

VISUAL CONFUSION LABEL TREE FOR IMAGE CLASSIFICATION

Yuntao Liu, Yong Dou, Ruochun Jin, Rongchun Li

National University of Defense Technology
National Laboratory for Parallel and Distributed Processing
Changsha, China, 410073

liuyuntao.me@gmail.com, {yongdou, jinruochun, rongchunli}@nudt.edu.cn

ABSTRACT

Convolution neural network models are widely used in image classification tasks. However, the running time of such models is so long that it is not conforming to the strict real-time requirement of mobile devices. In order to optimize models and meet the requirement mentioned above, we propose a method that replaces the fully-connected layers of convolution neural network models with a tree classifier. Specifically, we construct a Visual Confusion Label Tree based on the output of the convolution neural network models, and use a multi-kernel SVM plus classifier with hierarchical constraints to train the tree classifier. Focusing on those confusion subsets instead of the entire set of categories makes the tree classifier more discriminative and the replacement of the fully-connected layers reduces the original running time. Experiments show that our tree classifier obtains a significant improvement over the state-of-the-art tree classifier by 4.3% and 2.4% in terms of top-1 accuracy on *CIFAR-100* and *ImageNet* datasets respectively. Additionally, our method achieves $124\times$ and $115\times$ speedup ratio compared with fully-connected layers on AlexNet and VGG16 without accuracy decline.

Index Terms— Image classification, convolution neural network, confusion graph, label tree

1. INTRODUCTION

Deep convolution neural network(CNN) models are developing rapidly and it has evolved to the state-of-the-art technique [1] in image classification tasks. However, when applied to real-time applications on embedded devices where the power and storage are limited, CNN models can not meet the real-time demands because of its large amount of computation. Therefore, optimizing and accelerating these CNN models on embedded devices has become a challenge.

In view of this problem, researches have proposed a variety of compression and acceleration methods such as reducing the precision of multiplication and addition operations [2], setting the weights and inputs to binary codes [3], a skillful integration among several effective methods [4] and

structure changing [5, 6, 7]. However, the computation of a CNN model mainly induced by the computation of fully-connected(FC) layers [4]. The methods mentioned above mainly focus on the compression on the convolution layers and have not resolved the huge computation problem of FC layers.

Another classification method based on a tree classifier, which is an appropriate method for large-scale recognition problems with thousands of categories, has received extensive attention and substantial development. There are several methods to construct the structure of the tree classifier such as leveraging the semantic ontologies (taxonomies) [8, 9, 10], learning label trees [11] and probabilistic label trees [12], learning visual trees [13, 14] and enhanced visual tree [15]. Compared with the FC layers in the CNN models, the tree classifier has the advantage of small amount of calculation and the computation complexity of the tree classifier is $O(\log N)$ [11]. However, there has been no work that replaces FC layers with the tree classifier because most previous work construct the structure of the tree classifier by clustering directly from the image dataset. Previous methods do not utilize the information of FC layers so the accuracy is limited, which restricts the application of tree classification in accelerating depth CNN models. Moreover, this limitation also results in the separation of research on tree classification and deep CNN models and both of them can not benefit from each other.

[16] discovered that deep CNN models have visual confusions that is similar to human beings and we believe this characteristic can be used as the metric to construct the Label Tree. Therefore, we propose to use the community detection algorithm to construct the Label Tree, called **Visual Confusion Label Tree(VCLT)**. With this method, we can fully utilize the information of FC layers in the CNN models. Compared with previous Label Tree building methods, the advantage of the VCLT is that there is no need to manually set parameters and do clustering tasks during tree construction. In addition, our VCLT is constructed directly based on the features in deep CNN models so it has a more reasonable structure which is beneficial for improving the accuracy of the tree classifier. Moreover, to the best of our knowledge, VCLT is the first ef-

fort that connects the CNN model and the Label Tree directly so the tree structure fully inherits the information contained in the FC layers.

There are two main contributions in this paper as follows.

