## Robust Semi-supervised Learning by Wisely Leveraging Open-set Data

Yang Yang, Member, IEEE, Nan Jiang, Yi Xu, and De-Chuan Zhan

Index Terms—Semi-supervised Learning, OOD Detection, Open-set Data.

## REFERENCES

## APPENDIX A PROOF OF THE THEOREM 1

*Proof.* (a) Following the standard analysis, we have the following inequality in expectation

$$\begin{aligned}
& \mathbf{E}[\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_{t})] \\
& \leq \mathbf{E}[\langle \theta_{t+1} - \theta_{t}, \nabla \mathcal{L}(\theta_{t}) \rangle] + \frac{L}{2} \mathbf{E}[\|\theta_{t+1} - \theta_{t}\|^{2}] \\
& \stackrel{\text{(ii)}}{=} - \eta \mathbf{E}[\langle \widehat{g}(\theta_{t}), \nabla \mathcal{L}(\theta_{t}) \rangle] + \frac{\eta^{2}L}{2} \mathbf{E}[\|\widehat{g}(\theta_{t})\|^{2}] \\
& \stackrel{\text{(iii)}}{=} - \eta (1 - \frac{\eta L}{2}) \|\nabla \mathcal{L}(\theta_{t})\|^{2} + \frac{\eta^{2}L}{2} \mathbf{E}[\|\nabla \mathcal{L}(\theta_{t}) - \widehat{g}(\theta_{t})\|^{2}] \\
& \stackrel{\text{(iv)}}{\leq} - \frac{\eta}{2} \|\nabla \mathcal{L}(\theta_{t})\|^{2} + \frac{\eta^{2}L}{2} \mathbf{E}[\|\nabla \mathcal{L}(\theta_{t}) - \widehat{g}(\theta_{t})\|^{2}], \end{aligned} \tag{1}$$

where (i) is due to the smoothness of loss function  $\mathcal{L}$  in Assumption  $\ref{GD}$ ; (ii) is due to the update of SGD; (iii) is due to  $\mathrm{E}[\widehat{g}(\theta)] = \nabla \mathcal{L}(\theta)$ ; (iv) is due to  $\eta \leq 1/L$ . By the definition of  $\widehat{g}(\theta_t)$  and the convexity of norm  $\|\cdot\|^2$ , we know

$$E[\|\nabla \mathcal{L}(\theta_{t}) - \widehat{g}(\theta_{t})\|^{2}]$$

$$\leq \lambda E[\|\nabla \mathcal{L}(\theta_{t}) - g_{id}(\theta_{t})\|^{2}]$$

$$+ (1 - \lambda)(\tau E[\|\nabla \mathcal{L}(\theta_{t}) - g_{fr}(\theta_{t})\|^{2}$$

$$+ (1 - \tau)E[\|\nabla \mathcal{L}(\theta_{t}) - g_{uf}(\theta_{t})\|^{2}])$$

$$\leq \lambda \sigma^{2} + (1 - \lambda)[\tau(\frac{\epsilon}{2}\|\nabla \mathcal{L}(\theta_{t})\|^{2} + \sigma^{2})$$

$$+ (1 - \tau)(\frac{\nu}{2}\|\nabla \mathcal{L}(\theta_{t})\|^{2} + \sigma^{2})]$$

$$= \sigma^{2} + (1 - \lambda)\frac{\tau \epsilon + (1 - \tau)\nu}{2}\|\nabla \mathcal{L}(\theta_{t})\|^{2}, \qquad (2)$$

- Yang Yang is with the school of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing 210094, China. E-mail: yyang@njust.edu.cn
- Nan Jiang and De-Chuan Zhan are with the National Key Laboratory for Novel Software Technology, Nanjing University, and also with the School of Artificial Intelligence, Nanjing University, Nanjing 210023, China. E-mail: jiangn@lamda.nju.edu.cn, zhandc@nju.edu.cn
- Yi Xu is with the School of Control Science and Engineering, Dalian University of Technology, Dalian 116081, China.
   E-mail: yxu@dlut.edu.cn

Corresponding authors: Yi Xu and De-Chuan Zhan.

where the second inequality is due to Assumptions ?? and ??. Therefore, by (1) and (2) we have

$$E[\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_t)]$$

$$\leq -\frac{\eta}{2} \left( 1 - \eta L (1 - \lambda) \frac{\tau \epsilon + (1 - \tau)\nu}{2} \right) E[\|\nabla \mathcal{L}(\theta_t)\|^2] + \frac{\eta^2 L \sigma^2}{2}$$

$$\leq -\frac{\mu \eta}{2} E[\mathcal{L}(\theta_t) - \mathcal{L}(\theta_*)] + \frac{\eta^2 L \sigma^2}{2}, \tag{3}$$

where the last inequality is due to Assumption ?? and  $\eta = \frac{1}{(1-\lambda)(\tau\epsilon+(1-\tau)\nu)L}$ . Therefore, we have

$$E[\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_*)] \le \left(1 - \frac{\mu\eta}{2}\right) E[\mathcal{L}(\theta_t) - \mathcal{L}(\theta_*)] + \frac{\eta^2 L \sigma^2}{2},$$
(4)

which implies

$$E[\mathcal{L}(\theta_{n+m+m'+1}) - \mathcal{L}(\theta_*)]$$

$$\leq \exp(-\mu\eta(n+m+m')/2) \left(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*)\right) + \frac{\eta L\sigma^2}{\mu}, (5)$$

