

# IOFM: Using the Interpolation Technique on the Over-Fitted Models to Identify Clean-Annotated Samples

Dongha Kim, Yongchan Choi, Kunwoong Kim, Ilsang Ohn, Yongdai Kim  
AAAI 2024

Department of Statistics, Seoul National University  
Presented by Sangmoon Han

2025-12-23

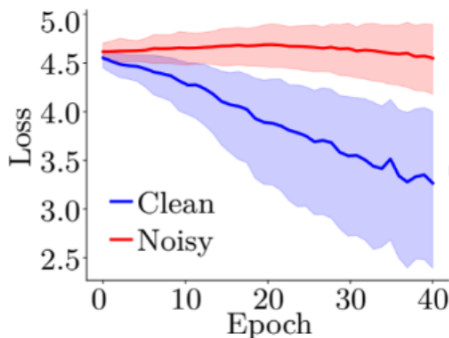
# Table of Contents

## 1. Overview

## 2. Proposed Method

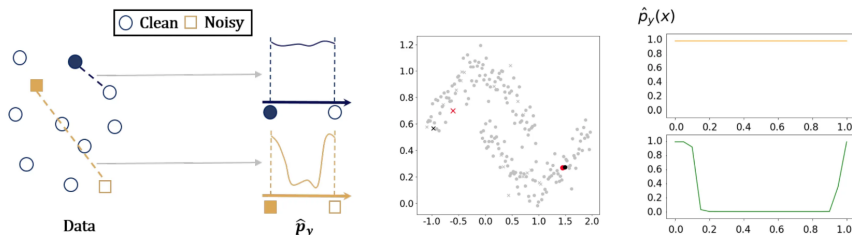
## 3. Experiments

# The memorization effect



- The memorization effect(**ME**) is a phenomenon that deep neural networks (DNNs) memorize clean data before noisy ones.
- Recent SOTA algorithms for handling noisy label problems are based on the **ME**.
- However, the **ME** may not occur in certain scenarios,
  1. labels have an imbalanced distribution.
  2. labels are heavily contaminated.

# Overview of IOFM



- Let's consider an over-fitted DNN and its feature space, i.e., the map of the highest hidden layer.
- For a given training sample  $(x_*, y_*)$ ,  $x_*$  is located close to other inputs sharing the label  $y_*$  regardless of anomalousness of  $y_*$  on the feature space.
- When they consider the original input space, the similarity between  $x_*$  and its neighbors, **chosen from the feature space**, would be quite different on the original input space depending on whether  $y_*$  is clean or noisy.
- Conceptually, The IOFM measures how similar the neighbors of a given data chosen from the feature space are on the input space and decides **the data as clean when the similarity is large**.

# Notation

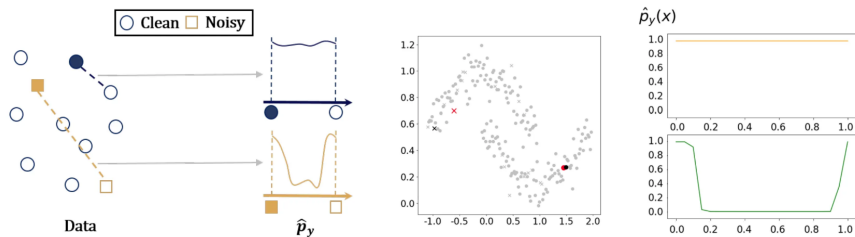
- $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ : Observed input vector.
- $y \in [K]$ : Observed label, where  $[K] = \{1, \dots, K\}$ .
- $y^{\text{gt}} \in [K]$ : Ground-truth label, where  $[K] = \{1, \dots, K\}$ .
- $(\mathbf{x}, y)$  is cleanly labeled if  $y = y^{\text{gt}}$  and noisily labeled if  $y \neq y^{\text{gt}}$ .
- Let  $\mathcal{D}^{\text{tr}} = \{(\mathbf{x}_i, y_i), i \in [n]\}$  be a training data set, and  $\mathcal{C}^{\text{tr}} = \{(\mathbf{x}, y) \in \mathcal{D}^{\text{tr}} : y = y^{\text{gt}}\}$  be the set of clean labeled samples.
- Our goal is to identify the clean labeled subset  $\mathcal{C}^{\text{tr}}$  from  $\mathcal{D}^{\text{tr}}$  accurately.

