



Normalized Loss Functions for Deep Learning with Noisy Labels

Xingjun Ma^{*1} Hanxun Huang^{*1} Yisen Wang² Simone Romano Sarah Erfani¹ James Bailey¹

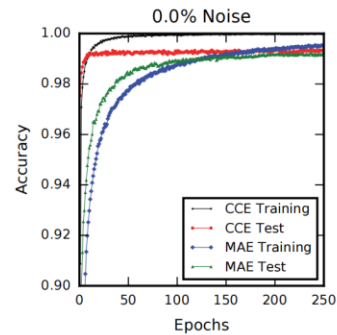
¹The University of Melbourne, Australia ²Shanghai Jiao Tong University, China.

ICML2020

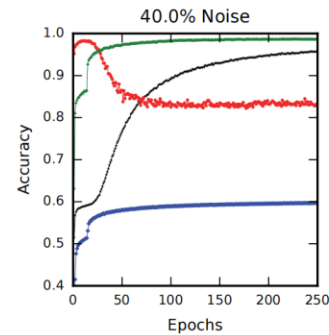
Robust loss functions under label noise for deep neural networks

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

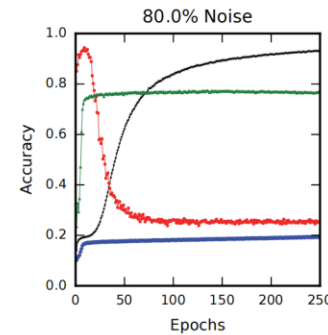
$$\sum_{i=1}^k L(f(\mathbf{x}), i) = C, \forall \mathbf{x} \in \mathcal{X}, \forall f$$



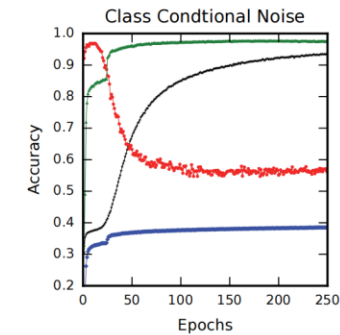
(a)



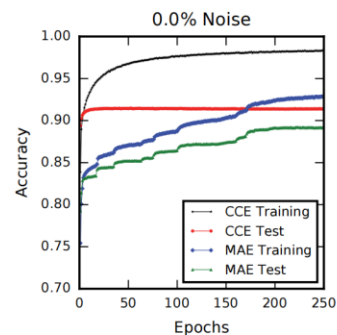
(b)



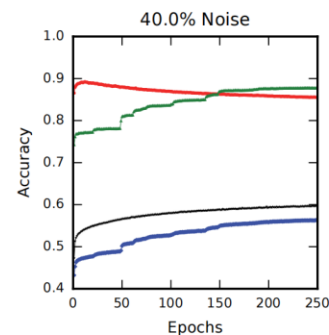
(c)



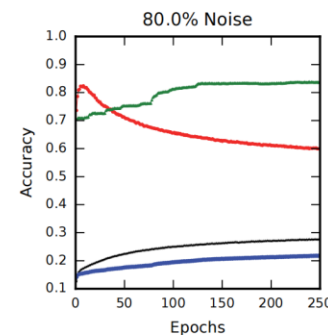
(d)



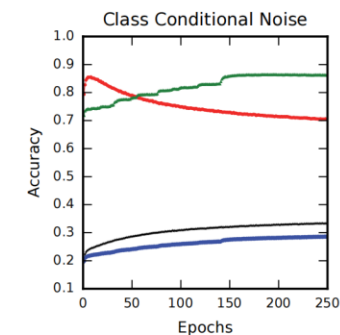
(e)



(f)



(g)



(h)

Proof 1 under symmetric noise

Theorem 1 In a multi-class classification problem, let loss function L satisfy Eq 2. Then L is noise tolerant under symmetric or uniform label noise if $\eta < \frac{k-1}{k}$.

Proof 1 Recall that for any f ,

risk \leftarrow $R_L(f) = \mathbb{E}_{\mathbf{x}, y_{\mathbf{x}}} L(f(\mathbf{x}), y_{\mathbf{x}})$ (3)

For uniform noise, we have, for any f ,¹

$$\begin{aligned} R_L^\eta(f) &= \mathbb{E}_{\mathbf{x}, \hat{y}_{\mathbf{x}}} L(f(\mathbf{x}), \hat{y}_{\mathbf{x}}) \\ &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{y_{\mathbf{x}}|\mathbf{x}} \mathbb{E}_{\hat{y}_{\mathbf{x}}|\mathbf{x}, y_{\mathbf{x}}} L(f(\mathbf{x}), \hat{y}_{\mathbf{x}}) \\ &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{y_{\mathbf{x}}|\mathbf{x}} \left[(1 - \eta) L(f(\mathbf{x}), y_{\mathbf{x}}) + \frac{\eta}{k-1} \sum_{i \neq y_{\mathbf{x}}} L(f(\mathbf{x}), i) \right] \\ &= (1 - \eta) R_L(f) + \frac{\eta}{k-1} (C - R_L(f)) \\ &= \frac{C\eta}{k-1} + \left(1 - \frac{\eta k}{k-1} \right) R_L(f). \end{aligned}$$

Thus, for any f ,

$$R_L^\eta(f^*) - R_L^\eta(f) = \left(1 - \frac{\eta k}{k-1} \right) (R_L(f^*) - R_L(f)) \leq 0$$

because $\eta < \frac{k-1}{k}$ and f^* is a minimizer of R_L . This proves f^* is also minimizer of risk under uniform noise.

