

# Stability and Generalization of Bilevel Programming in Hyperparameter Optimization: Appendix

Fan Bao\*, Guoqiang Wu\*<sup>†</sup>, Chongxuan Li\*<sup>†</sup>, Jun Zhu<sup>‡</sup>, Bo Zhang

Dept. of Comp. Sci. & Tech., Institute for AI, Tsinghua-Huawei Joint Center for AI  
BNRist Center, State Key Lab for Intell. Tech. & Sys., Tsinghua University, Beijing, China  
bf19@mails.tsinghua.edu.cn, {guoqiangwu90, chongxuanli1991}@gmail.com,  
{dcszj, dcszb}@tsinghua.edu.cn

## A Proofs of Main Theoretical Results

### A.1 Proof of Theorem 1

**Theorem 1** (Generalization bound of a uniformly stable algorithm). *Suppose a randomized HO algorithm  $\mathbf{A}$  is  $\beta$ -uniformly stable on validation in expectation, then*

$$|\mathbb{E}_{\mathbf{A}, S^{tr} \sim (D^{tr})^n, S^{val} \sim (D^{val})^m} [R(\mathbf{A}(S^{tr}, S^{val}), D^{val}) - \hat{R}^{val}(\mathbf{A}(S^{tr}, S^{val}), S^{val})]| \leq \beta.$$

*Proof.*

$$\begin{aligned} & |\mathbb{E}_{\mathbf{A}, S^{tr}, S^{val}} [R(\mathbf{A}(S^{tr}, S^{val}), D^{val}) - \hat{R}^{val}(\mathbf{A}(S^{tr}, S^{val}), S^{val})]| \\ &= |\mathbb{E}_{\mathbf{A}, S^{tr}, S^{val}, z \sim D^{val}} [\ell(\mathbf{A}(S^{tr}, S^{val}), z) - \ell(\mathbf{A}(S^{tr}, S^{val}), z_1^{val})]| \\ &= |\mathbb{E}_{\mathbf{A}, S^{tr}, S^{val}, z \sim D^{val}} [\ell(\mathbf{A}(S^{tr}, z, z_2^{val}, \dots, z_m^{val}), z_1^{val}) - \ell(\mathbf{A}(S^{tr}, S^{val}), z_1^{val})]| \\ &\leq \mathbb{E}_{S^{tr}, S^{val}, z \sim D^{val}} |\mathbb{E}_{\mathbf{A}} [\ell(\mathbf{A}(S^{tr}, z, z_2^{val}, \dots, z_m^{val}), z_1^{val}) - \ell(\mathbf{A}(S^{tr}, S^{val}), z_1^{val})]| \leq \beta, \end{aligned}$$

where the last inequality is due to the definition of stability.  $\square$

### A.2 Proof of Theorem 2

Here we prove a more general version of Theorem 2 in the full paper by considering SGD with weight decay in the outer level, i.e.,

$$\lambda_{t+1} = (1 - \alpha_{t+1}\mu)\lambda_t - \alpha_{t+1}\nabla_{\lambda_t} \ell(\lambda_t, \hat{\theta}(\lambda_t, S^{tr}), z_j^{val}), \quad (1)$$

where  $\alpha_t$  is the learning rate,  $\mu$  is the weight decay,  $j$  is randomly selected from  $\{1, \dots, m\}$  and  $\hat{\theta}$  is a random function. Theorem 2 in the full paper can be simply derived by letting  $\mu = 0$ .

**Theorem 2** (Uniform stability of algorithms with SGD in the outer level). *Suppose  $\hat{\theta}$  is a random function in a function space  $\mathcal{G}_{\hat{\theta}}$  and  $\forall S^{tr} \in Z^n, \forall z \in Z, \forall g \in \mathcal{G}_{\hat{\theta}}, \ell(\lambda, g(\lambda, S^{tr}), z)$  as a function of  $\lambda$  is  $L$ -Lipschitz continuous and  $\gamma$ -Lipschitz smooth, let  $c \leq \frac{s(\ell)}{2L^2}$ ,  $\mu \leq \min(\frac{1}{c}, (1-1/m)\gamma)$  and  $\kappa = \frac{c((1-1/m)\gamma - \mu)}{c((1-1/m)\gamma - \mu) + 1}$ . Then, solving Eq. (4) in the full paper with  $T$  steps SGD, learning rate  $\alpha_t \leq \frac{c}{t}$  and weight decay  $\mu$  in the outer level is  $\beta$ -uniformly stable on validation in expectation with*

$$\beta = \frac{2cL^2}{m} \left( \frac{1}{\kappa} \left( \left( \frac{Ts(\ell)}{2cL^2} \right)^\kappa - 1 \right) + 1 \right),$$

\*Equal contribution

<sup>†</sup>G. Wu is now at School of Software, Shandong University and C. Li is now at Gaoling School of AI, Renmin University of China. The work was done when they were at Tsinghua University.

<sup>‡</sup>Corresponding author.

which is increasing w.r.t.  $L$ ,  $\gamma$  and decreasing w.r.t.  $\mu$ .

*Proof.* Suppose  $S^{tr} \in Z^n$  and  $z \in Z$ , let  $f(\lambda, g) = \ell(\lambda, g(\lambda, S^{tr}), z)$ , where we omit the dependency on  $S^{tr}$  and  $z$  for simplicity, then  $f(\lambda, g)$  is as a function of  $\lambda$  is  $L$ -Lipschitz continuous and  $\gamma$ -Lipschitz smooth. Suppose  $S^{val}$  and  $S^{o\text{val}}$  differ in at most one point, let  $\{\lambda_t\}_{t \geq 0}$  and  $\{\lambda'_t\}_{t \geq 0}$  be the trace of Eq. (1) with  $S^{val}$  and  $S^{o\text{val}}$  respectively. Then the output of the HO algorithm  $\mathbf{A}$  with  $t$  steps SGD in the outer level is

$$\mathbf{A}(S^{tr}, S^{val}) = (\lambda_t, \hat{\theta}(\lambda_t, S^{tr})), \quad \mathbf{A}(S^{tr}, S^{o\text{val}}) = (\lambda'_t, \hat{\theta}(\lambda'_t, S^{tr})),$$

and

$$\begin{aligned} \ell(\mathbf{A}(S^{tr}, S^{val}), z) &= \ell(\lambda_t, \hat{\theta}(\lambda_t, S^{tr}), z) = f(\lambda_t, \hat{\theta}), \\ \ell(\mathbf{A}(S^{tr}, S^{o\text{val}}), z) &= \ell(\lambda'_t, \hat{\theta}(\lambda'_t, S^{tr}), z) = f(\lambda'_t, \hat{\theta}). \end{aligned}$$

Let  $\delta_t = \|\lambda_t - \lambda'_t\|$ . Suppose  $0 \leq t_0 \leq t$ , we have

$$\begin{aligned} \mathbf{E} \left[ |f(\lambda_t, \hat{\theta}) - f(\lambda'_t, \hat{\theta})| \right] &= \mathbf{E} \left[ |f(\lambda_t, \hat{\theta}) - f(\lambda'_t, \hat{\theta})| \cdot 1_{\delta_{t_0}=0} \right] \\ &\quad + \mathbf{E} \left[ |f(\lambda_t, \hat{\theta}) - f(\lambda'_t, \hat{\theta})| \cdot 1_{\delta_{t_0}>0} \right] \\ &\leq L \mathbf{E} [\delta_t \cdot 1_{\delta_{t_0}=0}] + P(\delta_{t_0} > 0) s(\ell). \end{aligned}$$

Without loss of generality, we assume  $S^{val}$  and  $S^{o\text{val}}$  at most differ in at the first point. If SGD doesn't select the first point for the first  $t_0$  iterations, then  $\delta_{t_0} = 0$ . As a result,

$$P(\delta_{t_0} = 0) \geq (1 - \frac{1}{m})^{t_0} \geq 1 - \frac{t_0}{m}.$$

Therefore,  $P(\delta_{t_0} > 0) \leq \frac{t_0}{m}$  and we have

$$\mathbf{E} \left[ |f(\lambda_t, \hat{\theta}) - f(\lambda'_t, \hat{\theta})| \right] \leq L \mathbf{E} [\delta_t \cdot 1_{\delta_{t_0}=0}] + \frac{t_0}{m} s(\ell). \quad (2)$$