# Proposed Method: IOFM

- Let  $p(\mathbf{x}; \theta) : \mathbb{R}^d \rightarrow \mathbb{R}^K$  be a discriminative DNN parametrized by  $\theta$ .
- Also, let  $h(\mathbf{x}; \theta)$  be the feature vector of  $p(\mathbf{x}; \theta)$ , the output of the DNN's highest hidden layer.
- Furthermore, let  $p(\mathbf{x}; \hat{\theta})(\hat{p}(\mathbf{x}))$  be a DNN that perfectly memorizes  $\mathcal{D}^{\text{tr}}$  and  $h(\mathbf{x}; \hat{\theta})(\hat{h}(\mathbf{x}))$  be its feature vector.
- For a given training sample  $(\mathbf{x}_*, y_*) \in \mathcal{D}^{\text{tr}}$ , let  $(\mathbf{x}_{\text{nbd}}, y_{\text{nbd}}) \in \mathcal{D}^{\text{tr}}$  be the nearest to  $(\mathbf{x}_*, y_*)$  on the feature space, defined by

$$\mathbf{x}_{\text{nbd}} = \underset{\mathbf{x} \in \mathcal{D}^{\text{tr}} \setminus \{\mathbf{x}_*\}}{\operatorname{argmin}} \left\| \hat{h}(\mathbf{x}) - \hat{h}(\mathbf{x}_*) \right\|_2.$$

# Proposed Method: IOFM



- The authors have noticed a distinct difference in the behavior of  $\hat{p}_{y_*}(\mathbf{x})$  depending on its cleanness, which will be utilized for developing the new score.
- That is, the area under  $\hat{p}_{y_*}(\mathbf{x})$  over the intervals between  $\mathbf{x}_*$  and  $\mathbf{x}_{\text{nb},l}$ , defined as

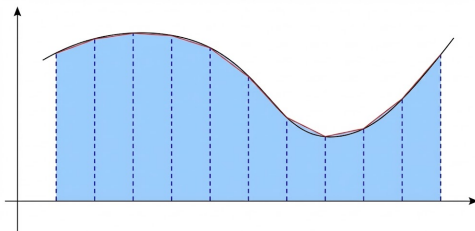
$$\int_0^1 \hat{p}_{y_*}(\alpha \mathbf{x}_* + (1 - \alpha) \mathbf{x}_{\text{nb},l}) d\alpha$$

is large when  $y_*$  is clean while it is relatively small when  $y_*$  is corrupted.

- Let  $\{\mathbf{x}_{\text{nb},l}\}_{l=1}^L \subset \{\mathbf{x}_i : y_i = y_*, i \in [n]\} \setminus \{\mathbf{x}_*\}$  be  $L$  neighborhood training inputs of  $\mathbf{x}_*$  in the feature space. Then, they proposed the following averaged score

$$\frac{1}{L} \sum_{l=1}^L \int_0^1 \hat{p}_{y_*}(\alpha \mathbf{x}_* + (1 - \alpha) \mathbf{x}_{\text{nb},l}) d\alpha. \quad (1)$$

# Proposed Method: IOFM



- Finally, The authors approximate the integration in (1) by the trapezoidal rule as follows:

$$s^{\text{IOFM}}(\mathbf{x}_*, y_*) = \frac{1}{L} \sum_{l=1}^L \sum_{h=1}^H \frac{1}{2H} (\hat{p}_y(\mathbf{x}_{l,h-1}) + \hat{p}_y(\mathbf{x}_{l,h})),$$

where  $\mathbf{x}_{l,h} = \frac{H-h}{H} \mathbf{x}_* + \frac{h}{H} \mathbf{x}_{\text{nbd},l}$  and  $H$  is the number of trapezoids, to have the IOFM score.