- *Visual Confusion Label Tree:* Our construction method is based on the hierarchical community detection algorithm. Using this algorithm on the output of FC layers we can construct a tree classifier whose structure is more suitable for deep CNN models. With this method we can improve the accuracy of the tree classifier compared with previous work by 2.4%–4.3% and we can prove in theory.
- *Replace the FC layers with the tree classifier:* After constructing the label tree, we replace the FC layers from the deep CNN models with our VCLT and we propose an effective algorithm to train our tree classifier. With this replacement we reduce the amount of computation in FC layers by $37\times-124\times$ without sacrificing the accuracy of original CNN-based methods.

2. LABEL TREE IN A NUTSHELL

The concept of the Label Tree was first proposed in [11] aiming at classification and a label tree is a tree $T = (N, E, F, L)$ with $n + 1$ indexed nodes $N = \{0, \dots, n\}$ where $E = \{(p_1, c_1), (p_{|E|}, c_{|E|})\}$ is a set of edges that are ordered pairs of parent and child node indices, $F = \{f_1, \dots, f_n\}$ are label predictors and label sets $L = \{l_0, \dots, l_n\}$ are associated to each node. Except the root of the tree, all other nodes have one parent and arbitrary number of children. The label sets indicate the set of labels to which a point should belong if it arrives at the given node. Classifying an example begins at the root node and for each edge leading to a child $(s, c) \in E$ one computes the score of the label predictor $f_c(x)$ which predicts whether the example x belongs to the set of labels l_c . One takes the most confident prediction, traverses to that child node, and then repeats the process. Classification is complete when one arrives at a node that identifies only a single label that is the predicted class. More details about Label Trees can be found in [11, 14, 17].

3. VISUAL CONFUSION LABEL TREE AND TRAINING

3.1. Definition of the Visual Confusion Label Tree

Definition 1. A *Visual Confusion Label Tree* is a tree $T = (\mathcal{N}, \mathcal{V}, \mathcal{E}, \mathcal{L})$ with k hierarchical layers $\mathcal{N} = \{n_1, \dots, n_k\}$ where n_{id} denotes the number of nodes in the id th layer, the node sets $\mathcal{V} = \{V_1, \dots, V_k\}$ where V_{id} is a set of nodes in the id th layer and $V_{id} = \{v_1, \dots, v_{n_{id}}\}$, branch edges $\mathcal{E} = \{(p_1, c_1), \dots, (p_{|\mathcal{E}|}, c_{|\mathcal{E}|})\}$ which are ordered pairs of parent and child node indices and labels sets

$\mathcal{L} = \{L_1, \dots, L_s\}$ where L_{id} is a label set of nodes in the id th layer and $L_{id} = \{l_1, \dots, l_{n_{id}}\}$ where l_s denotes the label set of the s th node in this layer.

We extend the notation in [11] to Definition 1. The number k in Definition 1, which is the number of the hierarchical layers of VCLT, is equal to the number of iterations of hierarchical community detection algorithm run on a related confusion graph(detailed in Section 3.2).

3.2. Visual Confusion Label Tree Establishing

Algorithm 1: Establish Visual Confusion Label Tree of a N -category classification

Input: A N -category classification model M ; A dataset D ; Top concern number $\tau (\tau \leq N)$;
Output: The VCLT $T = (\mathcal{N}, \mathcal{V}, \mathcal{E}, \mathcal{L})$;
1 $G \leftarrow \text{GenerateConfusionGraph}(M, D, N, \tau)$;
2 $\mathcal{P}, \mathcal{C} \leftarrow \text{HierarchicalCommunityDetect}(G)$;
3 $k \leftarrow \text{length}(\mathcal{P})$;
4 **for** i from 1 to k **do**
5 $\mathcal{N} = \mathcal{N} \cup \text{length}(\mathcal{P}[i])$;
6 $\mathcal{V} = \mathcal{V} \cup \mathcal{P}[i]$;
7 $\mathcal{L} = \mathcal{L} \cup \mathcal{C}[i]$;
8 **end**
9 **for** i from 2 to k **do**
10 **for** j from 1 to $\mathcal{N}[i]$ **do**
11 **foreach** $v_s \in \mathcal{P}[i][j]$ **do**
12 $\mathcal{E} = \mathcal{E} \cup e_{\mathcal{V}[i][j], v_s}$;
13 **end**
14 **end**
15 **end**
16 **return** T ;

Given a dataset D and its corresponding classification model M , Algorithm 1 establishes the VCLT T defined in Definition 1 using confusion graph generation and community detection algorithm. There are 3 main steps in Algorithm 1. The first is using the confusion graph generation algorithm to build a confusion graph. The second is using the hierarchical community detection algorithm to reveal communities in the confusion graph. The last is establishing a VCLT with the results of the second step.

