

A. Experimental Details

A.1. Data Sets and Network Architectures

The MNIST, CIFAR-10, CIFAR-100 data sets are obtained from PyTorch’s torchvision package.¹ A summary is shown in Table 3. The networks (MLP on MNIST, and CNN on CIFAR-10, CIFAR-100) used are shown in Table 4. Models 1 and 2 have been used in (Yu et al., 2019) on CIFAR-10 and CIFAR-100, respectively. Model 3 has been used in (Han et al., 2018).

Table 3. Data sets with artificial label noise.

	#tra	#val	#test	#classes
MNIST	60,000	5,000	5,000	10
CIFAR-10	50,000	5,000	5,000	10
CIFAR-100	50,000	5,000	5,000	100

A.2. Details for Figure 1

A.2.1. FIGURE 1(A)

We use the CIFAR-10 dataset (Table 3), and model 1 in Table 4. The number of training epochs is 200. We use the Adam optimizer (Kingma & Ba, 2014) with momentum 0.9 and batch size 128. The initial learning rate is 0.001, and is linearly decayed to zero from the 80th epoch. The 5 random $R(T)$ s (denoted “Random R(T)” 1-5) are generated by uniform sampling the corresponding hyperparameter $x = \{\alpha, \beta\}$.

Besides the test accuracies shown in Figure 1(a), we also show in Figure 8 the random $R(T)$ s, the original $R(T)$ used in Co-teaching (Han et al., 2018), the $R(T)$ obtained by $S2E$ (denoted “Searched”), and the implicit $R(T)$ corresponding to training on the whole noisy dataset (denoted “Baseline”).

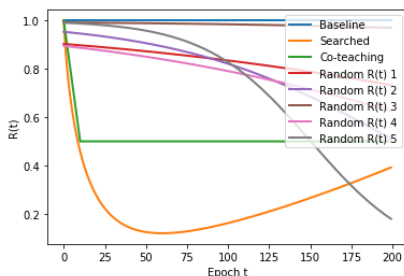


Figure 8. $R(t)$ used in Figure 1(a).

A.2.2. FIGURES 1(B)-1(C)

Experiment is performed on the MNIST/CIFAR-10/CIFAR-100 datasets (Table 3). The number of training epochs, batch size, and learning rate schedule are the same as that

¹<https://pytorch.org/docs/stable/torchvision/datasets.html>

in Figure 1(a).

A.2.3. FIGURE 1(D)

We use the CNN models 1-3 in Table 4. As CIFAR-100 has 100 outputs, we also change the number of outputs of model 1 to 100. The number of training epochs, batch size, and learning rate schedule are the same as that in Figure 1(a).

A.2.4. FIGURE 1(E)

We use model 1 in Table 4. For Adam, the learning rate schedule is the same as that in Figure 1(a). For SGD, the initial learning rate is 0.1, and decayed to 0.01 and 0.001 at the 500th and 750th epoch, respectively. Moreover, the number of training epochs is 1000 instead of 200. For RMSProp, the learning rate is fixed at 0.01.

A.2.5. FIGURE 1(F)

The number of training epochs, batch size, and learning rate schedule are the same as that in Figure 1(a). We only change the batch size and initial learning rate as shown in the figure of Figure 1(f). Moreover, to better demonstrate the memorization effect for small learning rates, the number of training epochs is set to 1000 instead of 200.

A.3. Additional Plots for Section 4.1.2

Figure 9 compares the label precisions of the various methods on CIFAR-10 and CIFAR-100.

B. Additional Experiments

B.1. Approximation to $R(\cdot)$ in Co-teaching

Recall that $R(t)$ in Co-teaching is generated from (1). As all basis functions in Table 1 are smooth, it is not possible for (4) to exactly subsume (1). However, $R(t)$ in (4) can well approximate (1). To illustrate this, we randomly generate three $R(t)$ ’s in Co-teaching’s search space by uniform sampling the corresponding hyperparameters $\tau \in (0, 1)$, $c \in \{0.5, 1, 2\}$ and $t_k \in (0, 200)$. Figure 10 shows the function in (4) that best approximates each of these $R(t)$ ’s with the least squared error.

