

A survey of Knowledge Distillation

Trung Dao
VinAI Research
14 Thuy Khue, Tay Ho, Ha Noi
v.trungdt21@vinai.io

Abstract

For the last few years, deep learning has proved its ability to solve complicated problems that were too cumbersome to handle. It is undeniable that key factors contributing to this success are mainly large-scale data and deep neural networks with billions of parameters. However, to deploy and smoothly run these enormous models on on-edge devices is not always a straightforward process. In order to overcome this, a variety of model compression and acceleration techniques have been developed. As one of the outstanding techniques, knowledge distillation (KD) seeks to transfer information learned from one large model (called teacher) to smaller one (called student). This survey firstly gives a brief overview of what KD is then quickly investigates some of the recent work about training schemes and teacher-student architectures.

1. Introduction

After achieving extraordinary establishment in various fields of Machine Learning using neural networks, scientists began to move interests into network compression and enhancement. Several worth-mentioned approaches which seek to create smaller model and cost-efficient are: parameter pruning/sharing, model quantization, low-rank factorization and knowledge distillation. A comprehensive overview on these approaches is outside of the scope of this survey, and the focus of this paper is solely about knowledge distillation, which has attracted many attentions from the deep learning research community in the recent years.

Comparing to other compression methods, KD might be considered superior since it can squeeze down a network disregarding the structural difference between the teacher and student network. The original idea of KD was proposed in [3], where the author focused on transferring the information from a large model or an ensemble of models into a smaller model without performance loss. This work is later generalized and popularized by Hinton et al. [11]. In their work, they found out an astounding observation that it is

simpler to train a smaller-scale network (classifier) using the soften predictions of another classifier as target value rather than the ground truth (one-hot) labels and later on called this procedure as *distillation* (short for knowledge distillation). They also stated that these probabilities offer richer information than labels alone, and empirically help the student network learn better. Knowledge distillation field since then has been developed and applied in various forms, but the main characteristic of KD is still remained the same: Teacher-Student framework, where the teacher model provides useful *knowledge* to the student model in order to improve its learning performance.

Thanks to its effective on neural network compression and acceleration, KD has been widely applied in different fields of artificial intelligence: visual/speech recognition, natural language processing (NLP). In visual recognition, KD was firstly applied on classification tasks [11, 15, 4, 2, 7] then later on expanded to other applications such as image segmentation [9], lane detection [12], facial landmark detection [26, 8]. In other field such as NLP, KD also proves its worth as being used to compress complex structures such as BERT [21, 14]. To briefly generalized, KD offers not only lightweight deep models, which allows deploying models to on-edge devices efficiently but recently also competitive performant ones.

The main contributions of this paper consist of:

1. Provide a comprehensive overview of Knowledge Distillation: problem definition, type of knowledge, and its recent progress.
2. Investigate how research community try to theoretically explain KD.
3. Offer several scenarios and its feasible solution based on other recent works.

2. Background: Knowledge Distillation

2.1. Problem statement

Neural network models have been successful in a myriad of fields including extremely complex problems. However,

these models are enormous in size and computational hungry, thus cannot be deployed to on-edge devices or simply not feasible in some situations. To compensate this situation, knowledge distillation was first proposed in [3] and properly formalize in [11]. Buciluă et al. [3], the knowledge is transferred from the larger model to smaller model while minimizing the logits difference produced by those two models respectively. This vanilla method open the first path of knowledge type that can be use to distilled, later on, there started to have some other works using the activations, feature maps of intermediate layers as the guide for the student network [19, 24, 10]. Some other methods to extract knowledge by comparing relationship between difference layers of teacher model [22, 25] are also worth mentioning but won't be discussed in this survey.