Proof 2 under class-conditional noise

Theorem 3 Suppose L satisfies Eq 2 and $0 \leq L(f(\mathbf{x}), i) \leq C/(k-1), \forall i \in [k]$. If $R_L(f^*) = 0$, then, L is noise tolerant under class conditional noise when $\bar{\eta}_{ij} < (1 - \eta_i), \forall j \neq i, \forall i, j \in [k]$.

Proof 3 For class-conditional noise, we have

$$\begin{aligned}
 R_L^\eta(f) &= \mathbb{E}_{\mathcal{D}}(1 - \eta_{y_{\mathbf{x}}})L(f(\mathbf{x}), y_{\mathbf{x}}) + \mathbb{E}_{\mathcal{D}} \sum_{i \neq y_{\mathbf{x}}} \bar{\eta}_{y_{\mathbf{x}}i} L(f(\mathbf{x}), i) \\
 &= \mathbb{E}_{\mathcal{D}}(1 - \eta_{y_{\mathbf{x}}})(C - \sum_{i \neq y_{\mathbf{x}}} L(f(\mathbf{x}), i)) \\
 &\quad + \mathbb{E}_{\mathcal{D}} \sum_{i \neq y_{\mathbf{x}}} \bar{\eta}_{y_{\mathbf{x}}i} L(f(\mathbf{x}), i) \\
 &= C\mathbb{E}_{\mathcal{D}}(1 - \eta_{y_{\mathbf{x}}}) - \mathbb{E}_{\mathcal{D}} \sum_{i \neq y_{\mathbf{x}}} (1 - \eta_{y_{\mathbf{x}}} - \bar{\eta}_{y_{\mathbf{x}}i}) L(f(\mathbf{x}), i)
 \end{aligned} \tag{6}$$

Since f_η^* is the minimizer of R_L^η , we have $R_L^\eta(f_\eta^*) - R_L^\eta(f^*) \leq 0$ and hence from Eq.(6) we have

$$\mathbb{E}_{\mathcal{D}} \sum_{i \neq y_{\mathbf{x}}} \underbrace{(1 - \eta_{y_{\mathbf{x}}} - \bar{\eta}_{y_{\mathbf{x}}i})}_{> 0} \underbrace{(L(f^*(\mathbf{x}), i) - L(f_\eta^*(\mathbf{x}), i))}_{\geq 0} \leq 0 \tag{7}$$

$$R_L(f^*) = L(f^*(\mathbf{x}), y_{\mathbf{x}}) = 0 \implies L(f^*(\mathbf{x}), i) = C/(k-1), \quad i \neq y_{\mathbf{x}}$$

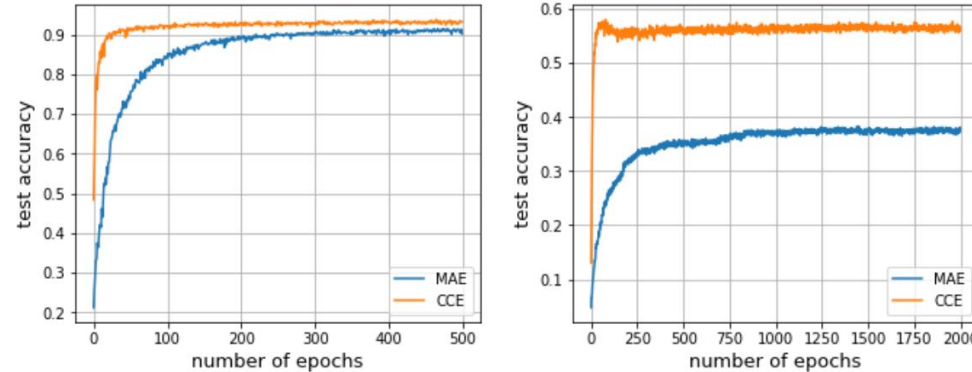
$$R_L(f_\eta^*) = L(f_\eta^*(\mathbf{x}), y_{\mathbf{x}}) \geq 0 \implies L(f_\eta^*(\mathbf{x}), i) \leq C/(k-1), \forall i$$

$$L(f_\eta^*(\mathbf{x}), i) = C/(k-1) \implies L(f_\eta^*(\mathbf{x}), y_{\mathbf{x}}) = 0$$

$$\Downarrow \\ f_\eta^*(\mathbf{x}) = f^*(\mathbf{x})$$

GCE-Generalized Cross Entropy Loss for Training Deep Neural Networks with Noisy Labels

MAE is robust but underperform Cross Entropy



CIFAR-10(clean)

CIFAR-100(clean)

$$\sum_{i=1}^n \frac{\partial \mathcal{L}(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)}{\partial \boldsymbol{\theta}} = \begin{cases} \sum_{i=1}^n \frac{1}{f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta}) & \text{for CCE} \\ \sum_{i=1}^n -\nabla_{\boldsymbol{\theta}} f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta}) & \text{for MAE/unhinged loss.} \end{cases}$$

$$\mathcal{L}_q(f(\mathbf{x}), \mathbf{e}_j) = \frac{(1 - f_j(\mathbf{x})^q)}{q} \quad \text{where } q \in (0, 1]$$

$\xleftarrow{q \rightarrow 0}$

$$\lim_{q \rightarrow 0} \mathcal{L}_q(f(\mathbf{x}), \mathbf{e}_j) = -\ln f_j(\mathbf{x})$$

$\xrightarrow{q=1}$

$$\mathcal{L}_q(f(\mathbf{x}), \mathbf{e}_j) = 1 - f_j(\mathbf{x})$$
$$\mathcal{L}_{MAE}(f(\mathbf{x}), \mathbf{e}_j) = \|\mathbf{e}_j - f(\mathbf{x})\|_1 = 2 - 2f_j(\mathbf{x})$$

