

Learning with Biased Labels

Cognitive Computing Lab, Baidu Research

Zheng He

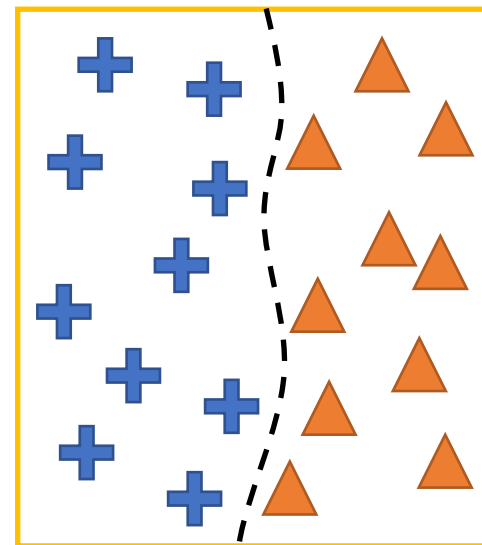
- Label noise robust learning
 - Background
 - Methods
 - New directions
- Quantized label learning
 - Motivation
 - Problem setting
 - Method
 - Experiment

Review of Label Noise Robust Learning

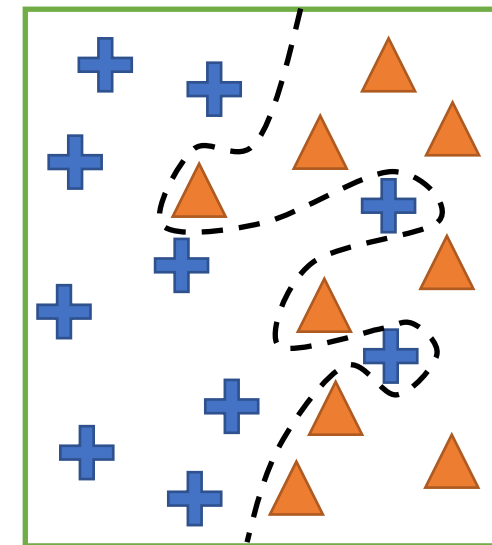
Problem definition

- Clean distribution: $p(x, y)$
- Corrupted distribution: $p(x, \bar{y})$, \bar{y} is a noisy label which may not be true
- Given the training dataset $\bar{D} = \{x_i, \bar{y}_i\}_{i=1}^N$, the goal is to minimize the risk on clean distribution (test data):

$$R(f_\theta) = \mathbb{E}_{(x,y) \sim p(x,y)} [L(f(x), y)]$$



Standard supervised learning



Learning with label noise

Label noise types

- Class-dependent noise: the true label is corrupted by a *noise transition matrix* $T \in [0,1]^{c \times c}$, where c is the number of classes.

$$T_{i,j} := p(\bar{y} = j | y = i)$$

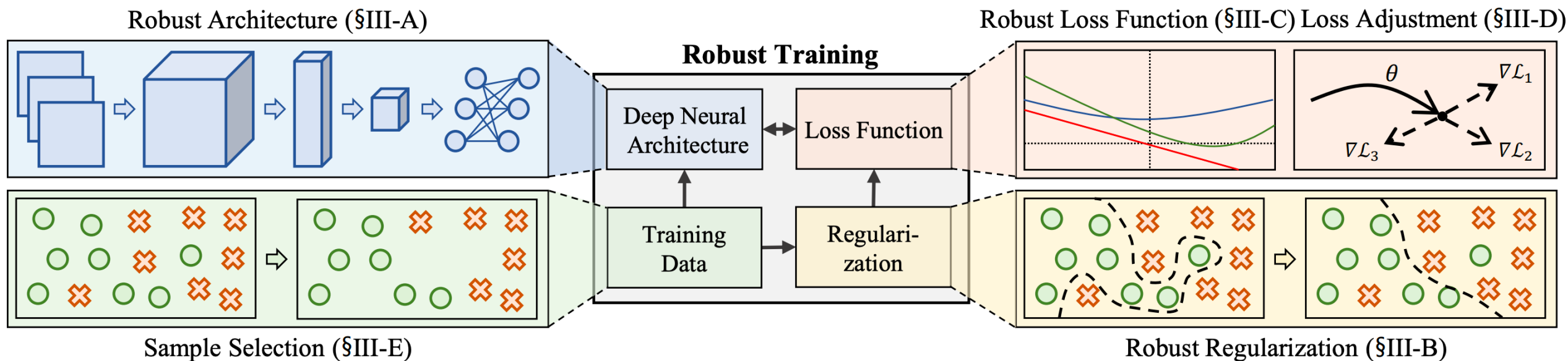
$$\begin{bmatrix} 1-\tau & \frac{\tau}{n-1} & \cdots & \frac{\tau}{n-1} \\ \frac{\tau}{n-1} & 1-\tau & & \frac{\tau}{n-1} \\ \vdots & & \ddots & \vdots \\ \frac{\tau}{n-1} & \frac{\tau}{n-1} & \cdots & 1-\tau \end{bmatrix}$$