2.2. Knowledge

2.2.1 Knowledge from logits

In this type of knowledge extraction, it usually refers to using the neural response of the final output layer, or also known as *logits* of the teacher model. As stated before, the first application of logits in KD is used in [3], but in many situations, given a highly confident teacher model, the output of softmax function gives more less the same information as the ground truth label (since this is the core idea of softmax function - maximizing the probabilities of one while minimizing the others). Tackling this problem, Hinton et al. [11] proposed the concept of 'soft labels' and declared that this type of label contains informative *dark knowledge*. Given the **logits** z from a network, the 'soft label' p_i of an image is defined as:

$$p_i = \frac{\exp(\frac{z_i}{\rho})}{\sum_j \exp(\frac{z_j}{\rho})}$$

where ρ is the temperature parameter, notice that when $\rho = 1$, we get the normal softmax function. Determine which value is optimal for ρ is still a debatable topics but it is argued that while increasing ρ , the label becomes softer and providing more information about which class is similar to the predicted label. Accordingly, the logits-based distillation loss function is defined as follow:

$$\mathcal{L}_{LB}(x; \theta) = \alpha * \mathcal{H}(y, \sigma(z_s)) + \beta * \mathcal{H}(\sigma(z_t; \rho), \sigma(z_s; \rho))$$

where x is the input, θ are the student model weights, $\mathcal{H}(\cdot)$ is the cross-entropy loss function, y is the ground truth label, $\sigma(\cdot)$ is the ρ -parameterized softmax function, α, β are the coefficients to tune in order to balance between both cross-entropies, and z_t, z_s are the logits of teacher and student model respectively.

Logits-based knowledge is fairly straightforward and easy to implement so it is often the first approach used when

need to apply knowledge distillation. As far as the survey goes, there are two main motives of using logits-based knowledge: 1) using soft labels 2) create/use noisy data to train. A brief table containing its description and several related work can be found at ?? . Nevertheless, since this type of knowledge only utilize the final output of the teacher model, it fails to offer information gathered in the intermediate layers, which later was proved to contain many informative result [19]. Other than that, due to its characteristic, logits-based knowledge is bounded with the supervised learning.

2.2.2 Knowledge from intermediate layers

What's extraordinary about deep neural networks is how they extract and abstractly represent features through the feature maps generated between each layer. Hence, feature-based knowledge would theoretically contain richer information than logits-based knowledge. Not only so, by combining both intermediate and last layer's output, we can consider feature-based knowledge is an extension of the previous knowledge type.

Romero et al. [19] introduced the term *hint* as the outputs of a teacher hidden layer, which are used to guide the student learning process. The core idea is to choose some intermediate layers in both teacher and student layer and force the student to mimics the result of the teacher's feature maps. Motivated by this work, there has been many studies to investigate which hint layer/ guided layer to choose and how to measure the distance between them. The feature-based distillation loss function is often defined as follow:

$$\mathcal{L}_{FB}(f_t, f_s) = \mathcal{D}(\Phi_t(f_t), \Phi_s(f_s))$$

where f_t, f_s are the selected hint and guided layers, Φ_t, Φ_s are the transformation functions for the mentioned layers respectively, which often are applied when the result of two layers are mismatched, \mathcal{D} is the distance function measuring the similarity between the hint and guided layer. L1, L2 is often used for this distance function, but for some works, more complex function might be used such as the Maximum Mean Discrepancy (MMD) metric [13] or Kullback-Leibler divergence [1].

Eventhough this knowledge type offers more favorable information for the student learning process, where to pick the hint/ guided layers or the transformation functions still need more investigations. For examples, the teacher's layer transformation function are directly critical for this process since there are risks of losing information in the process of transforming. Several past work has already encountered this problem since the feature map's dimension got down-scaled, causing loss of knowledge [24, 20]. To counter this situation, one could either come up with novel transformation function such as margin RELU[10] or not use any

transformation at all [19]. But as investigated by Heo et al. [10], the hints consist of both beneficial and adverse information so we should try to only gather the positive ones rather than distilling all of it. The same goes for the guided layers, sometimes researchers use the same transformation [24], which might lead to the same information loss. Some other works use 1×1 convolutional layer as a student transform [19, 10] in order to match the depth dimension with the hints, so no information would be lost during the process. A more detail of progressive work using feature-based knowledge can be found at figure ??.

