

Supplementary Material of Dirichlet-Based Prediction Calibration for Learning with Noisy Labels

1 Derivations

1.1 Equation 3

Given an example \mathbf{x}_i , the predicted probability obtained by the softmax function and our calibrated softmax function can be obtained as:

$$\boldsymbol{\rho}_i = [\rho_{i1}, \rho_{i2}, \dots, \rho_{iC}] = \left[\frac{e^{o_{i1}}}{\sum_{j=1}^C e^{o_{ij}}}, \frac{e^{o_{i2}}}{\sum_{j=1}^C e^{o_{ij}}}, \dots, \frac{e^{o_{iC}}}{\sum_{j=1}^C e^{o_{ij}}} \right], \quad (1)$$

and

$$\hat{\boldsymbol{\rho}}_i = [\hat{\rho}_{i1}, \hat{\rho}_{i2}, \dots, \hat{\rho}_{iC}] = \left[\frac{e^{o_{i1}} + \gamma}{\sum_{j=1}^C e^{o_{ij}} + \gamma}, \frac{e^{o_{i2}} + \gamma}{\sum_{j=1}^C e^{o_{ij}} + \gamma}, \dots, \frac{e^{o_{iC}} + \gamma}{\sum_{j=1}^C e^{o_{ij}} + \gamma} \right]. \quad (2)$$

Suppose the given label is y , the difference between the gradient of the softmax-based cross-entropy loss with respect to model output \mathbf{o}_i on the complementary label c and the gradient of the calibrated softmax-based cross-entropy loss can be expressed as:

$$\begin{aligned} \frac{\partial \mathcal{L}_{ce|\boldsymbol{\rho}_i}}{\partial o_{ic}} \Big|_{\mathbf{x}_i} - \frac{\partial \mathcal{L}_{ce|\hat{\boldsymbol{\rho}}_i}}{\partial o_{ic}} \Big|_{\mathbf{x}_i} &= \frac{\partial (-y \log \rho_{iy})}{\partial o_{ic}} - \frac{\partial (-y \log \hat{\rho}_{iy})}{\partial o_{ic}} \\ &= \frac{\partial \left(-y \log \frac{e^{o_{iy}}}{\sum_{j=1}^C e^{o_{ij}}} \right)}{\partial o_{ic}} - \frac{\partial \left(-y \log \frac{e^{o_{iy}} + \gamma}{\sum_{j=1}^C e^{o_{ij}} + \gamma} \right)}{\partial o_{ic}} \\ &= \frac{e^{o_{ic}}}{\sum_{j=1}^C e^{o_{ij}}} - \frac{e^{o_{ic}}}{\sum_{j=1}^C (e^{o_{ij}} + \gamma)} \\ &= \frac{\gamma C e^{o_{ic}}}{\sum_{j=1}^C e^{o_{ij}} \sum_{j=1}^C (e^{o_{ij}} + \gamma)} > 0, \end{aligned} \quad (3)$$

where the result is always greater than 0.

1.2 Equation 5

For example \mathbf{x}_i , the predicted probability for a given class c can be expressed as:

$$\begin{aligned}
P(y = c \mid \mathbf{x}_i, \boldsymbol{\theta}) &= \int p(y = c \mid \boldsymbol{\rho}_i) p(\boldsymbol{\rho}_i \mid \mathbf{x}_i, \boldsymbol{\theta}) d\boldsymbol{\rho}_i \\
&= \int \rho_{ic} p(\boldsymbol{\rho}_i \mid \mathbf{x}_i, \boldsymbol{\theta}) d\boldsymbol{\rho}_i \\
&= \int \cdots \int \cdots \int \rho_{ic} \cdot p(\rho_{i1}, \cdots, \rho_{ic}, \cdots, \rho_{iC} \mid \mathbf{x}_i, \boldsymbol{\theta}) d\rho_{i1} \cdots d\rho_{ic} \cdots d\rho_{iC} \\
&= \int \rho_c \left[\int \cdots \int \int \cdots \int p(\rho_{i1}, \cdots, \rho_{iC} \mid \mathbf{x}_i, \boldsymbol{\theta}) d\rho_{i1} \cdots d\rho_{i(c-1)} d\rho_{i(c+1)} \cdots d\rho_{iC} \right] d\rho_{ic} \\
&= \int \rho_{ic} p(\rho_{ic} \mid \mathbf{x}_i, \boldsymbol{\theta}) d\rho_{ic} \\
&= \int \rho_{ic} \left[\frac{1}{\mathcal{B}(\alpha_{ic}, \sum_{j=1, j \neq c}^C \alpha_{ij})} \rho_{ic}^{\alpha_{ic}-1} (1 - \rho_{ic})^{\sum_{j=1, j \neq c}^C \alpha_{ij}-1} \right] d\rho_{ic} \\
&= \frac{\mathcal{B}(\alpha_{ic} + 1, \sum_{j=1, j \neq c}^C \alpha_{ij})}{\mathcal{B}(\alpha_{ic}, \sum_{j=1, j \neq c}^C \alpha_{ij})} \int \frac{1}{\mathcal{B}(\alpha_{ic} + 1, \sum_{j=1, j \neq c}^C \alpha_{ij})} \rho_{ic}^{\alpha_{ic}} (1 - \rho_{ic})^{\sum_{j=1, j \neq c}^C \alpha_{ij}-1} d\rho_{ic} \\
&= \frac{\mathcal{B}(\alpha_{ic} + 1, \sum_{j=1, j \neq c}^C \alpha_{ij})}{\mathcal{B}(\alpha_{ic}, \sum_{j=1, j \neq c}^C \alpha_{ij})} \cdot 1 \\
&= \frac{\Gamma(\alpha_{ic} + 1) \Gamma(\sum_{j=1}^C \alpha_{ij})}{\Gamma(\sum_{j=1}^C \alpha_{ij} + 1) \Gamma(\alpha_{ic})} \\
&= \frac{\alpha_{ic}}{\sum_{j=1}^C \alpha_{ij}} \\
&= \frac{g(o_{ic}) + \gamma}{\sum_{j=1}^C (g(o_{ij}) + \gamma)},
\end{aligned} \tag{4}$$

where $\mathcal{B}(\cdot, \cdot)$ denotes the Beta function, and we have $p(\rho_{ic} \mid \mathbf{x}_i, \boldsymbol{\theta}) \sim \mathcal{B}(\rho_{ic} \mid \alpha_{ic}, \sum_{j=1, j \neq c}^C \alpha_{ij})$ and the following equation according to [Ng et al.(2011)Ng, Tian, and Tang]:

$$\begin{aligned}
p(\rho_{ic} \mid \mathbf{x}_i, \boldsymbol{\theta}) &= \frac{1}{\mathcal{B}(\alpha_{ic}, \sum_{j=1, j \neq c}^C \alpha_{ij})} \rho_{ic}^{\alpha_{ic}-1} (1 - \rho_{ic})^{\sum_{j=1, j \neq c}^C \alpha_{ij}-1} \\
&= \frac{\Gamma(\alpha_{ic} + \sum_{j=1}^C \alpha_{ij})}{\Gamma(\alpha_{ic}) \Gamma(\sum_{j=1}^C \alpha_{ij})} \rho_{ic}^{\alpha_{ic}-1} (1 - \rho_{ic})^{\sum_{j=1, j \neq c}^C \alpha_{ij}-1}.
\end{aligned} \tag{5}$$

1.3 Equation 6

According to Eq. 4, we can easily obtain \mathcal{L}_{nll} as:

$$\begin{aligned}
\mathcal{L}_{nll} &= -\frac{1}{N} \sum_{i=1}^N \log [p(y = c \mid \mathbf{x}_i, \boldsymbol{\theta})] = -\frac{1}{N} \sum_{i=1}^N \log \left(\frac{\alpha_{ic}}{\sum_{j=1}^c \alpha_{ij}} \right) = \frac{1}{N} \sum_{i=1}^N \log \left(\frac{\sum_{i=1}^c \alpha_{ij}}{\alpha_{ic}} \right) \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{ic} \left[\log \left(\sum_{j=1}^c \alpha_{ij} \right) - \log \alpha_{ic} \right].
\end{aligned} \tag{6}$$

1.4 Equation 7

The KL-divergence term \mathcal{L}_{kl} can be further derived as:

$$\mathcal{L}_{kl} = \frac{1}{NC} \sum_{i=1}^N D_{KL} (Dir(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i) \| Dir(\boldsymbol{\rho}_i | \mathbf{1})) = \frac{1}{NC} \sum_{i=1}^N \int p(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i) \log \frac{p(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i)}{p(\boldsymbol{\rho}_i | \mathbf{1})} d\boldsymbol{\rho}_i. \quad (7)$$