(a) Sym-flipping.

$$\begin{bmatrix} 1-\tau & \tau & 0 & 0 \\ 0 & 1-\tau & \tau & 0 \\ \vdots & & \ddots & \vdots \\ 0 & & & \tau \\ \tau & 0 & \cdots & 1-\tau \end{bmatrix}$$

(b) Pair-flipping.

➤ Deep learning approaches for dealing with label noise



Categorization of recent deep learning methods for overcoming noisy labels.

Robust loss function

- We call a loss function L symmetric if it satisfies, for some constant K

$$\sum_{j=1}^c L(f(x), j) = K, \forall x \in X, \forall f$$

- Symmetric loss is proved to be robust under class-dependent label noises.

Robust loss function

Proof. Given symmetric loss L , for any f , $R_L(f) = \mathbb{E}_{(x,y) \sim p(x,y)} [L(f(x), y)]$.

For symmetric label noise ($T_{i,j} = \begin{cases} 1 - \eta, & i = j \\ \frac{\eta}{c-1}, & i \neq j \end{cases}$), for any f

$$\begin{aligned} R_L^\eta(f) &= \mathbb{E}_{(x,\bar{y}) \sim p(x,\bar{y})} [L(f(x), \bar{y})] \\ &= \mathbb{E}_x \mathbb{E}_{y|x} \mathbb{E}_{\bar{y}|(x,y)} [L(f(x), \bar{y})] \\ &= \mathbb{E}_x \mathbb{E}_{y|x} \left[(1 - \eta)L(f(x), y) + \frac{\eta}{c-1} \sum_{i \neq y} L(f(x), i) \right] \\ &= \mathbb{E}_x \mathbb{E}_{y|x} \left[(1 - \eta)L(f(x), y) + \frac{\eta}{c-1} (K - L(f(x), y)) \right] \\ &= \frac{K\eta}{c-1} + \left(1 - \frac{\eta c}{c-1}\right) R_L(f) \end{aligned}$$

Thus, if $\eta < \frac{c-1}{c}$, the minimizer of R_L is also a minimizer of R_L^η .

Robust loss function

- Mean Absolute Error (MAE)
 - For DNNs with a softmax output layer, MAE can be computed as:

$$L_{MAE}(f(X), Y) = \sum_{i=1}^N \|e_{y_i} - f(x_i)\|_1 = \sum_{i=1}^N 2 - 2f_{y_i}(x_i)$$

where $e_{y_i} \in \{0,1\}^c$ refers to the one-hot encoding of y_i , and $f_{y_i}(x)$ denotes the y_i th element of $f(x)$

- ! MAE loss is symmetric but difficult to converge

Robust loss function

- Generalized Cross Entropy (GCE)

$$L_{\text{GCE}}(f(X), Y) = \sum_{i=1}^N \frac{1 - f_{y_i}^q(x_i)}{q}, q \in (0, 1]$$

- if $q \rightarrow 0$, GCE is equivalent to Cross Entropy; if $q \rightarrow 1$, GCE is equivalent to MAE.
- Robustness analysis. When using softmax output layer, L_{GCE} is bounded as

$$\frac{c - c^{1-q}}{q} \leq \sum_{j=1}^c \frac{1 - f_{y_j}^q(x)}{q} \leq \frac{c - 1}{q}$$

$$\frac{\eta(1 - c^{1-q})}{q(c - 1 - \eta c)} \leq R_{L_{\text{GCE}}}(f^*) - R_{L_{\text{GCE}}}^\eta(\hat{f}) \leq 0, \quad \text{if } \eta < \frac{c - 1}{c}$$

