Understanding the Dynamics of DNNs Using Graph Modularity

Yao Lu¹[®], Wen Yang¹[®], Yunzhe Zhang¹[®], Zuohui Chen¹[®], Jinyin Chen¹[®], Qi Xuan¹ [∞][®], Zhen Wang² [∞][®], and Xiaoniu Yang^{1,3}[®]

¹ Institute of Cyberspace Security, Zhejiang University of Technology, Hangzhou, 310023. China {yaolu.zjut,czuohui}@gmail.com, {chenjinyin,xuanqi}@zjut.edu.cn, wenyang.zjut@outlook.com, xsgxlz@live.cn
² School of Artificial Intelligence, Optics and Electronics (iOPEN), Northwestern Polytechnical University, Xi'an 710072. China

zhenwang0@gmail.com ³ Science and Technology on Communication Information Security Control Laboratory, Jiaxing 314033, China yxn2117@126.com

Abstract. There are good arguments to support the claim that deep neural networks (DNNs) capture better feature representations than the previous hand-crafted feature engineering, which leads to a significant performance improvement. In this paper, we move a tiny step towards understanding the dynamics of feature representations over layers. Specifically, we model the process of class separation of intermediate representations in pre-trained DNNs as the evolution of communities in dynamic graphs. Then, we introduce modularity, a generic metric in graph theory, to quantify the evolution of communities. In the preliminary experiment, we find that modularity roughly tends to increase as the layer goes deeper and the degradation and plateau arise when the model complexity is great relative to the dataset. Through an asymptotic analysis, we prove that modularity can be broadly used for different applications. For example, modularity provides new insights to quantify the difference between feature representations. More crucially, we demonstrate that the degradation and plateau in modularity curves represent redundant layers in DNNs and can be pruned with minimal impact on performance, which provides theoretical guidance for layer pruning. Our code is available at https://github.com/yaolu-zjut/Dynamic-Graphs-Construction.

Keywords: interpretability, modularity, layer pruning

1 Introduction

DNNs have gained remarkable achievements in many tasks, from computer vision [21,47,38] to natural language processing [54], which can arguably be attributed



Fig. 1. (Left): t-SNE outputs for CIFAR-10 testing data after each layer in VGG11. (Right): modularity curve of VGG11 on CIFAR-10. Best viewed in color.

to powerful feature representations learned from data [48,17]. Giving insight into DNNs' feature representations is helpful to better understand neural network behavior, which attracts much attention in recent years. Some works seek to characterize feature representations by measuring similarities between the representations of various layers and various trained models [25,43,46,39,59,56,14]. Others visualize the feature representations in intermediate layers for an intuitive understanding, revealing that the feature representations in the shallow layers are relatively general, while those in the deep layers are more specific [7,72,50,58,36]. These studies are insightful, but fundamentally limited, because they ignore the dynamics of DNNs or can only understand the dynamics of DNNs through qualitative visualization instead of quantitative study.

Hence, in this paper, we build upon previous studies and investigate the dynamics of intermediate layers. The left part of Fig. 1 shows t-SNE [34] outputs for 500 CIFAR-10 testing samples after each convolutional layer in VGG11, from which we are able to see how class separation in the feature representations progresses as the layer goes deeper. Inspired by this, we seek to quantify the process of class separation of intermediate representations for better understanding the dynamics of DNNs. Specifically, we treat each sample as a node, and there is an edge between two nodes if their feature representations are similar in the corresponding layer. Then we construct a series of graphs that share the same nodes, which can be modeled as a dynamic graph due to the feature continuity. In this way, we convert quantifying the process of class separation of intermediate representations to investigate the evolution of communities in dynamic graphs. Then, we introduce modularity to quantify the evolution of communities. As shown in the right part of Fig. 1, the value of modularity indeed grows with the depth, which is consistent with the process of class separation shown in the left part of Fig. 1. This indicates that modularity provides a quantifiable interpretation perspective for understanding the dynamics of DNNs.

Then we conduct systematic experiments on exploring how the modularity changes in different scenarios, including the training process, standard and adversarial scenarios. Through further analysis of the modularity, we provide two application scenarios for it: (i) representing the difference of various layers. (ii) providing theoretical guidance for layer pruning. To summarize, we make the following contributions:

- We model the class separation of feature representations from layer to layer as the evolution of communities in dynamic graphs, which provides a novel perspective for researchers to better understand the dynamics of DNNs.
- We leverage modularity to quantify the evolution of communities in dynamic graphs, which tends to increase as the layer goes deeper, but descends or reaches a plateau at particular layers. To preserve the generality of modularity, systematic experiments are conducted in various scenarios, e.g., standard scenarios, adversarial scenarios and training processes.
- Additional experiments show that modularity can also be utilized to represent the difference of various layers, which can provide insights on further theoretical analysis and empirical studies.
- Through further analysis on the degradation and plateau in the modularity curves, we demonstrate that the degradation and plateau reveal the redundancy of DNNs in depth, which provides a theoretical guideline for layer pruning. Extensive experiments show that layer pruning guided by modularity can achieve a considerable acceleration ratio with minimal impact on performance.

2 Related Work

Many researchers have proposed various techniques to analyze certain aspects of DNNs. Hence, in this section, we would like to provide a brief survey of the literature related to our current work.