Specifically, for the function “GenerateConfusionGraph”, we utilize the confusion graph establishing algorithm from [16]. This algorithm firstly normalizing the top- τ classification scores of each test sample and then accumulating each normalized score to the weight of the edge that connects the labeled category and the predicted category. For the function “HierarchicalCommunityDetect”, we use the algorithm from [18]. This algorithm is an iterative algorithm and it will continue running until the modularity is converged. This function outputs the \mathcal{P} and \mathcal{C} . The \mathcal{P} is a set of arrays which

node). All these methods make the tree classifier over the VCLT more discriminative.

In order to discriminate a given coarse-grained or fine-grained categories on a node c_j from its sibling nodes $\{s(c_j)\}$ under the same parent node $c_i = p(c_j)$, its multi-kernel SVM classifier is defined as:

$$f_{c_j}^l = K(\mathbf{x}, \mathbf{x}') + b \quad (1)$$

where l is the level of node c_j and $K(\mathbf{x}, \mathbf{x}')$ is the multi-kernel which is defined as:

$$K(\mathbf{x}, \mathbf{x}') = \sum_{m=1}^M d_m K_m(\mathbf{x}, \mathbf{x}') \quad (2)$$

with:

$$d_m \geq 0, \sum_{m=1}^M d_m = 1 \quad (3)$$

In our method, we use common kernels such as the linear kernel, the polynomial kernel, and the Gaussian kernel.

We train each classifier node by node from the root to leaf nodes and use the strategy of SVM Plus [21] to train the multi-label SVM classifiers. Specifically, there is a set of labeled training images for R sibling nodes $\{s(c_j)\}$ under the same parent node $c_i = p(c_j)$, $R \in [2, B]$ (B is the number of sibling nodes) and there are training samples $\Omega = \{(x_j^l, y_j^l) | c_j \in c_i\}$ (l is the level of the sibling nodes), training the multi-kernel SVM Plus classifiers for R sibling nodes is achieved by optimizing a objective function:

$$\min_{f_0, f_{c_j}, b, \xi, d} \frac{1}{2} \|f_0\|^2 + \frac{\lambda}{2} \sum_{j=1}^R \sum_{m=1}^M \frac{1}{d_m} \|f_{c_j}\|^2 + C \sum_{j=1}^R \xi_j \quad (4)$$

subject to:

$$\forall_{j=1}^R : y_j f_{c_j} = y_j \left(\sum_{m=1}^M K_{m,c_j}(\mathbf{x}_j, \mathbf{x}_j') + b \right) \geq 1 - \xi_j, \xi_j > 0 \quad (5)$$

where ξ_j indicates the slack variable, λ is the positive regularization parameters, C is the penalty term.

One problem of label tree is that the error propagation may have negative influence on the classification result. If a classification error happens on the parent node, the prediction of this sample will be wrong because the labels of leaf nodes under this misclassified parent node are all incorrect. In order to resolve this problem, we add an inter-level constraint to the SVM plus classifiers. Our strategy is that samples should be classified correctly at the parent level(level: $l+1$) if we want it to be further classified at the current node level(level: l). Thus, we add a constraint to guarantee that the score of current node classifier must be larger than the score of its parent node classifier, which can be denoted by Eq.(8). Therefore, we extend Eq.(4,5) to:

$$\min_{f_0, f_{c_j}, b, \xi, d} \frac{1}{2} \|f_0^l\|^2 + \frac{\lambda}{2} \sum_{j=1}^R \sum_{m=1}^M \frac{1}{d_m} \|f_{c_j}^l\|^2 + C \sum_{j=1}^R \xi_j \quad (6)$$

subject to:

$$\forall_{j=1}^R : y_j f_{c_j} = y_j \left(\sum_{m=1}^M K_{m,c_j}(\mathbf{x}_j, \mathbf{x}_j') + b \right) \geq 1 - \xi_j, \xi_j > 0 \quad (7)$$

$$\forall_{j=1}^R : f_{c_j}^l(\mathbf{x}_j) \geq f_{c_i}^{l+1}(\mathbf{x}_j), (\mathbf{x}_j, y_j) \in c_j \in c_i \quad (8)$$

