

Template for Preparation of Manuscripts for *Tsinghua Science and Technology*

This template is to be used for preparing manuscripts for submission to *Tsinghua Science and Technology*. Use of this template will save time in the review and production processes and will expedite publication. However, use of the template is not a requirement of submission. Do not modify the template in any way (delete spaces, modify font size/line height, etc.).

Attention to Difficult Pixels: A Cascade Coarse to Fine Network Architecture for Semantic Segmentation

Wang Zhenyang, Deng Zhidong*, and Wang Shiyao

Abstract: Scene labeling, based on semantic segmentation, is a fundamental topic in computer vision. The goal is to assign each pixel in the image a category label. Convolutional neural networks, especially the fully convolutional neural networks, have attracted increasing attention on semantic segmentation due to the powerful capabilities of extracting hierarchical features. Since it is required to learn to make dense predictions for each pixel, a simple network is hardly to obtain considerable performances on different scenes. In this paper, we propose a novel semantic segmentation network called HMNet, which aims to pay more attention to the hard examples. The network has three branches, where the first branch produces coarse output predictions, and the second branch selects the hard examples which will be fed to the last branch. All above branches focus on their own objectives and collaboratively learn to predict from coarse to fine inference. Since the semantic segmentation dataset contains a large number of relatively easy samples and some hard ones, HMNet is encouraged to select these hard examples to make further predictions which is help to improve the final performance. In order to evaluate predicting performance of the proposed HMNet, we conduct experiments on two public datasets including Sift Flow and Stanford Background Dataset. We show that the three branches can be trained in an end-to-end manner and the experimental results show that compared to all existing models, our HMNet consistently yields the best performance, with accuracy of 91.6% and 89.7%, respectively.

Key words: semantic segmentation; online hard example mining; convolutional neural network

1 Introduction

Semantic segmentation, also known as scene labeling, is one of the fundamental research topics in computer vision. The goal of semantic segmentation is to identify and assign each pixel in the image with a category. This requires a complete understanding of the semantic information of the entire image. That means, for the testing image, it needs to predict the label of each object, and it is also very important to determine

the boundary of each object in pixel level. Semantic segmentation has a strong application requirement in the field of environment perception and autonomous self driving car.

In recent years, deep learning has made breakthroughs in image classification, speech recognition, visual object recognition, object detection, and so on. For the task of image semantic segmentation, there are also many methods based on deep learning methods. Early research attempts to apply the convolution neural network designed for image recognition directly to semantic segmentation. Although good segmentation results are obtained, the prediction results are rough and the edges of objects can not be separated correctly. This is mainly caused by the losing of location information. The fully convolutional neural network is proposed by ??

• Wang Zhenyang, Deng Zhidong, Wang Shiyao are with the Department of Computer Science, Tsinghua University, Beijing 100084, China. E-mail: crazycry2010@gmail.com, michael@tsinghua.edu.cn, sy-wang14@mails.tsinghua.edu.cn

* To whom correspondence should be addressed.

Manuscript received: 2017-09-20; revised: year-month-day; accepted: year-month-day

to overcome this disadvantage, and become the most popular segmentation framework for semantic segmentation task.

In order to solve the problem of edge blur, conditional random field, Markov random field, Gaussian conditional random field, and other variations are proposed. But with the development of network architecture, especially with the presenting of ResNet ?? and batch normalization ??, an end-to-end training procedure can achieve a same even better segmentation results.

Another challenge is to simultaneously predict large and dense semantic labels, especially when the distribution of different classes is unbalanced. An extreme unbalanced example is the background and foreground pixels. In addition, the unbalance of easy and difficulty pixels can also affect the convergence of the network. Aiming at the problem of unbalanced distribution of different difficulty pixels, we propose a novel semantic segmentation network framework based on the idea of hard negative example mining.