Understanding feature representations. Understanding feature representations of a DNN can obtain more information about the interaction between machine learning algorithms and data than the loss function value alone. Previous works on understanding feature representations can be mainly divided into two categories. One category quantitatively calculates the similarities between the feature representations of different layers and models [25,43,46,39,59,56,14]. For example, Kornblith et al. [25] introduce centered kernel alignment (CKA) to measure the relationship between intermediate representations. Feng et al. [14] propose a new metric, termed as transferred discrepancy, to quantify the difference between two representations based on their downstream-task performance. Compared to previous studies which only utilize feature vectors, Tang et al. [56] leverage both feature vectors and gradients into designing the representations of DNNs. On the basis of [25], Nguyen et al. [43] utilize CKA to explore how varying depth and width affects model feature representations, and find that overparameterized models exhibit the block structure. Through further analysis, they show that some layers exhibit the block structure and can be pruned

with minimal impact on performance. Another category attempts to obtain an insightful understanding of feature representations through interpreting feature semantics [7,72,50,58,36]. Wang et al. [58] and Zeiler et al. [72] discover the hierarchical nature of the features in the neural networks. Specifically, shallow layers extract basic and general features while deep layers learn more specifically and globally. Yosinski et al. [68] quantify the degree to which a particular layer is general or specific. Besides, Donahue et al. [11] investigate the transfer of feature representations from the last few layers of a DNN to a novel generic task.

Modularity and community in DNNs. Previous empirical studies have explored modularity and community in neural networks. Some seek to investigate learned modularity and community structure at the neuron level[63,64,62,8,22,69] or at the subnetwork level [6,27]. Others train an explicitly modular architecture [2,19] or promote modularity via parameter isolation [24] or regularization [9] during training to develop more modular neural networks. Different from these existing works, in this paper, we explore the evolution of communities at the feature representation level, which provides a new perspective to characterize the dynamics of DNNs.

Layer pruning. State-of-the-art DNNs often involve very deep networks, which are bound to bring massive parameters and floating-point operations. Therefore, many efforts have been made to design compact models. Pruning is one stream among them, which can be roughly devided into three categories, namely, weight pruning [16,3,31], filter pruning [30,29,57] and layer pruning [66,13,73,5,60,61]. Weight pruning compresses over-parameterized models by dropping redundant individual weights, which has limited applications on general-purpose hardware. Filter pruning seeks to remove entire redundant filters or channels instead of individual weights. Compared to weight pruning and filter pruning, layer pruning removes the entire redundant layers, which is more suitable for general-purpose hardware. Existing layer pruning methods mainly differ in how to determine which layers need to be pruned. For example, Xu et al. [66] first introduce a trainable layer scaling factor to identify the redundant layers during the sparse training. And then they prune the layers with very small factors and retrain the model to recover the accuracy. Elkerdawy et al. [13] leverage imprinting to calculate a per-layer importance score in one-shot and then prune the least important layers and fine-tune the shallower model. Zhou et al. [73] leverage the ensemble view of block-wise DNNs and employ the multi-objective optimization paradigm to prune redundant blocks while avoiding performance degradation. Based on the observations of [1], Chen and Zhao [5], Wang et al. [60], and Wang et al. [61], respectively, utilize linear classifier probes to guide the layer pruning. Specifically, they prune the layers which provide minor contributions on boosting the performance of features.

Adversarial samples. Although DNNs have gained remarkable achievements in many tasks [21,47,38,54], they have been found vulnerable to adversarial examples, which are born of intentionally perturbing benign samples in a human-imperceptible fashion [55,18]. The vulnerability to adversarial examples hinders DNNs from being applied in safety-critical environments. Therefore, at-



Fig. 2. A dynamic graph with 4 snapshots. Nodes of the same color represent that they are in the same community and the thickness of the line represents the weight of the edge. we use Eq. (1) to calculate the modularity.

tacks and defenses on adversarial examples have attracted significant attention in machine learning. There has been a multitude of work studying methods to obtain adversarial examples [18,35,44,4,53,40,37,28] and to make DNNs robust against adversarial examples [51,65].

3 Methodology

3.1 Preliminary

We start by introducing some concepts in graph theory.

Dynamic graphs are the graphs that change with time [20]. In this paper, we focus on edge-dynamic graphs, i.e., edges may be added or deleted from the graph. Given a dynamic graph: $\mathcal{DG} = \{G_1, G_2, \dots, G_T\}$, where T is the number of snapshots. Repeatedly leveraging the static methods on each snapshot can collectively look insight into the graph's dynamics. **Communities**, which are quite common in many real networks [70,33,23], are defined as sets of nodes that are more densely connected internally than externally [41,42]. **Ground-truth communities** can be defined as the sets of nodes with common properties, e.g., common attribute, affiliation, role, or function [67]. **Modularity**, as a common measurement to quantify the quality of communities [15,42], carries advantages including intelligibility and adaptivity. In this paper, we adopt modularity Q as following:

$$Q = \frac{1}{2W} \sum_{ij} \left(a_{ij} - \frac{s_i s_j}{2W} \right) \delta(c_i, c_j), \tag{1}$$

where a_{ij} is the weight of the edge between node *i* and node *j*, $W = \sum_i \sum_j a_{ij}$ is the sum of the weights of all edges, which is used as a normalization factor. $s_i = \sum_j a_{ij}$ and $s_j = \sum_i a_{ij}$ are the strength of nodes *i* and *j*, c_i and c_j denote the community that nodes *i* and *j* belong to, respectively. $\delta(c_i, c_j)$ is 1 if node *i* and node *j* are in the same community and 0 otherwise. In this paper, each node represents an image with the corresponding label. Hence, these nodes can be divided into the corresponding ground-truth communities, which we utilize to calculate the modularity. Fig. 2 gives an example to intuitively understand the evolution of communities, from which we find that the communities with a high value of modularity tend to strengthen the intra-community connections and weaken the inter-community connections.