4. EXPERIMENT

4.1. Datasets and Experimental Settings

We use *CIFAR-100* [22] and *ILSVRC2012* [10] to evaluate the performance of the proposed classification method. *CIFAR-100* has 60000 images of 100 categories. Each category has 600 images in which 500 for training and 100 for validation. Divided into a training set and a validation set, *ILSVRC2012* has over 120 million images of 1000 categories and is commonly used to evaluate image classification algorithms. We use the training set for training and validation set for testing. The Mean Accuracy (MA) (%) [15] is used to capture the performance of each method. A PC with Intel Core i7 and 64GB memory is utilized to run all experiments.

4.2. Comparison of different tree classifiers

In this section, we compare the classification accuracy of our proposed VCLT classifier with those of other state-of-the-art tree classifiers. Trained and tested with *CIFAR-100* and *ImageNet* datasets, we compare the MA of each of the following tree classifiers: **semantic ontology** [8], **label tree** [11], **visual tree** [14] and the **enhanced visual tree** [15]. In order to train and test each model, we employ the *DeCAF* [23] features extracted from the FC6 layer of the AlexNet model (its first FC layer) and the classification accuracy of each tree classifier (quantified in MA) is demonstrated in Table 1.

Approaches	CIFAR-100	ImageNet
Semantic ontology	48.95%	44.47%
Label tree	52.04%	49.64%
Visual tree	51.30%	51.03%
Enhanced visual tree	54.33%	58.76%
Visual confusion label tree	58.63%	61.18%

Table 1. Classification accuracy of different tree classifiers.

From Table 1, we find the performance of Semantic ontology is the worst because its tree structure is constructed based on semantic space and the image classification process based on feature space. For another four methods based on feature space, the performance of Label tree is worse because of using OvR classifier to construct its tree structure, which is limited to sample imbalance and the performance of classifier. As for Visual tree, it uses the average features extracted directly from the dataset. Enhanced visual tree adopts the spectral clustering method that better reflects the diversity of

categories, so its performance is better than the Visual tree. Our VCLT constructs the tree structure based on the confusion of the CNN model, which makes sibling nodes as close as possible and the parent nodes as far as possible. So this tree structure is more proper and we obtain a significant improvement over the Enhanced visual tree by 4.3% and 2.4%.

4.3. Comparison between our tree classifier and CNN models

In this section, we compare the classification accuracy and test time of our VCLT with those of the corresponding CNN model. We choose AlexNet and VGG-Verydeep-16(VGG16) for this comparison. In the AlexNet-based experiment, we firstly train an AlexNet using the *CIFAR-100* dataset. Then we employ the *DeCAF* features extracted from the FC6 layer of AlexNet to train the corresponding VCLT classifier. The classification accuracy and test time of both are shown in Table 2. For a CNN model, the “test time” in Table 2 is the average running time of its FC layers when processing one image. For the VCLT classifier, the “test time” is the average running time of the whole tree classifier when processing one image. As for the VGG16-based experiment, we do the same thing except using the features extracted from the FC14 layer of VGG16 to train the VCLT classifier. Table 3 illustrates this comparison using the *ImageNet* dataset.

