
A Note on Knowledge Distillation Loss Function for Object Classification

Defang Chen
Zhejiang University
defchern@zju.edu.cn

Abstract

This research note provides a quick introduction to the knowledge distillation loss function used in object classification. In particular, we discuss its connection to a previously proposed logits matching loss function. We further treat knowledge distillation as a specific form of output regularization and demonstrate its connection to label smoothing and entropy-based regularization.

1 Introduction

Knowledge distillation (KD) was proposed as a model compression technique to distill the knowledge in a powerful yet cumbersome teacher model into a lightweight student model, hoping that it will enhance the student generalization ability [HVD15]. Since the advent of knowledge distillation, variant approaches have been developed to improve its effectiveness and widen its applications, including feature distillation [CMZ⁺21, CMZ⁺22], collaborative distillation [CMW⁺20], and diffusion-based distillation [CZM⁺23, ZCWC24, CZW⁺24]. In the vanilla KD [HVD15], class predictions between teacher and student models are aligned with a newly introduced hyperparameter—*temperature* in the softmax activation to control the softness of predicted distributions. In this way, Hinton and his collaborators [HVD15] argued that the pioneering model compression technique termed logits matching [BCNM06] is actually a special case of their proposed approach, provided that the temperature is much higher than the logits in order of magnitude and the logits are zero-mean normalized explicitly. In this note, we provide detailed derivations to review and deepen our understanding of this connection. We first prove that with a single *infinity temperature* condition, we could already build the connection between these two loss functions, although attached with an extra regularization. We further point out that *equal-mean normalization* for logits is enough to establish the exact equivalence. Finally, we discuss knowledge distillation from the output regularization perspective.

2 Preliminary

We first briefly recap the basic concept of object classification using deep neural networks, especially from the perspective of knowledge distillation. Then, a formal description of the standard knowledge distillation and logits matching loss functions are introduced with necessary notations.

2.1 Distilling Knowledge from Human Labellings

Given a training dataset $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})\}$ consisting of n objects from k categories, we denote the category label of the object \mathbf{x} as \mathbf{y} . Take one-class image classification problem as an example, we have $\mathbf{x} \in \mathbb{R}^{c \times h \times w}$, where c denotes the channel dimension, h and w denote the spatial dimensions. Besides, only one element in the vector $\mathbf{y} \in \mathbb{R}^k$ equals one and the other elements all equal zero

(e.g., $\mathbf{y}_i = 1$, and $\mathbf{y}_j = 0, \forall j \neq i$ if this object belongs to the i -th category). We then train a deep neural network with the parameter θ to learn a mapping from the object space to the category space.

Given the object \mathbf{x} , we denote the unnormalized prediction of a deep neural network as \mathbf{z} , which is also known as *logits*, and denote its softmax-based normalized version as \mathbf{p} , which is also known as *class predictions*. Mathematically, we have $\mathbf{p}_i = \exp(\mathbf{z}_i) / \sum_{j=1}^k \exp(\mathbf{z}_j)$, where \mathbf{z}_i and \mathbf{p}_i denote the i -th element of \mathbf{z} and \mathbf{p} , respectively. The model parameter θ is randomly initialized and updated using the following cross-entropy loss function:

$$\mathcal{L}_{\text{CE}}(\mathbf{y}, \mathbf{p}) = - \sum_{i=1}^k \mathbf{y}_i \log \mathbf{p}_i, \quad \text{or equivalently,} \quad \mathcal{L}_{\text{KL}}(\mathbf{y}, \mathbf{p}) = \sum_{i=1}^k \mathbf{y}_i \log \frac{\mathbf{y}_i}{\mathbf{p}_i}, \quad (1)$$

with the gradient as follows

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial \mathbf{p}_j} = -\frac{\mathbf{y}_j}{\mathbf{p}_j}, \quad \frac{\partial \mathbf{p}_j}{\partial \mathbf{z}_i} = \begin{cases} \mathbf{p}_j(1 - \mathbf{p}_j) & i = j \\ -\mathbf{p}_i \mathbf{p}_j & i \neq j. \end{cases} \quad (2)$$