Fig. 3. Pipeline for the dynamic graph construction and the application scenarios of the modularity metric. Best viewed in color, zoomed in for details.

3.2 Dynamic Graph Construction

Our proposed dynamic graph construction framework to understand the dynamics of a given DNN is visually summarized in Fig. 3.

Considering a typical image classification problem, the feature representation of the sample is transformed over the course of many layers, to be finally utilized by a classifier at the last layer. Therefore, we can model this process as follows:

$$\tilde{y} = f_{l+1}(f_l(\dots f_1(x)\dots)), \tag{2}$$

where $f_1(\cdot), f_2(\cdot), \cdots, f_l(\cdot)$ are *n* functions to extract feature representations of layers from bottom to top, $f_{l+1}(\cdot)$ and l denote the classifier and the number of layers, respectively. In this paper, we treat a bottleneck building block or a building block as a layer for ResNets [21] and take a sequence of consecutive layers, e.g., Conv-BN-ReLu, as a layer for VGGs [52]. We randomly sample Nsamples $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ with corresponding labels $\mathcal{Y} = \{y_1, y_2, \dots, y_N\}$ from the test set and feed them into a well-optimized DNN with fixed parameters to obtain an intermediate representation set $\mathcal{R} = \{r_1, r_2, \cdots, r_l\}$, where $r_i \in \mathbb{R}^{N \times C_i \times W_i \times H_i}$ denotes the feature representations in *i*-th layer, C_i , W_i and H_i are the number of channels, width and height of feature maps in *i*-th layer, respectively. Then we apply a flatten mapping $f : \mathbb{R}^{N \times C \times W \times H} \to \mathbb{R}^{N \times M}$ to each element in \mathcal{R} , where $M = C \times W \times H$. In order to capture the underlying relationship of samples in feature space, we construct a series of k-nearest neighbor (k-NN) graphs $G_i = (A_i, r_i)$, where A_i is the adjacency matrix of the k-NN graph in *i*-th layer. Specifically, we calculate the similarity matrix $S = \mathbb{R}^{N \times N}$ among N nodes using Eq. (3).

$$S_{jk}^{i} = \begin{cases} \frac{r_{i,j}^{T} r_{i,k}}{||r_{i,j}|| \ ||r_{i,k}||}, & \text{if } j \neq k \\ 0, & \text{if } j = k \end{cases},$$
(3)

where $r_{i,j}$ and $r_{i,k}$ are the feature representation vectors of samples j and k in *i*-th layer. According to the obtained similarity matrix S_{ik}^{i} , we choose top

k similar node pairs for each node to set edges and corresponding weights, so as to obtain the adjacency matrix A_i . Now, we obtain a series of k-NN graphs $\{G_i = (A_i, r_i) | i = 1, 2, \dots, l\}$, which reveal the internal relationship between feature representations of different samples in various layers. Due to the continuity of feature representations, i.e., feature representation of the current layer is obtained on the basis of the previous one. Hence, these k-NN graphs are relevant to each other and can be treated as multiple snapshots of the dynamic graph $\mathcal{DG} = \{G_1, G_2, \dots, G_l\}$ at different time intervals.

Then we revisit the dynamic graph for a better understanding of community evolution. Fig. 2 exhibits a demo, from which we can intuitively understand the process of community evolution in the dynamic graph. \mathcal{DG} consists of l snapshots, each snapshot shares same nodes (samples) and reveals the inherent correlations between samples in the corresponding layer. Therefore, our dynamic graph is actually an edge-dynamic graph. Due to the existence of ground-truth label of each sample, we can easily divide samples into K ground-truth communities. Hence, we can calculate the modularity of each snapshot in \mathcal{DG} using Eq. (1) together with the ground-truth communities. According to the obtained modularity of each snapshot, we finally obtain a modularity set $\mathcal{Q} = \{Q_1, Q_2, \dots, Q_l\}$, which reveals the evolution of communities in the dynamic graph.

4 Experiments

Our goal is to intuitively understand the dynamics in well-optimized DNNs. Reflecting this, our experimental setup consists of a family of VGGs [52] and ResNets [21] trained on standard image classification datasets CIFAR-10 [26], CIFAR-100 [26] and ImageNet [49] (For the sake of simplicity, we choose 50 classes in original ImageNet, termed as ImageNet50). Specifically, we leverage stochastic gradient descent algorithm with an initial learning rate of 0.01 to optimize the model. The batch size, weight decay, epoch and momentum are set to 256, 0.005, 150 and 0.9, respectively. The statistics of pre-trained models and ImageNet50 are shown in Appendix B. All experiments are conducted on two NVIDIA Tesla A100 GPUs. If not special specified, we set k = 3, N = 500 for CIFAR-10 and ImageNet50, k = 3, N = 1000 for CIFAR-100 to construct dynamic graphs.