Approaches	Accuracy	Test time (ms)	Speedup
AlexNet	54.02%	3.4215	-
VCLT_AlexNet	58.63%	0.0275	124×
VGG16	72.21%	4.3713	-
VCLT_VGG16	72.23%	0.0381	115×

Table 2. Experiment results on *CIFAR-100*

Approaches	Accuracy	Test time (ms)	Speedup
AlexNet	57.24%	5.5263	-
VCLT_AlexNet	61.18%	0.1493	37×
VGG16	71.53%	6.7981	-
VCLT_VGG16	66.71%	0.2427	28×

Table 3. Experiment results on *ImageNet*

From Table 2 and Table 3, we can see that the classification accuracy of VCLT with AlexNet is around 3% higher than that of the original AlexNet on both *CIFAR-100* and *ImageNet* datasets. In addition, on both datasets, the speedup ratios achieved by replacing the FC layers of AlexNet with our tree classifier are significant (124× on *CIFAR-100* and 37× on *ImageNet*). Though the accuracy improvement of VCLT with VGG16 is trivial on *CIFAR-100* and even negative on *ImageNet*, the speedup ratios are remarkable, achieving 115× and 28× respectively, which demonstrates VCLT’s promising potential to accelerate CNN-based applications.

Approaches	Accuracy
DeepCom	57.24%
VCLT_DeepCom	60.73%
BWN	57.08%
VCLT_BWN	57.11%
XNOR	42.37%
VCLT_XNOR	42.48%

Table 4. Results on *ImageNet*(compressed CNN models)

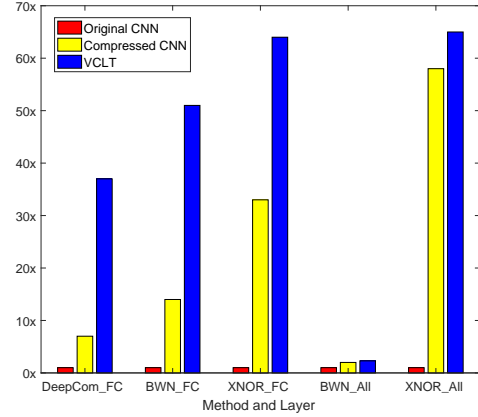


Fig. 3. The VCLT for *CIFAR-100* image set with 100 categories.

Additionally, we compare the classification accuracy and the speedup ratio of our VCLT with compressed CNN models. Here we choose the Binary-Weights-Net(BWN), XNOR-Net(XNOR) [3] and the DeepCompression Network(DeepCom) [4] for comparison. These compressed CNN models are based on AlexNet and we use their pre-trained models on *ImageNet* in our experiment. The accuracy results are shown in Table 4, we find out that our VCLT has no accuracy decline.

The speedup ratio results are shown in Fig. 3. For DeepCom, followed [4], comparisons of speedup mainly focus on FC layers and results are shown in *DeepCom_FC* in Fig. 3. We find the speedup ratio of our VCLT(37×) is 30× higher than DeepCom(7×) when compared with FC layers in the original AlexNet model. For BWN and XNOR, followed [3], comparisons of speedup ratio focus on both FC layers and the entire network. The results are shown respectively in *BWN_FC*, *XNOR_FC*, *BWN_All* and *XNOR_All*. We find the speedup ratio of our VCLT is 37× higher than BWN(14×) and 31× higher than XNOR(33×) when compared on FC layers in the original AlexNet model. As for the entire network model, our VCLT(2.3×) obtain an improvement over the BWN(2×) by 0.3× while VCLT(65×) over the XNOR(58×) by 7× in terms of speedup ratio.

5. CONCLUSION

In this paper, we propose a method of replacing the fully-connected layers in CNN models with a tree classifier in image classification applications. We utilize the community detection algorithm to construct a Visual Confusion Label Tree based on the confusion characteristics of CNN models. Then, we use the multi-kernel SVM plus classifier with hierarchical constraints to train the tree classifier on the Visual Confusion Label Tree. Finally, we use this tree classifier to replace fully-connected layers in CNN models. The experimental results on *CIFAR-100* and *ImageNet* demonstrated the advantages of the proposed method over other tree classifiers and original CNN models such as AlexNet and VGG16.

6. ACKNOWLEDGEMENTS

This work was supported by the Natural Science Foundation of China under the grant No. U1435219, No. 61402507 and No. 61303070.