Then, we take the partial derivative of \mathcal{L}_{CE} with respect to \mathbf{z}_i

$$\begin{aligned} \frac{\partial \mathcal{L}_{\text{CE}}}{\partial \mathbf{z}_i} &= \sum_j \frac{\partial \mathcal{L}_{\text{CE}}}{\partial \mathbf{p}_j} \frac{\partial \mathbf{p}_j}{\partial \mathbf{z}_i} = - \sum_j \frac{\mathbf{y}_j}{\mathbf{p}_j} \frac{\partial \mathbf{p}_j}{\partial \mathbf{z}_i} = - \left[\underbrace{\frac{\mathbf{y}_j}{\mathbf{p}_j} \mathbf{p}_j (1 - \mathbf{p}_j)}_{j=i} + \sum_{j \neq i} \underbrace{\frac{\mathbf{y}_j}{\mathbf{p}_j} (-\mathbf{p}_i \mathbf{p}_j)}_{j \neq i} \right] \\ &= - \left[\underbrace{\mathbf{y}_j (1 - \mathbf{p}_j)}_{j=i} + \sum_{j \neq i} \underbrace{\mathbf{y}_j (-\mathbf{p}_i)}_{j \neq i} \right] = - \left[\mathbf{y}_i - \mathbf{p}_i \sum_j \mathbf{y}_j \right] = \mathbf{p}_i - \mathbf{y}_i. \end{aligned} \quad (3)$$

Remark 1. The standard learning objective for object classification can be interpreted as distilling knowledge from human labellings, where the regular model acts as a “student” model and the ground-truth label for each object acts as the output of a(n) “teacher/oracle” model.

Label Smoothing: To alleviate the over-fitting issue in model training or the mis-labeling issue in data collection, label smoothing [SVI⁺16] technique was proposed. It replaces the hard label \mathbf{y} as the smoothed label \mathbf{y}' , with $\mathbf{y}' = (1 - \alpha)\mathbf{y} + \alpha\mathbf{u}$ and $\mathbf{u}_i = 1/k, \forall i$. In this case, the above loss function becomes¹

$$\mathcal{L}_{\text{LS}}(\mathbf{y}', \mathbf{p}) = (1 - \alpha)\mathcal{L}_{\text{CE}}(\mathbf{y}, \mathbf{p}) + \alpha\mathcal{L}_{\text{KL}}(\mathbf{u}, \mathbf{p}). \quad (4)$$

Confidence Penalty: Another similar output-based regularization reverses the KL term in the above equation, aiming to penalize a low-entropy model prediction [PTC⁺17]

$$\mathcal{L}_{\text{CP}}(\mathbf{y}', \mathbf{p}) = (1 - \alpha)\mathcal{L}_{\text{CE}}(\mathbf{y}, \mathbf{p}) + \alpha\mathcal{L}_{\text{KL}}(\mathbf{p}, \mathbf{u}). \quad (5)$$

The above two loss functions can be unified using skew-Jensen divergence [NB11, MSC20].

2.2 Distilling Knowledge from Neural Networks

We next denote class predictions of the teacher model with the parameter θ^t and the student model with the parameter θ^s as \mathbf{p}^t and \mathbf{p}^s , which are produced by normalized logits \mathbf{z}^t and \mathbf{z}^s , respectively.

A straightforward approach to distill knowledge from a pre-trained powerful teacher model to a simple student model is adopting the following logits matching loss function [BCNM06, BC14] to update the student model’s parameter θ^s :

$$\mathcal{L}_{\text{LM}}(\mathbf{z}^t, \mathbf{z}^s) = \frac{1}{2k} \|\mathbf{z}^t - \mathbf{z}^s\|_2^2 = \frac{1}{2k} \sum_{i=1}^k (\mathbf{z}_i^t - \mathbf{z}_i^s)^2, \quad (6)$$

with the gradient $\partial \mathcal{L}_{\text{LM}} / \partial \mathbf{z}_i^s = (\mathbf{z}_i^s - \mathbf{z}_i^t) / k$, and the gradient $\partial \mathcal{L}_{\text{LM}} / \partial \theta^s = \partial \mathcal{L}_{\text{LM}} / \partial \mathbf{z}^s \cdot \partial \mathbf{z}^s / \partial \theta^s$.

¹We omit the constant term irrelevant to the parameter optimization.

