



模式识别与神经计算研究组

PATtern Recognition and NEural Computing

---

# Detecting Corrupted Labels Without Training a Model to Predict

---

*Zhaowei Zhu Zihao Dong Yang Liu*

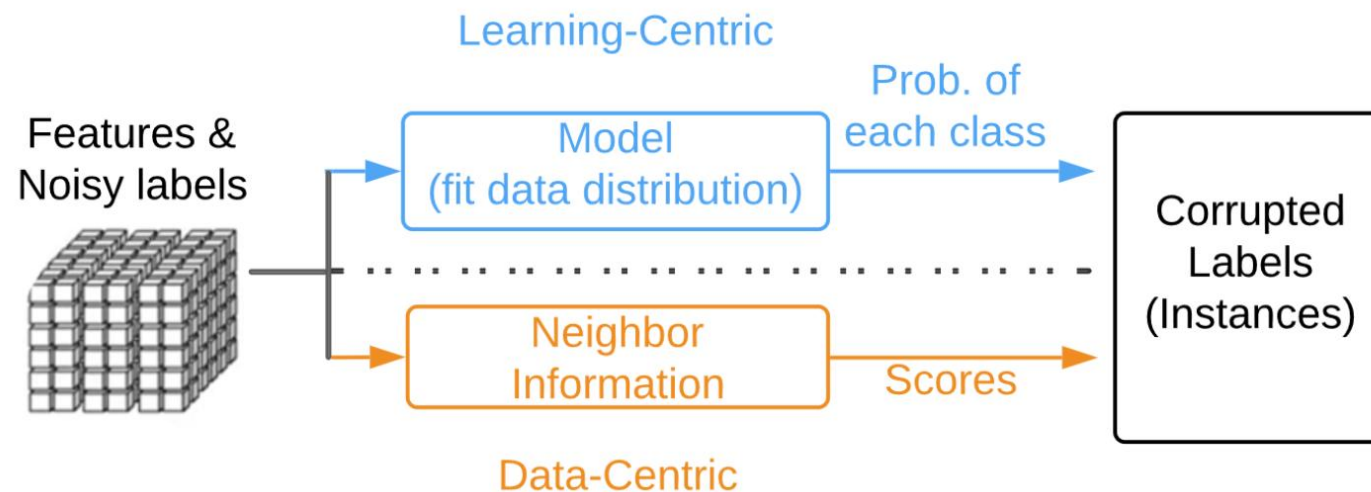
*ICML 2022*

# Background



- Label noise in real-world datasets encodes wrong correlation patterns and impairs the generalization of deep neural networks (DNNs).
- It is critical to find efficient ways to detect corrupted patterns.

- Learning-centric
  - prediction & compare
- Data-Centric
  - neighbor information



*Figure 1.* The existing learning-centric pipeline vs. our proposed data-centric pipeline. The inputs are features and the corresponding noisy labels, and the outputs are a set of corrupted labels. **Blue:** The learning-centric solution. **Orange:** The data-centric solution.

# Motivation

- Previous learning-centric methods
  - train DNNs with noisy supervisions
  - design robust loss functions
- Limitations of the learning-centric methods
  - task-specific and fine-tuning hyperparameters for different datasets/noise
  - as long as the model is trained with noisy supervisions, the memorization of corrupted instances exists.

The model will “subjectively” and wrongly treat the memorized/overfitted corrupted instances as clean.

- Solution
  - drop the dependency on the noisy supervision
  - design a training-free method to find label errors.

**Definition 2.1** ( $(k, \delta_k)$  label clusterability). A dataset  $D$  satisfies  $(k, \delta_k)$  label clusterability if:  $\forall n \in [N]$ , the feature  $x_n$  and its  $k$ -Nearest-Neighbors ( $k$ -NN)  $x_{n_1}, \dots, x_{n_k}$  belong to the same true class with probability at least  $1 - \delta_k$ .

- a probabilistic  $Y$  given  $X$
  - the quality of features and the value of  $k$
- 
- The  $(k, 0)$  label clusterability is also known as  $k$ -NN label clusterability

# Method

- share the same noise transition matrix

$$\mathbb{P}(\tilde{Y} = \tilde{y}_n \mid X = x_n, Y = y_n) = \mathbb{P}(\tilde{Y} = \tilde{y}_n \mid X = x_{n'}, Y = y_n), \forall x_{n'} \in \{x_{n_1}, \dots, x_{n_k}\}$$

- using appropriate features may be better than model logits/predictions when the dataset is noisy.