Now we bound  $\mathbf{E} [\delta_t \cdot 1_{\delta_{t_0}=0}]$ . Let  $\gamma' = (1 - 1/m)\gamma - \mu$  and let  $j$  be the index selected by SGD at the  $t + 1$  iteration, then we have

$$\begin{aligned} \mathbf{E} [\delta_{t+1} \cdot 1_{\delta_{t_0}=0}] &\leq \mathbf{E} [\delta_{t+1} \cdot 1_{j=1} \cdot 1_{\delta_{t_0}=0}] + \mathbf{E} [\delta_{t+1} \cdot 1_{j>1} \cdot 1_{\delta_{t_0}=0}] \\ &\leq \frac{1}{m} (|1 - \alpha_{t+1}\mu| \cdot \mathbf{E}[\delta_t \cdot 1_{\delta_{t_0}=0}] + 2\alpha_{t+1}L) \\ &\quad + \frac{m-1}{m} (|1 - \alpha_{t+1}\mu| + \alpha_{t+1}\gamma) \mathbf{E}[\delta_t \cdot 1_{\delta_{t_0}=0}] \\ &= (1 + \alpha_{t+1}\gamma') \mathbf{E}[\delta_t \cdot 1_{\delta_{t_0}=0}] + \frac{2\alpha_{t+1}L}{m} \\ &\leq \exp(\alpha_{t+1}\gamma') \mathbf{E}[\delta_t \cdot 1_{\delta_{t_0}=0}] + \frac{2\alpha_{t+1}L}{m} \\ &\leq \exp(\frac{c}{t+1}\gamma') \mathbf{E}[\delta_t \cdot 1_{\delta_{t_0}=0}] + \frac{2cL}{(t+1)m}. \end{aligned}$$

As a result,

$$\begin{aligned} \mathbf{E}[\delta_t \cdot 1_{\delta_{t_0}=0}] &\leq \sum_{j=t_0+1}^t \frac{2cL}{jm} \prod_{k=j+1}^t \exp(\frac{c\gamma'}{k}) = \sum_{j=t_0+1}^t \frac{2cL}{jm} \exp(c\gamma' \sum_{k=j+1}^t \frac{1}{k}) \\ &\leq \sum_{j=t_0+1}^t \frac{2cL}{jm} \exp(c\gamma' \ln \frac{t}{j}) = \sum_{j=t_0+1}^t \frac{2cL}{jm} \left(\frac{t}{j}\right)^{c\gamma'} \\ &= \frac{2cLt^{c\gamma'}}{m} \sum_{j=t_0+1}^t \left(\frac{1}{j}\right)^{1+c\gamma'} \leq \frac{2cLt^{c\gamma'}}{m} \frac{t^{-c\gamma'} - t_0^{-c\gamma'}}{-c\gamma'} \\ &= \frac{2L}{m\gamma'} \left( \left(\frac{t}{t_0}\right)^{c\gamma'} - 1 \right). \end{aligned}$$

Combining with Eq. (2), we have

$$\mathbf{E} \left[ |f(\lambda_T, \hat{\theta}) - f(\lambda'_T, \hat{\theta})| \right] \leq \inf_{0 \leq t_0 \leq T} \frac{2L^2}{m\gamma'} \left( \left( \frac{T}{t_0} \right)^{c\gamma'} - 1 \right) + \frac{t_0}{m} s(\ell). \quad (3)$$

The right hand side is approximately minimized when

$$t_0 = \left( \frac{2cL^2}{s(\ell)} \right)^{\frac{1}{c\gamma'+1}} T^{\frac{c\gamma'}{c\gamma'+1}} \leq T,$$

which gives

$$\mathbf{E} \left[ |f(\lambda_T, \hat{\theta}) - f(\lambda'_T, \hat{\theta})| \right] \leq \frac{1 + 1/c\gamma'}{m} (2cL^2)^{\frac{1}{c\gamma'+1}} T^{\frac{c\gamma'}{c\gamma'+1}} (s(\ell))^{\frac{c\gamma'}{c\gamma'+1}} - \frac{2L^2}{m\gamma'} =: \beta.$$

Let  $\kappa = \frac{c\gamma'}{c\gamma'+1} = \frac{c((1-1/m)\gamma-\mu)}{c((1-1/m)\gamma-\mu)+1}$ , then  $\beta$  can be written as

$$\beta = \frac{2cL^2}{m} \left( \frac{1}{\kappa} \left( \left( \frac{Ts(\ell)}{2cL^2} \right)^\kappa - 1 \right) + 1 \right).$$

Since the r.h.s. of Eq. (3) is increasing w.r.t.  $L$  and  $\gamma'$ , where  $\gamma'$  is further increasing w.r.t.  $\gamma$  and decreasing w.r.t.  $\mu$ , we can conclude  $\beta$  is increasing w.r.t.  $L, \gamma$  and decreasing w.r.t.  $\mu$ .  $\square$

### A.3 Proof of Theorem 3

**Definition 1.** (Lipschitz continuous) Suppose  $(X, d_X), (Y, d_Y)$  are two metric spaces and  $f : X \rightarrow Y$ . We define  $f$  is  $L$  Lipschitz continuous iff  $\forall a, b \in X, d_Y(f(a), f(b)) \leq Ld_X(a, b)$ .

**Definition 2.** (Lipschitz smooth) Suppose  $X, Y$  are subsets of two real normed vector spaces and  $f : X \rightarrow Y$  is differentiable. We define  $f$  is  $\gamma$  Lipschitz smooth iff  $f'$  is  $\gamma$  Lipschitz continuous.

**Definition 3.** (Lipschitz norm) Suppose  $(X, d_X), (Y, d_Y)$  are two metric spaces,  $f : X \rightarrow Y$ , we define  $\|f\|_{Lip} = \inf\{L \in [0, \infty] : \forall a, b \in X, d_Y(f(a), f(b)) \leq Ld_X(a, b)\}$ , i.e., the minimum  $L$  such that  $f$  is  $L$  Lipschitz continuous.

**Definition 4.** Given a function  $f(\lambda, \theta)$ , we use  $\|f(\lambda, \theta)\|_{\lambda \in \Lambda, Lip}$  and  $\|f(\lambda, \theta)\|_{\theta \in \Theta, Lip}$  to explicitly denote the Lipschitz norm of  $f$  w.r.t.  $\lambda \in \Lambda$  and  $\theta \in \Theta$  respectively.

**Definition 5.** (Vector norm) Suppose  $a \in \mathbf{R}^m$ , we use  $\|a\|$  to denote the  $l_2$  norm of  $a$ .

**Definition 6.** (Matrix norm) Suppose  $A \in \mathbf{R}^{m \times n}$ , we define  $\|A\| = \sup_{0 \neq a \in \mathbf{R}^n} \frac{\|Aa\|}{\|a\|}$ , i.e., the norm of the linear operator induced by  $A$ .

**Lemma 1.** Suppose  $X, Y$  are two real normed vector spaces,  $\Omega$  is an open set of  $X$ ,  $f : \Omega \rightarrow Y$  is continuously differentiable,  $S \subset \Omega$  is convex and has non-empty interior, then  $\|f|_S\|_{Lip} = \sup_{c \in S} \|f'(c)\|$ .

*Proof.* Suppose  $a, b \in S$ , according to the mean value theorem, there is a  $c$  lies in the segment determined by  $a$  and  $b$ , s.t.,  $\|f(b) - f(a)\| \leq \|f'(c)(b - a)\|$ . Furthermore, we have

$$\|f'(c)(b - a)\| \leq \|f'(c)\| \cdot \|b - a\| \leq \sup_{c \in S} \|f'(c)\| \cdot \|b - a\|.$$

Thereby,  $f|_S$  is  $\sup_{c \in S} \|f'(c)\|$  Lipschitz continuous and  $\|f|_S\|_{Lip} \leq \sup_{c \in S} \|f'(c)\|$ .