4.1 Understanding the Evolution of Communities

We systematically investigate the modularity curves of various models and repeat each experiment 5 times to obtain the mean and variance of modularity curves. From the results reported in Fig. 4, we have the following observations:

- The modularity roughly tends to increase as the layer goes deeper.
- On the same dataset, the frequency of degradation and plateau existing in the modularity curve gradually increases as models get deeper. Specifically, the modularity curve gradually reaches a plateau on CIFAR-10 and CIFAR-100



Fig. 4. Modularity curves of different models. C10, C100 and I50 denote CIFAR-10, CIFAR-100 and ImageNet50, respectively. Shaded regions indicate standard deviation.



Fig. 5. Modularity curves of ResNets on ImageNet.

at the deep layer as VGG gets deeper. Compared to VGGs, the modularity curve of ResNets descends, reaches a plateau, or rises very slowly mostly happening in the repeatable layers.

• According to the modularity curves of VGG16 and VGG19, we can see that as the complexity of the dataset increases, the plateau gradually disappears.

Since previous works [58,72] have shown that shallow layers extract general features while deep layers learn more specifically. Hence, feature representations of the same category are more similar in deep layers than in shallow layers, which can be seen in Fig. 1. Therefore, the samples in the same community (category) tend to connect with each other, i.e., the modularity increases as the layer goes deeper. In this sense, the growth of modularity quantitatively reflects the process of class separation in feature representations inside the DNNs. Besides, the tendency of modularity is consistent with the observation of [1], which utilizes linear classifier probes to measure how suitable the feature representations at every layer are for classification. Compared to [1] that requires training a linear classifier for each layer, our modularity provides a more convenient and effective tool to understand the dynamics of DNNs.

According to the second and third findings, we can draw a conclusion that the degradation and plateau arise when model complexity is great relative to the dataset. With the relative complexity of the model getting greater, a little bit



Fig. 6. Modularity curves in different scenarios.



(a) The influence of N on the modularity. (b) The influence of k on the modularity.

Fig. 7. Hyperparameter sensitivity analysis.

of performance improvement costs nearly doubling the number of layers, e.g., ResNet18 and ResNet34 on CIFAR-10. Many successive layers are essentially performing the same operation over many times, refining feature representations just a little more each time instead of making a more fundamental change of feature representations. Hence, these features representations exhibit a similar degree of class separation, which explains the existence of the plateau and degradation in the modularity curves. To further explore the relationship between the degradation as well as plateau and model relative complexity, we use pre-trained ResNets in torchvision¹ to conduct experiments on ImageNet. From Fig. 5 we find that the tendency of modularity curves is almost consistent with the observation in Fig. 4(b). The main difference lies in the modularity curves of ResNet101 and ResNet152 trained on ImageNet50 exhibit the degradation and plateau in the middle layer, while those trained on ImageNet do not, which further confirms the conclusion we make. Additional experiment on exploring the evolution of communities is presented in Appendix A.1.

4.2 Modularity Curves of Adversarial Samples

In Section 4.1, we discuss the modularity curves of normal samples. Here, we would like to explore the modularity curves of adversarial samples. Specifically,

¹ https://pytorch.org/vision/stable/index.html

we first utilize FGSM [18], PGD [35], JSMA [44], CW [4], OnePixel [53] and Local Search [40] to attack the pre-trained VGG16 and ResNet18 on CIFAR-10 for generating adversarial samples. To make it easier for others to reproduce our results, we utilize default parameter settings in AdverTorch [10] to obtain untargeted adversarial samples. The attack success rate is 100% for each attack mode. Then we randomly choose 500 adversarial samples in each attack mode, set k = 3 to construct dynamic graphs and plot modularity curves. As shown in Fig. 6(a), we can find that the modularity curve of adversarial samples can reach a smaller peak than normal samples, which can be interpreted as adversarial attacks blur the distinctions among various categories.

4.3 Modularity during Training Time

During the training process, the model is considered to learn valid feature representations. In order to explore the dynamics of DNNs in this process, we conduct experiments on ResNet18 and VGG16. Specifically, in the first 30 training epochs, we save the model file and test accuracy of every epoch. Then we construct a dynamic graph for each model and calculate the modularity of each snapshot in this dynamic graph. Fig. 6(b) shows the modularity and accuracy curves on ResNet18 and VGG16, from which we find that the values of modularity in shallow layers nearly keep constant or small fluctuations. We believe that is because feature representations in shallow layers are general [58,72], samples in the same category do not cluster together. Hence, despite learning effective feature representations, the value of modularity does not increase significantly. Compared to shallow layers, feature representations in deep layers are more specific. With the improvement of model performance, the valid feature representations gradually learned by the model make intra-community connections tighter, which explains the modularity curves of deep layers are almost proportional to the accuracy curve. These phenomena may give some information about how the training evolves inside the DNNs and guide the intuition of researchers. Besides, we provide the experiment on randomly initialized models in Appendix A.2.

4.4 Ablation Study

To provide further insight into the modularity, we conduct ablation studies to evaluate the hyperparameter sensitivity of the modularity. The batch size N and the number of edges k that each node connects with are two hyperparameters.