Robust loss function

- Generalized Cross Entropy (GCE)

$$L_{\text{GCE}}(f(X), Y) = \sum_{i=1}^N \frac{1 - f_{y_i}^q(x_i)}{q}, q \in (0,1]$$

- Gradient analysis:
 - GCE behaves like a weighted MAE

$$\frac{\partial L(f(x_i; \theta), y_i)}{\partial \theta} = \begin{cases} -\nabla_{\theta} f_{y_i}(x_i; \theta), & \text{MAE} \\ -\frac{1}{f_{y_i}(x_i; \theta)} \nabla_{\theta} f_{y_i}(x_i; \theta), & \text{CE} \\ -f_{y_i}(x_i; \theta)^{q-1} \nabla_{\theta} f_{y_i}(x_i; \theta), & \text{GCE} \end{cases}$$

Robust loss function

- Symmetric Cross Entropy (SCE)

- KL divergence: $KL(q||p) = H(q, p) - H(q)$

- $CE = H(q, p)$, Reverse Cross Entropy $RCE = H(p, q)$

$$L_{RCE}(f(X), Y) = - \sum_{i=1}^N f^{\top}(x_i) \log e_i$$

- Define $\log 0 = A$, RCE term is symmetric. if $A=-2$, RCE is exact MAE:

$$L_{RCE}(f(X), Y) = - \sum_{i=1}^N A \left(1 - f_{y_i}(x_i) \right)$$

- $SCE = \alpha CE + \beta RCE$, ensure both converge and robustness

Loss adjustment

➤ Correcting loss with the estimated noise transition matrix \hat{T}

- Backward correction [1] : is unbiased if $\hat{T} = T$

$$L_{Backward} = \hat{T}^{-1} \langle l(f(x), 1), l(f(x), 2), \dots, l(f(x), c) \rangle^{\top}$$

- Forward correction [1]:

$$L_{Forward} = L(\hat{T}^{\top} f(X)^{\top}, \bar{Y})$$

Robust regularization

➤ Reduce model complexity may prevent overfitting

- Widely used regularization techniques: data augmentation, weight decay, dropout, label smoothing ...
- More advanced techniques:
 - Restrict the distance to initialization of parameters [1]

$$L_{\lambda}^{\text{RDI}}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - \tilde{y}_i)^2 + \frac{\lambda^2}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}(0)\|^2$$

- Reduce trainable parameters [2]

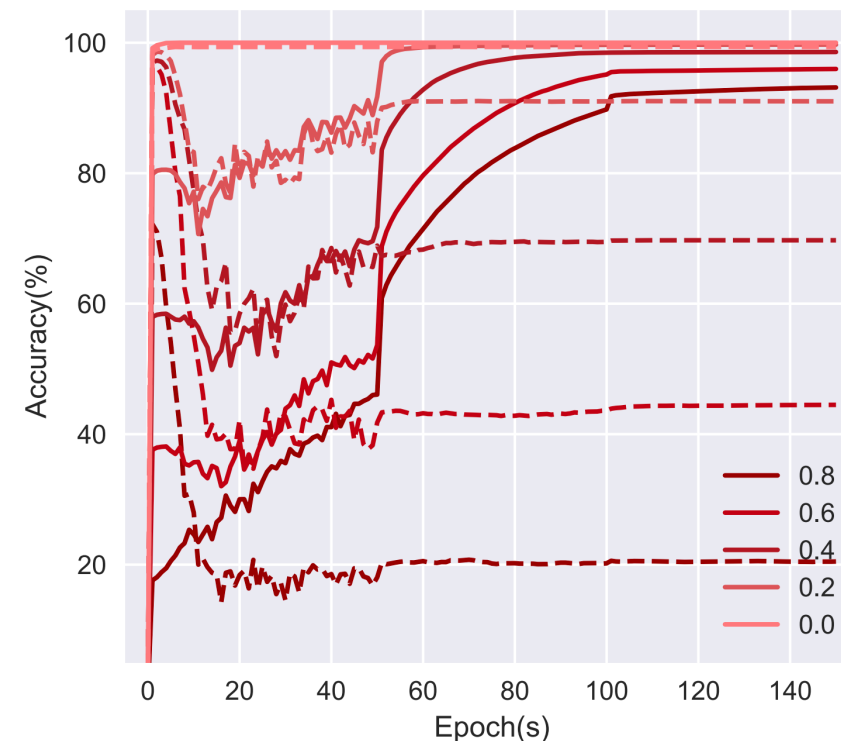
[1] Hu W, Li Z, Yu D. Simple and effective regularization methods for training on noisily labeled data with generalization guarantee. ICLR 2020.