Suppose  $c \in S^\circ$ , where  $S^\circ$  is the interior of  $S$  and  $u \in X$  with  $\|u\| = 1$ , then

$$\lim_{\epsilon \rightarrow 0} \frac{f(c + \epsilon u) - f(c)}{\epsilon} = f'(c)u.$$

Thereby,

$$\|f|_S\|_{Lip} \geq \lim_{\epsilon \rightarrow 0} \left\| \frac{f(c + \epsilon u) - f(c)}{\epsilon} \right\| = \|f'(c)u\|.$$

Since  $u$  is arbitrary, we have  $\|f'(c)\| = \sup_{u \in X, \|u\|=1} \|f'(c)u\| \leq \|f\|_S \|Lip\|$ .

Since  $S$  has non-empty interior, we have  $S \subset \overline{S^\circ}$  by the property of convex sets. Suppose  $c \in S$ , then  $c \in \overline{S^\circ}$  and there is a sequence  $c_n \in S^\circ$ , s.t.,  $c_n \rightarrow c$ . Since  $c_n \in S^\circ$ , we have  $\|f'(c_n)\| \leq \|f\|_S \|Lip\|$ . Let  $n \rightarrow \infty$ , by the continuity of  $f'$ , we have  $\|f'(c)\| \leq \|f\|_S \|Lip\|$ . Since  $c \in S$  is arbitrary, we have  $\sup_{c \in S} \|f'(c)\| \leq \|f\|_S \|Lip\|$ . Finally, we have  $\sup_{c \in S} \|f'(c)\| = \|f\|_S \|Lip\|$ .  $\square$

**Lemma 2.** Suppose  $\Lambda$  and  $\Theta$  are convex and compact with non-empty interiors,  $Z$  is compact,  $\Lambda \times \Theta \times Z$  is included in an open set  $\Omega$  and  $f(\lambda, \theta, z) \in C^k(\Omega)$ , then for all  $i \leq k-1$  order partial differential  $h(\lambda, \theta, z)$  of  $f(\lambda, \theta, z)$ , we have  $\sup_{\theta \in \Theta, z \in Z} \|h(\lambda, \theta, z)\|_{\lambda \in \Lambda, Lip} < \infty$  and

$$\sup_{\lambda \in \Lambda, z \in Z} \|h(\lambda, \theta, z)\|_{\theta \in \Theta, Lip} < \infty.$$

*Proof.* Suppose  $h(\lambda, \theta, z)$  is a  $i \leq k-1$  order partial differential of  $f(\lambda, \theta, z)$ , then  $h(\lambda, \theta, z) \in C^1(\Omega)$  and  $\nabla_\lambda h(\lambda, \theta, z) \in C(\Omega)$ . Since  $\Lambda \times \Theta \times Z$  is compact,  $\nabla_\lambda h(\lambda, \theta, z)$  is bounded in  $\Lambda \times \Theta \times Z$ . According to Lemma 1, we have

$$\sup_{\theta \in \Theta, z \in Z} \|h(\lambda, \theta, z)\|_{\lambda \in \Lambda, Lip} = \sup_{\theta \in \Theta, z \in Z} \sup_{\lambda \in \Lambda} \|\nabla_\lambda h(\lambda, \theta, z)\| < \infty.$$

Similarly, we can derive  $\sup_{\lambda \in \Lambda, z \in Z} \|h(\lambda, \theta, z)\|_{\theta \in \Theta, Lip} < \infty$ .  $\square$

**Lemma 3.** Suppose (1)  $\forall 1 \leq k \leq K, \forall \lambda \in \Lambda, G_{\lambda, k}(\theta)$  is a mapping from  $\Theta$  to  $\Theta$ , i.e.,  $G_{\lambda, k} : \Theta \rightarrow \Theta$ , (2)  $\forall 1 \leq k \leq K, \forall \theta \in \Theta, G_{\lambda, k}(\theta)$  as a function of  $\lambda$  is  $L_1^G < \infty$  Lipschitz continuous, (3)  $\forall 1 \leq k \leq K, \forall \lambda \in \Lambda, G_{\lambda, k}(\theta)$  as a function of  $\theta$  is  $L_2^G < \infty$  Lipschitz continuous. Let  $\hat{\theta}(\lambda) = G_{\lambda, K}(G_{\lambda, K-1}(\cdots(G_{\lambda, 1}(\theta_0))))$ , then  $\hat{\theta}(\lambda)$  is  $L^\theta$  Lipschitz continuous with

$$L^\theta = \begin{cases} L_1^G \frac{(L_2^G)^{K-1}}{L_2^{G-1}} & L_2^G \neq 1 \\ KL_1^G & L_2^G = 1 \end{cases}.$$

*Proof.* We use  $\theta_K(\lambda)$  to denote  $G_{\lambda, K}(G_{\lambda, K-1}(\cdots(G_{\lambda, 1}(\theta_0))))$ . Suppose  $\lambda, \lambda' \in \Lambda$  and  $K \geq 1$ , we have

$$\begin{aligned} \|\theta_K(\lambda) - \theta_K(\lambda')\| &= \|G_{\lambda, K}(\theta_{K-1}(\lambda)) - G_{\lambda', K}(\theta_{K-1}(\lambda'))\| \\ &\leq \|G_{\lambda, K}(\theta_{K-1}(\lambda)) - G_{\lambda', K}(\theta_{K-1}(\lambda))\| + \|G_{\lambda', K}(\theta_{K-1}(\lambda)) - G_{\lambda', K}(\theta_{K-1}(\lambda'))\| \\ &\leq L_1^G \|\lambda - \lambda'\| + L_2^G \|\theta_{K-1}(\lambda) - \theta_{K-1}(\lambda')\|. \end{aligned}$$

If  $L_2^G \neq 1$ , we have  $\|\theta_K(\lambda) - \theta_K(\lambda')\| \leq \frac{(L_2^G)^{K-1}}{L_2^{G-1}} L_1^G \|\lambda - \lambda'\|$ .

If  $L_2^G = 1$ , we have  $\|\theta_K(\lambda) - \theta_K(\lambda')\| \leq KL_1^G \|\lambda - \lambda'\|$ .  $\square$

**Lemma 4.** Suppose (1)  $\forall 1 \leq k \leq K, \forall \lambda \in \Lambda, G_{\lambda, k}(\theta)$  is a mapping from  $\Theta$  to  $\Theta$ , i.e.,  $G_{\lambda, k} : \Theta \rightarrow \Theta$ , (2)  $\forall 1 \leq k \leq K, \forall \theta \in \Theta, G_{\lambda, k}(\theta)$  and  $\frac{\partial}{\partial \lambda} G_{\lambda, k}(\theta)$  as a function of  $\lambda$  is  $L_1^G$  and  $\gamma_1^G$  Lipschitz continuous respectively, (3)  $\forall 1 \leq k \leq K, \forall \lambda \in \Lambda, G_{\lambda, k}(\theta)$  and  $\frac{\partial}{\partial \theta} G_{\lambda, k}(\theta)$  as a function of  $\theta$  is  $L_2^G$  and  $\gamma_2^G$  Lipschitz continuous respectively, (4)  $\forall 1 \leq k \leq K, \forall \theta \in \Theta, \frac{\partial}{\partial \theta} G_{\lambda, k}(\theta)$  as a function of  $\lambda$  is  $\gamma_3^G \geq 0$  Lipschitz continuous, (5)  $\forall 1 \leq k \leq K, \forall \lambda \in \Lambda, \frac{\partial}{\partial \lambda} G_{\lambda, k}(\theta)$  as a function of  $\theta$  is  $\gamma_4^G \geq 0$  Lipschitz continuous. Let  $\hat{\theta}(\lambda) = G_{\lambda, K}(G_{\lambda, K-1}(\cdots(G_{\lambda, 1}(\theta_0))))$ , then  $\hat{\theta}(\lambda)$  is  $\gamma^\theta$  Lipschitz smooth with

$$\gamma^\theta = \begin{cases} \mathcal{O}((L_2^G)^{2K}) & L_2^G > 1 \\ \mathcal{O}(K^3) & L_2^G = 1, L_1^G > 0 \\ \mathcal{O}(K) & L_2^G = 1, L_1^G = 0 \\ \mathcal{O}(1) & L_2^G < 1 \end{cases},$$

and  $\gamma^\theta$  is determined by  $L_1^G, L_2^G, \gamma_1^G, \gamma_2^G, \gamma_3^G, \gamma_4^G, K$ .

