

Noisy Label Learning with Instance-Dependent Outliers: Identifiability via Crowd Wisdom

Tri Nguyen¹ Shahana Ibrahim² Xiao Fu¹

¹School of EECS
Oregon State University

²Department of ECE
University of Central Florida

NeurIPS 2024



Noisy Label Learning Problem

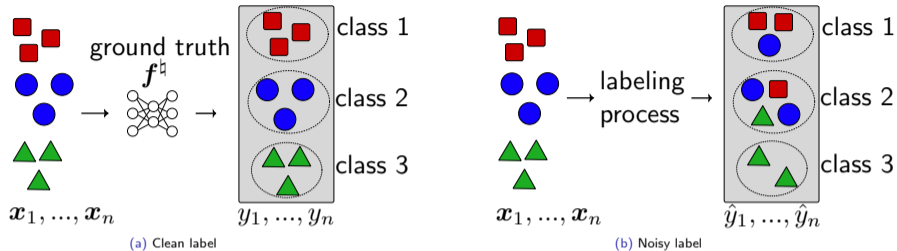


Figure: Source: internet/chatgpt.

Noisy Label Learning Problem

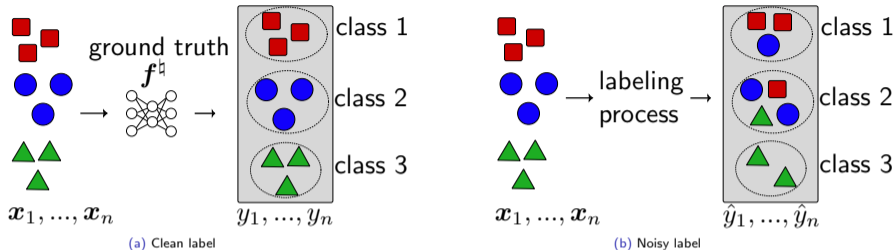


Figure: Source: internet/chatgpt.

Goal: Recover the ground truth (GT) classifier f^* given $(x_1, \hat{y}_1), \dots, (x_N, \hat{y}_N)$.

Noise Generation Modeling Approach



- ▶ Among approaches including noisy label filtering [1–4], robust noise design [5–9], the **noise generation modeling** [10–18] approach is the most popular.

- ▶ Among approaches including noisy label filtering [1–4], robust noise design [5–9], the **noise generation modeling** [10–18] approach is the most popular.
- ▶ Let \mathbf{x}, y, \hat{y} be the sample's feature, its true label, and the obtained noisy label,

$$\underbrace{\begin{bmatrix} \Pr(\hat{y} = 1 | \mathbf{x}) \\ \dots \\ \Pr(\hat{y} = K | \mathbf{x}) \end{bmatrix}}_{\mathbf{g}^{\natural}(\mathbf{x})} = \underbrace{\begin{bmatrix} \Pr(\hat{y} = 1 | y = 1, \mathbf{x}) & \dots & \Pr(\hat{y} = 1 | y = K, \mathbf{x}) \\ \dots & \dots & \dots \\ \Pr(\hat{y} = K | y = 1, \mathbf{x}) & \dots & \Pr(\hat{y} = K | y = K, \mathbf{x}) \end{bmatrix}}_{\mathbf{T}^{\natural}(\mathbf{x})} \underbrace{\begin{bmatrix} \Pr(y = 1 | \mathbf{x}) \\ \dots \\ \Pr(y = K | \mathbf{x}) \end{bmatrix}}_{\mathbf{f}^{\natural}(\mathbf{x})}$$

Noise Generation Modeling Approach

- ▶ Among approaches including noisy label filtering [1–4], robust noise design [5–9], the **noise generation modeling** [10–18] approach is the most popular.
- ▶ Let \mathbf{x}, y, \hat{y} be the sample's feature, its true label, and the obtained noisy label,

$$\underbrace{\begin{bmatrix} \Pr(\hat{y} = 1 | \mathbf{x}) \\ \dots \\ \Pr(\hat{y} = K | \mathbf{x}) \end{bmatrix}}_{\mathbf{g}^{\natural}(\mathbf{x})} = \underbrace{\begin{bmatrix} \Pr(\hat{y} = 1 | y = 1, \mathbf{x}) & \dots & \Pr(\hat{y} = 1 | y = K, \mathbf{x}) \\ \dots & \dots & \dots \\ \Pr(\hat{y} = K | y = 1, \mathbf{x}) & \dots & \Pr(\hat{y} = K | y = K, \mathbf{x}) \end{bmatrix}}_{\mathbf{T}^{\natural}(\mathbf{x})} \underbrace{\begin{bmatrix} \Pr(y = 1 | \mathbf{x}) \\ \dots \\ \Pr(y = K | \mathbf{x}) \end{bmatrix}}_{\mathbf{f}^{\natural}(\mathbf{x})}$$