Similarly, we can also adopt the vanilla knowledge distillation loss function [HVD15] for knowledge transfer. This approach introduces temperature τ as a hyper-parameter to soften the model-predicted probability distributions, and modifies the relationship between logits \mathbf{z} and predictions \mathbf{p} as follows

$$\mathbf{p}_i^t = \frac{\exp(\mathbf{z}_i^t/\tau)}{\sum_{j=1}^k \exp(\mathbf{z}_j^t/\tau)}, \quad \mathbf{p}_i^s = \frac{\exp(\mathbf{z}_i^s/\tau)}{\sum_{j=1}^k \exp(\mathbf{z}_j^s/\tau)}. \quad (7)$$

The vanilla knowledge distillation loss function is

$$\mathcal{L}_{\text{KD}}(\mathbf{p}^t, \mathbf{p}^s) = \sum_{i=1}^k \mathbf{p}_i^t \log \frac{\mathbf{p}_i^t}{\mathbf{p}_i^s} = - \sum_{i=1}^k \mathbf{p}_i^t \log \mathbf{p}_i^s + \sum_{i=1}^k \mathbf{p}_i^t \log \mathbf{p}_i^t. \quad (8)$$

The second term in Equation (8) is a negative entropy of \mathbf{p}^t , which is irrelevant to the update of student's parameters. Similarly, we take the partial derivative of \mathcal{L}_{KD} w.r.t. \mathbf{z}_i^s :

$$\frac{\partial \mathcal{L}_{\text{KD}}}{\partial \mathbf{p}_j^s} = -\frac{\mathbf{p}_j^t}{\mathbf{p}_j^s}, \quad \frac{\partial \mathbf{p}_j^s}{\partial \mathbf{z}_i^s} = \begin{cases} \frac{1}{\tau} \mathbf{p}_j^s (1 - \mathbf{p}_j^s) & i = j \\ -\frac{1}{\tau} \mathbf{p}_i^s \mathbf{p}_j^s & i \neq j. \end{cases} \quad (9)$$

$$\frac{\partial \mathcal{L}_{\text{KD}}}{\partial \mathbf{z}_i^s} = \sum_j \frac{\partial \mathcal{L}_{\text{KD}}}{\partial \mathbf{p}_j^s} \frac{\partial \mathbf{p}_j^s}{\partial \mathbf{z}_i^s} = - \sum_j \frac{\mathbf{p}_j^t}{\mathbf{p}_j^s} \frac{\partial \mathbf{p}_j^s}{\partial \mathbf{z}_i^s} = \frac{1}{\tau} (\mathbf{p}_i^s - \mathbf{p}_i^t). \quad (10)$$

Remark 2. The standard knowledge distillation loss function (Equation (8)) multiplying τ^2 acts as a regularized logits matching loss function (Equation (6)) under the infinity temperature.

Proof. We then prove the equivalence between $\tau^2 \mathcal{L}_{\text{KD}}$ and \mathcal{L}_{LM} [HVD15, KOK+21] below.

$$\begin{aligned} \lim_{\tau \rightarrow \infty} \tau^2 \frac{\partial \mathcal{L}_{\text{KD}}}{\partial \mathbf{z}_i^s} &\stackrel{\textcircled{1}}{=} \lim_{\tau \rightarrow \infty} \tau (\mathbf{p}_i^s - \mathbf{p}_i^t) = \lim_{\tau \rightarrow \infty} \tau \left(\frac{\exp(\mathbf{z}_i^s/\tau)}{\sum_{j=1}^k \exp(\mathbf{z}_j^s/\tau)} - \frac{\exp(\mathbf{z}_i^t/\tau)}{\sum_{j=1}^k \exp(\mathbf{z}_j^t/\tau)} \right) \\ &\stackrel{\textcircled{2}}{=} \lim_{\tau \rightarrow \infty} \left(\frac{\sum_{j=1}^k \tau (\exp((\mathbf{z}_j^t - \mathbf{z}_i^t)/\tau) - 1) - \sum_{j=1}^k \tau (\exp((\mathbf{z}_j^s - \mathbf{z}_i^s)/\tau) - 1)}{\left(\sum_{j=1}^k \exp((\mathbf{z}_j^t - \mathbf{z}_i^t)/\tau) \right) \left(\sum_{j=1}^k \exp((\mathbf{z}_j^s - \mathbf{z}_i^s)/\tau) \right)} \right) \\ &= \frac{1}{k} (\mathbf{z}_i^s - \mathbf{z}_i^t) - \frac{1}{k^2} \sum_{j=1}^k (\mathbf{z}_j^s - \mathbf{z}_j^t) \end{aligned} \quad (11)$$