SL/SCE-Symmetric Cross Entropy for Robust Learning with Noisy Labels

$$\ell_{ce} = - \sum_{k=1}^K q(k | \mathbf{x}) \log p(k | \mathbf{x})$$

Clean targets (pointing to $p(k | \mathbf{x})$)
dist of true labels (pointing to $q(k | \mathbf{x})$)

If given labels are noisy ? (q is not clean)

$$\ell_{rce} = - \sum_{k=1}^K p(k | \mathbf{x}) \log q(k | \mathbf{x})$$

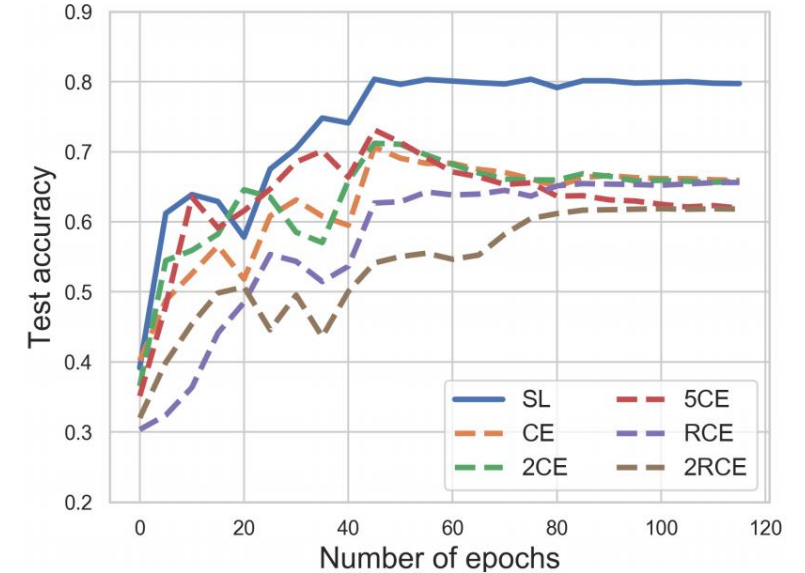
$$\log(q(k \neq y | \mathbf{x})) = A \text{ when } q(k \neq y | \mathbf{x}) = 0$$

$$\ell_{rce} = - \sum_{k=1}^K p(k | \mathbf{x}) \log q(k | \mathbf{x})$$

$$= -p(y | \mathbf{x}) \log 1 - \sum_{k \neq y} p(k | \mathbf{x}) A = -A \sum_{k \neq y} p(k | \mathbf{x})$$

$$= -A(1 - p(y | \mathbf{x}))$$

$$\sum_{k=1}^K \ell_{rce}(f(x), j) = A(1 - K)$$



CIFAR-10 with 60% symmetric label noise.

$$\ell = \alpha \ell_{ce} + \beta \ell_{rce}$$

Motivations

1. If a loss function satisfies **symmetric condition**, it is noise-robust under symmetric and class-conditional condition
2. CE is non-robust but better for convergence, robust loss functions (at least MAE and RCE) suffer from underfitting
3. Combine CE and a robust loss function can make them benefit from each other

$$\mathcal{L}_{\text{norm}} = \frac{\mathcal{L}(f(\mathbf{x}), y)}{\sum_{j=1}^K \mathcal{L}(f(\mathbf{x}), j)}.$$

Is normalized loss term still robust?

Robust loss functions suffer from underfitting

$$NCE = \frac{-\sum_{k=1}^K \mathbf{q}(k|\mathbf{x}) \log \mathbf{p}(k|\mathbf{x})}{-\sum_{j=1}^K \sum_{k=1}^K \mathbf{q}(y=j|\mathbf{x}) \log \mathbf{p}(k|\mathbf{x})}$$

$$= \log_{\prod_k^K \mathbf{p}(k|\mathbf{x})} \mathbf{p}(y|\mathbf{x}),$$

$$NMAE = \frac{\sum_{k=1}^K |\mathbf{p}(k|\mathbf{x}) - \mathbf{q}(k|\mathbf{x})|}{\sum_{j=1}^K \sum_{k=1}^K |\mathbf{p}(k|\mathbf{x}) - \mathbf{q}(y=j|\mathbf{x})|}$$

$$= \frac{1}{K-1} (1 - \mathbf{p}(y|\mathbf{x})) = \frac{1}{2(K-1)} \cdot MAE.$$

$$NRCE = \frac{-\sum_{k=1}^K \mathbf{p}(k|\mathbf{x}) \log \mathbf{q}(k|\mathbf{x})}{-\sum_{j=1}^K \sum_{k=1}^K \mathbf{p}(k|\mathbf{x}) \log \mathbf{q}(y=j|\mathbf{x})}$$

$$= \frac{1}{K-1} (1 - \mathbf{p}(y|\mathbf{x})) = \frac{1}{A(K-1)} \cdot RCE$$

$$NFL = \frac{-\sum_{k=1}^K \mathbf{q}(k|\mathbf{x}) (1 - \mathbf{p}(k|\mathbf{x}))^\gamma \log \mathbf{p}(k|\mathbf{x})}{-\sum_{j=1}^K \sum_{k=1}^K \mathbf{q}(y=j|\mathbf{x}) (1 - \mathbf{p}(k|\mathbf{x}))^\gamma \log \mathbf{p}(k|\mathbf{x})}$$

$$= \log_{\prod_k^K (1 - \mathbf{p}(k|\mathbf{x}))^\gamma \mathbf{p}(k|\mathbf{x})} (1 - \mathbf{p}(y|\mathbf{x}))^\gamma \mathbf{p}(y|\mathbf{x}).$$