*Proof.* Suppose  $1 \leq k \leq K$ , we use  $\theta_k(\lambda)$  to denote  $G_{\lambda,k}(G_{\lambda,k-1}(\cdots(G_{\lambda,1}(\theta_0))))$ . According to Lemma 3,  $\theta_k(\lambda)$  is  $L^{\hat{\theta},k} = \begin{cases} L_1^G \frac{(L_2^G)^{k-1}}{L_2^{G-1}} & L_2^G \neq 1 \\ kL_1^G & L_2^G = 1 \end{cases}$  Lipschitz continuous. Taking gradient to  $\theta_k(\lambda)$  w.r.t.  $\lambda$ , we have

$$\frac{\partial}{\partial \lambda} \theta_k(\lambda) = \frac{\partial}{\partial \lambda} G_{\lambda,k}(\theta_{k-1}(\lambda)) = \left[ \frac{\partial}{\partial \lambda} G_{\lambda,k}(\theta) \right] \Big|_{\theta=\theta_{k-1}(\lambda)} + \left[ \frac{\partial}{\partial \theta} G_{\lambda,k}(\theta) \right] \Big|_{\theta=\theta_{k-1}(\lambda)} \left[ \frac{\partial}{\partial \lambda} \theta_{k-1}(\lambda) \right].$$

Taking the Lipschitz constant w.r.t.  $\lambda$ , we have

$$\begin{aligned} & \left\| \left[ \frac{\partial}{\partial \lambda} G_{\lambda,k}(\theta) \right] \Big|_{\theta=\theta_{k-1}(\lambda)} \right\|_{\lambda, Lip} \leq \gamma_1^G + \gamma_4^G L^{\hat{\theta},k-1}, \\ & \left\| \left[ \frac{\partial}{\partial \theta} G_{\lambda,k}(\theta) \right] \Big|_{\theta=\theta_{k-1}(\lambda)} \right\|_{\lambda, Lip} \leq \gamma_3^G + \gamma_2^G L^{\hat{\theta},k-1}, \\ & \left\| \frac{\partial}{\partial \lambda} \theta_k(\lambda) \right\|_{\lambda, Lip} \leq \left\| \left[ \frac{\partial}{\partial \lambda} G_{\lambda,k}(\theta) \right] \Big|_{\theta=\theta_{k-1}(\lambda)} \right\|_{\lambda, Lip} \\ & \quad + \left\| \left[ \frac{\partial}{\partial \theta} G_{\lambda,k}(\theta) \right] \Big|_{\theta=\theta_{k-1}(\lambda)} \right\|_{\lambda, Lip} \sup_{\lambda \in \Lambda} \left\| \frac{\partial}{\partial \lambda} \theta_{k-1}(\lambda) \right\| \\ & \quad + \sup_{\lambda \in \Lambda, \theta \in \Theta} \left\| \frac{\partial}{\partial \theta} G_{\lambda,k}(\theta) \right\| \left\| \frac{\partial}{\partial \lambda} \theta_{k-1}(\lambda) \right\|_{\lambda, Lip} \\ & \leq \gamma_1^G + \gamma_4^G L^{\hat{\theta},k-1} + (\gamma_3^G + \gamma_2^G L^{\hat{\theta},k-1}) L^{\hat{\theta},k-1} + L_2^G \left\| \frac{\partial}{\partial \lambda} \theta_{k-1}(\lambda) \right\|_{\lambda, Lip} \\ & = \gamma_2^G (L^{\hat{\theta},k-1})^2 + (\gamma_3^G + \gamma_4^G) L^{\hat{\theta},k-1} + \gamma_1^G + L_2^G \left\| \frac{\partial}{\partial \lambda} \theta_{k-1}(\lambda) \right\|_{\lambda, Lip}. \end{aligned}$$

As for  $\theta_0$ , we have

$$\left\| \frac{\partial}{\partial \lambda} \theta_0(\lambda) \right\|_{\lambda, Lip} = 0.$$

Let  $\gamma^{\hat{\theta}}$  be the  $K$ th term of the sequence defined by

$$a_k = \gamma_2^G (L^{\hat{\theta},k-1})^2 + (\gamma_3^G + \gamma_4^G) L^{\hat{\theta},k-1} + \gamma_1^G + L_2^G a_{k-1}, \quad a_0 = 0,$$

which is determined by  $L_1^G, L_2^G, \gamma_1^G, \gamma_2^G, \gamma_3^G, \gamma_4^G$ , then  $\left\| \frac{\partial}{\partial \lambda} \theta_K(\lambda) \right\|_{\lambda, Lip} \leq \gamma^{\hat{\theta}}$  and  $\hat{\theta}(\lambda) = \theta_K(\lambda)$  is  $\gamma^{\hat{\theta}}$  Lipschitz smooth. Finally, we analyze the order of  $\gamma^{\hat{\theta}}$ . If  $L_2^G > 1$ , then  $L^{\hat{\theta},K} = \mathcal{O}((L_2^G)^K)$  and  $\gamma^{\hat{\theta}} = \mathcal{O}((L_2^G)^{2K})$ . If  $L_2^G = 1$ , then  $L^{\hat{\theta},K} = KL_1^G + L^{\theta_0}$  and  $\gamma^{\hat{\theta}} = \begin{cases} \mathcal{O}(K) & L_1^G = 0 \\ \mathcal{O}(K^3) & L_1^G > 0 \end{cases}$ . If  $L_2^G < 1$ , then  $L^{\hat{\theta},K} = \mathcal{O}(1)$  and  $\gamma^{\hat{\theta}} = \mathcal{O}(1)$ .  $\square$

**Assumption 1.**  $\Lambda$  and  $\Theta$  are compact and convex with non-empty interiors, and  $Z$  is compact.

**Assumption 2.**  $\ell(\lambda, \theta, z) \in C^2(\Omega)$ , where  $\Omega$  is an open set including  $\Lambda \times \Theta \times Z$  (i.e.,  $\ell$  is second order continuously differentiable on  $\Omega$ ).

**Assumption 3.**  $\varphi_i(\lambda, \theta, z) \in C^3(\Omega)$ , where  $\Omega$  is an open set including  $\Lambda \times \Theta \times Z$  (i.e.,  $\varphi_i$  is third order continuously differentiable on  $\Omega$ ).

**Assumption 4.**  $\varphi_i(\lambda, \theta, z)$  is  $\gamma_\varphi$ -Lipschitz smooth as a function of  $\theta$  for all  $1 \leq i \leq n$ ,  $z \in Z$  and  $\lambda \in \Lambda$  (Assumption 3 implies such a constant  $\gamma_\varphi$  exists).

Here we prove a more general version of Theorem 3 in the full paper by considering SGD or GD with weight decay  $\nu$  in the inner level. Theorem 3 in the full paper can be simply derived by letting  $\nu = 0$ .

**Theorem 3.** Suppose Assumption 1,2,3,4 hold and the inner level problem is solved with  $K$  steps SGD or GD with learning rate  $\eta$  and weight decay  $\nu$ , then  $\forall S^{tr} \in Z^n, \forall z \in Z, \forall g \in \mathcal{G}_{\hat{\theta}}, \ell(\lambda, g(\lambda, S^{tr}), z)$  as a function of  $\lambda$  is  $L = \mathcal{O}((1 + \eta(\gamma_\varphi - \nu))^K)$  Lipschitz continuous and  $\gamma = \mathcal{O}((1 + \eta(\gamma_\varphi - \nu))^2 K)$  Lipschitz smooth.