3. Explaining Knowledge Distillation

As discussed in the previous sections, knowledge distillation has been successfully applied in various fields and applications. However, clarying how and why KD works in general and its models are usually superior than those trained from raw data still remains an work in progress. Not only so, there is no universally accepted hypothesis as to how knowledge is transmitted, making it impossible to accurately test emperically scientific results and to design new approaches in a more structured fashion. Most of past works seldom go beyond qualitative assertions, such as suggesting that learning from soft labels may be better than learning from hard labels, or that the teacher’s performance in a multi-class environment offers details on how close different classes are to one another [3, 11, 10].

Recently, some scientific researchers start to focus on finding out what are the real factors making KD works. In the case of deep linear classifiers, Phuong et al. [18] came up with a theoretical explanation for a generalization bound for fast convergence of learning distilled student networks. This justification clarifies what the student would learn and how quickly they do it, as well as the factors that influence distillation performance, to be more precise, the authors stated there are three main factors: data geometry, optimization bias and srtong monotocity of the student classifier.

This survey is going to focus on another work, where the author - Chen et al. quantified the distilled knowledge using the term visual concepts from the intermediate layers of a deep neural network (DNN) in order to explain KD [5]. What separates this paper from other work is that the author tried to interpret KD from a different perspective of information theory by quantifying, analyzing and comparing encoded information in the intermediate layers between DNNs learned by KD and normal DNNs, which not only investigate the success of KD but also what DNNs have learned during training in order to perform effective inference.

3.1. Setup

Given the teacher network is a pretrained classification model then distill to the student one. The main focus is to compare and analyse the difference between the student network with the *baseline network* (DNN learned from raw data) (both three models have the same architecture for fair comparison). Let $x \in R^n$ be the input image, $f_t(x), f_s(x) \in R^L$ be the intermediate-layer features and $y_t = g_t(f_t(x)), y_s = g_s(f_s(x))$ be the classification results of teacher and student, respectively. As the effect of KD, $f_s(x)$ is forced to be approximately equivanlent to $f_t(x)$

The author used the method proposed by [17] in order to quantify the discarded information of the input (consider as the conditional entropy $H(X')$) given the intermediate-layer feature $f^* = f(x)$ as follows:

$$H(X') \text{ s.t. } \forall x' \in X', \|f(x') - f^*\|^2 \leq \tau$$

where X' denotes the set of images having a specific object instance (represented by a small range of feature $\|f(x') - f^*\|^2 \leq \tau$). With a proper condition of x' ($x' \sim \mathcal{N}(x, \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2))$, where n is the number of images), the entropy $H(X')$ of the entire image can be formed by the entropies of every pixel $\{H_i\}$:

$$H(X') = \sum_{i=1}^n H_i$$

where $H_i = \log \sigma_i + \frac{1}{2} \log(2\pi e)$ measures the discarded information of the i -th pixel.

3.2. Proposed hypotheses

Hypothesis 1: Knowledge distillation makes the DNN learn more reliable visual concepts then learning from raw data.

The paper aims to compare the number of visual concepts that are encoded in the baseline network and those in the student network. They assumed that since those pixels having high H_i would not contribute much to the prediction result of the model, they would be considered as task-irrelevant feature, while those having smaller value are called *visual concepts*: image region whose information is significantly less discarded and mainly used by the DNN. Using the same concept of information theory, this is how the authors quantify the knowledge:

$$\begin{aligned} N_{concept}^{bg}(x) &= \sum_{i \in \Lambda_{bg}^{w.r.t.x}} \mathbb{1}(\bar{H} - H_i > b), \\ N_{concept}^{fg}(x) &= \sum_{i \in \Lambda_{fg}^{w.r.t.x}} \mathbb{1}(\bar{H} - H_i > b) \\ \lambda &= \mathbb{E}_{x \in \mathbf{I}} \left[\frac{N_{concept}^{fg}(x)}{N_{concept}^{fg}(x) + N_{concept}^{bg}(x)} \right] \end{aligned}$$

where $N_{concept}^{bg}(x)$, $N_{concept}^{fg}(x)$ are the number of visual concepts encoded on the background and foreground respectively, Λ_{bg} , Λ_{fg} are the sets of pixel of background and foreground w.r.t. image x , $\bar{H} = E_{i \in \Lambda_{bg}}[H_i]$ is the average entropy value of the background pixels which can be used as a baseline entropy. The pixels that having smaller entropy value than this baseline by an small amount b are considered task-relevant visual concepts. Finally, λ is proposed as the metric used to measure how effective the model's feature extraction process. The reason why λ can do such thing is argued by the author that statistically the foreground contains more informative features than the background, so a well-trained DNN model should focus on the visual concepts in the foreground.