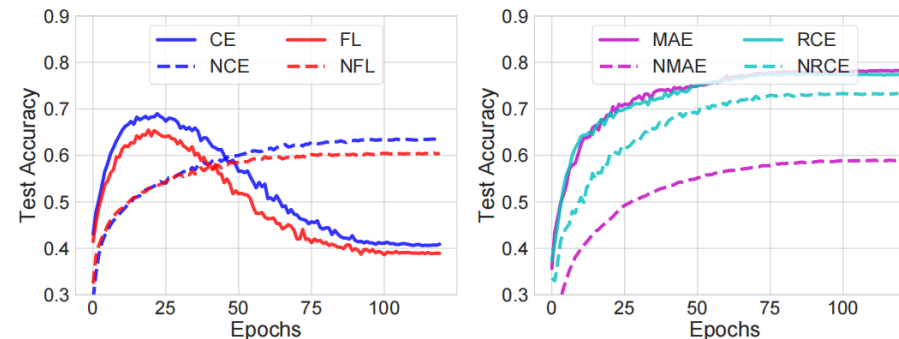


Figure 2. Test accuracies of unnormalized versus normalized loss functions on CIFAR-10 under 0.6 symmetric noise.

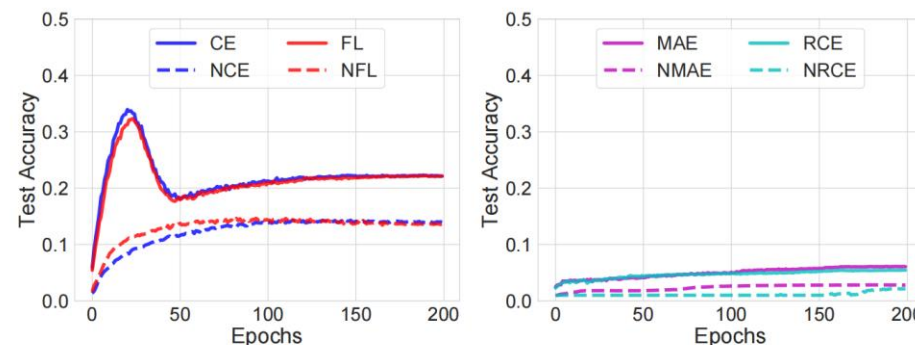


Figure 1. Test accuracies of unnormalized versus normalized loss functions on CIFAR-100 under 0.6 symmetric noise.

Active Loss and Passive Loss

Definition 1. (Active loss function) \mathcal{L}_{Active} is an active loss function if $\forall (\mathbf{x}, y) \in \mathcal{D} \forall k \neq y \ell(f(\mathbf{x}), k) = 0$.

Definition 2. (Passive loss function) $\mathcal{L}_{Passive}$ is a passive loss function if $\forall (\mathbf{x}, y) \in \mathcal{D} \exists k \neq y \ell(f(\mathbf{x}), k) \neq 0$.

Table 1. Examples of active and passive loss functions.

Loss Type	Active	Passive
Examples	CE, NCE, FL, NFL	MAE, NMAE, RCE, NRCE

$$\mathcal{L}_{APL} = \alpha \cdot \mathcal{L}_{Active} + \beta \cdot \mathcal{L}_{Passive}$$

A robust loss function with two complementary loss

Q1: complementary is necessary?

Q2: robust L_{Active} is necessary?

Only explicitly optimize the target class

Besides explicitly optimize other classes

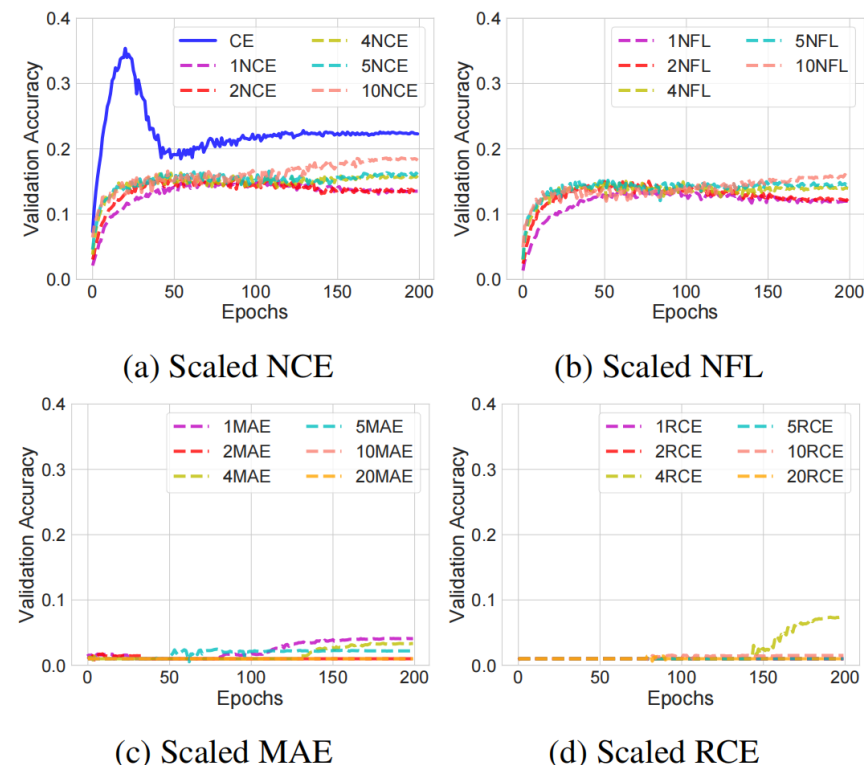
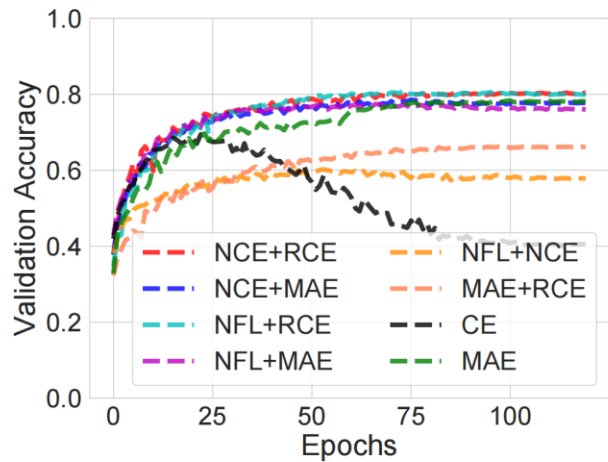
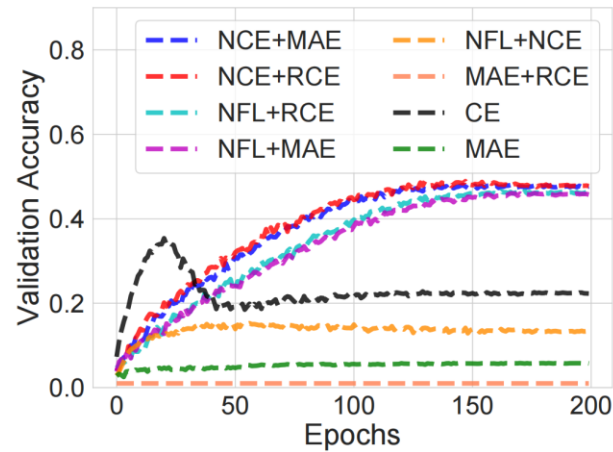