$$\underbrace{\mathbf{g}^{\natural}(\mathbf{x})}_{\text{noisy label probability vector}} = \underbrace{\mathbf{T}^{\natural}(\mathbf{x})}_{\text{confusion matrix}} \underbrace{\mathbf{f}^{\natural}(\mathbf{x})}_{\text{true label probability vector}}$$

Instance-Dependent Noise Model

Noise generation model:

$$\mathbf{g}^{\natural}(\mathbf{x}) = \mathbf{T}^{\natural}(\mathbf{x})\mathbf{f}^{\natural}(\mathbf{x}).$$

- ▶ Learning under instance-dependent confusion matrices is an ill-posed problem.

Instance-Dependent Noise Model

Noise generation model:

$$\mathbf{g}^{\natural}(\mathbf{x}) = \mathbf{T}^{\natural}(\mathbf{x}) \mathbf{f}^{\natural}(\mathbf{x}).$$

- ▶ Learning under instance-dependent confusion matrices is an ill-posed problem.
- ▶ Most existing works [12, 14, 15, 17, 19, 20] resort to a simplification: $\mathbf{T}^{\natural}(\mathbf{x}) = \mathbf{A}^{\natural}$, $\forall \mathbf{x}$.

Instance-Dependent Noise Model

Noise generation model:

$$\mathbf{g}^{\natural}(\mathbf{x}) = \mathbf{T}^{\natural}(\mathbf{x}) \mathbf{f}^{\natural}(\mathbf{x}).$$

- ▶ Learning under instance-dependent confusion matrices is an ill-posed problem.
- ▶ Most existing works [12, 14, 15, 17, 19, 20] resort to a simplification: $\mathbf{T}^{\natural}(\mathbf{x}) = \mathbf{A}^{\natural}$, $\forall \mathbf{x}$.
- ▶ However, real data exhibits a more complex confusion matrix.



Figure: Nominal images (left) exhibits similar labeling difficulty, whereas special/outlier images (right) display a wide range of labeling challenges.

Instance-Dependent Noise Model

Noise generation model:

$$\mathbf{g}^{\natural}(\mathbf{x}) = \mathbf{T}^{\natural}(\mathbf{x})\mathbf{f}^{\natural}(\mathbf{x}).$$

- ▶ Learning under instance-dependent confusion matrices is an ill-posed problem.
- ▶ Most existing works [12, 14, 15, 17, 19, 20] resort to a simplification: $\mathbf{T}^{\natural}(\mathbf{x}) = \mathbf{A}^{\natural}$, $\forall \mathbf{x}$.
- ▶ However, real data exhibits a more complex confusion matrix.



Figure: Nominal images (left) exhibits similar labeling difficulty, whereas special/outlier images (right) display a wide range of labeling challenges.

- ▶ We consider an instance-dependent noise model

Nominal samples: $\mathbf{T}^{\natural}(\mathbf{x}_n) = \mathbf{A}^{\natural}$, if $n \in \mathcal{O} \subseteq [N]$

Outlier samples: $\mathbf{T}^{\natural}(\mathbf{x}_n) = \mathbf{A}^{\natural}(\mathbf{x}_n)$ for some $\mathbf{A}^{\natural}(\cdot)$, otherwise.

Identifiability Guarantee

Proposed criterion:

$$\underset{\{\mathbf{A}_m \in \mathcal{A}\}, \{\mathbf{e}_n^{(m)} \in \mathcal{E}\}, \mathbf{f} \in \mathcal{F}}{\text{minimize}} \quad \mathbf{L}_{\text{ce}} \triangleq -\frac{1}{S} \sum_{(m,n) \in \mathcal{S}} \sum_{k=1}^K \mathbb{1}[\hat{y}_n^{(m)} = k] \log \left[\mathbf{A}_m \mathbf{f}(\mathbf{x}_n) + \mathbf{e}_n^{(m)} \right]_k, \quad (1a)$$

$$\text{subject to} \quad \sum_{n=1}^N \mathbb{1} \left\{ \sum_{m=1}^M \|\mathbf{e}_n^{(m)}\|_2 > 0 \right\} \leq C, \quad (1b)$$