Since  $\nu$  is large enough, then we know that  $\eta':=\frac{1}{(1-\lambda)(\tau\epsilon+(1-\tau)\nu)L}$  is small enough such that  $\exp\left(-\mu\eta'(n+m+m')/2\right)$   $\gg \frac{L\sigma^2}{(n+m+m')\mu^2(\mathcal{L}(\theta_0)-\mathcal{L}(\theta_*))}$ . Thus, we have

$$E[\mathcal{L}(\theta_{n+m+m'+1}) - \mathcal{L}(\theta_*)]$$

$$\leq \exp(-\mu \eta'(n+m+m')/2) \left(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*)\right) + \frac{\eta' L \sigma^2}{\mu}$$
(6)

Let's consider a special case of  $n+m+m'=\frac{2(1-\lambda)(\tau\epsilon+(1-\tau)\nu)L}{\mu}$ , then  $\eta'=\frac{2}{(n+m+m')\mu}$ , implying

$$E[\mathcal{L}(\theta_{n+m+m'+1}) - \mathcal{L}(\theta_*)] \le O\left(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*)\right)$$
 (7)

(b) Next, let's consider the case without either far OOD data or near OOD data, i.e.,  $\lambda=1$ . Then, following the similar analysis in (a), we have

$$E[\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_t)]$$

$$\leq -\frac{\eta}{2} E[\|\nabla \mathcal{L}(\theta_t)\|^2] + \frac{\eta^2 L \sigma^2}{2}$$

$$\leq -\frac{\mu \eta}{2} E[\mathcal{L}(\theta_t) - \mathcal{L}(\theta_*)] + \frac{\eta^2 L \sigma^2}{2}, \tag{8}$$

where the last inequality is due to Assumption ??. Therefore, we have

$$E[\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_*)] \le \left(1 - \frac{\mu\eta}{2}\right) E[\mathcal{L}(\theta_t) - \mathcal{L}(\theta_*)] + \frac{\eta^2 L \sigma^2}{2},$$
(9)

which implies

$$E[\mathcal{L}(\theta_{n+1}) - \mathcal{L}(\theta_*)] \le \exp(-\mu \eta n/2) \left(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*)\right) + \frac{\eta L \sigma^2}{\mu},$$
(10)

By setting 
$$\eta = \frac{2}{n\mu} \log \left( \frac{n\mu^2(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*))}{\sigma^2 L} \right)$$
,
$$\mathbb{E}[\mathcal{L}(\theta_{n+1}) - \mathcal{L}(\theta_*)]$$

$$\leq \frac{L\sigma^2}{n\mu^2} + \frac{2L\sigma^2}{n\mu^2} \log \left( \frac{n\mu^2(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*))}{\sigma^2 L} \right)$$

$$\leq O\left( \frac{\log(n)}{n} \right) \tag{11}$$

(c) Finally, let's consider the case without far OOD data (but with near OOD data), i.e.,  $\tau=1$ . Then, following the similar analysis in (a), we have

$$E[\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_t)]$$

$$\leq -\frac{\eta}{2} \left( 1 - \frac{(1-\lambda)\epsilon\eta L}{2} \right) E[\|\nabla \mathcal{L}(\theta_t)\|^2] + \frac{\eta^2 L\sigma^2}{2}$$

$$\leq -\frac{\mu\eta}{2} E[\mathcal{L}(\theta_t) - \mathcal{L}(\theta_*)] + \frac{\eta^2 L\sigma^2}{2}, \tag{12}$$

where the last inequality is due to Assumption ?? and  $\eta \leq \frac{1}{(1-\lambda)\epsilon\nu L}$ . Therefore, we have

$$E[\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_*)] \le \left(1 - \frac{\mu\eta}{2}\right) E[\mathcal{L}(\theta_t) - \mathcal{L}(\theta_*)] + \frac{\eta^2 L \sigma^2}{2},$$
(13)

which implies

$$E[\mathcal{L}(\theta_{n+m+1}) - \mathcal{L}(\theta_*)]$$

$$\leq \exp\left(-\mu\eta(n+m)/2\right) \left(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*)\right) + \frac{\eta L\sigma^2}{\mu}, \quad (14)$$
By setting  $\eta = \frac{2}{(n+m)\mu} \log\left(\frac{(n+m)\mu^2(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*))}{\sigma^2 L}\right),$ 

$$E[\mathcal{L}(\theta_{n+m+1}) - \mathcal{L}(\theta_*)]$$

$$\leq \frac{L\sigma^2}{(n+m)\mu^2} + \frac{2L\sigma^2}{(n+m)\mu^2} \log\left(\frac{(n+m)\mu^2(\mathcal{L}(\theta_0) - \mathcal{L}(\theta_*))}{\sigma^2 L}\right)$$

$$\leq O\left(\frac{\log(n+m)}{(n+m)}\right) \tag{15}$$