science and Center for OPTical IMagery Analysis and Learning,
Northwestern Polytechnical University, Xi'an 710072, Shaanxi, P. R. China.

² Unmanned System Research Institute,
Northwestern Polytechnical University, Xi'an 710072, Shaanxi, P. R. China.

³ Xi'an Institute of Optics and Precision Mechanics, Chinese Academy of Sciences,
Xi'an 710119, Shaanxi, P. R. China,
and University of Chinese Academy of Sciences, Beijing 100049, P. R. China.

ABSTRACT

Scene classification of very high resolution remote sensing images is becoming more and more important because of its wide range of applications. However, previous works are mainly based on handcrafted features which do not have enough adaptability and expression ability. In this paper, inspired by the attention mechanism of human visual system, we propose a novel attention based network (AttNet) for scene classification. It can focus selectively on some key areas of images so that it can abandon redundant information. Essentially, AttNet gives a way to readjust the signal of supervision, and it is one of the first successful attempts on visual attention for remote sensing scene classification. Our method is evaluated on the UC Merced Land-Use Dataset, in comparison with some state-of-the-art methods. The experimental result shows that the proposed method makes a great improvement on both convergence speed and classification accuracy, and it also shows the effectiveness of visual attention for this task.

Index Terms— Scene classification, remote sensing, visual attention, deep learning, convolutional neural networks, long short-term memory

1. INTRODUCTION

Recently, scene classification of very high resolution (VHR) remote sensing images is becoming more and more important because of its wide range of applications, such as natural disaster detection, land-cover/land-use classification, geographic space object detection, geographic image retrieval, urban planning, and environment monitoring. In the early works, handcrafted features are the most widely used in this task and have been intensively investigated, such as color histograms, scale-invariant feature transform (SIFT), and histogram of oriented gradients (HOG). These methods rely heavily on professional skills and domain expertise to design

various features so that their adaptability and expression ability are not strong enough. In the meanwhile, these methods usually need mid-level encoders as an auxiliary, such as the famous bag-of-visual-words (BoVW) [1], fisher vector (FV) coding and spatial pyramid matching (SPM) [2]. However, as this task becomes more difficult, the above-mentioned methods can not meet our requirements.

When our human go to understand a scene, we tend to select some key areas first and then combine these areas to generate the understanding of the entire scene, which is called visual attention. This mechanism will help us calculate faster and more accurately. So if it can be introduced into the domain of scene classification, the accuracy will be greatly improved in theory. Besides, in recent years, deep learning methods almost replace other traditional methods rely on its strong expressive capacity and the multi-level features. Among them, the most widely used are convolutional neural networks (CNN) [3] and recurrent neural networks (RNN). Especially for CNN, its convolution operation makes it very suitable for image classification. In the domain of remote sensing scene classification, many excellent methods have been proposed and have achieved surprising performance compared to the earlier state-of-the-art methods [1, 4, 5, 6, 7]. However, deep learning methods are data-driven so that they need huge data to learn good fitting [8, 9].

As the advantages above, we are committed to bringing this attention mechanism to the field of remote sensing, while fully integrating the strengths of CNN and RNN. In this paper, a novel attention based network (AttNet) is proposed. To the best of our knowledge, this method is one of the first successful attempts on visual attention for scene classification of VHR remote sensing images. First, we construct a CNN as our high-level feature extractor, which is trained on other huge datasets to solve the problem of insufficient data. Second, a attention mechanism with mask matrix is proposed to focus on the key areas of high-level feature map, which is the core of our network. Finally, in order to update the mask ma-