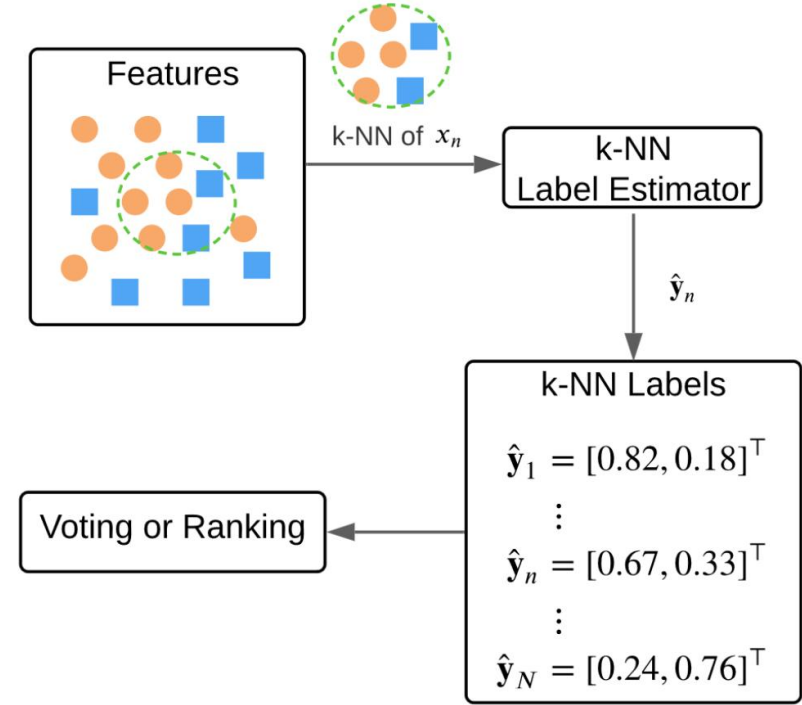


Figure 2. Detect corrupted labels with similar features. **Orange circle**: instance with noisy label 1. **Blue square**: instance with noisy label 2. **Green dashed circle**: A  $k$ -NN example.

- The F1 score on corrupted instances is sensitive to the case when the noise rate is mild to low, which is typically the case in practice.

$$F_1 = 2/(\text{Precision}^{-1} + \text{Recall}^{-1}).$$

$$\text{Precision} = \frac{\sum_{n \in [N]} \mathbb{1}(v_n = 1, \tilde{y}_n \neq y_n)}{\sum_{n \in [N]} \mathbb{1}(v_n = 1)},$$

$$\text{Recall} = \frac{\sum_{n \in [N]} \mathbb{1}(v_n = 1, \tilde{y}_n \neq y_n)}{\sum_{n \in [N]} \mathbb{1}(\tilde{y}_n \neq y_n)}.$$



- Voting-Based Local Detection

$$y_n^{\text{vote}} = \arg \max_{i \in [K]} \hat{\mathbf{y}}_n[i]$$

$$v_n := \mathbf{1}(y_n^{\text{vote}} \neq \tilde{y}_n)$$

- Ranking-Based Global Detection

- scoring function
- threshold

*Property 3.1* (Relative score). Within the same instance, the score of the majority class is higher than the others, i.e.,  $\forall j \neq y_n^{\text{vote}}, j \in [K], \forall n \in [N] : \text{Score}(\hat{\mathbf{y}}_n, y_n^{\text{vote}}) > \text{Score}(\hat{\mathbf{y}}_n, j)$ .

*Property 3.2* (Absolute score).  $\text{Score}(\hat{\mathbf{y}}_n, j)$  is jointly determined by both  $\hat{\mathbf{y}}_n[j]$  and  $\hat{\mathbf{y}}_n[j']$ ,  $\forall j' \neq j$ .

$$\text{Score}(\hat{\mathbf{y}}_n, j) = \frac{\hat{\mathbf{y}}_n^\top \mathbf{e}_j}{\|\hat{\mathbf{y}}_n\|_2 \|\mathbf{e}_j\|_2},$$

$$\mathcal{I} = \text{argsort}\{\text{Score}(\hat{\mathbf{y}}_n, j)\}_{n \in \mathcal{N}_j},$$

$$v_n = \mathbf{1}(\text{Loc}(n, \mathcal{I}) \leq \tilde{N}_j),$$

$$\tilde{N}_j = \mathbb{P}(Y \neq j \mid \tilde{Y} = j) \cdot N_j$$

$$\mathbb{P}(Y \neq j \mid \tilde{Y} = j) = 1 - \mathbb{P}(Y = j \mid \tilde{Y} = j),$$