Identifiability Guarantee

Proposed criterion:

$$\underset{\{\mathbf{A}_m \in \mathcal{A}\}, \{\mathbf{e}_n^{(m)} \in \mathcal{E}\}, \mathbf{f} \in \mathcal{F}}{\text{minimize}} \quad \mathcal{L}_{\text{ce}} \triangleq -\frac{1}{S} \sum_{(m,n) \in \mathcal{S}} \sum_{k=1}^K \mathbb{1}[\hat{y}_n^{(m)} = k] \log \left[\mathbf{A}_m \mathbf{f}(\mathbf{x}_n) + \mathbf{e}_n^{(m)} \right]_k, \quad (1a)$$

$$\text{subject to} \quad \sum_{n=1}^N \mathbb{1} \left\{ \sum_{m=1}^M \|\mathbf{e}_n^{(m)}\|_2 > 0 \right\} \leq C, \quad (1b)$$

Theorem (Identifiability and Generalization)

Let $(\{\hat{\mathbf{A}}_m\}, \{\hat{\mathbf{e}}_n^{(m)}\}, \hat{\mathbf{f}})$ be any optimal solution of (1). The following result holds with probability greater than $1 - 2/S - K/T^\alpha$:

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\mathbf{x}}} \left[\min_{\mathbf{\Pi}} \|\hat{\mathbf{f}}(\mathbf{x}) - \mathbf{\Pi}^\top \mathbf{f}^\natural(\mathbf{x})\|_2^2 \right] \leq K(\eta + \xi_1 + \xi_2),$$

$$\min_{\mathbf{\Pi}} \|\hat{\mathbf{A}}_m - \mathbf{A}_m^\natural \mathbf{\Pi}\|_F^2 = K\sigma^2(\eta + \xi_1 + \xi_2), \quad \forall m,$$

where $\eta^2 = \mathcal{O} \left(\beta M T^\alpha \sqrt{S} (\sqrt{M \log S} + (\|\mathbf{X}\| R_{\mathcal{F}})^{0.25}) \right)$, $\mathbf{\Pi}$ a permutation matrix, and $T = N - |\mathcal{I}|$. In addition, we have exact outlier detection, i.e., $\hat{\mathcal{I}} = \mathcal{I}$.

Identifiability Guarantee

Proposed criterion:

$$\underset{\{\mathbf{A}_m \in \mathcal{A}\}, \{\mathbf{e}_n^{(m)} \in \mathcal{E}\}, \mathbf{f} \in \mathcal{F}}{\text{minimize}} \quad \mathcal{L}_{\text{ce}} \triangleq -\frac{1}{S} \sum_{(m,n) \in \mathcal{S}} \sum_{k=1}^K \mathbb{1}[\hat{y}_n^{(m)} = k] \log \left[\mathbf{A}_m \mathbf{f}(\mathbf{x}_n) + \mathbf{e}_n^{(m)} \right]_k, \quad (1a)$$

$$\text{subject to} \quad \sum_{n=1}^N \mathbb{1} \left\{ \sum_{m=1}^M \|\mathbf{e}_n^{(m)}\|_2 > 0 \right\} \leq C, \quad (1b)$$

Theorem (Identifiability and Generalization)

Let $(\{\hat{\mathbf{A}}_m\}, \{\hat{\mathbf{e}}_n^{(m)}\}, \hat{\mathbf{f}})$ be any optimal solution of (1). The following result holds with probability greater than $1 - 2/S - K/T^\alpha$:

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\mathbf{x}}} \left[\min_{\mathbf{\Pi}} \|\hat{\mathbf{f}}(\mathbf{x}) - \mathbf{\Pi}^\top \mathbf{f}^\natural(\mathbf{x})\|_2^2 \right] \leq K(\eta + \xi_1 + \xi_2),$$

$$\min_{\mathbf{\Pi}} \|\hat{\mathbf{A}}_m - \mathbf{A}_m^\natural \mathbf{\Pi}\|_F^2 = K\sigma^2(\eta + \xi_1 + \xi_2), \quad \forall m,$$

where $\eta^2 = \mathcal{O} \left(\beta M T^\alpha \sqrt{S} (\sqrt{M \log S} + (\|\mathbf{X}\| R_{\mathcal{F}})^{0.25}) \right)$, $\mathbf{\Pi}$ a permutation matrix, and $T = N - |\mathcal{I}|$. In addition, we have exact outlier detection, i.e., $\hat{\mathcal{I}} = \mathcal{I}$.