7. REFERENCES

- [1] SQ. Ren J. Sun KM. He, XY. Zhang, “Deep residual learning for image recognition,” pp. 770–778, 2015.
- [2] JP. David M. Courbariaux, Y. Bengio, “Training deep neural networks with low precision multiplications,” *Computer Science*, 2014.
- [3] J. Redmon A. Farhadi M. Rastegari, V. Ordonez, “Xnet: Imagenet classification using binary convolutional neural networks,” pp. 525–542, 2016.
- [4] WJ. Dally S. Han, H. Mao, “Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding,” in *ICLR*, 2016.
- [5] S. Yan M. Lin, Q. Chen, “Network in network,” *Computer Science*, 2013.
- [6] T. Brox M. Riedmiller JT. Springenberg, A. Dosovitskiy, “Striving for simplicity: The all convolutional net,” *Eprint Arxiv*, 2014.
- [7] MW. Moskewicz K. Ashraf WJ. Dally K. Keutzer FN. Iandola, H. Song, “Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size,” *arXiv preprint arXiv:1602.07360*, 2016.
- [8] Y. Lim DM. Blei FF. Li LJ. Li, C. Wang, “Building and using a semantivisual image hierarchy,” in *CVPR 2010*. IEEE, 2010, pp. 3336–3343.
- [9] EP. Xing B. Zhao, FF. Li, “Large-scale category structure aware image categorization,” in *Advances in Neural Information Processing Systems*, 2011, pp. 1251–1259.
- [10] R. Socher LJ. Li K. Li FF. Li J. Deng, W. Dong, “Imagenet: A large-scale hierarchical image database,” in *CVPR 2009*. IEEE, 2009, pp. 248–255.
- [11] D. Grangier S. Bengio, J. Weston, “Label embedding trees for large multi-class tasks,” in *Advances in Neural Information Processing Systems*, 2010, pp. 163–171.
- [12] M. Tappen O. Shamir C. Liu B. Liu, F. Sadeghi, “Probabilistic label trees for efficient large scale image classification,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 843–850.
- [13] N. Zhou J. Peng R. Jain J. Fan, X. He, “Quantitative characterization of semantic gaps for learning complexity estimation and inference model selection,” *IEEE Transactions on Multimedia*, vol. 14, no. 5, pp. 1414–1428, 2012.
- [14] J. Peng L. Gao J. Fan, N. Zhou, “Hierarchical learning of tree classifiers for large-scale plant species identification,” *IEEE Transactions on Image Processing*, vol. 24, no. 11, pp. 4172–4184, 2015.
- [15] J. Zhang X. Gao Y. Zheng, J. Fan, “Hierarchical learning of multi-task sparse metrics for large-scale image classification,” *Pattern Recognition*, vol. 67, pp. 97–109, 2017.
- [16] Y Wang X Niu R Jin, Y Dou, “Confusion graph: Detecting confusion communities in large scale image classification,” .
- [17] AC. Berg FF. Li J. Deng, S. Satheesh, “Fast and balanced: Efficient label tree learning for large scale object recognition,” in *Advances in Neural Information Processing Systems*, 2011, pp. 567–575.
- [18] R. Lambiotte E. Lefebvre VD. Blondel, JL. Guillaume, “Fast unfolding of community hierarchies in large network, 2008,” *J. Stat. Mech. P*, vol. 1008.
- [19] Yueqing Wang, Xinwang Liu, Yong Dou, Qi Lv, and Yao Lu, “Multiple kernel learning with hybrid kernel alignment maximization,” *Pattern Recognition*, vol. 70, pp. 104–111, 2017.
- [20] Qiang Wang, Yong Dou, Xinwang Liu, Qi Lv, and Shijie Li, “Multi-view clustering with extreme learning machine,” *Neurocomputing*, vol. 214, pp. 483–494, 2016.
- [21] Alexander J. Smola, Peter Bartlett, Bernhard Schölkopf, and Dale Schuurmans, “Probabilities for sv machines,” in *Advances in Large Margin Classifiers*, 2000, pp. 61–74.
- [22] A. Krizhevsky, “Learning multiple layers of features from tiny images,” 2009.

- [23] J. Donahue S. Karayev J. Long Y. Jia, E. Shelhamer,
“Caffe: Convolutional architecture for fast feature embedding,” pp. 675–678, 2014.