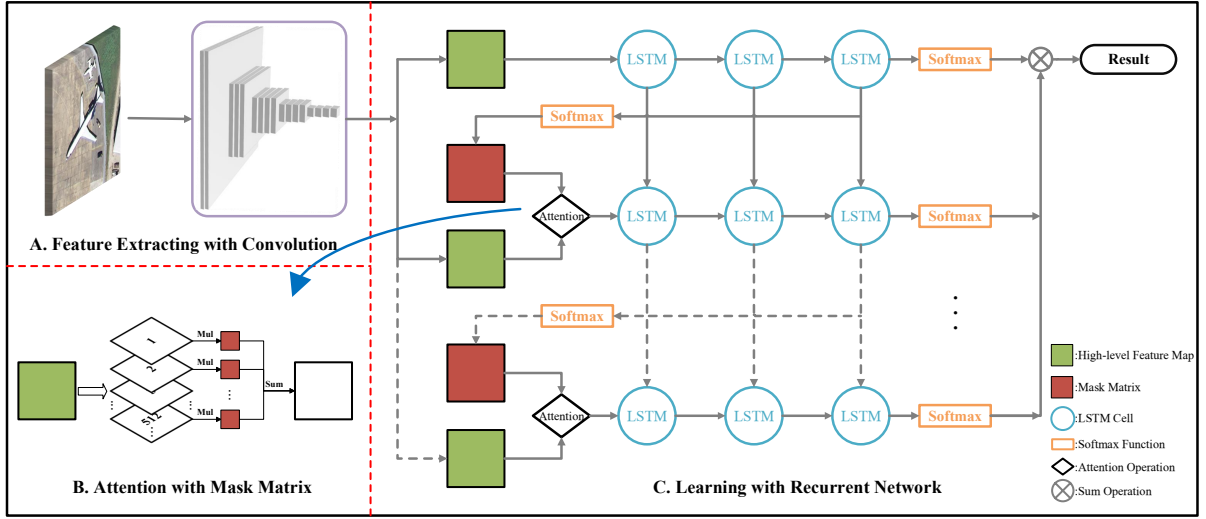


Fig. 1. The overall architecture of the proposed *attention based network*.

trix circularly, we design a recurrent section based on long short-term memory (LSTM) which is one type of RNN.

The paper is structured as follows: Section 1 provides background, motivation and overview of proposed AttNet. Section 2 introduce the details of proposed method. Section 3 shows the experimental results of our method in comparison with some state-of-the-art methods, and the experimental dataset is UC Merced Land-Use Dataset which contains 21 classes.

2. PROPOSED ATTENTION BASED NETWORK

The overall architecture of the proposed AttNet is shown in Fig. 1, and it is constructed around the attention mechanism. In order to facilitate the introduction, we divide it into the following three main parts.

2.1. Feature Extracting with Convolution

For better performance, scene classification usually operates at the feature level. So how to extract strong features is particularly important. Nowadays, with the popularity of deep learning, CNN becomes the most powerful feature extractor due to its deep structure. With millions of parameters need to be adjusted, the dataset of remote sensing is far from satisfied which only contains thousands of images. So we construct a CNN and train it on other huge dataset. After that we transfer the parameters to our network directly. Our implementation is based on VGGNet-16, and the pre-training dataset is Places205.

2.2. Attention with Mask Matrix

Visual attention is the core idea of this paper and there are many ways to implement it. Through experiments and analysis, we design a mask matrix based attention mechanism as can be seen in Fig. 1. Specifically, it is to set a weight at pixel level to represent the degree of importance, and the collection of these weights is the mask matrix. However, the high-level feature map is usually high-dimensional and our network uses the size of $7 \times 7 \times 512$. So if we set independent mask matrix for each dimension, the computational complexity will be greatly increased or even the model will not converge. Under this circumstances, we share the mask matrix of one feature map which can make this attention mechanism workable without any significant loss.

2.3. Learning with Recurrent Network

After we have designed the attention mechanism for scene classification, we need to find an effective way to learn the mask matrix. Given that there are often multiple key areas in the same scene, we decided to design a recurrent network by ourselves, as shown in Fig. 1. In this network, we stack several layers LSTM to generate the mask matrix circularly while also give the final classification prediction. More specifically, we use the hidden state of the last layer of LSTM combined with softmax function to calculate the mask matrix of next recurrence. The benefit of this design is that it can form an end-to-end network which is very conducive for training and can get better performance.