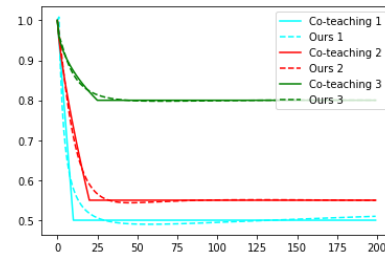


Figure 10. $R(t)$ in Co-teaching and the best approximation from the proposed search space.

Table 4. MLP and CNN models used in the experiments.

MLP on MNIST	CNN on CIFAR-10 (Model 1)	CNN on CIFAR-100 (Model 2)	CNN (Model 3)
28×28 gray image	32×32 RGB image	32×32 RGB image	32×32 RGB image
Dense 28×28→256 ReLU	5×5 Conv, 6 ReLU 2×2 Max-pool	3×3 Conv, 64 BN, ReLU 3×3 Conv, 64 BN, ReLU 2×2 Max-pool	3×3 Conv, 128 BN, LReLU 3×3 Conv, 128 BN, LReLU 3×3 Conv, 128 BN, LReLU 2×2 Max-pool, stride 2 Dropout, p=0.25
	5×5 Conv, 16 ReLU 2×2 Max-pool	3×3 Conv, 128 BN, ReLU 3×3 Conv, 128 BN, ReLU 2×2 Max-pool	3×3 Conv, 256 BN, LReLU 3×3 Conv, 256 BN, LReLU 3×3 Conv, 256 BN, LReLU 2×2 Max-pool, stride 2 Dropout, p=0.25
	Dense 16×5×5→120 ReLU Dense 120→84 ReLU	3×3 Conv, 196 BN, ReLU 3×3 Conv, 196 BN, ReLU 2×2 Max-pool	3×3 Conv, 512 BN, LReLU 3×3 Conv, 256 BN, LReLU 3×3 Conv, 128 BN, LReLU Avg-pool
Dense 256→10	Dense 84→10	Dense 256→100	Dense 128→10

B.2. Comparison with Weight Sharing

Weight sharing (Pham et al., 2018; Liu et al., 2019) is a popular method to speed up the search in NAS. In this experiment, we study if weight-sharing is also beneficial to the search of $R(\cdot)$. We compare *S2E* with *ASNG* (Akimoto et al., 2019), which is a weight-sharing version of NG. Specifically, *ASNG* optimizes

$$\min_{\theta, w} \mathcal{G}(\theta, w) \equiv \int_{\mathbf{x} \in \mathcal{F}} \mathcal{L}_{\text{val}}(f(\mathbf{w}; R(\mathbf{x})), \mathcal{D}_{\text{val}}) p_{\theta}(\mathbf{x}) d\mathbf{x},$$

by alternating the updates of w (using gradient descent) and θ (using natural gradient descent). Unlike *S2E* in (6), in which each θ has its own optimal w^* , *ASNG* only uses one w that is shared by all θ .

Table 5 compare the test accuracies of *S2E* and *ASNG*. As can be seen, the $R(\cdot)$ obtained by *ASNG* is much worse than that from *S2E*, indicating weight-sharing is not a good choice here. Recently, the problem of weight sharing is also discussed in (Sciuto et al., 2020), which shows that it is not useful in NAS for convolutional and recurrent neural works.

 Table 5. Testing accuracies (%) obtained on CIFAR-10 by *ASNG* and *S2E*.

	sym-20%	sym-50%	pair-45%
ASNG	57.82	47.34	41.46
S2E	58.73	50.82	47.58

C. Proofs

C.1. Proposition 1

Proof. By definition,

$$\begin{aligned} \nabla^2 \mathcal{J}(\theta) &= \int \bar{f}(\mathbf{x}) \nabla^2 p_{\theta}(\mathbf{x}) d\mathbf{x} \\ &= \mathbb{E}_{p_{\theta}} \left[\bar{f}(\mathbf{x}) \frac{\nabla^2 p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right]. \end{aligned} \quad (10)$$

Now,

$$\begin{aligned} \nabla^2 \log p_{\theta}(\mathbf{x}) &= \nabla \left(\frac{\nabla p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \\ &= \frac{\nabla^2 p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{x})} - \frac{\nabla p_{\theta}(\mathbf{x}) \nabla p_{\theta}(\mathbf{x})^{\top}}{p_{\theta}^2(\mathbf{x})}. \end{aligned}$$