[2] Xia X, Liu T, Han B, et al. Robust early-learning: Hindering the memorization of noisy labels. ICLR 2021.

Sample selection

➤ Memorization effect

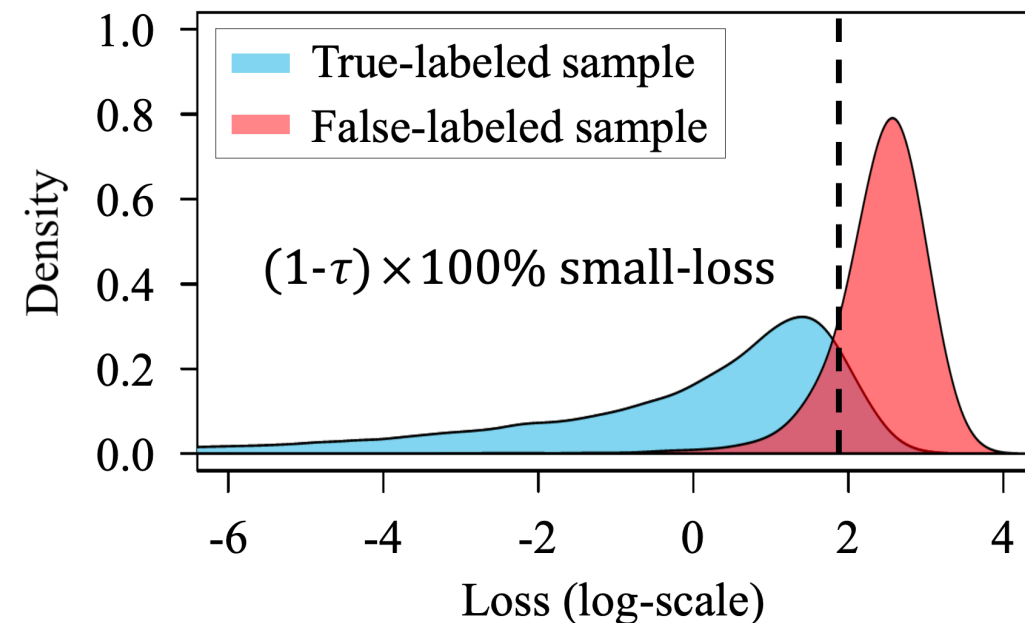
- Deep networks tend to fit easy (clean) patterns first, and gradually over-fit hard (noisy) patterns
- Memorization effect could be used to identify clean samples



Training and test accuracy curves.
Solid lines: training accuracy;
dashed lines: test accuracy.

Sample selection

- Select small-loss training examples as true-labeled examples.
 - Only train with small-loss examples [1];
 - Assigned more weights to small-loss examples [2].



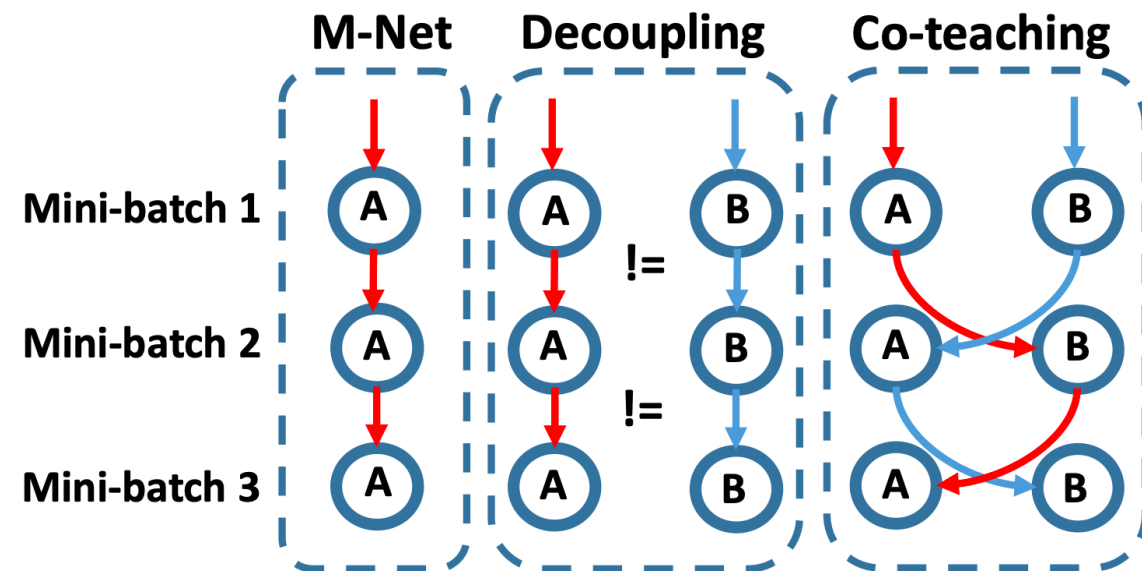
Loss distribution of training examples at the training accuracy of 50% on noisy CIFAR-100.