# Improvement of IOFM

- **MixUp Loss Function:** To enhance the smoothness of a trained DNN more while keeping memorizing training data, it would be helpful to use [2, Mixup], whose loss function is given as:

$$\mathbb{E}_{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \sim \mathcal{D}^{\text{tr}}} \mathbb{E}_{\lambda \sim B(\alpha, \alpha)} \text{CE}(\text{Mix}_{\lambda}(y_1, y_2), p(\text{Mix}_{\lambda}(\mathbf{x}_1, \mathbf{x}_2); \theta))$$

- **Use of Multiple IOFM Scores:** During the training process of the DNN until over-fitting, the authors calculate the IOFM scores at different training epochs and take the average for the final score.
- **Computation Time Reduction in Searching Neighbors:** This issue can be addressed by only restricting the neighbor search among a small subset of randomly sampled training data that share the same label.

# Experiment

Dataset	CIFAR100					
Imbalance type	Step			Long-tail		
Noise type	symm.		Asymm.	symm.		Asymm.
Noise rate ( $r$ )	0.3	0.5	0.2	0.3	0.5	0.2
Loss	0.958(0.698)	0.938(0.632)	0.826(0.674)	0.939(0.670)	0.896(0.609)	0.765(0.667)
Ens-Loss	0.969(0.945)	0.951(0.912)	0.829(0.816)	0.954(0.918)	0.915(0.873)	0.796(0.791)
Margin	0.952(0.689)	0.927(0.626)	0.847(0.684)	0.929(0.660)	0.887(0.603)	0.793(0.676)
AUM	0.966(0.851)	0.947(0.779)	0.871(0.802)	0.951(0.798)	0.911(0.727)	0.827(0.775)
sinIOFM	0.965(0.902)	0.949(0.874)	0.862(0.766)	0.953(0.872)	0.909(0.833)	0.826(0.723)
IOFM	<b>0.973(0.958)</b>	<b>0.958(0.936)</b>	<b>0.911(0.910)</b>	<b>0.961(0.942)</b>	<b>0.927(0.905)</b>	<b>0.875(0.874)</b>

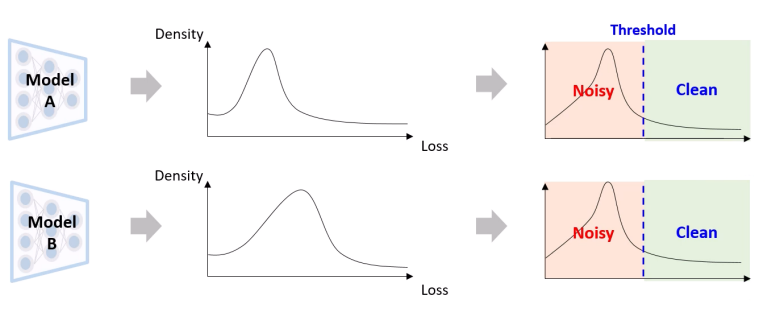
Table 1: Comparison of the clean/noisy classification AUC values on the imbalanced CIFAR100. We list the best and final (in the parentheses) results.

Dataset	CIFAR10			CIFAR100		
Noise type	Symm.		Asymm.	Symm.		Asymm.
Noise rate ( $r$ )	0.8	0.9	0.4	0.8	0.9	0.4
Loss	0.919(0.569)	0.836(0.556)	0.933(0.753)	0.847(0.557)	0.708(0.521)	0.621(0.552)
Ens-Loss	0.947(0.895)	0.867(0.763)	0.908(0.901)	0.873(0.704)	0.733(0.660)	0.642(0.638)
Margin	0.918(0.567)	0.834(0.554)	<b>0.940(0.756)</b>	0.833(0.553)	0.704(0.520)	0.634(0.553)
AUM	0.948(0.809)	0.869(0.687)	0.927(0.874)	0.874(0.693)	0.734(0.588)	0.677(0.637)
sinIOFM	0.924(0.699)	0.842(0.622)	0.912(0.825)	0.859(0.684)	0.714(0.567)	0.664(0.607)
IOFM	<b>0.954(0.907)</b>	<b>0.887(0.811)</b>	0.934( <b>0.921</b> )	<b>0.890(0.806)</b>	<b>0.746(0.653)</b>	<b>0.713(0.712)</b>

Table 2: Clean/noisy classification AUC results on heavily noisy CIFAR10&100. The best and final (in the parentheses) results are listed.