Given the probability density function $p(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i) \sim Dir(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i)$, we can obtain the following equation:

$$p(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i) = \frac{1}{\mathbb{B}(\tilde{\boldsymbol{\alpha}}_i)} \prod_{j=1}^C \rho_{ij}^{\tilde{\alpha}_{ij}-1} = \frac{\Gamma(\sum_{j=1}^C \tilde{\alpha}_{ij})}{\prod_{j=1}^C \Gamma(\tilde{\alpha}_{ij})}, \quad (8)$$

where $\mathbb{B}(\cdot)$ indicates the multivariate Beta function. Therefore, we have:

$$\begin{aligned} & p(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i) \log \frac{p(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i)}{p(\boldsymbol{\rho}_i | \mathbf{1})} d\boldsymbol{\rho}_i \\ &= \int \left(\frac{1}{\mathbb{B}(\tilde{\boldsymbol{\alpha}}_i)} \prod_{j=1}^C \rho_{ij}^{\tilde{\alpha}_{ij}-1} \right) \log \left(\frac{\mathbb{B}(\mathbf{1})}{\mathbb{B}(\tilde{\boldsymbol{\alpha}}_i)} \prod_{j=1}^C \rho_{ij}^{\tilde{\alpha}_{ij}-1} \right) d\boldsymbol{\rho}_i \\ &= \log \frac{\mathbb{B}(\mathbf{1})}{\mathbb{B}(\tilde{\boldsymbol{\alpha}}_i)} \int \frac{1}{\mathbb{B}(\tilde{\boldsymbol{\alpha}}_i)} \prod_{j=1}^C \rho_{ij}^{\tilde{\alpha}_{ij}-1} d\boldsymbol{\rho}_i + \int \left(\log \prod_{j=1}^C \rho_{ij}^{\tilde{\alpha}_{ij}-1} \right) \left(\frac{1}{\mathbb{B}(\tilde{\boldsymbol{\alpha}}_i)} \prod_{j=1}^C \rho_{ij}^{\tilde{\alpha}_{ij}-1} \right) d\boldsymbol{\rho}_i \\ &= \log \frac{\mathbb{B}(\mathbf{1})}{\mathbb{B}(\tilde{\boldsymbol{\alpha}}_i)} + \mathbb{E}_{\boldsymbol{\rho}_i \sim Dir(\boldsymbol{\rho}_i | \tilde{\boldsymbol{\alpha}}_i)} \left[\log \prod_{j=1}^C \rho_{ij}^{\tilde{\alpha}_{ij}-1} \right] \\ &= \log \frac{\mathbb{B}(\mathbf{1})}{\mathbb{B}(\tilde{\boldsymbol{\alpha}}_i)} + \sum_{c=1}^C (\tilde{\alpha}_{ic} - 1) \mathbb{E}_{\boldsymbol{\rho}_c \sim \mathcal{B}(\boldsymbol{\rho}_c | \tilde{\boldsymbol{\alpha}}_{ic}, \sum_{j=1, j \neq c}^C \tilde{\alpha}_{ij})} [\log \rho_{ic}] \\ &= \log \left[\frac{\Gamma(\sum_{j=1}^C \tilde{\alpha}_{ij})}{\Gamma(C) \prod_{j=1}^C \Gamma(\tilde{\alpha}_{ij})} \right] + \sum_{c=1}^C (\tilde{\alpha}_{ic} - 1) \left[\psi(\tilde{\alpha}_{ic}) - \psi \left(\sum_{j=1}^C \tilde{\alpha}_{ij} \right) \right]. \end{aligned} \quad (9)$$

Eventually, \mathcal{L}_{kl} can be simplify as:

$$\mathcal{L}_{kl} = \frac{1}{NC} \sum_{i=1}^N \log \left[\frac{\Gamma(\sum_{j=1}^C \tilde{\alpha}_{ij})}{\Gamma(C) \prod_{j=1}^C \Gamma(\tilde{\alpha}_{ij})} \right] + \sum_{c=1}^C (\tilde{\alpha}_{ic} - 1) \left[\psi(\tilde{\alpha}_{ic}) - \psi \left(\sum_{j=1}^C \tilde{\alpha}_{ij} \right) \right]. \quad (10)$$

2 Pseudo-Code

The training procedure of DPC is shown in Algorithm 1.

3 Additional Experimental Settings

For all CIFAR experiments, we use the same set of hyperparameters: the number of augmentations $M = 2$, sharpening temperature $T = 0.5$, the clean probability threshold $\tau = 0.5$, and Beta distribution parameter $\alpha = 4$. Similar to [Li et al.(2020)Li, Socher, and Hoi], the λ_u is chosen from $\{0, 25, 50, 150\}$ using a small validation set and is shown in Table 1 and 2 for detail. The balancing weight of \mathcal{L}_{con} is set to 1 for CIFAR-10 and 3 for CIFAR-100. For the WebVision dataset, we set $M = 2$, $T = 0.5$, $\tau = 0.5$, $\lambda_u = 0$, $\alpha = 0.5$, and the balancing weight of \mathcal{L}_{con} to 1.