Thus,

$$\begin{aligned} \frac{\nabla^2 p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{x})} &= \nabla^2 \log p_{\theta}(\mathbf{x}) + \frac{\nabla p_{\theta}(\mathbf{x}) \nabla p_{\theta}(\mathbf{x})^{\top}}{p_{\theta}^2(\mathbf{x})} \\ &= \nabla^2 \log p_{\theta}(\mathbf{x}) + \left(\frac{\nabla p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \left(\frac{\nabla p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right)^{\top} \\ &= \nabla^2 \log p_{\theta}(\mathbf{x}) + \bar{p}_{\theta} \bar{p}_{\theta}^{\top}, \end{aligned}$$

Result follows on substituting this into (10). \square

C.2. Proposition 2

Proof. First, we introduce the following Lemma 1 which results from Assumption 2.

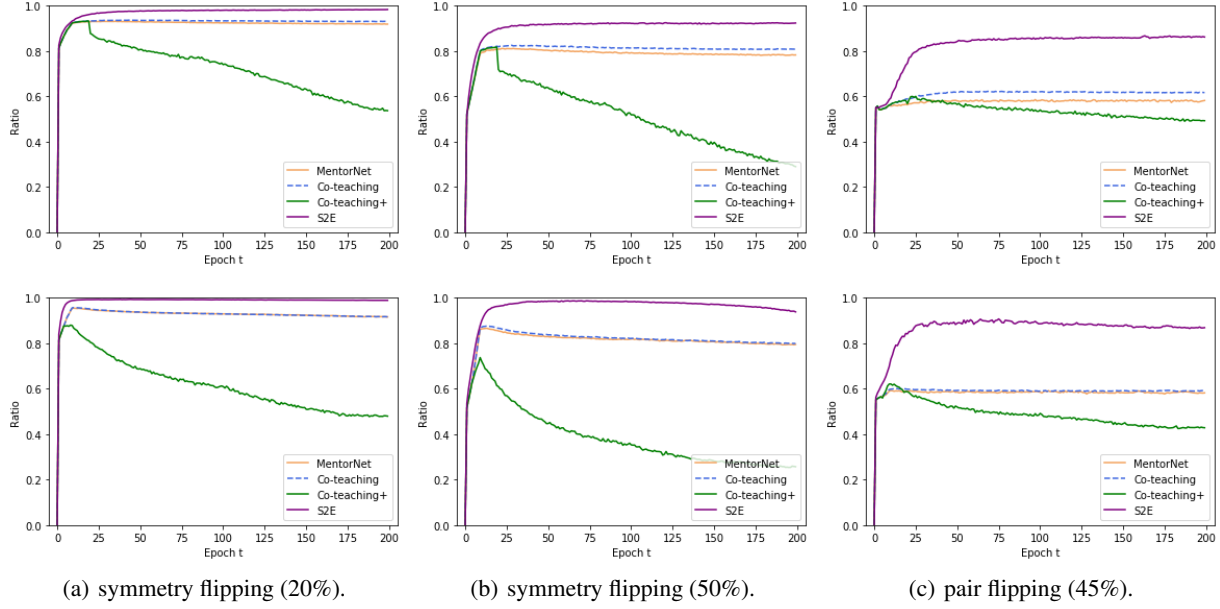


Figure 9. Label precisions of *MentorNet*, *Co-teaching*, *Co-teaching+* and *S2E* on CIFAR-10 (top) and CIFAR-100 (bottom).

Lemma 1 ((Rockafellar, 1970)). *Since \mathcal{J} is L -Lipschitz smooth, we have $\mathcal{J}(\mathbf{y}) \leq \mathcal{J}(\mathbf{x}) + \langle \nabla \mathcal{J}(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2$ for any \mathbf{x} and \mathbf{y} .*

Define a function g as

$$g(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}, \mathbf{H}) = (\boldsymbol{\theta} - \mathbf{y})^\top \mathbf{z} + \frac{1}{2} (\boldsymbol{\theta} - \mathbf{y})^\top \mathbf{H} (\boldsymbol{\theta} - \mathbf{y}).$$

Due to (9), we can express $\boldsymbol{\theta}^{m+1}$ as

$$\boldsymbol{\theta}^{m+1} = \arg \min_{\boldsymbol{\theta}} g(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}, \mathbf{H}), \quad (11)$$

where

$$\mathbf{y} = \boldsymbol{\theta}^m, \mathbf{z} = \nabla \mathcal{J}(\boldsymbol{\theta}^m) - \mathbf{e}^m \text{ and } \mathbf{H} = \boldsymbol{\Delta}^m. \quad (12)$$