Figure 3. Test accuracies of scaled loss functions on CIFAR-100 with 0.6 symmetric noise.

Experiments



(a) CIFAR-10



(b) CIFAR-100

Under 60% symmetric noise

Active + Passive:
NCE + RCE, NCE + MAE, NFL + RCE, NFL + MAE

Active + Active :
NFL + NCE

Passive + Passive :
MAE + RCE

Baselines:
CE, MAE

Experiments-Symmetric Noise

Table 2. Test accuracies (%) of different methods on benchmark datasets with clean or symmetric label noise ($\eta \in [0.2, 0.8]$). The results (mean \pm std) are reported over 3 random runs and the top 2 best results are **boldfaced**.

Datasets	Methods	Clean ($\eta=0.0$)	Symmetric Noise Rate (η)			
			0.2	0.4	0.6	0.8
MNIST	CE	99.25 \pm 0.08	97.42 \pm 0.06	94.21 \pm 0.54	86.00 \pm 1.48	47.08 \pm 1.15
	FL	99.30 \pm 0.02	97.45 \pm 0.19	94.71 \pm 0.25	85.76 \pm 1.85	49.77 \pm 2.26
	GCE	99.27 \pm 0.01	99.18 \pm 0.06	98.72 \pm 0.05	97.43 \pm 0.23	12.77 \pm 2.00
	NLNL	99.27 \pm 0.02	97.49 \pm 0.30	96.64 \pm 0.52	97.22 \pm 0.06	10.32 \pm 0.73
	SCE	99.24 \pm 0.08	99.15 \pm 0.04	98.78 \pm 0.09	97.45 \pm 0.29	73.70 \pm 0.84
	NFL+MAE	99.39 \pm 0.04	99.12 \pm 0.06	98.74 \pm 0.14	96.91 \pm 0.09	74.98 \pm 1.99
	NFL+RCE	99.38 \pm 0.02	99.19 \pm 0.06	98.79 \pm 0.10	97.46 \pm 0.03	74.59 \pm 2.23
	NCE+MAE	99.37 \pm 0.02	99.14 \pm 0.05	98.78 \pm 0.00	96.76 \pm 0.34	74.66 \pm 1.11
	NCE+RCE	99.37 \pm 0.02	99.20 \pm 0.04	98.79 \pm 0.12	97.48 \pm 0.13	75.18 \pm 1.19
CIFAR-10	CE	90.36 \pm 0.03	75.90 \pm 0.28	60.28 \pm 0.27	40.90 \pm 0.35	19.65 \pm 0.46
	FL	89.63 \pm 0.25	74.59 \pm 0.49	57.55 \pm 0.39	38.91 \pm 0.62	19.43 \pm 0.27
	GCE	89.38 \pm 0.23	87.27 \pm 0.21	83.33 \pm 0.39	72.00 \pm 0.37	29.08 \pm 0.80
	NLNL	91.93 \pm 0.20	83.98 \pm 0.18	76.58 \pm 0.44	72.85 \pm 0.39	51.41 \pm 0.85
	SCE	91.30 \pm 0.22	88.05 \pm 0.26	82.06 \pm 0.24	66.08 \pm 0.25	30.69 \pm 0.63
	NFL+MAE	89.25 \pm 0.19	87.33 \pm 0.14	83.81 \pm 0.06	76.36 \pm 0.31	45.23 \pm 0.52
	NFL+RCE	90.91 \pm 0.02	89.14 \pm 0.13	86.05 \pm 0.12	79.78 \pm 0.13	55.06 \pm 1.08
	NCE+MAE	88.83 \pm 0.34	87.12 \pm 0.21	84.19 \pm 0.43	77.61 \pm 0.05	49.62 \pm 0.72
	NCE+RCE	90.76 \pm 0.22	89.22 \pm 0.27	86.02 \pm 0.09	79.78 \pm 0.50	52.71 \pm 1.90
CIFAR-100	CE	70.89 \pm 0.22	56.99 \pm 0.41	41.40 \pm 0.36	22.15 \pm 0.40	7.58 \pm 0.44
	FL	70.61 \pm 0.44	56.10 \pm 0.48	40.77 \pm 0.62	22.14 \pm 1.00	7.21 \pm 0.25
	GCE	69.00 \pm 0.56	65.24 \pm 0.56	58.94 \pm 0.50	45.18 \pm 0.93	16.18 \pm 0.46
	NLNL	68.72 \pm 0.60	46.99 \pm 0.91	30.29 \pm 1.64	16.60 \pm 0.90	11.01 \pm 2.48
	SCE	70.38 \pm 0.45	55.39 \pm 0.18	39.99 \pm 0.59	22.35 \pm 0.65	7.57 \pm 0.28
	NFL+MAE	67.98 \pm 0.52	63.58 \pm 0.09	58.18 \pm 0.08	46.10 \pm 0.50	24.78 \pm 0.82
	NFL+RCE	68.23 \pm 0.62	64.52 \pm 0.35	58.20 \pm 0.31	46.30 \pm 0.45	25.16 \pm 0.55
	NCE+MAE	68.75 \pm 0.54	65.25 \pm 0.62	59.22 \pm 0.36	48.06 \pm 0.34	25.50 \pm 0.76
	NCE+RCE	69.02 \pm 0.11	65.31 \pm 0.07	59.48 \pm 0.56	47.12 \pm 0.62	25.80 \pm 1.12