Appendices

A. THEORETICAL ANALYSIS OF THE EFFECT ON ACCURACY WHEN TREE STRUCTURE CHANGED

We consider a problem of three categories classification. We assume that the categories are A, B and C and we can construct four different tree structures, which are shown in Fig.1. According to the definition of the Label Tree, tree structure T_1 is more reasonable than others for the classification task.

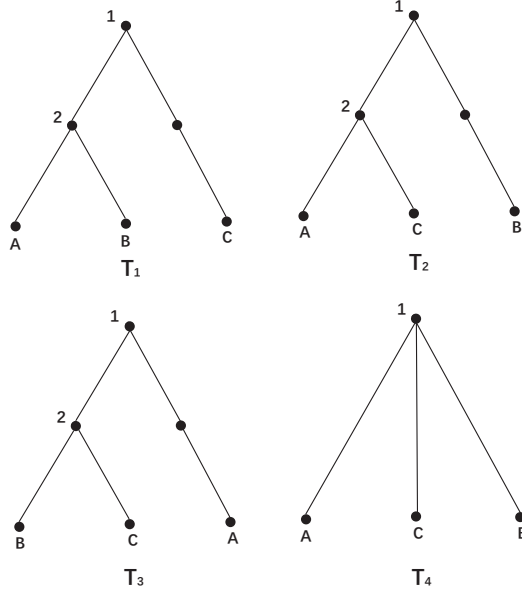


Fig. 1. Four tree structures.

In order to prove it, we assume that we have trained classifiers on these tree structures. At each node of these tree structures, we use SVM for classifier and the SVM classifiers on these nodes are same. Here we assume the distances between every pair of three categories are d_{AB} , d_{AC} , d_{BC} and $d_{BC} > d_{AC} \gg d_{AB}$, which means category A is similar to B while C is different from both of them. Ideally, the distances in high feature space which is projected by classifier (distance of SVM) among these categories is shown in Fig. 2.

We will prove the tree structure T_1 is the most ideal tree structure. We choose one of T_2 and T_3 to be compared with T_1 because T_2 and T_3 are actually the same.

Proposition A.1. *Tree structure T_1 is better than T_2 and T_4 .*

Proof. For T_1 , we assume that P_{AB}^1 denotes the probability of classifying samples in A and B from all the samples correctly at the SVM on the node 1. Then we define P_C^1 , P_A^2 , P_B^2 and so on. In addition, we define P_A , P_B and P_C as the probability of classifying the three categories from all the samples correctly. Here we know:

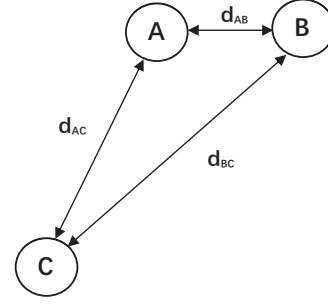


Fig. 2. The distances among category A, B and C.

$$\begin{aligned} P_A &= P_{AB}^1 \times P_A^2 \\ P_B &= P_{AB}^1 \times P_B^2 \\ P_C &= P_C^1 \end{aligned} \quad (1)$$

Because the probability that SVM makes a correct classification is in proportion to the distance of SVM. So we get:

$$\begin{aligned} P_A &\propto d_{AC} \times d_{AB} \\ P_B &\propto d_{AC} \times d_{AB} \\ P_C &\propto d_{AC} \end{aligned} \quad (2)$$

Similarly, for T_2 , we get:

$$\begin{aligned} P_A &\propto d_{AB} \times d_{AC} \\ P_B &\propto d_{AB} \\ P_C &\propto d_{AB} \times d_{AC} \end{aligned} \quad (3)$$

For T_4 , we get:

$$\begin{aligned} P_A &\propto d_{AB} \\ P_B &\propto d_{AB} \\ P_C &\propto d_{AC} \end{aligned} \quad (4)$$