Table 1. Overall accuracy (%) of different methods with the UC Merced Land-Use Dataset

Methods	80% samples for training	50% samples for training
AttNet	99.12 \pm 0.40	96.81 \pm 0.14
Combing Scenarios I and II [4]	98.49	-
Fusion by Addition [10]	97.42 \pm 1.79	-
CNN-NN [5]	97.19	-
Fine-tuning GoogLeNet [1]	97.10	-
GoogLeNet [6]	94.31 \pm 0.89	92.70 \pm 0.60
CaffeNet [6]	95.02 \pm 0.81	93.98 \pm 0.67
VGG-VD-16 [6]	95.21 \pm 1.20	94.14 \pm 0.69
MS-CLBP+FV [11]	93.00 \pm 1.20	88.76 \pm 0.76
Gradient Boosting Random CNNs [7]	94.53	-
Pyramid of Spatial Relations [2]	89.10	-
BoVW [1]	76.81	-

3. EXPERIMENTS

3.1. Experimental Datasets

The UC Merced Land-Use (UCM) dataset [12] is one of the first ground truth datasets derived from a publicly available high resolution overhead image. It was manually extracted from aerial orthoimagery and downloaded from the United States Geological Survey (USGS) National Map. This dataset contains 21 typical land-use scene categories. Each of categories consists of 100 images measuring 256×256 pixels with a pixel resolution of 30 cm in the red-green-blue color space. The classification of UCM dataset is challenging because of the high inter-class similarity among categories such as medium residential and dense residential areas.

3.2. Experimental Details

The setting of specific parameters is through theoretical analysis and several experiments. First of all, making attention operation with large feature map can result in expensive computation, so the size of high-level feature map is chosen as $7 \times 7 \times 512$. Next, we decide to stack 3 layers LSTM networks, because experiments show that the small number of layer will lead to the decline of classification accuracy and large number of it will increase the difficulty of convergence and computational complexity. In the meanwhile, the hidden size of each LSTM cell is set to 256 to ensure that it has sufficient expression ability. Finally, in the process of single image, our network needs to constantly update the mask matrix and the parameter of recurrence number determines the number of loops. Through several experiments, we find that the larger recurrence number will bring higher classification accuracy, but the convergence rate will also decrease heavily. In this case, we choose 20 as the recurrence number.

In normal supervised learning methods, we usually divide the experimental dataset into two parts, one part as a training

set and the other as a validation set. In general, the larger the proportion of training set, the higher the classification accuracy. For the UCM dataset, we choose two common allocation ratios by referring to previous works, one using 80% samples for training and the other using 50% samples. Besides, we choose overall accuracy (OA) and confusion matrix (CM) as the evaluation metrics, and all implementations are based on PyTorch with the NVIDIA Titan X.

3.3. Experimental Results

The proposed *attention based network* is compared with some state-of-the-art methods, and the overall accuracy is shown in Table 1. As can be seen in this table, the proposed *attention based network* outperforms other comparison approaches under all training ratios. This proves the superiority of our method. In the meanwhile, our network is based on the high-level feature extracted from VGGNet-16 as same as the method called VGG-VD-16 [6], but the performance of overall accuracy is further boosted by 3.91% under 80% ratio and 2.67% under 50% ratio. This result confirms that the visual attention has a positive effect on the improvement of classification accuracy.

The confusion matrix under the training ratio of 80% is shown in Fig. 2. This figure shows the best result of our experiment. We can see that only two images are misclassified. One image of freeway is divided into overpass, and the other wrongly divide the medium residential into sparse residential. These are indeed very confusing categories, but can be solved by using a more powerful CNN or augmenting the remote sensing dataset.