Experiments-Class-conditional Noise

Table 3. Test accuracy (%) of different methods on benchmark datasets with clean or asymmetric label noise ($\eta \in [0.1, 0.4]$). The results (mean \pm std) are reported over 3 random runs and the top 2 best results are **boldfaced**.

Datasets	Methods	Asymmetric Noise Rate (η)			
		0.1	0.2	0.3	0.4
MNIST	CE	98.53 \pm 0.11	96.75 \pm 0.31	92.98 \pm 1.41	85.74 \pm 2.70
	FL	98.97 \pm 0.10	98.35 \pm 0.17	96.57 \pm 0.36	91.18 \pm 2.02
	GCE	99.25 \pm 0.03	99.11 \pm 0.04	96.99 \pm 0.53	88.56 \pm 2.40
	NLNL	98.38 \pm 0.17	95.98 \pm 0.58	91.52 \pm 1.14	86.36 \pm 0.40
	SCE	99.15 \pm 0.07	99.05 \pm 0.05	97.96 \pm 0.40	91.89 \pm 3.32
	NFL+MAE	99.31 \pm 0.05	99.09 \pm 0.12	97.88 \pm 0.16	93.52 \pm 0.19
	NFL+RCE	99.33 \pm 0.06	99.13 \pm 0.01	97.99 \pm 0.05	93.59 \pm 0.82
	NCE+MAE	99.26 \pm 0.02	99.21 \pm 0.04	98.99 \pm 0.03	93.40 \pm 1.28
	NCE+RCE	99.34 \pm 0.06	99.17 \pm 0.02	97.94 \pm 0.21	93.12 \pm 1.17
CIFAR-10	CE	87.38 \pm 0.16	83.62 \pm 0.15	79.38 \pm 0.28	75.00 \pm 0.50
	FL	86.35 \pm 0.30	82.97 \pm 0.14	79.48 \pm 0.21	74.60 \pm 0.15
	GCE	88.42 \pm 0.07	86.07 \pm 0.31	80.78 \pm 0.21	74.98 \pm 0.32
	NLNL	88.54 \pm 0.25	84.74 \pm 0.08	81.26 \pm 0.43	76.97 \pm 0.52
	SCE	88.13 \pm 0.21	83.92 \pm 0.07	79.70 \pm 0.27	78.20 \pm 0.03
	NFL+MAE	88.46 \pm 0.20	86.81 \pm 0.32	83.91 \pm 0.34	77.16 \pm 0.10
	NFL+RCE	90.20 \pm 0.15	88.73 \pm 0.29	85.74 \pm 0.22	79.27 \pm 0.43
	NCE+MAE	88.25 \pm 0.09	86.44 \pm 0.23	83.98 \pm 0.52	78.23 \pm 0.42
	NCE+RCE	89.95 \pm 0.20	88.56 \pm 0.17	85.58 \pm 0.44	79.59 \pm 0.40
CIFAR-100	CE	65.42 \pm 0.22	58.45 \pm 0.45	51.09 \pm 0.29	41.68 \pm 0.45
	FL	64.79 \pm 0.18	58.59 \pm 0.81	51.26 \pm 0.18	42.15 \pm 0.44
	GCE	61.98 \pm 0.81	59.99 \pm 0.83	53.99 \pm 0.29	41.49 \pm 0.79
	NLNL	59.55 \pm 1.22	50.19 \pm 0.56	42.81 \pm 1.13	35.10 \pm 0.20
	SCE	64.15 \pm 0.61	58.22 \pm 0.47	49.85 \pm 0.91	42.19 \pm 0.19
	NFL+MAE	66.06 \pm 0.23	63.10 \pm 0.22	56.19 \pm 0.61	43.51 \pm 0.42
	NFL+RCE	66.13 \pm 0.31	63.12 \pm 0.41	54.72 \pm 0.38	42.97 \pm 1.03
	NCE+MAE	65.71 \pm 0.34	62.38 \pm 0.60	58.02 \pm 0.48	47.22 \pm 0.30
	NCE+RCE	65.68 \pm 0.25	62.68 \pm 0.79	57.82 \pm 0.41	46.79 \pm 0.96

Experiments-NGCE & real-world Noise

Table 4. Test accuracy (%) of APL losses NGCE+MAE and NGCE+RCE on CIFAR-10 under both symmetric and asymmetric noise. The top-2 best results are in **bold**.