# DivideMix in semi-supervised learning

- The IOFM can be applied to learn deep classification models with noisy labeled data.
- In this paper, the authors consider a combination of the IOFM with the [1, DivideMix], one of the state-of-the-art methods for learning classification models in the presence of noisy labels.
- To combine the IOFM and DivideMix, the authors simply substitute the per-sample losses used in the DivideMix with the IOFM scores.

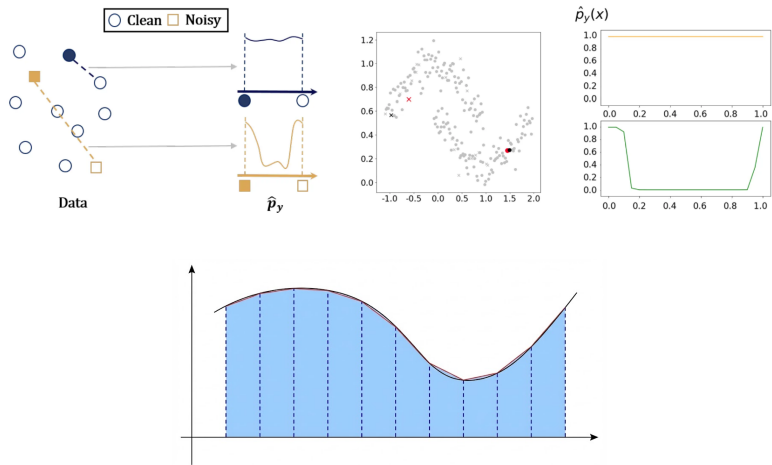


# Experiment

Dataset	CIFAR10			
Imbalance type	Step			
Noise type	symm.		Asymm.	
Noise rate ( $r$ )	0.3	0.5	0.2	0.4
Cross-Entropy	72.72	68.99	81.19	74.95
Mixup	74.81	65.63	81.97	75.93
DivideMix	85.76	85.16	87.15	78.35
IOFM+DivideMix	<b>88.19</b>	<b>88.07</b>	<b>87.35</b>	<b>78.73</b>

Table 3: Comparison of the best test accuracies(%) of various methods on the imbalanced CIFAR10.

# Appendix: Proposal



We consider two versions of the IOFM: 1) the original IOFM (IOFM) based on the ensemble of multiple IOFM scores from multiple epochs, and 2) the IOFM based on the single score at the last epoch (sinIOFM). We consider

## Appendix: Algorithm

---

### Algorithm 1: IOFM

In practice, we set  $(T_1, T_2, L, H) = (150, 10, 10, 3)$ .

---

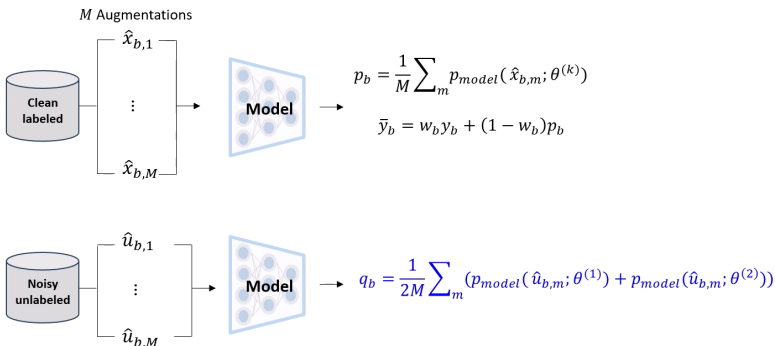
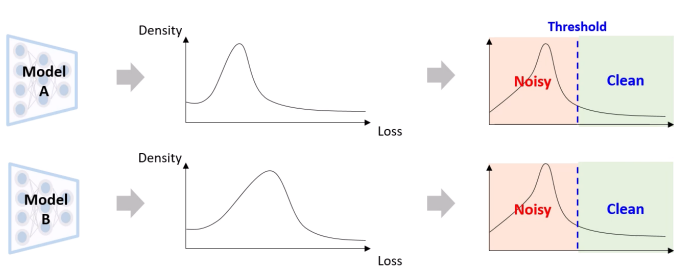
**input** Training data:  $\mathcal{D}^{\text{tr}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , a prediction model and its feature function:  $p(\cdot; \theta)$  and  $h(\cdot; \theta)$ , an optimizer  $\mathcal{O}$ , four integers:  $T_1, T_2, L$ , and  $H$ .