In fact, hard sample mining is not a new concept. As early as 1994, Sung and Poggio ?? proposed bootstrapping method in their face detection algorithm. The main purpose is to enhance the algorithm's capacity by changing the distribution of difficult samples. Inspired by the similar ideas, a cascade coarse to fine network architecture(CCFNet) is proposed. CCFNet is composed by three branches which share a same feature extraction network. The first branch is a coarse segmentation network which can predict correctly for most pixel. The second branch is an attention network used to predict a degree of difficulty per pixel. For the difficult pixels, the third segmentation branch is proposed to refine the final results.

For semantic segmentation tasks, each pixel should be assigned a prediction label. So each training image contains a large number of optimization targets, which makes a single segmentation network be difficult to meet all the optimization goals. CCFNet uses a multi-branch network structure, integrating the prediction from two different branches with an attention model to optimize the final segmentation results.

In conclusion, we propose a cascade coarse to fine network architecture for semantic segmentation. Compared with the traditional semantic segmentation networks, CCFNet has the following characteristics: 1. Propose a cascade segmentation network architecture to refine the final segmentation results. 2. Incorporate

the idea of hard mining into an attention module.

3. Validate CCFNet on both Sift Flow Dataset and Stanford Background Dataset.

2 Related Works

3 Method

Inspired by online hard example mining(OHEM) algorithm, we propose a cascade coarse to fine network architecture CCFNet. Section 3.1 illustrates how to turn a ImageNet pre-trained model to a cascade coarse to fine semantic segmentation network CCFNet. Taking the characteristics of semantic segmentation task into consideration, section ?? introduces a hard instance mining method to learn the attention about the difficulty of each instance. Section ?? details the training and testing process of CCFNet.

3.1 Cascade Coarse to Fine Network Architecture

We choose ResNet-50 pre-trained on ImageNet as our baseline model. ResNet originally is designed for image classification which win ILSRVC 2015 competition and surpass the human performance on ImageNet dataset. ResNet-50 is one of the version provided in experiments, which is faster than VGG-16 and more accurate than VGG-19. Compared with ResNet-101, ResNet-50 is cheaper in computation resources and memory consumption, but can achieve a comparable accuracy. Figure ??a visualizes the network architecture of ResNet-50, which is composed by five stages with different configurations of layers and a classification stage. We treat the ResNet-50 as a common feature extraction part of CCFNet by discarding the classification stage.

The architecture of CCFNet is shown as Fig. ??b, which is composed by three cascade network branches: a coarse segmentation branch as a baseline result, an attention branch to predict the difficulty of labeling each pixel instance, and a refine segmentation branch to refine the final segmentation results. These three branches share a common feature extraction network.

The feature extraction network is a fundamental convolutional neural network. A modified version of ResNet-50 is adopted by this paper. For semantic segmentation task, the context is important to predict the correct label of each pixel instance. But it is difficult to determine the boundary of each pixel's context, since different objects may have different sizes. And the problem gets more complicated when considering the

various perspective of each images. A simple yet effective method to solve this problem is to integrate multi-scale features for label predicting. The residual error model itself has the property of extracting and integrating multi-scale features, which can be seen from Fig. ???. From the unravelled view by Veit et al. ??, a two-unit ResNet is equivalent to an ensemble of four sub-networks with different receptive fields, as illustrated in Fig. ???. So the whole ResNet-50 can be expanded as a linearly growing ensemble of sub-networks, which can extract and integrate multi-scale features.

Besides, there are two improvements adopted by ResNet-50 to make it more suitable for segmentation task. First, we only keep the first three pooling layers to preserve the resolution. So the final resolution of prediction is 1/8 of the original input image. Secondly, we replace the convolutional layer in the last two stages with the dilated convolutions. It can help to enlarge the reception field of predicted feature maps. The modified ResNet-50 is used as feature extraction network in CCFNet.