Note that $\boldsymbol{\Delta}^m$ is a positive definite matrix, thus g is a convex function on $\boldsymbol{\theta}$. Consider the directional derivative of g w.r.t. $\boldsymbol{\theta}$ at the optimal point $\boldsymbol{\theta} = \boldsymbol{\theta}^{m+1}$, and using the fact that g is a convex function, we have

$$\langle \mathbf{z} + \mathbf{H}(\boldsymbol{\theta}^{m+1} - \mathbf{y}), \mathbf{w}^m \rangle \geq 0 \quad (13)$$

for any direction \mathbf{w} .

Let $\mathbf{w} = \boldsymbol{\theta}^m - \boldsymbol{\theta}^{m+1}$. Combining (12) and (13), we have

$$\langle \nabla \mathcal{J}(\boldsymbol{\theta}^m) - \mathbf{e}^m, \boldsymbol{\gamma}^m \rangle \leq -(\boldsymbol{\gamma}^m)^\top \boldsymbol{\Delta}^m \boldsymbol{\gamma}^m. \quad (14)$$

Next, using Lemma 1, we have

$$\begin{aligned} \mathcal{J}(\boldsymbol{\theta}^{m+1}) &\leq \mathcal{J}(\boldsymbol{\theta}^m) + \langle \nabla \mathcal{J}(\boldsymbol{\theta}^m), \boldsymbol{\theta}^{m+1} - \boldsymbol{\theta}^m \rangle \\ &\quad + \frac{L}{2} \|\boldsymbol{\theta}^{m+1} - \boldsymbol{\theta}^m\|^2. \end{aligned} \quad (15)$$

Now, we add the error term \mathbf{e}^m in (15), i.e.,

$$\begin{aligned} \mathcal{J}(\boldsymbol{\theta}^{m+1}) &\leq \mathcal{J}(\boldsymbol{\theta}^m) + \langle \nabla \mathcal{J}(\boldsymbol{\theta}^m) - \mathbf{e}^m, \boldsymbol{\theta}^{m+1} - \boldsymbol{\theta}^m \rangle \\ &\quad + \frac{L}{2} \|\boldsymbol{\theta}^{m+1} - \boldsymbol{\theta}^m\|^2 + \langle \mathbf{e}^m, \boldsymbol{\theta}^{m+1} - \boldsymbol{\theta}^m \rangle, \\ &\leq \mathcal{J}(\boldsymbol{\theta}^m) - (\boldsymbol{\gamma}^m)^\top \boldsymbol{\Delta}^m \boldsymbol{\gamma}^m + \frac{L}{2} \|\boldsymbol{\theta}^{m+1} - \boldsymbol{\theta}^m\|^2 \\ &\quad + \langle \mathbf{e}^m, \boldsymbol{\theta}^{m+1} - \boldsymbol{\theta}^m \rangle, \end{aligned} \quad (16)$$

$$\begin{aligned} &= \mathcal{J}(\boldsymbol{\theta}^m) - (\boldsymbol{\gamma}^m)^\top \boldsymbol{\Delta}^m \boldsymbol{\gamma}^m + \frac{L}{2} \|\boldsymbol{\gamma}^m\|^2 + \langle \mathbf{e}^m, \boldsymbol{\gamma}^m \rangle, \\ &\leq \mathcal{J}(\boldsymbol{\theta}^m) - (\boldsymbol{\gamma}^m)^\top \boldsymbol{\Delta}^m \boldsymbol{\gamma}^m + \frac{L}{2} \|\boldsymbol{\gamma}^m\|^2 \\ &\quad + \|\mathbf{e}^m\|_2 \|\boldsymbol{\gamma}^m\|_2, \end{aligned} \quad (17)$$

$$\begin{aligned} &\leq \mathcal{J}(\boldsymbol{\theta}^m) - \frac{1}{\eta} \|\boldsymbol{\gamma}^m\|^2 + \frac{L}{2} \|\boldsymbol{\gamma}^m\|^2 \\ &\quad + \|\mathbf{e}^m\|_2 \|\boldsymbol{\gamma}^m\|_2, \end{aligned} \quad (18)$$

$$= \mathcal{J}(\boldsymbol{\theta}^m) + \frac{\eta L - 2\eta}{2\eta} \|\boldsymbol{\gamma}^m\|^2 + \|\mathbf{e}^m\|_2 \|\boldsymbol{\gamma}^m\|_2, \quad (19)$$

where (16) is from (14), (17) is from inequality $\langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle \leq \|\boldsymbol{\alpha}\| \|\boldsymbol{\beta}\|$, (18) results from Assumption 3, i.e., the smallest eigen value of $\boldsymbol{\Delta}^m$ is not smaller than η . Finally, rearranging terms in (19), we will obtain the Proposition. \square

C.3. Theorem 1

Before proving this Theorem 1, we first introduce the following Lemma 1.