Whether batch size N has a non-negligible impact on the modularity curve? we set N = 200, 500, 1000, 2000, 5000, k = 3 to repeat experiments on CIFAR-10 with pre-trained VGG16 and ResNet18. In Fig. 7(a), we show that smaller N has relatively lower modularity in early layers but N has less influence on modularity in final layers. Generally speaking, different values of N have the same tendency on modularity curves.

Whether k has a non-negligible impact on the modularity curve? we set N = 500, k = 3, 5, 7, 9, 11 to repeat the above experiments. Fig. 7(b) shows that the different selections of k have a negligible impact on modularity



Fig. 8. Emergence of the block structure with different methods.

because the modularity curves almost overlap together. We observe the similar tendency even when k is large (see Appendix A.3 for details)).

Hence, we can conclude that the modularity is reliable for hyperparameters.

5 Application Scenarios of Modularity

In Section 4, we investigate the dynamics of DNNs, from which we gain some insights. On this basis, we provide two application scenarios for the modularity.

5.1 Representing the Difference of Various Layers

In addition to quantifying the degree of class separation of intermediate representations, modularity can also be used to represent the difference of various layers like previous works [25,43,46,39,59,56,14]. Instead of directly calculating the similarity between two feature representations, we compute the difference of modularity of two layers because our modularity reflects the global attribute of the corresponding layer, i.e., the degree of class separation. Specifically, we calculate the difference matrix D, with its element defined as $D_{ij} = |Q_i - Q_j|$. Q_i and Q_j denote the value of modularity in *i*-th and *j*-th layer, respectively. We visualize the result as a heatmap, with the x and y axes representing the layers of the model. As shown in Fig. 8(a), the heatmap shows a block structure in representational difference, which arises because representations after residual connections are more different from representations inside ResNet blocks than other post-residual representations. Moreover, we also reproduce the result of linear CKA [25] in Fig. 8(b), from which we can see the same block structure. Note that we measure the difference of various layers while [25] focuses on the similarity, so the darker the color in Fig. 8(b), the less similar.

5.2 Guiding Layer Pruning with Modularity

In Section 4.1, we find that the degradation and plateau arise when model complexity is great relative to the dataset. Previous works have demonstrated that DNNs are redundant in depth [71] and overparameterized models exist many consecutive hidden layers that have highly similar feature representations [43].



Fig. 9. The modularity curves of VGGs and ResNets of different depths.



Fig. 10. Experiments on pruning the irredundant layers.

Informed by these conclusions, we wonder if the degradation and plateau offer an intuitive instruction in identifying the redundant layers. Hence, we assume that the plateau and degradation make no contribution or negative contribution to the model. In other words, these layers we consider are redundant and can be pruned with acceptable loss. To verify our assumption, we conduct systematic experiments on CIFAR-10, CIFAR-100 and ImageNet50.

Pruning the redundant layers. Since the plateau mostly appears in the last few layers of VGG19, we remove the Conv-BN-ReLu in VGG19 one by one from back to front to obtain a series of variants. As for ResNet152, we remove the bottleneck building blocks from back to front in stage 3 (the degradation and plateau mostly emerge in stage 3) to get variant models. Detailed structures are shown in Appendix B. For example, VGG19_1 denotes VGG19 prunes 1 layer. Then we finetune these variant models with the same hyperparameter setting as section 4.1. The left part of Fig. 9 shows the results of VGGs on CIFAR-10, from which we find that the modularity curves of different VGG almost overlap together. The only difference between these modularity curves is whether the plateau emerges in the last few layers. Specifically, the modularity curve of VGG19_6 does not have the plateau, while VGG19_1 has the obvious plateau. Moreover, the plateau gradually disappears as the layer is removed one by one from back to front. Note that these variant models have similar performance, which proves that the plateau is indeed redundant. The middle part and right

part of Fig. 9 show the modularity curves of different ResNet on CIFAR-10 and ImageNet50. These modularity curves almost coincide in the shallow layers, while the plateau gradually narrows with the continuous removal of the middle layers. Consequently, this strongly proves that the plateau can be pruned with minimal impact on performance.

Pruning the irredundant layers. In the previous paragraph, we verify that removing redundant layers will not affect the accuracy. Here, we wonder whether removing irredundant layers will result in a significant performance drop. Hence, we conduct further experiments on CIFAR-10 and CIFAR-100 with VGG11 (According to the modularity curves in Fig. 4, we think VGG11 is relatively irredundant on CIFAR-10 and CIFAR-100). We remove the Conv-BN-ReLu in VGG11 one by one from back to front to obtain a series of variants and finetune them. Fig. 10(a) exhibits the modularity curves of those variants, from which we can see that they can finally reach almost the same peak. Next we calculate the corresponding variation of accuracy brought about by pruning the layer. Specifically, we calculate the variation of modularity ΔQ of the final layer with $\Delta Q = Q_i - Q_{i-1}$, where Q_i denotes the modularity of *i*-th layer that we want to prune. Finally, we calculate the variation of accuracy ΔAcc using $\Delta Acc = Acc_i - Acc_j$, where Acc_i and Acc_j represent the accuracy of the original model and pruned model, respectively. Fig. 10(b) shows the results, from which we can find that on CIFAR-10, removing the layer that modularity increases 0.01 results in nearly no influence (0.01%) on model performance, while pruning the layer that has a 0.16 increment on modularity results in 0.57% degradation. With the complexity of the dataset increasing, this gap becomes more obvious. On CIFAR-100, pruning the layer that modularity goes up 0.08 leads to a 0.37%drop in accuracy, while a variation of 0.26 in modularity causes a variation of 2.11% in accuracy. Hence, we draw a conclusion that the variation of accuracy is proportional to the variation of modularity, which means pruning irredundant layers will result in a more significant drop in performance than removing redundant layers. Besides, This phenomenon becomes more obvious as the complexity of the dataset increases.