- Borrow the results from the HOC , where the noise transition probability and the marginal distribution of clean label can be estimated with only features and the corresponding noisy labels.
- Bayes' rule:

$$\mathbb{P}(Y = j \mid \tilde{Y} = j) = \mathbb{P}(\tilde{Y} = j \mid Y = j) \cdot \mathbb{P}(Y = j) / \mathbb{P}(\tilde{Y} = j),$$



---

**Algorithm 1** Detection with **Similar Features** (The SimiFeat Detector)

---

```

1: Input: Number of epochs:  $M$ .  $k$ -NN parameter:  $k$ . Noisy dataset:  $\tilde{D} = \{(x_n, \tilde{y}_n)\}_{n \in [N]}$ . Feature extractor:  $g(\cdot)$ .
   Method: Vote or Rank. Epoch counter  $m = 0$ .
2: repeat
3:    $x'_n \leftarrow \text{RandPreProcess}(x_n), \forall n;$  # Initialize & Standard data augmentations
4:    $x_n \leftarrow g(x'_n), \forall n;$  # For tasks with rarely clusterable features, extract features with  $g(\cdot)$ 
5:    $\hat{y}_n \leftarrow \text{kNNLabel}(\{x_n\}_{n \in [N]}, k)$  # Get soft labels. One can weight instances by the similarity to the center instance.
6:   if Vote then
7:      $y_n^{\text{vote}} \leftarrow \arg \max_{i \in [K]} \hat{y}_n[i];$  # Apply local majority vote
8:      $v_n \leftarrow \mathbb{1}(y_n^{\text{vote}} \neq \tilde{y}_n), \forall n \in [N];$  # Treat as corrupted if majority votes disagree with noisy labels
9:   else
10:     $\mathbb{P}(Y), \mathbb{P}(\tilde{Y}|Y) \leftarrow \text{HOC}(\{(x_n, \tilde{y}_n)\}_{n \in [N]});$  # Estimate clean priors  $\mathbb{P}(Y)$  and noise transitions  $\mathbb{P}(\tilde{Y}|Y)$  by HOC
11:     $\mathbb{P}(Y|\tilde{Y}) = \mathbb{P}(\tilde{Y}|Y) \cdot \mathbb{P}(Y) / \mathbb{P}(\tilde{Y});$  # Estimate thresholds by Bayes' rule
12:    for  $j$  in  $[K]$  do
13:       $\mathcal{N}_j := \{n | \tilde{y}_n = j\};$  # Detect corrupted labels in each set  $\mathcal{N}_j$ 
14:       $\mathcal{I} \leftarrow \text{argsort}\{\text{Score}(\hat{y}_n, j)\}_{n \in \mathcal{N}_j};$  #  $\mathcal{I}$  records the raw index of each sorted value
15:       $v_n \leftarrow \mathbb{1}(\text{Loc}(n, \mathcal{I}) \leq \lfloor (1 - \mathbb{P}(Y = j | \tilde{Y} = j)) \cdot N_j \rfloor);$  # Select low-score (head) instances as corrupted ones
16:    end for
17:  end if
18:   $\mathcal{V}_m = \{v_n\}_{n \in [N]};$  # Record detection results in the  $m$ -th epoch
19: until  $M$  times
20:  $\mathcal{V} = \text{Vote}(\mathcal{V}_m, \forall m \in [M]);$  # Do majority vote based on results from  $M$  epochs
21: Output:  $[N] \setminus \mathcal{V}$ .

```

---

# Experiments



- Fitting Noisy Distributions May Not Be Necessary

Table 1. Comparisons of  $F_1$ -scores (%). CORES, CL, TracIn: Train with noisy supervisions. SimiFeat-V and SimiFeat-R: Get  $g(\cdot)$  without any supervision. Top 2 are **bold**.