Lemma 2. *Define*

$$\boldsymbol{\varepsilon}^m = \left\| (\boldsymbol{\Delta}^m)^{-1} \mathbf{e}^m \right\| \text{ and } \rho^m = (\boldsymbol{\Delta}^m)^{-1} \mathcal{J}(\boldsymbol{\theta}^m).$$

We have

$$c^m \leq \|\gamma^m\| \leq \|\rho^m\| + \varepsilon^m,$$

where $c^m = \max(\|\rho^m\| - \varepsilon^m, \varepsilon^m - \|\rho^m\|)$.

Proof. Since θ^{m+1} is generated by (9), thus

$$\begin{aligned} \|\gamma^m\| &= \left\| (\Delta^m)^{-1} (\nabla \mathcal{J}(\theta^m) - e^m) \right\|, \\ &= \left\| (\Delta^m)^{-1} e^m + \rho^m \right\|. \end{aligned}$$

Then, the Lemma follows from Cauchy inequality. \square

Now, we start to prove Theorem 1.

Proof of Theorem 1. Since the eigen values of Δ^m are in $[\eta, L]$ (by Assumption 3), we have

$$\frac{1}{L} \|e^m\| \leq \left\| (\Delta^m)^{-1} e^m \right\| \leq \frac{1}{\eta} \|e^m\|. \quad (20)$$

Combining (20) and Proposition 2, we obtain

$$\begin{aligned} \mathcal{J}(\theta^m) - \mathcal{J}(\theta^{m+1}) &\geq \frac{2-L\eta}{2\eta} \|\gamma^m\|^2 - \|e^m\| \|\gamma^m\|, \\ &\geq \frac{2-L\eta}{2\eta} \|\gamma^m\|^2 - \eta(\varepsilon^m)^2 \|\gamma^m\|. \end{aligned} \quad (21)$$

Next, using Lemma 2 in (21), we have

$$\begin{aligned} \mathcal{J}(\theta^m) - \mathcal{J}(\theta^{m+1}) &\geq \frac{2-L\eta}{2\eta} (\|\rho^m\| - \varepsilon^m)^2 \\ &\quad - \eta(\varepsilon^m)^2 (\|\rho^m\| + \varepsilon^m). \end{aligned}$$

Rearranging teams in the above inequality, we have where

$$\begin{aligned} b_1 &= \frac{2-L\eta}{2\eta}, \\ b_2 &= \frac{2-L\eta + \eta^2}{\eta}, \\ \text{and } b_3 &= \frac{2\eta^2 + L\eta - 2}{2\eta}. \end{aligned}$$

First Assertion. Define the following auxiliary function

$$\psi(\theta^m) = b_1 \|\rho^m\|^2 - b_2 \|\rho^m\| \varepsilon - b_3 (\varepsilon^m)^2.$$

With this definition, we have

$$\mathcal{J}(\theta^m) - \mathcal{J}(\theta^{m+1}) \geq \psi(\theta^m). \quad (22)$$

It is easy to see that since $\|\rho^m\|$ and ε are continuous, b_1, b_2 and b_3 are non-negative, then $\psi(\theta^m)$ is lower semi-continuous. Let the sub-level set of ψ be

$$\mathcal{L}(\psi, a) \equiv \{\theta \mid \psi(\theta) \leq a\}, \quad a \geq 0.$$

Note that the sub-level set of $\mathcal{L}(\psi, a)$ is closed for any $a \geq 0$ (see Theorem 7.1 in (Rockafellar, 1970)). Denote $u = \|\rho^m\|$, and resolving the quadratic inequality in u :

$$b_1 u^2 - b_2 \varepsilon^m u - b_3 (\varepsilon^m)^2 - t \leq 0,$$

we conclude

$$u \leq \frac{b_2 \varepsilon^m}{2b_1} + \frac{1}{2b_1} \sqrt{(b_2^2 + 4b_1 b_3) (\varepsilon^m)^2 + 4b_1 a}.$$