The coarse segmentation branch is a baseline model for image semantic segmentation. This branch is shown as the red part in Fig. ?? b, we adopt a fully convolutional network(FCN) with two convolutional layers to predict the semantic classes for their regions. Since the resolution is 1/8 of the original input image, the feature maps are up-sampled by bilinear interpolation. Finally, a pixel-wise softmax loss is adopted to predict the probabilities of each pixel.

The attention branch is proposed to learn a soft-attention, which is a one-channel feature map with the same resolution as the input image. It is mainly used to indicate the segment difficulty of each pixel. The idea behind is simple yet effective. The segmentation datasets contain a large number of easy pixel instances and a small number of difficulty pixel instances. Paying more attention on these difficulty pixel instances can make the training process converges faster and efficiently. The attention branch shares the same feature extraction network as the coarse segmentation branch, and has a similar network structure. The major different is that the attention branch is a two-category semantic segmentation network which is only used to predict the segment difficulty. Just the same as the coarse segmentation branch, softmax layer is adopted again to generate a soft-attention weighting coefficient.

The refine segmentation branch refines the

segmentation results as the final network output. This branch is more complicated compared with the first two branches. The coarse segmentation branch is hard to segment all the pixels correctly. So the pixel instances can be divided into two groups by the coarse segmentation branch. The pixels which can be segment correctly by the coarse segmentation branch is denote as easy pixel instances, while the others are difficult ones. A fine segmentation network is introduced to reclassification the difficult pixel instances. Inspired by the PSPNet ??, pyramid pooling is adopted by the fine segmentation network to extract multi-scale features. And the refine segmentation branch is a weighted summation of the coarse segmentation branch and the fine segmentation network, with the weighting coefficient predicted from the attention branch. So an end-to-end learning branch is proposed to learn the final segmentation result directly.

The three branches are cascaded one by one, and constitute an end-to-end learning network with multiple loss functions.

3.2 Attention to Difficulty Pixels

Hard example mining is one of the commonly used training techniques for machine learning. The traditional implement is a continuous iterative process which could be divided into two steps. Firstly, the training model is fixed to screen out the difficult examples, and the training set is updated by adding a certain rate of difficult examples. Secondly, in the fixed training set, the training model is re-trained.

In this paper, the two-step process of hard example mining is optimised to an end-to-end learning network framework. For semantic segmentation task, each pixel should be assigned a category label. So a single image contain enough training samples for hard example mining. The attention branch is used to predict each pixel is easy or difficult for the coarse segmentation branch. It makes the end-to-end learning possible by heuristically filtering out the hard examples online. The final segmentation results relay more on the coarse segmentation branch if the pixel is predicted as an easy one. Otherwise, the fine segmentation network takes up a larger proportion. Inspired by online hard example mining, the attention branch with a heuristic strategy is introduced to CCFNet to predict the difficult of each pixel. And the final segmentation results are promoted by hard sample selection.

During the training process, the attention branch is

supervised by a label map with 0/1 values, indicating easy or difficult for the pixel in corresponding position. The attention branch is cascaded behind the coarse segmentation branch, so the 0/1 label map can be generated by a comparison between the prediction of the coarse segmentation branch and the segmentation ground truth. 0 indicates the pixel is misclassified by the coarse segmentation branch, while 1 represents a correctly prediction. The 0/1 label map is used as the ground truth of the attention branch, supervising the attention branch to learn the difficulty of each pixel.

4 Experimental results

4.1 Network Configuration

4.2 Data Sets

4.3 Comparison Results

4.4 Results Visualization

5 Conclusions

For example: The parallelization of cutoff pair interactions is mature on CPUs, and typically employs a voxel-based method.

References to the literature are cited by *number in square brackets* at appropriate locations (*before* a period, comma, etc.) in the text.

Examples:

- Negotiation research spans many disciplines [3].
- This result was later contradicted by Becker and Seligman [5, 6], who
- This effect has been widely studied [1-3, 7].
-achieved until rather recently [11, 21, 22], with.....
-stage of cap formation (see Fig. 5 in Ref. [14]).