4. CONCLUSIONS

In this paper, we propose a novel attention based network (AttNet) for remote sensing scene classification, which can focus

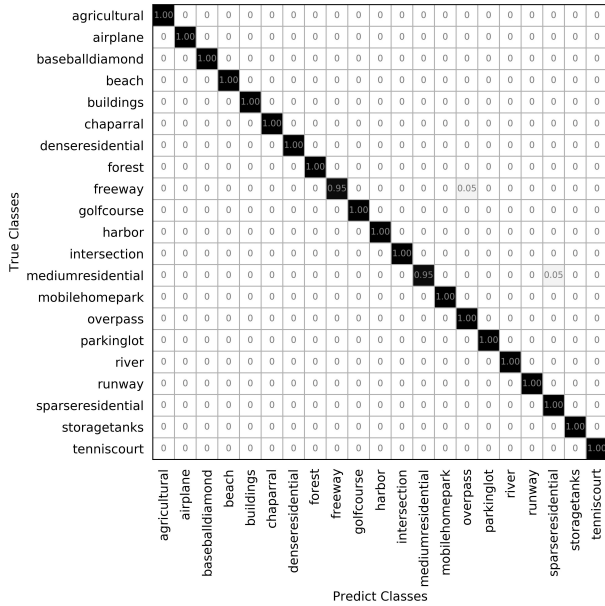


Fig. 2. The confusion matrix with UC Merced Land Use Dataset under the training ratio of 80%.

selectively on the key areas of images so that it can abandon redundant information. Taking the visual attention as the core, our network fully exploits the advantages of CNN as well as RNN, and outperforms other state-of-the-art methods with the UC Merced Land-Use Dataset. The experimental results demonstrate the superiority of our method and the effectiveness of visual attention.

5. ACKNOWLEDGMENT

This work was supported by the National Key R&D Program of China under Grant 2017YFB1002202, National Natural Science Foundation of China under Grant 61773316, Fundamental Research Funds for the Central Universities under Grant 3102017AX010, and the Open Research Fund of Key Laboratory of Spectral Imaging Technology, Chinese Academy of Sciences.

6. REFERENCES

- [1] Marco Castelluccio, Giovanni Poggi, Carlo Sansone, and Luisa Verdoliva, "Land use classification in remote sensing images by convolutional neural networks," *arXiv preprint arXiv:1508.00092*, 2015.
- [2] Shizhi Chen and YingLi Tian, "Pyramid of spatial relations for scene-level land use classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 4, pp. 1947–1957, 2015.
- [3] Qi Wang, Junyu Gao, and Yuan Yuan, "Embedding structured contour and location prior in siamesed fully convolutional networks for road detection," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 1, pp. 230–241, 2018.
- [4] Fan Hu, Gui-Song Xia, Jingwen Hu, and Liangpei Zhang, "Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery," *Remote Sensing*, vol. 7, no. 11, pp. 14680–14707, 2015.
- [5] Esam Othman, Yakoub Bazi, Naif Alajlan, Haikel Al-hichri, and Farid Melgani, "Using convolutional features and a sparse autoencoder for land-use scene classification," *International Journal of Remote Sensing*, vol. 37, no. 10, pp. 2149–2167, 2016.
- [6] Gui-Song Xia, Jingwen Hu, Fan Hu, Baoguang Shi, Xiang Bai, Yanfei Zhong, Liangpei Zhang, and Xiaoqiang Lu, "Aid: A benchmark data set for performance evaluation of aerial scene classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 7, pp. 3965–3981, 2017.
- [7] Fan Zhang, Bo Du, and Liangpei Zhang, "Scene classification via a gradient boosting random convolutional network framework," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 3, pp. 1793–1802, 2016.
- [8] Qi Wang, Jia Wan, and Yuan Yuan, "Locality constraint distance metric learning for traffic congestion detection," *Pattern Recognition*, vol. 75, pp. 272–281, 2018.
- [9] Qi Wang, Jia Wan, and Yuan Yuan, "Deep metric learning for crowdedness regression," *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2017.2703920, 2018.
- [10] Souleyman Chaib, Huan Liu, Yanfeng Gu, and Hongxun Yao, "Deep feature fusion for vhr remote sensing scene classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 8, pp. 4775–4784, 2017.
- [11] Longhui Huang, Chen Chen, Wei Li, and Qian Du, "Remote sensing image scene classification using multi-scale completed local binary patterns and fisher vectors," *Remote Sensing*, vol. 8, no. 6, pp. 483, 2016.
- [12] Yi Yang and Shawn Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proceedings of the 18th SIGSPATIAL international conference on advances in geographic information systems*. ACM, 2010, pp. 270–279.