Proof. The \$k\$th updating step of SGD can be written as

$$G_{:,k}(\cdot) = (1 - \eta) G_{:,k}(\cdot) + \eta \nabla_{j_k} \ell(\cdot; Z_{j_k}^{\text{tr}}) = \left(1 - \frac{\eta}{2}\right) G_{:,k}(\cdot) + \frac{\eta}{2} \nabla_{j_k} \ell(\cdot; Z_{j_k}^{\text{tr}});$$

where \$j\_k\$ is randomly selected from \$\{1, 2, \dots, n\}\$. The output of \$K\$ steps SGD is \$\hat{G}(\cdot; S^{\text{tr}}) = G\_{:,K}(\cdot) (G\_{:,K-1}(\cdot) (\dots (G\_{:,1}(\cdot) (G\_{:,0}(\cdot))))\$ and \$\hat{G}\$ is formed by iterates over \$\{j\_1, j\_2, \dots, j\_K\} \subset \{1, 2, \dots, n\}^K\$.

According to Lemma 2 and Assumption 3, we have

$$L_1^G, \sup_{k; j_k; S^{\text{tr}}} \|G_{:,k}(\cdot)\|_{2, \text{Lip}} = \sup_{i; Z} \left\| \left(1 - \frac{\eta}{2}\right) G_{:,k}(\cdot) + \frac{\eta}{2} \nabla_{j_k} \ell(\cdot; Z) \right\|_{2, \text{Lip}} < 1;$$

Similarly, we have

$$\begin{aligned} L_1^G, \sup_{k; j_k; S^{\text{tr}}} \|G_{:,k}(\cdot)\|_{2, \text{Lip}} < 1; \quad L_2^G, \sup_{k; j_k; S^{\text{tr}}} \left\| \frac{\partial}{\partial z} G_{:,k}(\cdot) \right\|_{2, \text{Lip}} < 1; \\ L_3^G, \sup_{k; j_k; S^{\text{tr}}} \left\| \frac{\partial^2}{\partial z^2} G_{:,k}(\cdot) \right\|_{2, \text{Lip}} < 1; \quad L_4^G, \sup_{k; j_k; S^{\text{tr}}} \left\| \frac{\partial^3}{\partial z^3} G_{:,k}(\cdot) \right\|_{2, \text{Lip}} < 1; \end{aligned}$$

According to Assumption 4, we have

$$\sup_{k; j_k; S^{\text{tr}}} \|G_{:,k}(\cdot)\|_{2, \text{Lip}} \leq 1 + \eta = 1 + \left(\frac{\eta}{2}\right), \quad L_2^G < 1;$$

According to Lemma 3 and Lemma 4, \$\hat{G}(\cdot; S^{\text{tr}})\$ is \$L^\wedge = L\_1^G \frac{(L\_2^G)^K}{L\_2^G - 1} - 1\$ Lipschitz continuous and \$\hat{G}^\wedge = O((L\_2^G)^{2K})\$ Lipschitz smooth as a function of \$z\$. By definition, \$L^\wedge\$ and \$\hat{G}^\wedge\$ are independent of the training dataset \$S^{\text{tr}}\$ and the random indices \$\{j\_1, j\_2, \dots, j\_K\}\$ and thereby the randomness of \$\hat{G}\$.

According to Lemma 2 and Assumption 2, we have

$$L_1^\wedge = \sup_{z; Z} \|\hat{G}^\wedge(\cdot; z)\|_{2, \text{Lip}} < 1; \quad L_2^\wedge = \sup_{z; Z} \|\hat{G}^\wedge(\cdot; z)\|_{2, \text{Lip}} < 1;$$

Similarly, we have

$$\begin{aligned} L_1^\wedge, \sup_z \left\| \frac{\partial}{\partial z} \hat{G}^\wedge(\cdot; z) \right\|_{2, \text{Lip}} < 1; \quad L_2^\wedge, \sup_z \left\| \frac{\partial^2}{\partial z^2} \hat{G}^\wedge(\cdot; z) \right\|_{2, \text{Lip}} < 1; \\ L_3^\wedge, \sup_z \left\| \frac{\partial^3}{\partial z^3} \hat{G}^\wedge(\cdot; z) \right\|_{2, \text{Lip}} < 1; \quad L_4^\wedge, \sup_z \left\| \frac{\partial^4}{\partial z^4} \hat{G}^\wedge(\cdot; z) \right\|_{2, \text{Lip}} < 1; \end{aligned}$$

Suppose \$z \in Z\$, firstly we consider the Lipschitz continuity of \$\hat{G}(\cdot; S^{\text{tr}}; z)\$:

$$\begin{aligned} & \|\hat{G}(\cdot; S^{\text{tr}}; z)\|_{2, \text{Lip}} \\ & \leq \sup_{z; Z} \|\hat{G}^\wedge(\cdot; z)\|_{2, \text{Lip}} + \sup_{z; Z} \|\hat{G}^\wedge(\cdot; z)\|_{2, \text{Lip}} \|\hat{G}^\wedge(\cdot; S^{\text{tr}})\|_{2, \text{Lip}} \\ & \leq L_1^\wedge + L_2^\wedge L^\wedge, \quad L; \end{aligned} \tag{4}$$

Then we consider the Lipschitz continuity of \$\frac{\partial}{\partial z} \hat{G}(\cdot; S^{\text{tr}}; z)\$, which can be expanded as

$$\frac{\partial}{\partial z} \hat{G}(\cdot; S^{\text{tr}}; z) = \frac{\partial}{\partial z} \hat{G}^\wedge(\cdot; z) \Big|_{= \hat{G}^\wedge(S^{\text{tr}})} + \frac{\partial}{\partial z} \hat{G}^\wedge(\cdot; z) \Big|_{= \hat{G}^\wedge(S^{\text{tr}})} \frac{\partial}{\partial z} \hat{G}^\wedge(\cdot; S^{\text{tr}});$$

Taking the Lipschitz norm w.r.t. \$z\$, we have

$$\left\| \frac{\partial}{\partial z} \hat{G}(\cdot; S^{\text{tr}}; z) \right\|_{2, \text{Lip}} \Big|_{= \hat{G}^\wedge(S^{\text{tr}})} \leq L_1^\wedge + L_4^\wedge L^\wedge;$$

$$\| \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{z}) - \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{S}^{\text{tr}}) \|_{\text{Lip}} \leq \frac{1}{3} + \frac{1}{2} L^{\wedge};$$

which yields

$$\begin{aligned} & \| \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{z}) - \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{S}^{\text{tr}}) \|_{\text{Lip}} \\ & \| \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{z}) - \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{S}^{\text{tr}}) \|_{\text{Lip}} \\ & + \| \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{z}) - \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{S}^{\text{tr}}) \|_{\text{Lip}} L^{\wedge} + L_2 \| \frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathbf{S}^{\text{tr}}) \|_{\text{Lip}} \\ & \frac{1}{3} + \frac{1}{4} L^{\wedge} + (\frac{1}{3} + \frac{1}{2} L^{\wedge}) L^{\wedge} + L_2 L^{\wedge}, \end{aligned} \quad (5)$$

With Eq. (4) and Eq. (5), we can conclude  $\mathcal{L}(\theta; \mathbf{z})$  as a function of  $\theta$  is  $L = O((1 + \frac{1}{3}))^K)$  Lipschitz continuous and  $= O((1 + \frac{1}{3}))^{2K})$  Lipschitz smooth. By definition  $L$  and  $\frac{1}{3}$  are independent of the training data set,  $\mathbf{z}$ , the random indices  $j_1, j_2, \dots, j_K$  and thereby the randomness of  $\mathbf{S}^{\text{tr}}$ . Thereby, we have  $\mathcal{L}(\theta; \mathbf{z})$  as a function of  $\theta$  is  $L = O((1 + \frac{1}{3}))^K)$  Lipschitz continuous and  $= O((1 + \frac{1}{3}))^{2K})$  Lipschitz smooth. Similarly, the result also holds for GD.  $\square$