[1] Co-teaching: Robust training of deep neural networks with extremely noisy labels. NeurIPS 2018.

[2] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in ICML, 2018.

Sample selection

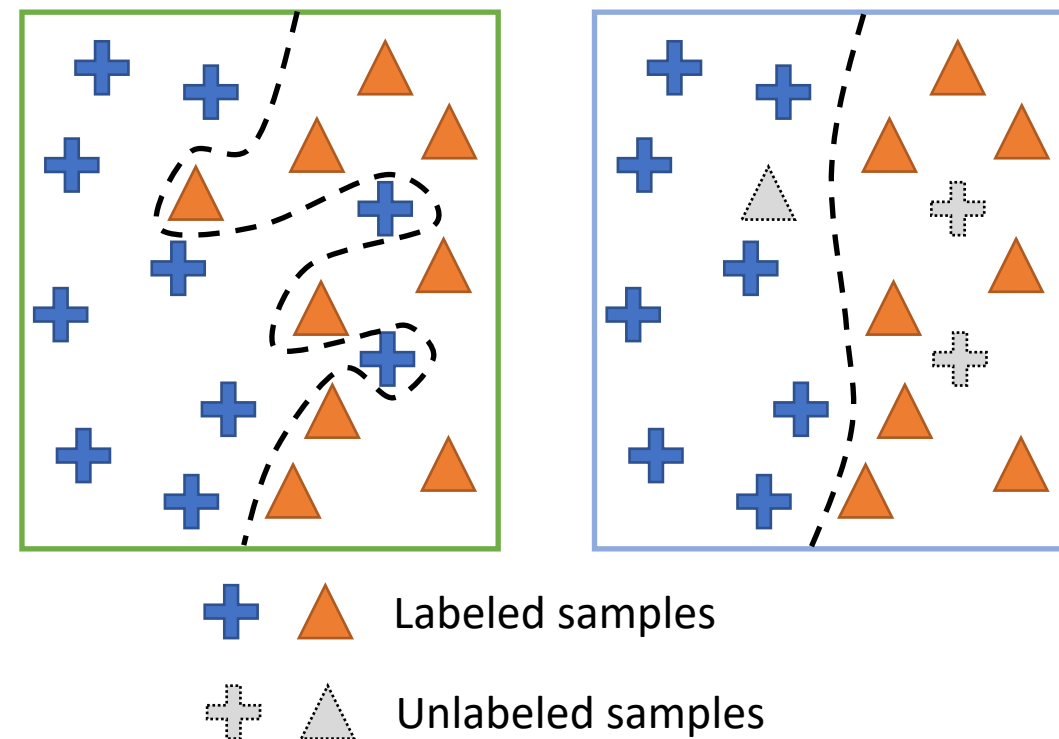
- Co-teaching
 - Train two networks simultaneously
 - In each mini-batch data, each network samples its small-loss instances, and teaches such useful instances to its peer network.
 - Two networks can attenuate each others' error.



Hybrid approaches

➤ Combine with semi-supervised learning methods.

- Relabeling [1]
- Consistency regularization
 - Dividemix [2]
 - ELR [3]



[1] Zhou T, Wang S, Bilmes J. Robust curriculum learning: from clean label detection to noisy label self-correction. ICLR 2020.