In addition, the CIFAR experiments are conducted on 4-card 2080 servers, while the WebVision experiments are conducted on 4-card 3090-Ti servers.

Algorithm 1 The DPC algorithm.

Input: Training dataset \mathcal{D} , two networks $f(\cdot; \theta^{(1)})$ and $f(\cdot; \theta^{(2)})$, calibrated softmax function $\hat{\rho}$, evidential deep learning loss \mathcal{L}_{edl} loss with balancing factor β , unsupervised loss \mathcal{L}_{uns} with weight λ_{uns} , clean probability threshold τ , number of augmentations M , and sharpening temperature T .

Output: Learned model parameters $\theta^{(1)}$ and $\theta^{(2)}$.

```

1:  $\theta^{(1)}, \theta^{(2)} = \text{WarmUp}(\mathcal{D}, \mathcal{L}_{edl}, \theta^{(1)}, \theta^{(2)})$            # warmup with evidential deep learning loss  $\mathcal{L}_{edl}$ 
2: for  $e = 1 : \text{MaxEpoch}$  do
3:    $\mathcal{M}^{(1)} = f(\mathcal{D}, \theta^{(1)})$                                      # calculate per-example margin with  $\theta^{(1)}$ 
4:    $\mathcal{M}^{(2)} = f(\mathcal{D}, \theta^{(2)})$                                      # calculate per-example margin with  $\theta^{(2)}$ 
5:   # Kurtosis  $\leq 1.2$  means that the distribution exhibits a multimodal tendency and is suitable for GMM modeling.
6:   if Kurtosis  $\leq 1.2$  then
7:      $\mathcal{P}^{(1)} = \text{GMM}(\mathcal{D}, \theta^{(1)})$            # model per-example margin to obtain clean probability for  $\theta^{(2)}$ 
8:      $\mathcal{P}^{(2)} = \text{GMM}(\mathcal{D}, \theta^{(2)})$            # model per-example margin to obtain clean probability for  $\theta^{(1)}$ 
9:   else
10:     $\mathcal{P}^{(1)} = \text{Norm}(\mathcal{D}, \theta^{(1)})$            # Normalize margins to directly obtain clean probability for  $\theta^{(2)}$ 
11:     $\mathcal{P}^{(2)} = \text{Norm}(\mathcal{D}, \theta^{(2)})$            # Normalize margins to directly obtain clean probability for  $\theta^{(1)}$ 
12:   end if
13:   for  $k = 1, 2$  do
14:      $\mathcal{X}^{(k)} = \{(\mathbf{x}_i, \mathbf{y}_i, p_i) \mid p_i \geq \tau, \forall (\mathbf{x}_i, \mathbf{y}_i, p_i) \in (\mathcal{D}, \mathcal{P}^{(k)})\}$ 
15:      $\mathcal{U}^{(k)} = \{\mathbf{x}_i \mid p_i < \tau, \forall (\mathbf{x}_i, p_i) \in (\mathcal{D}, \mathcal{M}^{(k)})\}$ 
16:     for  $i = 1 : \text{MaxIter}$  do
17:       From  $\mathcal{X}^{(k)}$ , draw a mini-batch  $\{(\mathbf{x}_b, \mathbf{y}_b, p_b); b \in (1, \dots, B)\}$ 
18:       From  $\mathcal{U}^{(k)}$ , draw a mini-batch  $\{\mathbf{x}_b; b \in (1, \dots, B)\}$ 
19:       for  $b = 1 : B$  do
20:          $\{\hat{\mathbf{x}}_{b,1}, \hat{\mathbf{x}}_{b,2}, \dots, \hat{\mathbf{x}}_{b,M}\} = \text{Augment}(\{\mathbf{x}_{b,1}, \mathbf{x}_{b,2}, \dots, \mathbf{x}_{b,M}\})$ 
21:          $\{\hat{\mathbf{u}}_{b,1}, \hat{\mathbf{u}}_{b,2}, \dots, \hat{\mathbf{u}}_{b,M}\} = \text{Augment}(\{\mathbf{u}_{b,1}, \mathbf{u}_{b,2}, \dots, \mathbf{u}_{b,M}\})$ 
22:          $\hat{\rho}_b^{(k)} = \frac{1}{M} \sum_m \hat{\rho}(f(\hat{\mathbf{x}}_{b,m}; \theta^{(k)}))$            # use the proposed calibrated softmax function
23:          $\hat{\rho}_b^{(k)} = \text{Sharpen}(p_b \mathbf{y}_b^{(k)} + (1 - p_b) \hat{\rho}_b^{(k)}, T)$ 
24:          $\hat{\rho}_b = \frac{1}{2M} \sum_m [\hat{\rho}(f(\hat{\mathbf{u}}_{b,m}; \theta^{(1)}) + \hat{\rho}(f(\hat{\mathbf{u}}_{b,m}; \theta^{(2)}))]$            # with calibrated softmax
25:          $\hat{\rho}_b = \text{Sharpen}(\hat{\rho}_b, T)$ 
26:       end for
27:        $\hat{\mathcal{X}} = \{(\hat{\mathbf{x}}_{b,m}, \mathbf{y}_b, \hat{\rho}_b^{(k)}); b \in (1, 2, \dots, B), m \in (1, 2, \dots, M)\}$ 
28:        $\hat{\mathcal{U}} = \{(\hat{\mathbf{u}}_{b,m}, \hat{\rho}_b); b \in (1, 2, \dots, B), m \in (1, 2, \dots, M)\}$ 
29:        $\mathcal{X}', \mathcal{U}' = \text{MixMatch}^{(1)}(\hat{\mathcal{X}}, \hat{\mathcal{U}})$            # MixMatch stage 1: produce mix-up examples
30:       # MixMatch stage 2: calculate the training loss with our proposed  $\mathcal{L}_{edl}$  and  $\mathcal{L}_{uns}$ 
31:        $\mathcal{L}_{total} = \text{MixMatch}^{(2)}(\mathcal{X}', \mathcal{U}', \mathcal{L}_{edl}, \beta, \mathcal{L}_{uns}, \lambda_{uns})$ 
32:        $\theta^{(k)} = \text{SGD}(\mathcal{D}, \mathcal{L}_{total}, \theta^{(k)})$ 
33:     end for
34:   end for
35: end for