Hypothesis 2: Knowledge distillation ensures that the DNN is prone to learning various concepts simultaneously, in contrast of DNN learning from raw data which learns these sequentially.

The author proposed another metrics based on the number of learned foreground visual concepts along each training epochs and take measurement with respect to the epoch that the model learn those the most: $\hat{m} = \arg\max_k N_k^{fg}(I)$ and "weight distance" function to estimate the learning effect $\sum_{k=1}^{\hat{m}} \frac{||w_k - w_{k-1}||}{||w_0||}$ (where w_i is the weight of the model at epoch i -th):

$$D_{\text{mean}} = \mathbb{E}_{I \in \mathcal{I}} \left[\sum_{k=1}^{\hat{m}} \frac{||w_k - w_{k-1}||}{||w_0||} \right]$$

$$D_{\text{var}} = \text{Var}_{I \in \mathcal{I}} \left[\sum_{k=1}^{\hat{m}} \frac{||w_k - w_{k-1}||}{||w_0||} \right]$$

where D_{mean} is the average weight distance where the DNN usually extracts the most task-relevant visual concepts, the smaller the value the faster the DNN learns the visual concepts during training. While D_{var} was considered inversely proportional to how simultaneously the model learns different visual concepts during training.

Hypothesis 3: Comparing to learning from raw data, knowledge distillation produces more consistent optimization directions.

The final metrics is computed by comparing the amount of learned and chosen foreground visual concepts in the final model $||S_M(I)||$ with the union of those learned in each epoch $||\bigcup_{j=1}^M S_j(I)||$. The proper form of the metric is defined as follow:

$$\rho = \frac{||S_M(I)||}{||\bigcup_{j=1}^M S_j(I)||}$$

The higher the value ρ indicates the less detours and more stable training process.

3.3. Discussion

The paper offers a novel and understandable interpretation of knowledge distillation from the perspective of information theory, to be more precise, measuring the information encoded in a DNN's intermediate layers. For each hypothesis, the author proposed a metrics which allows to partly theoretical explain and empirically prove that corresponding theory. But there are some confusing details might need further investigation: 1) this explaining framework only focuses on classification problem, which is one of the simplest form, in order to apply this framework to another applications might need extra information 2) The mathematical concept of presenting object and assumption about distribution of the image set: given a diversified training dataset, whether these metrics function similarly? 3) The estimated epoch \hat{m} used for the second metrics is not a precise estimation.

4. Regularizing Knowledge distillation

As one of the most beneficial characteristic of KD, any student can learn from any teacher disregarding the structural difference. But it is empirically proved that well-trained large DNN doesn't often make good teachers due to the mismatched capacity, which makes the student unable to mimic it [6].

To tackle this problem, multiple work shares the same idea of regularizing the teacher model [6, 16]. In [6], Cho et al. proposed a new process ESKD (Early-stopped knowledge distillation) after empirically prove sequential knowledge distillation is also not that efficient since it can only outperform one model training from scratch but not ensemble of those. They argued that the found solution space of the teacher is not accessible from the student, which means to find a teacher whose discovered solution should be discoverable by the student. Based on another works, the author assumed early stopping allows large model to behave as a small network while still having better search space than smaller ones. Different from Cho et al., Lukasik et al. [16] investigates the denoising effect of label smoothing on noisy data then applies it on the distillation process in order to test its effectiveness. The ending result is that applying label smoothing on the teacher significantly enhances over vanilla distillation, while applying the same on the student has mixed results.