Thus,

$$\begin{aligned} \mathcal{L}(\psi, a) &= \left\{ \theta \mid \|\rho^m\| \leq \frac{b_2 \varepsilon^m}{2b_1} \right. \\ &\quad \left. + \frac{1}{2b_1} \sqrt{(b_2^2 + 4b_1 b_3) (\varepsilon^m)^2 + 4b_1 a} \right\}, \end{aligned}$$

In particular,

$$\mathcal{L}(\psi, 0) = \left\{ \theta \mid \|\rho^m\| \leq \frac{b_2 + \sqrt{(b_2^2 + 4b_1 b_3)} \varepsilon}{2b_1} \right\}.$$

Define

$$d_1 = \frac{b_2 + \sqrt{b_2^2 + 4b_1 b_3}}{2b_1} \quad \text{and} \quad d_2 = b_1^{-1}.$$

We conclude

$$\mathcal{L}(\psi, a) \subseteq \left\{ \theta \mid \|\rho^m\| \leq d_1 \varepsilon^m + d_2 a^{\frac{1}{2}} \right\},$$

and

$$\mathcal{L}(\psi, 0) \subseteq \{\theta \mid \|\rho^m\| \leq d_1 \varepsilon^m\}. \quad (23)$$

Next, we prove that there exists a limit point $\bar{\theta}$ of $\{\theta^m\}$ such that $\bar{\theta} \in \mathcal{L}(\psi, 0)$. Suppose the opposite holds. By (22), we have

$$\mathcal{J}(\theta^m) - \mathcal{J}(\theta^{m+1}) \geq \psi(\theta^m), \quad \forall m \geq m_1. \quad (24)$$

By Assumption $\lim_{m \rightarrow \infty} \{\theta^m\} \cap \mathcal{L}(\psi, 0) = \emptyset$. Then, since ψ is lower semi-continuous, we have

$$\psi(\theta^m) \geq c > 0,$$

when $m \geq m_2$ for some sufficiently large m_2 and a positive constant c .

Denote $k = \max\{m_1, m_2\}$, for any $m \geq k$, we have

$$\begin{aligned} \mathcal{J}(\theta^k) - \mathcal{J}(\theta^m) &= \sum_{j=m}^k (\mathcal{J}(\theta^j) - \mathcal{J}(\theta^{j+1})), \\ &\geq (k-m)c. \end{aligned}$$

Let $k \rightarrow \infty$, we have $\lim_{m \rightarrow \infty} \mathcal{J}(\theta^m) = -\infty$, which contradicts with Assumption 2, i.e., $\inf \mathcal{J} > -\infty$. Thus, from (23), for every limit point $\bar{\theta}$ of $\{\theta^m\}$, we must have

$$\|\rho^m\| \leq d_1 \varepsilon^m.$$

Recall the definition of ρ^m in Lemma 2, and by Assumption 3 that the error ε^m on gradient is upper-bounded by $\bar{\varepsilon}$, we obtain the first assertion.

Second Assertion. If the sequence $\{\boldsymbol{\theta}^m\}$ converges, then for every sub-sequence $\{\boldsymbol{\theta}^{m_i}\}$ of $\{\boldsymbol{\theta}^m\}$ it follows

$$\limsup_{i \rightarrow \infty} \mathcal{J}(\boldsymbol{\theta}^{m_i}) = \liminf_{m \rightarrow \infty} \mathcal{J}(\boldsymbol{\theta}^m),$$

Thus,

$$\lim_{m \rightarrow \infty} \boldsymbol{\theta}^m \subseteq \mathcal{L}(\psi, 0), \quad (25)$$

where $\mathcal{L}(\psi, 0)$ is in (23). Combing (25) with the first assertion, we then have

$$\lim_{m \rightarrow \infty} \|\varepsilon^m\| \leq c_1 \bar{\varepsilon}.$$

Finally, by (i) in Assumption 2 and definition of ε^m in Lemma 2, we have

$$\lim_{m \rightarrow \infty} \|\mathbf{e}^m\| \leq c_2 \bar{\varepsilon},$$

for a positive constant c_2 , which proves the second assertion. \square