```

Table 1: Unsupervised loss weight λ_u for CIFAR-10 and CIFAR-100 with synthetic noise.

Hyperparameter	CIFAR-10						CIFAR-100					
	Symmetric			Asymmetric			Symmetric			Asymmetric		
	20%	50%	80%	10%	30%	40%	20%	50%	80%	10%	30%	40%
λ_u	0	25	25	0	0	0	25	150	150	25	25	150

Table 2: Unsupervised loss weight λ_u for CIFAR-N with real-world noise.

Hyperparameter	CIFAR-10N					CIFAR-100N	
	aggre	rand1	rand2	rand3	worst	noisy100	
λ_u	0	0	0	0	25	25	

Table 3: Additional ablation studies on CIFAR-N with real-world noise.

Hyperparameter	CIFAR-10N					CIFAR-100N	
	aggre	rand1	rand2	rand3	worst	noisy100	
MLNT	91.37	88.65	88.54	89.13	84.15	58.24	
MLNT w \mathcal{L}_{edl}	91.45	89.37	88.56	89.50	83.99	60.35	
Improve	↑ 0.08	↑ 0.72	↑ 0.02	↑ 0.37	↓ 0.16	↑ 2.11	
SOP	95.61	95.28	95.31	95.39	93.24	67.81	
SOP w \mathcal{L}_{edl}	95.87	95.75	95.35	95.65	93.29	68.25	
Improve	↑ 0.26	↑ 0.47	↑ 0.04	↑ 0.26	↑ 0.05	↑ 0.44	

4 Additional Experimental Results

To further verify the generalizability of the Dirichlet-based prediction calibration method, we also provide the results of another two methods, MLNT [Li et al.(2019)Li, Wong, Zhao, and Kankanhalli] and SOP [Liu et al.(2022)Liu, Zhu, Qu, and You], integrated with \mathcal{L}_{edl} in Table 3 on CIFAR-N with real-world noise. The proposed method continues to be effective, indicating its general applicability.

References

- [Li et al.(2019)Li, Wong, Zhao, and Kankanhalli] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. Learning to learn from noisy labeled data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5051–5059, 2019.
- [Li et al.(2020)Li, Socher, and Hoi] Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. *arXiv preprint arXiv:2002.07394*, 2020.
- [Liu et al.(2022)Liu, Zhu, Qu, and You] Sheng Liu, Zhihui Zhu, Qing Qu, and Chong You. Robust training under label noise by over-parameterization. In *International Conference on Machine Learning*, pages 14153–14172. PMLR, 2022.
- [Ng et al.(2011)Ng, Tian, and Tang] Kai Wang Ng, Guo-Liang Tian, and Man-Lai Tang. Dirichlet and related distributions: Theory, methods and applications. 2011.