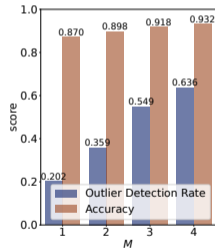


Figure: Performance of the proposal on CIFAR-10 with synthetic labels against different number of annotators.

Experiments using Real Datasets

Data.

- ▶ CIFAR10-N [21]. $N = 60000, K = 10, M = 3$. The error rates of annotators are 17.23%, 18.12%, and 17.64%.
- ▶ LabelMe [22, 23]. $N = 2688, K = 8, M = 59$. The average error rate is 25.95%.
- ▶ ImageNet-15N: we acquire noisy annotations by asking AMT workers to annotate some images from ImageNet. $K = 15, N = 2,514, M = 100$. The average error rate of the annotators is 42.68%.

Table: Average classification accuracy on CIFAR-10N, LabelMe, and ImageNet-15N datasets, labeled by human annotators. **Bold black** represents the best and **blue** represents the second best.

Method/Dataset	CIFAR-10N	LabelMe	ImageNet-15N
PTD	89.52 ± 0.24	84.18 ± 1.36	65.53 ± 0.18
BLTM	75.68 ± 0.47	82.10 ± 0.56	66.57 ± 0.76
VolMinNet	86.58 ± 0.21	79.97 ± 0.16	63.11 ± 1.08
Reweight	89.56 ± 0.30	84.51 ± 0.50	65.85 ± 2.93
GCE	78.01 ± 7.23	83.41 ± 0.59	64.71 ± 1.38
MEIDTM	68.69 ± 0.31	83.53 ± 0.21	72.66 ± 0.58
CrowdLayer	87.38 ± 0.43	82.80 ± 0.90	61.36 ± 2.73
TraceReg	86.57 ± 0.24	82.83 ± 0.23	68.43 ± 0.12
MaxMIG	90.11 ± 0.09	83.73 ± 0.84	81.13 ± 1.42
GeoCrowdNet (F)	87.19 ± 0.37	87.21 ± 0.39	80.45 ± 1.77
GeoCrowdNet (W)	86.43 ± 0.44	82.83 ± 0.75	68.79 ± 0.27
COINNet (Ours)	92.09 ± 0.47	87.60 ± 0.54	93.71 ± 3.26

Qualitative Results

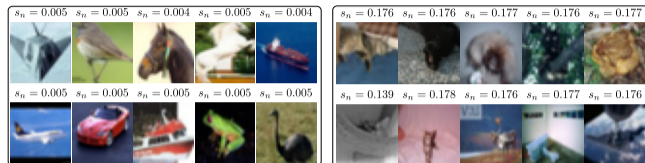


Figure: Examples from CIFAR-10N with low (left) and high (right) $s_n = \sum_{m=1}^M \|\hat{\mathbf{e}}_n^{(m)}\|^2$.

Qualitative Results

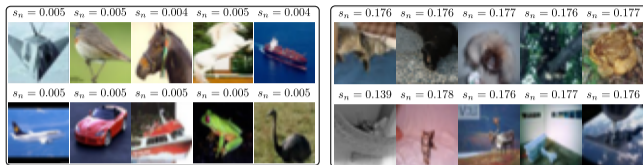


Figure: Examples from CIFAR-10N with low (left) and high (right) $s_n = \sum_{m=1}^M \|\hat{\mathbf{e}}_n^{(m)}\|^2$.

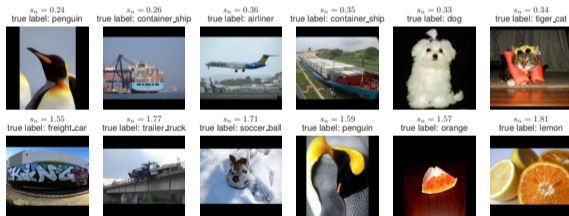


Figure: Examples from ImageNet-15N with low (top) and high (bottom) $s_n = \sum_{m=1}^M \|\hat{\mathbf{e}}_n^{(m)}\|^2$.