And we know the probability of correct classification for a tree structure T can be defined as:

$$P_T = P_A + P_B + P_C \quad (5)$$

The probabilities for T_1 , T_2 and T_4 are:

$$\begin{aligned} P_{T_1} &\propto (2d_{AC} \times d_{AB} + d_{AC}) \\ P_{T_2} &\propto (2d_{AC} \times d_{AB} + d_{AB}) \\ P_{T_3} &\propto (2d_{AB} + d_{AC}) \end{aligned} \quad (6)$$

The SVM classifiers on the nodes of a tree structure are same so the probabilities in Eq.(6) have a same proportion, we denote it as k :

$$\begin{aligned} P_{T_1} &= k(2d_{AC} \times d_{AB} + d_{AC}) \\ P_{T_2} &= k(2d_{AC} \times d_{AB} + d_{AB}) \\ P_{T_3} &= k(2d_{AB} + d_{AC}) \end{aligned} \quad (7)$$

For P_{T_1} and P_{T_2} :

$$\begin{aligned}
P_{T_1} - P_{T_2} &= k(d_{AC} - d_{AB}) \\
\because d_{AC} &> d_{AB} \\
\therefore P_{T_1} &> P_{T_2}
\end{aligned} \tag{8}$$

For P_{T_1} and P_{T_4} :

$$P_{T_1} - P_{T_4} = k(2d_{AB}(d_{AC} - 1)) \tag{9}$$

If $d_{AC} > 1$ then $P_{T_1} > P_{T_4}$, which demands d_{AC} is very large. From the assumption we know that category C is far from category B and C in the distance of SVM, so $P_{T_1} > P_{T_4}$.

In summary, Tree structure T_1 is better than T_2 , T_3 and T_4 . \square

B. THEORETICAL ANALYSIS OF THE SPEEDUP RATIO ON REPLACING FULLY-CONNECTED LAYERS WITH THE TREE CLASSIFIER

In this section we talk about the speedup ratio of our method in theory. Here we take AlexNet on *CIFAR-100* and *ImageNet* datasets as the analysis object and we replace fully-connected(FC) 7 and FC8 layers with the tree classifier. As we all know, each layer of the FC layers is actually a vector inner product process, which can be defined as:

$$Out_m = \sum_{i=1}^N C(w_m, f_i) + b_m, m = 1, 2, \dots, M \tag{10}$$

with:

$$C(w, f) = \sum_{i=1}^n \sum_{j=1}^n w_{ij} f_{ij} \tag{11}$$

where Out is the output of the FC layer, w is the convolution kernel, f is the input of the FC layer, b is the bias, their subscript m denotes that it is the output of the m th channel and there are M channels in total. C means convolution operation and $C(w, f)$ denotes a convolution operation between kernel w with a size of $n \times n$ and a feature map f . As for FC7 and FC8, the inputs both are feature maps with a size of $1 \times 1 \times 4096$ and the outputs are a feature map with the size of $1 \times 1 \times 4096$ and a score vector with the size of $1 \times 1 \times 1000$. We can calculate the computation using the equation mentioned above.

If we replace the FC layers with our tree classifier, then we should calculate the computation of the tree classifier. The tree classifiers have a hierarchical structure. For one node there is a classifier on each of the child nodes under it. Each classifier under the same parent node computes a result of one test image and the parent node selects which branch to go by comparing all these results. We repeat this process from root node to leaf nodes. Finally the category on the selected leaf node is the classification result. Therefore, the number(N) of classifiers involved in the classification process is equal to

Classifier	CIFAR-100	ImageNet
FC	34	42
Tree Classifier	0.14	0.52
Speedup	233×	78×

Table 1. Comparison of the computation on different datasets in theory(unit:million mult-adds)

the sum of the number of child nodes under all the nodes in the path from the root node to a specific leaf node. Here we find that the worst situation is the path with the most child nodes and we find this specific number of CIFAR-100 and ImageNet is 18 and 65. The dimension(d) of the features which is used for the classifier is 4096. Therefore, we can calculate the multi-adds computation of the tree classifier on *CIFAR-100* dataset is:

$$2Nd = 2 \times 18 \times 4096 = 147456 \tag{12}$$

Similarly, we can calculate the computation on *ImageNet* dataset is 524288.

We make a statistic comparison in Table 1.