- 1:  $\mathcal{S}^{\text{ens}} \leftarrow \emptyset$  // Ensemble IOFM score set
- 2: **for** ( $ep = 1$  **to**  $T_1$ ) **do**
- 3:   *MixUp*( $f(\cdot; \theta), \mathcal{D}^{\text{tr}}, \mathcal{O}$ ) // train  $p(\cdot; \theta)$  using MixUp
- 4:   **if** ( $ep \bmod T_2 = 0$ ) **then**
- 5:      $\mathcal{S}^{\text{tmp}} \leftarrow \emptyset$
- 6:     **for** ( $i = 1$  **to**  $n$ ) **do**
- 7:        $s_i \leftarrow s^{\text{IOFM}}(\mathbf{x}_i, y_i)$  //IOFM score of  $(\mathbf{x}_i, y_i)$
- 8:        $\mathcal{S}^{\text{tmp}} \leftarrow \text{append}(\mathcal{S}^{\text{tmp}}, s_i)$  // append  $s_i$  to  $\mathcal{S}^{\text{tmp}}$
- 9:     **end for**
- 10:     $\mathcal{S}^{\text{ens}} \leftarrow \mathcal{S}^{\text{ens}} + \mathcal{S}^{\text{tmp}}$
- 11:   **end if**
- 12: **end for**

**output**  $\mathcal{S}^{\text{ens}}$

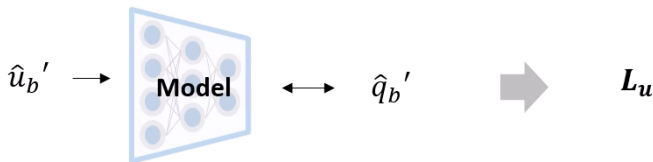
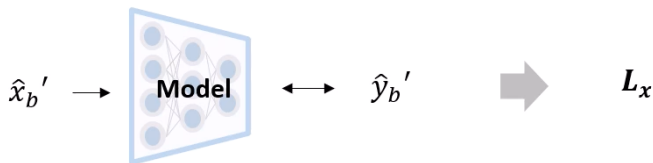
---

# Appendix: Dividemix & Mixmatch





## Appendix: Dividemix & Mixmatch

$$\begin{aligned}\hat{x}_b' &= \lambda' \hat{x}_1' + (1 - \lambda') \hat{x}_2' \\ \hat{y}_b' &= \lambda' \hat{y}_1' + (1 - \lambda') \hat{y}_2' \\ \hat{u}_b' &= \lambda' \hat{u}_1' + (1 - \lambda') \hat{u}_2' \\ \hat{q}_b' &= \lambda' \hat{q}_1' + (1 - \lambda') \hat{q}_2'\end{aligned}$$





# References I

-  Junnan Li, Richard Socher, and Steven C. H. Hoi. *DivideMix: Learning with Noisy Labels as Semi-supervised Learning*. 2020. arXiv: 2002.07394 [cs.CV]. URL: <https://arxiv.org/abs/2002.07394>.
-  Hongyi Zhang et al. *mixup: Beyond Empirical Risk Minimization*. 2018. arXiv: 1710.09412 [cs.LG]. URL: <https://arxiv.org/abs/1710.09412>.