[2] J. Li, R. Socher, and S. C. Hoi, "DivideMix: Learning with noisy labels as semi-supervised learning," in Proc. ICLR, 2020

[3] Liu S, Niles-Weed J, Razavian N, et al. Early-learning regularization prevents memorization of noisy labels. NeurIPS 2020

Hybrid approaches

- Consistency regularization
 - Assumption: If two points x_1, x_2 reside in a high-density region are close, then so should be their corresponding outputs y_1, y_2
 - If a realistic perturbation was applied to the unlabeled data points, the prediction should not change significantly.

$$L_{reg} = \mathbb{E}_{x \in D_u} [d(f(x; \theta), f(\hat{x}; \theta))]$$

where $d(\cdot)$ is a distance measure

Hybrid approaches

- Dividemix:

- Use Gaussian Mixture Model to select clean samples

- Apply M rounds of data augmentation

```
for m = 1 to M do
    |  $\hat{x}_{b,m} = \text{Augment}(x_b)$ 
    |  $\hat{u}_{b,m} = \text{Augment}(u_b)$ 
end
```

- Label refinement: guided by model predictions of labeled samples x_b

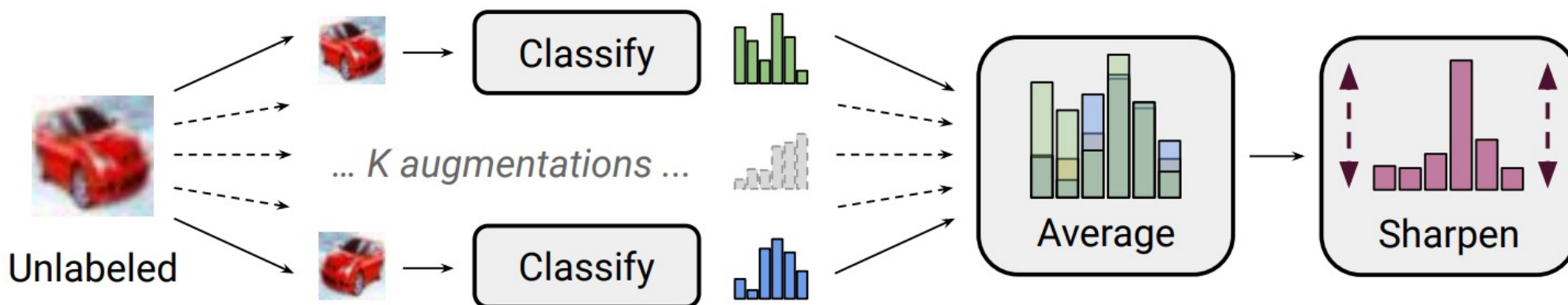
$$p_b = \frac{1}{M} \sum_m \text{P}_{\text{model}}(\hat{x}_{b,m}; \theta^{(k)})$$
$$\bar{y}_b = w_b y_b + (1 - w_b) p_b$$

// refine ground-truth label gui

$$\hat{y}_b = \text{Sharpen}(\bar{y}_b, T)$$

Hybrid approaches

- Dividemix:
 - Co-guessing: average the predictions from both networks of unlabeled sample u_b



$$\mathcal{L}_x = -\frac{1}{|\mathcal{X}'|} \sum_{x, p \in \mathcal{X}'} \sum_c p_c \log(p_{\text{model}}^c(x; \theta)), \quad \mathcal{L}_u = \frac{1}{|\mathcal{U}'|} \sum_{x, p \in \mathcal{U}'} \|p - p_{\text{model}}(x; \theta)\|_2^2.$$

Hybrid approaches

- ELR:

- $L_{ELR} = L_{CE} + \lambda L_{reg}$

$$L_{reg}(f(X), Y) = \sum_{i=1}^N \log(1 - \langle f(x_i), t_i \rangle)$$

where t_i is moving average of model historical prediction.

- Use early-learning stage of model prediction to hinder memorization.
- ELR+ tricks: train two models; use weight averaging; mixup data

augmentation

Summary of methods

- Loss function:
 - Utilize symmetric loss functions
 - Estimate noise transition matrix
- Regularization:
 - Reduce effective model capacity
- Sample selection:
 - Design criteria to identify noisy

data.

- Hybrid approaches:
 - Relabeling
 - Consistency regularization

Instance-dependent label noise

- The label corruption probability is assumed to be dependent on both the data features and class labels. The corruption probability is

$$T_{i,j}(x) := p(\bar{y} = j | y = i, x)$$

- Difficulties:
 - How to estimate the noise transition matrix? The size of T is very huge.
 - How to identify noisy samples? The loss distribution of true-labeled and false-labeled samples may heavily overlap.

Thanks for listening!