Methods	Symmetric noise		Asymmetric noise
	0.4	0.8	0.4
GCE	83.33 \pm 0.39	29.08 \pm 0.80	74.98 \pm 0.32
NGCE+MAE	84.14 \pm 0.15	50.55 \pm 1.08	76.55 \pm 0.48
NGCE+RCE	85.76 \pm 0.26	44.69 \pm 4.93	71.65 \pm 0.68

Table 5. Top-1 validation accuracies (%) on clean ILSVRC12 validation set of ResNet-50 models trained on WebVision using different loss functions, under the Mini setting in (Jiang et al., 2018). The top-2 best results are in **bold**.

Loss	CE	GCE	SCE	NCE+MAE	NCE+RCE
Acc	58.88	53.68	61.76	62.36	62.64



Learning from Noisy Labels with Complementary Loss Functions

Deng-Bao Wang,^{1,2} Yong Wen,³ Lujia Pan,^{4,3} Min-Ling Zhang^{1,2,5*}

¹School of Computer Science and Engineering, Southeast University, Nanjing 210096, China

²Key Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education, China

³Noah's Ark Lab, Huawei Technologies

⁴NSKEYLAB, Xi'an Jiaotong University

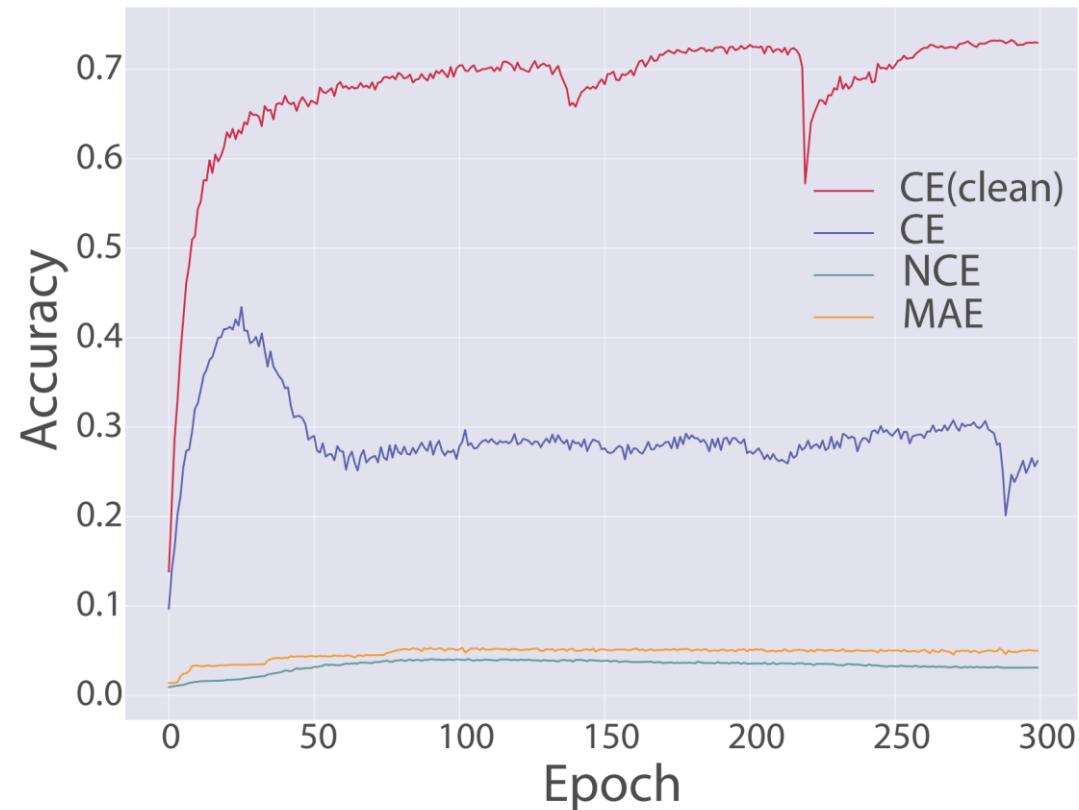
⁵Collaborative Innovation Center of Wireless Communications Technology, China

wangdb@seu.edu.cn, {wenyong4, panlujia}@huawei.com, zhangml@seu.edu.cn

AAAI 2021

Motivations

1. Cross entropy achieving good convergence but not robust (clean instances / easy instances)
2. Robust loss function suffers from underfitting (noisy instances / hard instances)



CIFAR-100 under 0.5 symmetric noise.

Learning with Complementary Losses

Algorithm 1: Learning with Complementary Losses.

Input: Training set \mathcal{D} , robust loss ℓ_{ROB} , T , T_{warm} , α , γ_1 , γ_2 , batch size B and Optimizer \mathcal{O} .

Output: Model parameters θ .

```

1 for  $t = 1; t \leq T$  do
2   for  $b = 1; b \leq \frac{|\mathcal{D}|}{B}$  do
3     Fetch mini-batch  $\mathcal{B}$  from  $\mathcal{D}$ ;
4     if  $t < T_{warm}$  then ▷ Warm up
5        $\mathcal{L} = \sum_{(\mathbf{x}, \bar{y}) \in \mathcal{B}} \ell_{CE}(f(\mathbf{x}), \bar{y})$ ;
6       Update  $\theta = \mathcal{O}(\mathcal{L}, \theta)$ ;
7     end
8     else ▷ Main train
9       Calculate the average prediction for each
10        sample using Eq. (7);
11        Obtain  $\tilde{\mathcal{B}}$  using Eq. (8);
12        Obtain  $\tilde{\mathcal{B}}'$  using Eq. (9);
13        Obtain  $\bar{\mathcal{B}}$  using Eq. (10);
14        Calculate  $\mathcal{L}_{CE} = \sum_{(\mathbf{x}, \tilde{y}) \in \tilde{\mathcal{B}}} \ell_{CE}(f(\mathbf{x}), \tilde{y})$ ;
15        Calculate  $\mathcal{L}_{ROB} = \sum_{(\mathbf{x}, \bar{y}) \in \bar{\mathcal{B}}} \ell_{ROB}(f(\mathbf{x}), \bar{y})$ ;
16         $\mathcal{L}_{CL} = \mathcal{L}_{CE} + \alpha \mathcal{L}_{ROB}$ ;
17        Update  $\theta = \mathcal{O}(\mathcal{L}_{CL}, \theta)$ 
18     end
19   Preserve the model outputs of current epoch;
20 end
21 return  $\theta$ .