METHOD	CIFAR10				CIFAR100			
	Human	Symm. 0.6	Asym. 0.3	Inst. 0.4	Human	Symm. 0.6	Asym. 0.3	Inst. 0.4
CORES	65.00	92.94	7.68	<b>87.43</b>	3.52	<b>92.34</b>	0.02	9.67
CL	55.85	80.59	76.45	62.89	64.58	78.98	52.96	50.08
TRACIN	55.02	76.94	73.47	58.85	61.75	76.74	48.42	49.89
DEEP $k$ -NN	56.21	82.35	75.24	63.08	57.40	70.69	56.75	63.85
SIMIFEAT-V	<b>82.30</b>	<b>93.21</b>	<b>82.52</b>	81.09	<b>73.19</b>	84.48	<b>65.42</b>	<b>74.26</b>
SIMIFEAT-R	<b>83.28</b>	<b>95.56</b>	<b>83.58</b>	<b>82.26</b>	<b>74.67</b>	<b>88.68</b>	<b>62.89</b>	<b>73.53</b>

- Features May Be Better Than Model Predictions

Table 2. Comparisons of  $F_1$ -scores (%). CORES, CL, TracIn: Use logit layers. SimiFeat-V/R: Use only representations. **All methods use the same fixed extractor from CLIP**. Top 2 are **bold**.

METHOD	CIFAR10				CIFAR100			
	Human	Symm. 0.6	Asym. 0.3	Inst. 0.4	Human	Symm. 0.6	Asym. 0.3	Inst. 0.4
CE SIEVE	67.21	94.56	5.24	8.41	16.24	88.55	2.6	1.63
CORES	83.18	<b>96.94</b>	12.05	<b>88.89</b>	38.52	<b>92.33</b>	7.02	<b>85.52</b>
CL	69.76	95.03	77.14	62.91	67.64	85.67	62.58	61.53
TRACIN	81.85	95.96	80.75	64.97	<b>79.32</b>	<b>91.03</b>	63.12	64.31
DEEP $k$ -NN	82.98	87.47	76.96	77.42	72.33	82.95	64.96	74.25
SIMIFEAT-V	<b>87.43</b>	96.44	<b>88.97</b>	87.11	76.26	86.88	<b>73.50</b>	<b>80.03</b>
SIMIFEAT-R	<b>87.45</b>	<b>96.74</b>	<b>89.04</b>	<b>91.14</b>	<b>79.21</b>	90.54	<b>68.14</b>	77.37



# Experiments

- The Effect of the Quality of Features

Table 3. Comparisons of  $F_1$ -scores (%) using  $g(\cdot)$  with different  $\delta_k$  (%). Model names are the same as Figure 3.

PRE-TRAINED MODEL	CIFAR10			CIFAR100		
	$1 - \delta_k$	Human	Inst. 0.4	$1 - \delta_k$	Human	Inst. 0.4
R18-IMG	35.73	75.40	80.22	11.30	74.91	71.99
R34-IMG	48.13	79.52	82.43	16.17	76.88	74.00
R50-IMG	45.77	78.40	82.06	15.81	76.55	73.51
ViT-B/32-CLIP	64.12	87.45	91.14	19.94	79.21	77.37
R34-C10-SSL	69.31	83.28	85.26	2.59	68.03	65.94
R34-C10-CLEAN	99.41	98.39	98.59	0.22	60.90	60.73
R34-C100-SSL	18.59	59.96	74.99	22.46	74.67	73.53
R34-C100-CLEAN	18.58	60.17	76.41	89.07	92.87	95.29

Table 4. Experiments on Clothing1M. None: Standard training with 1M noisy data. R50-Img (or ViT-B/32-CLIP, R50-Img Warmup-1): Apply our method with ResNet50 pre-trained on ImageNet (or ViT-B/32 pre-trained by CLIP, R50-Img with 1-epoch warmup). The clean test accuracy on the best epoch, the last 10 epochs, and the last epoch, are listed. Top-1 is **bold**.

DATA SELECTION	# TRAINING	BEST EPOCH	LAST 10	LAST
NONE	1M (100%)	70.32	$69.44 \pm 0.13$	69.53
R50-IMG	770K (77.0%)	72.37	$71.95 \pm 0.08$	71.89
ViT-B/32-CLIP	700K (70.0%)	72.54	$72.23 \pm 0.17$	72.11
R50-IMG WARMUP-1	767K (76.7%)	<b>73.64</b>	<b><math>73.28 \pm 0.18</math></b>	<b>73.41</b>



Thanks