Practicality of layer pruning by modularity. According to the above experimental results, we are able to conclude that modularity can be used to provide effective theoretical guidance for layer pruning. Here, we would like to evaluate its practicality. Specifically, we first plot the modularity curve of the original model, then we prune the layer where the curve drops, reaches a plateau or grows slowly, finally we finetune the new model. We adapt number of parameters and required Float Points Operations (denoted as FLOPs), to evaluate model size and computational requirement. We leverage a package in pytorch[45], which terms thop² to calculate FLOPs and parameters. Table 1 shows the performance of different layer pruning methods [5,61] on ResNet56 for CIFAR. Compared with Chen et al, our method provides considerable better parameters and FLOPs reductions (43.00% vs. 42.30%, 60.30% vs. 34.80%), while yielding a higher accuracy (93.38% vs. 93.29%). Compared to DBP-0.5,

² https://github.com/Lyken17/pytorch-OpCounter

Method	Top-1%	Params(PR)	$\mathrm{FLOPs}(\mathrm{PR})$
ResNet56	93.27	0%	0%
Chen et al $[5]$	93.29	42.30%	34.80%
DBP-0.5 [61]	93.39	/	53.41%
Ours	93.38	43.00%	60.30%

Table 1. Pruning results of ResNet56 on CIFAR-10. PR is the pruning rate.

our method shows more advantages in FLOPs reduction (60.30% vs. 53.41%), while maintaining a competitive accuracy (93.38% vs. 93.39%).

According to the above experiments, we demonstrate the effectiveness and efficiency of layer pruning guided by modularity.

6 Conclusion and Future Work

In this study, through modeling the process of class separation from layer to layer as the evolution of communities in dynamic graphs, we provide a graph perspective for researchers to better understand the dynamics of DNNs. Then we develop modularity as a conceptual tool and apply it to various scenarios, e.g., the training process, standard and adversarial scenarios, to gain insights into the dynamics of DNNs. Extensive experiments show that modularity tends to rise as the layer goes deeper, which quantitatively reveals the process of class separation in intermediate layers. Moreover, the degradation and plateau arise when model complexity is great relative to the dataset. Through further analysis on the degradation and plateau at particular layers, we demonstrate that modularity can provide theoretical guidance for layer pruning. In addition to guiding layer pruning, modularity can also be used to represent the difference of various layers.

We hope the simplicity of our dynamic graph construction approach could facilitate more research ideas in interpreting DNNs from a graph perspective. Besides, we wish that the modularity presented in this paper can make a tiny step forward in the direction of neural network structure design, layer pruning and other potential applications. Recent work has shown that Vision Transformers can achieve superior performance on image classification tasks [32,12]. In the future, we will further explore the dynamics of Visual Transformers.

Acknowledgements. This work was supported in part by the Key R&D Program of Zhejiang under Grant 2022C01018, by the National Natural Science Foundation of China under Grants U21B2001, 61973273, 62072406, 11931015, U1803263, by the Zhejiang Provincial Natural Science Foundation of China under Grant LR19F030001, by the National Science Fund for Distinguished Young Scholars under Grant 62025602, by the Fok Ying-Tong Education Foundation, China under Grant 171105, and by the Tencent Foundation and XPLORER PRIZE. We also sincerely thank Jinhuan Wang, Zhuangzhi Chen and Shengbo Gong for their excellent suggestions.