#### A.4 Proof of Theorem 4

**Theorem 4 (Expectation bound of CV)** Suppose  $\mathbf{S}^{\text{tr}} \in (\mathcal{D}^{\text{tr}})^n$ ,  $\mathbf{S}^{\text{val}} \in (\mathcal{D}^{\text{val}})^m$  and  $\mathbf{S}^{\text{tr}}$  and  $\mathbf{S}^{\text{val}}$  are independent, and  $\mathcal{L}^{\text{cv}}(\mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}})$  denote the results of CV as shown in Algorithm 2, then

$$\mathbb{E} \left[ \mathcal{R}(\mathcal{L}^{\text{cv}}(\mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}}); \mathcal{D}^{\text{val}}) - \mathcal{R}^{\text{val}}(\mathcal{L}^{\text{cv}}(\mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}}); \mathbf{S}^{\text{val}}) \right] \leq \frac{\log T}{2m}.$$

**Proof.** Let  $\theta_t \in \Theta$  be the  $t$ th hyperparameter, which is a random vector taking value  $\theta_t$  on  $\Theta$ . Let  $\mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}})$  be the random function corresponding to the optimization in the inner level, then  $\mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}})$  is the output hypothesis given hyperparameter  $\theta_t$  and training data set  $\mathbf{S}^{\text{tr}}$ . Let  $t^*$  be the index of the best hyperparameter, i.e.,

$$t^* = \arg \min_{1 \leq t \leq T} \mathcal{R}^{\text{val}}(\theta_t; \mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}}); \mathbf{S}^{\text{val}});$$

then the output of CV is  $\mathcal{L}^{\text{cv}}(\mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}}) = (\theta_{t^*}; \mathcal{L}(\theta_{t^*}; \mathbf{S}^{\text{tr}}))$ .

Let  $X_t = \mathcal{R}(\theta_t; \mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}}); \mathcal{D}^{\text{val}}) - \mathcal{R}^{\text{val}}(\theta_t; \mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}}); \mathbf{S}^{\text{val}})$ , then we have

$$\begin{aligned} & \mathcal{R}(\mathcal{L}^{\text{cv}}(\mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}}); \mathcal{D}^{\text{val}}) - \mathcal{R}^{\text{val}}(\mathcal{L}^{\text{cv}}(\mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}}); \mathbf{S}^{\text{val}}) \\ & = \mathcal{R}(\theta_{t^*}; \mathcal{L}(\theta_{t^*}; \mathbf{S}^{\text{tr}}); \mathcal{D}^{\text{val}}) - \mathcal{R}^{\text{val}}(\theta_{t^*}; \mathcal{L}(\theta_{t^*}; \mathbf{S}^{\text{tr}}); \mathbf{S}^{\text{val}}) = X_{t^*}. \end{aligned}$$

By Hoeffding's lemma, we have for any  $\gamma > 0$

$$\begin{aligned} \mathbb{E} e^{\gamma X_{t^*}} &= \mathbb{E}_{\theta_t; \mathbf{S}^{\text{tr}}} \mathbb{E}_{\mathbf{S}^{\text{val}}} \exp \left( \frac{\gamma}{m} \sum_{k=1}^n \mathcal{R}(\theta_t; \mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}}); \mathcal{D}^{\text{val}}) - \mathcal{R}(\theta_t; \mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}}); \mathbf{z}_k^{\text{val}}) \right) \\ &= \mathbb{E}_{\theta_t; \mathbf{S}^{\text{tr}}} \mathbb{E}_{\mathbf{z}_k^{\text{val}}} \exp \left( \frac{\gamma}{m} \mathcal{R}(\theta_t; \mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}}); \mathcal{D}^{\text{val}}) - \mathcal{R}(\theta_t; \mathcal{L}(\theta_t; \mathbf{S}^{\text{tr}}); \mathbf{z}_k^{\text{val}}) \right) \\ &= \mathbb{E}_{\theta_t; \mathbf{S}^{\text{tr}}} \exp \left( \frac{\gamma^2}{m^2} \frac{s(\cdot)^2}{8} \right) = \exp \left( \frac{\gamma^2}{m} \frac{s(\cdot)^2}{8} \right); \end{aligned}$$

Then we have

$$\begin{aligned} \mathbb{E} X_t &= \mathbb{E} \max_{1 \leq t \leq T} X_t = \frac{1}{s} \mathbb{E} \log \exp(s \max_{1 \leq t \leq T} X_t) = \frac{1}{s} \log \mathbb{E} \exp(s \max_{1 \leq t \leq T} X_t) \\ &= \frac{1}{s} \log \mathbb{E} \max_{1 \leq t \leq T} \exp(s X_t) = \frac{1}{s} \log \int \mathbb{E} \exp(s X_t) \\ &= \frac{1}{s} \log T \exp\left(\frac{s^2 s(\cdot)^2}{m} \frac{1}{8}\right) = \frac{\log T}{s} + \frac{s s(\cdot)^2}{8m}. \end{aligned}$$

Taking  $s = \frac{q \sqrt{\frac{8m \log T}{s(\cdot)^2}}}{q}$ , we have  $\mathbb{E} X_t = s(\cdot) \frac{q \sqrt{\frac{\log T}{2m}}}{q}$ . Similarly, we have  $\mathbb{E} X_t = s(\cdot) \frac{q \sqrt{\frac{\log T}{2m}}}{q}$ .

Finally,  $\mathbb{E} X_t \leq s(\cdot) \frac{q \sqrt{\frac{\log T}{2m}}}{q}$ .  $\square$

## B Construct a Worst Case for Theorem 3

We construct a worst case where the Lipschitz constant of Theorem 3 increases at least exponentially w.r.t.  $K$ . It is a feature learning example with a small neural network. The model has one parameter and one hyperparameter and uses squared activation function. We use the squared loss. The data distribution is any distribution in the support  $\mathcal{Z} = \{(x, y) : \frac{1}{2} \leq x \leq 1, 1 - y \leq 2g\}$ . The parameter space and hyperparameter space are  $\Theta = [0, 1]$  and  $\Lambda = [0, \frac{1}{4}]$  respectively. Formally, the loss function is  $\ell(\theta; z) = (y - \theta(x))^2$ . The inner loop is solved by SGD with a learning rate  $\eta$ . We formalize the result in Proposition 1.

Proposition 1. Suppose  $\ell(\theta; z) = (y - \theta(x))^2$ ,  $\Theta = [0, \frac{1}{4}]$ ,  $\Lambda = [0, 1]$ ,  $\mathcal{Z} = \{(x, y) : \frac{1}{2} \leq x \leq 1, 1 - y \leq 2g\}$  and the inner level problem is solved with steps SGD with learning rate, then  $8S^{\text{tr}} \leq 2Z^n$ ,  $8z \leq 2Z$ ,  $8g \leq 2G$ ,  $\ell(\theta; S^{\text{tr}}; z)$  as a function of  $\theta$  is at least  $L = ((1 + \frac{3}{16})^K)$  Lipschitz continuous.

Proof. We use  $z = (x, y) \in \mathcal{Z}$  to denote the data point used in one step of SGD, where we omit the index of the data point for simplicity. Firstly, the gradient of the loss function is  $\nabla_{\theta} \ell(\theta; z) = 2(y - \theta(x))^2(-x^2) = -4yx^2 - 2x^4$  and one step SGD satisfies

$$\begin{aligned} r_{k+1}(\theta; z) &= \theta + 4(yx^2 - 2x^4) = (1 + 4yx^2 - 4x^4)\theta + 4x^4 \\ (1 + 4yx^2 - 4x^4)\theta + 4x^4 &= (1 + 4yx^2 - 4x^4)\theta + (1 + 3x^2)\theta + (1 + \frac{3}{4})x^4 \end{aligned}$$

Let  $\hat{\theta}_k(\theta_0)$  be the trajectory of SGD, then we have  $\hat{\theta}_k(\theta_0) = (1 + \frac{3}{4})^k \theta_0$ .