Having the same idea, Li et al. [23] investigates and compares KD with label smoothing regularizer and later on proposed a novel Teacher-free Knowledge distillation (TfKD) framework. It started with the observation that using either poorly trained teacher to distil student or student to teach teacher model still can improve and enhanced the guided models, which suggests the author to consider KD as a regularization term (strongly related to label smoothing

regularizer). Having that in mind, the paper consider replacing the class distribution predictions of teacher model with a simpler one, implemented in the TfKD framework. This framework is particularly useful in circumstances where a more efficient teacher model is inaccessible or where only minimal computing resources are available. There are two methods of distilling in TfKD framework: 1) self-training distillation and 2) Combine KD with Label Smoothing Regularizer to create a 100% accuracy teacher. Both of the methods are very simple yet effective and also empirically proved to be performant.

5. Conclusion

The main technical information and applications of knowledge distillation have been covered in this survey. We also briefly review the taxonomy methods for current KD approaches and include description of the problem. We then investigate how KD is perceived and explained in the current past work. Finally some methods of regularizing KD while applying in real-world application are also discussed.

References

- [1] Gustavo Aguilar, Yuan Ling, Yu Zhang, Benjamin Yao, Xing Fan, and Chenlei Guo. Knowledge distillation from internal representations, 2020.
- [2] Hessam Bagherinezhad, Maxwell Horton, Mohammad Rastegari, and Ali Farhadi. Label refinery: Improving imagenet classification through label progression, 2018.
- [3] Cristian Bucilua, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '06, page 535–541, New York, NY, USA, 2006. Association for Computing Machinery.
- [4] Wei-Chun Chen, Chia-Che Chang, Chien-Yu Lu, and Che-Rung Lee. Knowledge distillation with feature maps for image classification, 2018.
- [5] Xu Cheng, Zhefan Rao, Yilan Chen, and Quanshi Zhang. Explaining knowledge distillation by quantifying the knowledge, 2020.
- [6] Jang Hyun Cho and Bharath Hariharan. On the efficacy of knowledge distillation, 2019.
- [7] Joseph DiPalma, Arief A. Suriawinata, Laura J. Tafe, Lorenzo Torresani, and Saeed Hassanpour. Resolution-based distillation for efficient histology image classification, 2021.
- [8] Shiming Ge, Shengwei Zhao, Chenyu Li, and Jia Li. Low-resolution face recognition in the wild via selective knowledge distillation. *IEEE Transactions on Image Processing*, 28(4):2051–2062, Apr 2019.
- [9] Tong He, Chunhua Shen, Zhi Tian, Dong Gong, Changming Sun, and Youliang Yan. Knowledge adaptation for efficient semantic segmentation, 2019.
- [10] Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive overhaul of feature distillation, 2019.
- [11] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015.
- [12] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. Learning lightweight lane detection cnns by self attention distillation, 2019.
- [13] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer, 2017.
- [14] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding, 2020.
- [15] Zhizhong Li and Derek Hoiem. Learning without forgetting, 2017.
- [16] Michal Lukasik, Srinadh Bhojanapalli, Aditya Krishna Menon, and Sanjiv Kumar. Does label smoothing mitigate label noise?, 2020.
- [17] Haotian Ma, Yinqing Zhang, Fan Zhou, and Quanshi Zhang. Quantifying layerwise information discarding of neural networks, 2019.
- [18] Mary Phuong and Christoph Lampert. Towards understanding knowledge distillation. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5142–5151. PMLR, 09–15 Jun 2019.
- [19] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets, 2015.
- [20] Zhiqiang Shen, Zhankui He, and Xiangyang Xue. Meal: Multi-model ensemble via adversarial learning, 2019.
- [21] Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. Patient knowledge distillation for BERT model compression. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4323–4332, Hong Kong, China, Nov. 2019. Association for Computational Linguistics.
- [22] J. Yim, D. Joo, J. Bae, and J. Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7130–7138, 2017.
- [23] Li Yuan, Francis E. H. Tay, Guilin Li, Tao Wang, and Jiashi Feng. Revisiting knowledge distillation via label smoothing regularization, 2021.
- [24] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer, 2017.
- [25] Chenrui Zhang and Yuxin Peng. Better and faster: Knowledge transfer from multiple self-supervised learning tasks via graph distillation for video classification, 2018.
- [26] Yang Zhao, Yifan Liu, Chunhua Shen, Yongsheng Gao, and Shengwu Xiong. Mobilefan: Transferring deep hidden representation for face alignment. *Pattern Recognition*, 100:107114, 2020.