For $\textcircled{1}$, the detailed derivation is provided in Equation (10), and we substitute the Equation (7) into \mathbf{p}_i^s and \mathbf{p}_i^t ; for $\textcircled{2}$, the detailed derivation is provided in Section A.2 of [KOK+21], and we require τ goes to positive infinity to leverage the Taylor approximation of exponential function. Therefore, the standard knowledge distillation loss function acts as a regularized logits matching loss function:

$$\mathcal{L}_{\text{LM}_r} = \frac{1}{2k} \|\mathbf{z}^t - \mathbf{z}^s\|_2^2 - \frac{1}{2k^2} \left(\sum_{j=1}^k \mathbf{z}_j^s - \sum_{j=1}^k \mathbf{z}_j^t \right)^2 + \text{Constant}. \quad (12)$$

The above loss function also implies that when the summation of logits predicted by teacher and student models are equal, i.e., $\sum_j^k \mathbf{z}_j^s = \sum_j^k \mathbf{z}_j^t$, the extra regularization term will just disappear. \square

Remark 3. Based on the infinity large temperature and equal-mean normalization for logits (i.e., $\sum_j^k \mathbf{z}_j^s = \sum_j^k \mathbf{z}_j^t$), the gradient of $\tau^2 \mathcal{L}_{\text{KD}}$ equals to that of \mathcal{L}_{LM} and thus we can conclude that the effects of these two losses are exactly the same.

2.3 Knowledge Distillation as Output Regularization

In practice, the knowledge distillation loss function is generally used together with the original cross-entropy loss function. That is to say, we adopt the following loss function to optimize parameters of a student model:

$$\mathcal{L}_{\text{KD}'} = (1 - \alpha) \mathcal{L}_{\text{CE}}(\mathbf{y}, \mathbf{p}^s) + \alpha \mathcal{L}_{\text{KD}}(\mathbf{p}^t, \mathbf{p}^s). \quad (13)$$

Comparing the loss function in Equation (4), we conclude that class predictions of the teacher model act as adaptive label smoothing to prevent the student output being over-confident [YTL+20].

Remark 4. The knowledge distillation loss function in Equation (13) acts as an adaptive label smoothing loss function in Equation (4).

References

- [BC14] Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *Advances in Neural Information Processing Systems*, pages 2654–2662, 2014.
- [BCNM06] 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*, pages 535–541, 2006.
- [CMW⁺20] Defang Chen, Jian-Ping Mei, Can Wang, Yan Feng, and Chun Chen. Online knowledge distillation with diverse peers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3430–3437, 2020.
- [CMZ⁺21] Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen. Cross-layer distillation with semantic calibration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7028–7036, 2021.
- [CMZ⁺22] Defang Chen, Jian-Ping Mei, Hailin Zhang, Can Wang, Yan Feng, and Chun Chen. Knowledge distillation with the reused teacher classifier. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11933–11942, 2022.
- [CZM⁺23] Defang Chen, Zhenyu Zhou, Jian-Ping Mei, Chunhua Shen, Chun Chen, and Can Wang. A geometric perspective on diffusion models. *arXiv preprint arXiv:2305.19947*, 2023.
- [CZW⁺24] Defang Chen, Zhenyu Zhou, Can Wang, Chunhua Shen, and Siwei Lyu. On the trajectory regularity of ode-based diffusion sampling. pages 7905–7934, 2024.
- [HVD15] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [KOK⁺21] Taehyeon Kim, Jaehoon Oh, NakYil Kim, Sangwook Cho, and Se-Young Yun. Comparing kullback-leibler divergence and mean squared error loss in knowledge distillation. *arXiv preprint arXiv:2105.08919*, 2021.
- [MSC20] Clara Meister, Elizabeth Salesky, and Ryan Cotterell. Generalized entropy regularization or: There’s nothing special about label smoothing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6870–6886, 2020.
- [NB11] Frank Nielsen and Sylvain Boltz. The burbea-rao and bhattacharyya centroids. *IEEE Transactions on Information Theory*, 57(8):5455–5466, 2011.
- [PTC⁺17] Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing neural networks by penalizing confident output distributions. *arXiv preprint arXiv:1701.06548*, 2017.
- [SVI⁺16] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016.
- [YTL⁺20] Li Yuan, Francis EH Tay, Guilin Li, Tao Wang, and Jiashi Feng. Revisiting knowledge distillation via label smoothing regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3903–3911, 2020.
- [ZCWC24] Zhenyu Zhou, Defang Chen, Can Wang, and Chun Chen. Fast ode-based sampling for diffusion models in around 5 steps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7777–7786, 2024.