References

- Alain, G., Bengio, Y.: Understanding intermediate layers using linear classifier probes. arXiv preprint arXiv:1610.01644 (2016)
- Alet, F., Lozano-Pérez, T., Kaelbling, L.P.: Modular meta-learning. In: Conference on Robot Learning. pp. 856–868 (2018)
- Azarian, K., Bhalgat, Y., Lee, J., Blankevoort, T.: Learned threshold pruning. arXiv preprint arXiv:2003.00075 (2020)
- 4. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: 2017 ieee symposium on security and privacy (sp). pp. 39–57 (2017)
- Chen, S., Zhao, Q.: Shallowing deep networks: layer-wise pruning based on feature representations. IEEE Transactions on Pattern Analysis and Machine Intelligence 41(12), 3048–3056 (2018)
- Csordás, R., van Steenkiste, S., Schmidhuber, J.: Are neural nets modular? inspecting their functionality through differentiable weight masks. In: International Conference on Learning Representations (2021)
- Das, A., Rad, P.: Opportunities and challenges in explainable artificial intelligence (xai): A survey. arXiv preprint arXiv:2006.11371 (2020)
- Davis, B., Bhatt, U., Bhardwaj, K., Marculescu, R., Moura, J.M.: On network science and mutual information for explaining deep neural networks. In: IEEE International Conference on Acoustics, Speech and Signal Processing. pp. 8399– 8403 (2020)
- Delange, M., Aljundi, R., Masana, M., Parisot, S., Jia, X., Leonardis, A., Slabaugh, G., Tuytelaars, T.: A continual learning survey: Defying forgetting in classification tasks. IEEE Transactions on Pattern Analysis and Machine Intelligence (2021)
- Ding, G.W., Wang, L., Jin, X.: AdverTorch v0.1: An adversarial robustness toolbox based on pytorch. arXiv preprint arXiv:1902.07623 (2019)
- Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., Darrell, T.: Decaf: A deep convolutional activation feature for generic visual recognition. In: International conference on machine learning. pp. 647–655. PMLR (2014)
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (2021)
- Elkerdawy, S., Elhoushi, M., Singh, A., Zhang, H., Ray, N.: To filter prune, or to layer prune, that is the question. In: Proceedings of the Asian Conference on Computer Vision (2020)
- 14. Feng, Y., Zhai, R., He, D., Wang, L., Dong, B.: Transferred discrepancy: Quantifying the difference between representations. arXiv preprint arXiv:2007.12446 (2020)
- 15. Fortunato, S.: Community detection in graphs. CoRR abs/0906.0612 (2009)
- 16. Frankle, J., Carbin, M.: The lottery ticket hypothesis: Finding sparse, trainable neural networks. In: International Conference on Learning Representations (2019)
- Goh, G., Cammarata, N., Voss, C., Carter, S., Petrov, M., Schubert, L., Radford, A., Olah, C.: Multimodal neurons in artificial neural networks. Distill 6(3), e30 (2021)
- 18. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. In: International Conference on Learning Representations (2015)
- Goyal, A., Lamb, A., Hoffmann, J., Sodhani, S., Levine, S., Bengio, Y., Schölkopf, B.: Recurrent independent mechanisms. In: International Conference on Learning Representations (2021)

- 16 Y. Lu et al.
- Harary, F., Gupta, G.: Dynamic graph models. Mathematical and Computer Modelling 25(7), 79–87 (1997)
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 770–778 (2016)
- 22. Hod, S., Casper, S., Filan, D., Wild, C., Critch, A., Russell, S.: Detecting modularity in deep neural networks. arXiv preprint arXiv:2110.08058 (2021)
- Jonsson, P.F., Cavanna, T., Zicha, D., Bates, P.A.: Cluster analysis of networks generated through homology: automatic identification of important protein communities involved in cancer metastasis. BMC bioinformatics 7(1), 1–13 (2006)
- Kirsch, L., Kunze, J., Barber, D.: Modular networks: Learning to decompose neural computation. In: Advances in Neural Information Processing Systems. pp. 2414– 2423 (2018)
- Kornblith, S., Norouzi, M., Lee, H., Hinton, G.: Similarity of neural network representations revisited. In: International Conference on Machine Learning. pp. 3519– 3529 (2019)
- Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009)
- 27. Lake, B.M., Ullman, T.D., Tenenbaum, J.B., Gershman, S.J.: Building machines that learn and think like people. Behavioral and Brain Sciences 40 (2017)
- Li, J., Ji, R., Chen, P., Zhang, B., Hong, X., Zhang, R., Li, S., Li, J., Huang, F., Wu, Y.: Aha! adaptive history-driven attack for decision-based black-box models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 16168–16177 (2021)
- Lin, M., Ji, R., Wang, Y., Zhang, Y., Zhang, B., Tian, Y., Shao, L.: Hrank: Filter pruning using high-rank feature map. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1529–1538 (2020)
- 30. Lin, S., Ji, R., Yan, C., Zhang, B., Cao, L., Ye, Q., Huang, F., Doermann, D.: Towards optimal structured cnn pruning via generative adversarial learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2790–2799 (2019)
- 31. Lin, T., Stich, S.U., Barba, L., Dmitriev, D., Jaggi, M.: Dynamic model pruning with feedback. In: International Conference on Learning Representations (2020)
- 32. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10012–10022 (2021)
- Lusseau, D.: The emergent properties of a dolphin social network. Proceedings of the Royal Society of London. Series B: Biological Sciences 270(suppl_2), S186–S188 (2003)
- Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. Journal of Machine Learning Research 9(11) (2008)
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: International Conference on Learning Representations (2018)
- Mahendran, A., Vedaldi, A.: Understanding deep image representations by inverting them. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5188–5196 (2015)
- 37. Maho, T., Furon, T., Le Merrer, E.: Surfree: a fast surrogate-free black-box attack. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10430–10439 (2021)