Taking gradient of  $\hat{\theta}_k(\theta_0)$  w.r.t.  $\theta_0$ , we have

$$\begin{aligned} r_{k+1}(\theta_0) &= 4yx^2 \hat{\theta}_k(\theta_0) + (1 + 4yx^2 - 4x^4)r_{k+1}(\theta_0) - 4x^4(2\hat{\theta}_k(\theta_0)^3 + 23\hat{\theta}_k(\theta_0)^2 r_{k+1}(\theta_0)) \\ &= 4yx^2 \hat{\theta}_k(\theta_0) + (1 + 4yx^2 - 4x^4)r_{k+1}(\theta_0) - 8x^4 \hat{\theta}_k(\theta_0)^3 - 12x^4 \hat{\theta}_k(\theta_0)^2 r_{k+1}(\theta_0) \\ &= 4x^2 \hat{\theta}_k(\theta_0)(y - 2x^2 \hat{\theta}_k(\theta_0)^2) + (1 + 4yx^2 - 12x^4 \hat{\theta}_k(\theta_0)^2)r_{k+1}(\theta_0) \end{aligned}$$

As for the first term, we have  $4x^2 \hat{\theta}_k(\theta_0)(y - 2x^2 \hat{\theta}_k(\theta_0)^2) - 2x^2 \hat{\theta}_k(\theta_0) \geq 0$ . As for the coefficient of the second term, we have  $1 + 4yx^2 - 12x^4 \hat{\theta}_k(\theta_0)^2 \geq 1 + x^2 - 12x^4 = 4 - 10x^2 \geq 0$ . Besides,  $r_{k+1}(\theta_0) = 4yx^2 \theta_0 - 8x^4 \theta_0^3 - 2x^2 \theta_0 - \frac{1}{2} \theta_0$ . Thereby,  $r_{k+1}(\theta_0) \geq 0$  and furthermore

$$r_{k+1}(\theta_0) \geq (1 + 4yx^2 - 4x^4)r_{k+1}(\theta_0) = (1 + 4yx^2 - 4x^4)^k r_{k+1}(\theta_0) = \frac{1}{2}(1 + 4yx^2 - 4x^4)^k \theta_0$$

Then, we consider  $\ell(\hat{\theta}_k(\theta_0); z) = (y - \hat{\theta}_k(\theta_0)x)^2$ . Its gradient w.r.t.  $\theta_0$  is

$$r_{k+1}(\theta_0; \hat{\theta}_k(\theta_0); z) = 2(y - \hat{\theta}_k(\theta_0)x)(-\hat{\theta}_k(\theta_0)x^2 - 2x^2 \hat{\theta}_k(\theta_0)) r_{k+1}(\theta_0)$$

Thereby,

$$\begin{aligned} |r_{k+1}(\theta_0; \hat{\theta}_k(\theta_0); z)| &= 2|y - \hat{\theta}_k(\theta_0)x|^2 |j(\hat{\theta}_k(\theta_0)x)^2 + 2x^2 \hat{\theta}_k(\theta_0)| |r_{k+1}(\theta_0)| \\ &= 2|y - \hat{\theta}_k(\theta_0)x|^2 |j(\hat{\theta}_k(\theta_0)x)^2 + 2x^2 \hat{\theta}_k(\theta_0)| |r_{k+1}(\theta_0)| \end{aligned}$$



Since  $\|j_k(\cdot; z) - j_k(\cdot; z')\|_2 \leq \frac{3}{4} \|z - z'\|_2$  and  $\|j_k(\cdot; z)\|_2 \leq (1 + \frac{3}{4})^K \|z\|_2$ , we have

$$\|j_k(\cdot; z) - j_k(\cdot; z')\|_2 \leq \frac{3}{4} (1 + \frac{3}{4})^{K-1} \|z - z'\|_2 \leq \frac{3}{4} (1 + \frac{3}{4})^K \|z - z'\|_2$$

$$= \frac{3}{8} (1 + \frac{3}{4})^{K-1} (1 + \frac{3}{4})^K \|z - z'\|_2$$

Finally,

$$\|j_k(\cdot; z)\|_{\text{Lip}} \leq \frac{3}{8} (1 + \frac{3}{4})^{K-1} (1 + \frac{3}{4})^K \|z\|_2$$

$$\frac{3}{32} (1 + \frac{3}{4})^{K-1} (1 + \frac{3}{4})^K := L;$$

and  $\|j_k(\cdot; z)\|_{\text{Lip}} \leq L = ((1 + \frac{3}{16})^K)$ . □

### C Improve Theorem 3 under Stronger Assumptions

When the inner loss  $\ell_i$  is convex or strongly convex, we can get tighter bounds  $L$  and  $\mu$  in Theorem 3. In Proposition 2, we show that  $L = O(K)$  and  $\mu = O(K^{-3})$  when the inner loss  $\ell_i$  is convex. In this case, the dependence of the generalization gap (i.e., in Theorem 2) is  $O(K^{-2})$ . In Proposition 3, we show that  $L = O(1)$  and  $\mu = O(1)$  w.r.t.  $K$  when the inner loss  $\ell_i$  is strongly convex. In this case, the dependence of the generalization gap is  $O(1)$ . We get these tighter results by deriving tighter Lipschitz constants for updating functions of SGD using the (strongly) convex properties of  $\ell_i$ . Other parts of the proof is the same as Theorem 3.

Notice that Theorem 3 implies that the learning rate  $\eta$  should be of the order  $\frac{1}{\sqrt{K}}$  for a moderate  $\epsilon$  and  $\delta$ . Therefore,  $\eta$  will be very small when  $K$  is very large, and the algorithm will converge slow in practice. However, Proposition 2 and Proposition 3 imply that if we use a (strongly) convex inner loss,  $\eta$  will not affect the order of  $L$  and  $\mu$ , and thereby we can use a large  $\eta$  in practice in this case.

**Proposition 2.** Suppose Assumption 1,2,3,4 hold.  $g(\cdot; z)$  as a function of  $\cdot$  is convex for all  $1 \leq i \leq n, z \in Z$  and  $\mu > 0$ , and the inner level problem is solved with steps SGD or GD with learning rate  $\frac{\eta}{2}$ , then  $G_{i,k}(\cdot; z)$  as a function of  $\cdot$  is  $L = O(K)$  Lipschitz continuous and  $\mu = O(K^{-3})$  Lipschitz smooth.

**Proof.** The  $k$ th updating step of SGD can be written as

$$G_{i,k}(\cdot) = \frac{1}{2} \left( j_k(\cdot; z_k^{\text{tr}}) + j_k(\cdot; z_k^{\text{tr}}) \right);$$

where  $j_k$  is randomly selected from  $\{j_1, \dots, j_n\}$ . The output of  $K$  steps SGD is  $\hat{G}_i(\cdot; S^{\text{tr}}) = G_{i,K}(G_{i,K-1}(\dots(G_{i,1}(\cdot; S^{\text{tr}}))))$  and  $G_i$  is formed by iterates over  $\{j_1, \dots, j_n\}$   $2 \leq k \leq K$ .

According to Lemma 2 and Assumption 3, we have

$$L_1^G, \sup_{k,j_k; S^{\text{tr}}} \|G_{i,k}(\cdot)\|_{2; \text{Lip}} \leq \sup_{i,z} \|j_k(\cdot; z)\|_{2; \text{Lip}} < 1;$$

Similarly, we have

$$G_1, \sup_{k,j_k; S^{\text{tr}}} \|j_k(\cdot; z)\|_{2; \text{Lip}} < 1; G_2, \sup_{k,j_k; S^{\text{tr}}} \|j_k(\cdot; z)\|_{2; \text{Lip}} < 1;$$

$$G_3, \sup_{k,j_k; S^{\text{tr}}} \|j_k(\cdot; z)\|_{2; \text{Lip}} < 1; G_4, \sup_{k,j_k; S^{\text{tr}}} \|j_k(\cdot; z)\|_{2; \text{Lip}} < 1;$$