```

$$p(\mathbf{x}) = \frac{1}{w} \sum_{i=k-w}^{k-1} f^{(i)}(\mathbf{x}) \quad \text{Model ensembling}$$

$$\tilde{\mathcal{B}} = \{(\mathbf{x}, \tilde{y}) \mid (\mathbf{x}, \bar{y}) \in \mathcal{B}\} \quad \text{where } \tilde{y} = \arg \max_{j \in [c]} p_j \quad \text{Pseudo labels}$$

$$\tilde{\mathcal{B}}' = \arg \min_{\mathcal{B}': |\mathcal{B}'| \geq \gamma_1 |\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \tilde{\mathcal{B}}'} H(p(\mathbf{x})) \quad \text{Easy instances}$$

$$\bar{\mathcal{B}} = \arg \max_{\mathcal{B}': |\mathcal{B}'| \geq \gamma_2 |\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \mathcal{B}'} E(p(\mathbf{x})) \quad \text{Hard instances}$$

Experiments(CIFAR-10)

Noise type	Uniform			Class-conditional		
Noise ratio	0.2	0.4	0.6	0.2	0.3	0.4
Standard CE	83.78±0.32	67.73±0.72	47.79±0.82	86.54±0.52	81.86±0.70	76.02±0.93
GCE (2018)	90.44±0.08	88.08±0.13	81.13±0.21	90.30±0.16	88.68±0.10	84.77±0.87
SCE (2019)	91.68±0.17	87.54±0.47	78.88±0.69	89.91±0.58	84.86±0.77	76.52±1.33
APL (2020)	87.79±0.48	79.13±0.75	66.51±1.38	90.14±0.34	83.70±1.08	76.02±1.20
Ours (CE+MAE)	93.49±0.14	92.04±0.20	89.37±0.14	94.10±0.17	93.01±0.20	91.52±0.21
Ours (CE+APL)	92.20±0.54	90.47±0.87	88.01±0.55	94.08±0.14	93.39±0.17	91.89±0.19
M-correction (2019)	93.69±0.32	93.18±0.10	90.59±0.33	—	—	—
DivideMix (2020)	95.23±0.22	93.62±0.15	92.84±0.19	93.20±0.28	92.38±0.24	91.15±0.69
Ours* (CE)	93.62±0.21	93.21±0.17	92.33±0.24	93.83±0.18	92.82±0.22	91.39±0.26
Ours* (CE+MAE)	95.37±0.09●	94.79±0.11●	93.59±0.19●	94.98±0.14●	94.33±0.24●	92.18±0.42●
Ours* (CE+APL)	94.70±0.13	93.80±0.14	92.68±0.27	94.34±0.14	93.38±0.28	91.10±0.40

Experiments(CIFAR-100)

Noise type	Uniform			Class-conditional		
Noise ratio	0.2	0.4	0.6	0.2	0.3	0.4
Standard CE	62.82±0.27	49.10±0.37	31.11±0.36	63.43±0.20	54.37±0.30	44.89±0.40
GCE (2018)	69.43±0.20	64.44±0.14	55.96±0.47	68.96±0.20	64.68±0.22	50.67±0.44
SCE (2019)	61.60±0.31	46.21±0.40	30.05±0.73	62.18±0.21	53.67±0.30	44.29±0.39
APL (2020)	70.74±0.33	61.92±0.52	48.64±0.83	63.11±0.45	52.91±0.46	42.48±0.68
Ours (CE+MAE)	71.63±0.26	68.61±0.35	62.52±0.22	69.99±0.23	68.32±0.20	65.72±0.28
Ours (CE+APL)	71.26±0.36	68.50±0.28	62.30±0.23	69.54±0.17	67.71±0.23	65.39±0.52
M-correction (2019)	68.95±0.53	65.43±0.30	59.43±0.36	—	—	—
DivideMix (2020)	74.80±0.28	72.92±0.20	69.38±0.26	74.90±0.41	72.14±0.45	50.79±0.62
Ours* (CE)	71.04±0.40	69.20±0.27	65.09±0.33	67.80±0.37	66.12±0.36	65.19±0.29
Ours* (CE+MAE)	77.61±0.22●	75.81±0.37●	72.17±0.30●	76.74±0.53●	75.32±0.39●	68.07±0.79
Ours* (CE+APL)	77.21±0.39	75.62±0.27	71.34±0.60	76.50±0.41	74.47±0.50	68.18±1.16●

Experiments

Noise ratio	Uniform (0.5)	Class-cond. (0.3)
Standard CE	23.16±0.53	43.67±0.22
SCE	23.89±0.32	38.92±0.44
APL	5.43±0.14	Fail
Ours (CE+MAE)	50.59±0.33	51.99±0.29
Ours (CE+APL)	50.67±0.29	51.31±0.19
Ours* (CE)	48.98±0.32	45.88±0.44
Ours* (CE+MAE)	56.73±0.35	57.27±0.40
Ours* (CE+APL)	55.24±0.39	56.64±0.29

TinyImageNet

	CE	GCE	SCE	M-corr.	D-mix.	Ours
Acc.	68.80	69.75	71.02	71.00	74.76	73.59

clothing1M