- Maqueda, A.I., Loquercio, A., Gallego, G., García, N., Scaramuzza, D.: Eventbased vision meets deep learning on steering prediction for self-driving cars. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5419–5427 (2018)
- Morcos, A.S., Raghu, M., Bengio, S.: Insights on representational similarity in neural networks with canonical correlation. In: Advances in Neural Information Processing Systems. pp. 5732–5741 (2018)
- Narodytska, N., Kasiviswanathan, S.P.: Simple black-box adversarial perturbations for deep networks. arXiv preprint arXiv:1612.06299 (2016)
- Newman, M.E.: Modularity and community structure in networks. Proceedings of the National Academy of Sciences 103(23), 8577–8582 (2006)
- Newman, M.E., Girvan, M.: Finding and evaluating community structure in networks. Physical Review E 69(2), 026113 (2004)
- 43. Nguyen, T., Raghu, M., Kornblith, S.: Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth. In: International Conference on Learning Representations (2021)
- 44. Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z.B., Swami, A.: The limitations of deep learning in adversarial settings. In: 2016 IEEE European symposium on security and privacy (EuroS&P). pp. 372–387 (2016)
- 45. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)
- 46. Raghu, M., Gilmer, J., Yosinski, J., Sohl-Dickstein, J.: SVCCA: singular vector canonical correlation analysis for deep learning dynamics and interpretability. In: Advances in Neural Information Processing Systems. pp. 6076–6085 (2017)
- Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, realtime object detection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 779–788 (2016)
- Rumelhart, D.E., Hinton, G.E., Williams, R.J.: Learning internal representations by error propagation. Tech. rep., California Univ San Diego La Jolla Inst for Cognitive Science (1985)
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M.S., Berg, A.C., Fei-Fei, L.: Imagenet large scale visual recognition challenge. International Journal of Computer Vision 115(3), 211–252 (2015)
- Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Gradcam: Visual explanations from deep networks via gradient-based localization. In: International Conference on Computer Vision. pp. 618–626 (2017)
- Shafahi, A., Najibi, M., Ghiasi, M.A., Xu, Z., Dickerson, J., Studer, C., Davis, L.S., Taylor, G., Goldstein, T.: Adversarial training for free! Advances in Neural Information Processing Systems **32** (2019)
- 52. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: International Conference on Learning Representations (2015)
- Su, J., Vargas, D.V., Sakurai, K.: One pixel attack for fooling deep neural networks. IEEE Transactions on Evolutionary Computation 23(5), 828–841 (2019)
- Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: Advances in Neural Information Processing Systems. pp. 3104–3112 (2014)
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I.J., Fergus, R.: Intriguing properties of neural networks. In: International Conference on Learning Representations (2014)

- 18 Y. Lu et al.
- 56. Tang, S., Maddox, W.J., Dickens, C., Diethe, T., Damianou, A.: Similarity of neural networks with gradients. arXiv preprint arXiv:2003.11498 (2020)
- 57. Tang, Y., Wang, Y., Xu, Y., Deng, Y., Xu, C., Tao, D., Xu, C.: Manifold regularized dynamic network pruning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5018–5028 (2021)
- Wang, F., Liu, H., Cheng, J.: Visualizing deep neural network by alternately image blurring and deblurring. Neural Networks 97, 162–172 (2018)
- Wang, L., Hu, L., Gu, J., Hu, Z., Wu, Y., He, K., Hopcroft, J.E.: Towards understanding learning representations: To what extent do different neural networks learn the same representation. In: Advances in Neural Information Processing Systems. pp. 9607–9616 (2018)
- 60. Wang, W., Chen, M., Zhao, S., Chen, L., Hu, J., Liu, H., Cai, D., He, X., Liu, W.: Accelerate cnns from three dimensions: A comprehensive pruning framework. In: International Conference on Machine Learning. pp. 10717–10726 (2021)
- Wang, W., Zhao, S., Chen, M., Hu, J., Cai, D., Liu, H.: Dbp: Discrimination based block-level pruning for deep model acceleration. arXiv preprint arXiv:1912.10178 (2019)
- Watanabe, C.: Interpreting layered neural networks via hierarchical modular representation. In: International Conference on Neural Information Processing. pp. 376–388 (2019)
- Watanabe, C., Hiramatsu, K., Kashino, K.: Modular representation of layered neural networks. Neural Networks 97, 62–73 (2018)
- Watanabe, C., Hiramatsu, K., Kashino, K.: Understanding community structure in layered neural networks. Neurocomputing 367, 84–102 (2019)
- 65. Wong, E., Rice, L., Kolter, J.Z.: Fast is better than free: Revisiting adversarial training. In: International Conference on Learning Representations (2020)
- Xu, P., Cao, J., Shang, F., Sun, W., Li, P.: Layer pruning via fusible residual convolutional block for deep neural networks. arXiv preprint arXiv:2011.14356 (2020)
- Yang, J., Leskovec, J.: Defining and evaluating network communities based on ground-truth. Knowledge and Information Systems 42(1), 181–213 (2015)
- Yosinski, J., Clune, J., Bengio, Y., Lipson, H.: How transferable are features in deep neural networks? In: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N.D., Weinberger, K.Q. (eds.) Advances in Neural Information Processing Systems. pp. 3320–3328 (2014)
- You, J., Leskovec, J., He, K., Xie, S.: Graph structure of neural networks. In: International Conference on Machine Learning. pp. 10881–10891 (2020)
- Zachary, W.W.: An information flow model for conflict and fission in small groups. Journal of anthropological research 33(4), 452–473 (1977)
- Zagoruyko, S., Komodakis, N.: Wide residual networks. In: Proceedings of the British Machine Vision Conference (2016)
- Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional networks. In: European Conference on Computer Vision. pp. 818–833 (2014)
- 73. Zhou, Y., Yen, G.G., Yi, Z.: Evolutionary shallowing deep neural networks at block levels. IEEE Transactions on Neural Networks and Learning Systems (2021)