Acknowledgements

This work was supported in part by the National Science Foundation of China (NSFC) under Grant Nos. 91420106, 90820305, and 60775040, and by the National High-Tech R&D Program of China under Grant No. 2012AA041402.

References

The font is Times New Roman \zihao{5-}. This part is placed at the end of the manuscript. References should be numbered sequentially as they appear throughout the text. Only one publication should be given for each number. The list of references should only include papers that have been published or accepted by a named publication or recognized preprint server. Authors should ensure the accuracy and completeness of all references before submission. Please ensure references are given in the correct style as shown below in order to avoid delays in typesetting your article

References

- [1] M. Abrahams and M. Kattenfeld, The role of turbidity as a constraint on predator-prey interactions in aquatic environments, *Behav. Ecol. Sociobiol.*, vol. 40, no. 3, pp.169-174, Mar. 1997. (**Journal style**)
- [2] R. C. Calfee and R. R. Valencia, *APA Guide to Preparing Manuscripts for Journal Publication*. Washington, DC: American Psychological Association, 1991. (**Authored book style**)
- [3] J. M. O'Neill and J. Egan, Mens and womens gender role journeys: metaphor for healing, transition, and transformation, in *Gender Issues Across the Life Cycle*, B. R. Wainrib, Ed. New York: Springer, 1992, pp. 107-123. (**Book chapter style**)
- [4] J. C. Phelan, B. G. Link, A. Stueve, and B. A. Pescosolido, Have public conceptions of mental health changed in the past half century? Does it matter? presented at the 124th Annu. Meeting American Public Health Association, New York, USA, 1969. (**Paper presented at conferences style**)
- [5] A. Mahdavi and B. Spasojevic, Incorporating simulation into building systems control logic, In *Proc. 10th Int. Building Performance Simulation Association Conf.*, Beijing, China, 2007, pp. 1175-1181. (**Paper in proceedings style**)
- [6] X. Yang, Study of building material emissions and indoor air quality, Ph.D. dissertation, Dept. Arch., MIT, MA, USA, 1999. (**Dissertation style**)
- [7] J. Dong, S. Martin, and P. Waldo, Method and system for dynamically presenting cluster analysis results, US Patent US6380937B1, March 30, 2002. (**Patent style**)
- [8] Y. Jiang, Liquid desiccant air-conditioning system and its applications, (in Chinese), *Heating Ventilating & Air Conditioning*, vol. 34, no. 11, pp. 88-97, 2004. (**Non-English publication style**)
- [9] M. K. Slifka and J. L. Whitton, Clinical implications of dysregulated cytokine production, *J. Mol. Med.*, doi:10.1007/s001090000086. (**Article by DOI style**)
- [10] J. Doe. The dictionary of substances and their effects, <http://www.rsc.org/dose/title>, 1999, Jan. 15. (**Online document style**)



First A. Author Photo. Biographies should be limited to one paragraph consisting of the following: sequentially ordered list of degrees, including years achieved; sequentially ordered places of employ concluding with current employment; association with any official journals or conferences; major professional and/or academic achievements, i.e., best paper awards, research grants, etc.; any publication information (number of papers and titles of books published); current research interests; association with any professional associations.



Second B. Author Photo. Biographies should be limited to one paragraph consisting of the following: sequentially ordered list of degrees, including years achieved; sequentially ordered places of employ concluding with current employment; association with any official journals or

conferences; major professional and/or academic achievements, i.e., best paper awards, research grants, etc.; any publication information (number of papers and titles of books published); current research interests; association with any professional associations.



Third C. Author Photo. Biographies should be limited to one paragraph consisting of the following: sequentially ordered list of degrees, including years achieved; sequentially ordered places of employ concluding with current employment; association with any official journals or conferences; major professional and/or academic achievements, i.e., best paper awards, research grants, etc.; any publication information (number of papers and titles of books published); current research interests; association with any professional associations.