Then we consider  $\|G_{\cdot;k}(\cdot)\|_{2;\text{Lip}}$ . According to the co-coercivity of  $r'_{j_k}(\cdot; z_{j_k}^{\text{tr}})$ , we have

$$\begin{aligned} \|G_{\cdot;k}(\cdot)\|_{2;\text{Lip}}^2 &= \mathbb{E} \left[ \sum_{j \in \mathcal{J}_k} \left( \sum_{j' \in \mathcal{J}_k} r'_{j'}(\cdot; z_{j'}^{\text{tr}}) r'_{j'}(\cdot; z_{j'}^{\text{tr}}) \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_{j \in \mathcal{J}_k} \left( \sum_{j' \in \mathcal{J}_k} r'_{j'}(\cdot; z_{j'}^{\text{tr}}) r'_{j'}(\cdot; z_{j'}^{\text{tr}}) \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_{j \in \mathcal{J}_k} \left( \sum_{j' \in \mathcal{J}_k} r'_{j'}(\cdot; z_{j'}^{\text{tr}}) r'_{j'}(\cdot; z_{j'}^{\text{tr}}) \right)^2 \right] \end{aligned}$$

Thereby,  $\|G_{\cdot;k}(\cdot)\|_{2;\text{Lip}} \leq 1$  and  $\sup_{k \in \mathcal{K}; S^{\text{tr}}} \|G_{\cdot;k}(\cdot)\|_{2;\text{Lip}} \leq 1$ ,  $L_2^G$ . According to

Lemma 3 and Lemma 4,  $\hat{\gamma}(\cdot; S^{\text{tr}})$  is  $L^G = KL_1^G$  Lipschitz continuous and  $\hat{\gamma} = O(K^3)$  Lipschitz smooth as a function of  $S^{\text{tr}}$ . By definition,  $L^G$  and  $\hat{\gamma}$  are independent of the training data set and the random indices  $(\mathcal{J}_1; \mathcal{J}_2; \dots; \mathcal{J}_K)$  and thereby the randomness of

According to Lemma 2 and Assumption 2, we have

$$L_1 = \sup_{z \in \mathcal{Z}} \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} < 1; \quad L_2 = \sup_{z \in \mathcal{Z}} \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} < 1;$$

Similarly, we have

$$\begin{aligned} \hat{L}_1, \sup_z \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} < 1; \quad \hat{L}_2, \sup_z \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} < 1; \\ \hat{L}_3, \sup_z \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} < 1; \quad \hat{L}_4, \sup_z \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} < 1; \end{aligned}$$

Suppose  $\alpha \geq 2$ , firstly we consider the Lipschitz continuity of  $\hat{\gamma}(\cdot; S^{\text{tr}}; z)$ :

$$\begin{aligned} &\|\hat{\gamma}(\cdot; S^{\text{tr}}; z)\|_{2;\text{Lip}} \\ &= \sup_{z \in \mathcal{Z}} \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} + \sup_{z \in \mathcal{Z}} \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} \|\hat{\gamma}(\cdot; S^{\text{tr}})\|_{2;\text{Lip}} \\ &= L_1 + L_2 L^G, \quad L: \end{aligned} \tag{6}$$

Then we consider the Lipschitz continuity of  $\hat{\gamma}(\cdot; S^{\text{tr}}; z)$ , which can be expanded as

$$\frac{\partial}{\partial z} \hat{\gamma}(\cdot; S^{\text{tr}}; z) = \frac{\partial}{\partial z} \hat{\gamma}(\cdot; z) \Big|_{S^{\text{tr}}} + \frac{\partial}{\partial z} \hat{\gamma}(\cdot; z) \Big|_{S^{\text{tr}}} \frac{\partial}{\partial S^{\text{tr}}} \hat{\gamma}(\cdot; S^{\text{tr}});$$

Taking the Lipschitz norm w.r.t.  $S^{\text{tr}}$ , we have

$$\begin{aligned} \|\hat{\gamma}(\cdot; S^{\text{tr}}; z)\|_{2;\text{Lip}} &= \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} \|\hat{\gamma}(\cdot; S^{\text{tr}})\|_{2;\text{Lip}} \\ \|\hat{\gamma}(\cdot; S^{\text{tr}}; z)\|_{2;\text{Lip}} &= \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} \|\hat{\gamma}(\cdot; S^{\text{tr}})\|_{2;\text{Lip}} \end{aligned}$$

which yields

$$\begin{aligned} &\|\hat{\gamma}(\cdot; S^{\text{tr}}; z)\|_{2;\text{Lip}} \\ &= \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} \|\hat{\gamma}(\cdot; S^{\text{tr}})\|_{2;\text{Lip}} \\ &+ \|\hat{\gamma}(\cdot; z)\|_{2;\text{Lip}} L^G + L_2 \|\hat{\gamma}(\cdot; S^{\text{tr}})\|_{2;\text{Lip}} \\ &= \hat{L}_1 + \hat{L}_4 L^G + (\hat{L}_3 + \hat{L}_2 L^G) L^G + L_2, \quad : \end{aligned} \tag{7}$$

With Eq. (6) and Eq. (7), we can conclude  $\hat{g}(\cdot; S^{tr}; z)$  as a function of  $\theta$  is  $L = O(K)$  Lipschitz continuous and  $\hat{g} = O(K^3)$  Lipschitz smooth. By definition  $\theta$  and  $z$  are independent of the training dataset  $S^{tr}$ ,  $z$ , the random indices  $(j_1; j_2; \dots; j_K)$  and thereby the randomness of  $\hat{g}$ . Thereby, we have  $\hat{S}^{tr} \in \mathbb{Z}^n$ ,  $\hat{z} \in \mathbb{Z}$ ,  $\hat{g} \in \mathcal{G}_\lambda(\cdot; \hat{g}(\cdot; S^{tr}); z)$  as a function of  $\theta$  is  $L = O(K)$  Lipschitz continuous and  $\hat{g} = O(K^3)$  Lipschitz smooth. Similarly, the result also holds for GD.  $\square$

**Proposition 3.** Suppose Assumption 1,2,3,4 hold  $\hat{g}(\cdot; z)$  as a function of  $\theta$  is  $\mu$ -strongly convex for all  $1 \leq i \leq n$ ,  $z \in \mathbb{Z}$  and  $\mu > 0$ , and the inner level problem is solved with steps SGD or GD with learning rate  $\frac{1}{2}$ , then  $\hat{S}^{tr} \in \mathbb{Z}^n$ ,  $\hat{z} \in \mathbb{Z}$ ,  $\hat{g} \in \mathcal{G}_\lambda(\cdot; \hat{g}(\cdot; S^{tr}); z)$  as a function of  $\theta$  is  $L = O(1)$  Lipschitz continuous and  $\hat{g} = O(1)$  Lipschitz smooth w.r.t  $K$ .

**Proof.** The  $k$ th updating step of SGD can be written as

$$G_{:,k}(\theta) = \theta - r \cdot j_k(\cdot; z_j^{tr}) = \theta - \frac{1}{2} j j^2 \cdot j_k(\cdot; z_j^{tr});$$

where  $j_k$  is randomly selected from  $\{1; 2; \dots; n\}$ . The output of  $K$  steps SGD is  $\hat{g}(\cdot; S^{tr}) = G_{:,K}(G_{:,K-1}(\dots(G_{:,1}(\theta))))$  and  $\hat{g}$  is formed by iterates over  $(j_1; j_2; \dots; j_K) \in \{1; 2; \dots; n\}^K$ .

According to Lemma 2 and Assumption 3, we have

$$L_1^G, \sup_{k; j_k; S^{tr}} j j G_{:,k}(\theta) j j^2_{2;Lip} = \sup_{i; z} j j r \cdot \frac{1}{2} j j j^2 \cdot j_i(\cdot; z) j j^2_{2;Lip} < 1;$$

Similarly, we have

$$\begin{aligned} G_1^G, \sup_{k; j_k; S^{tr}} j j \frac{\partial}{\partial \theta} G_{:,k}(\theta) j j^2_{2;Lip} < 1; \quad G_2^G, \sup_{k; j_k; S^{tr}} j j \frac{\partial^2}{\partial \theta^2} G_{:,k}(\theta) j j^2_{2;Lip} < 1; \\ G_3^G, \sup_{k; j_k; S^{tr}} j j \frac{\partial^3}{\partial \theta^3} G_{:,k}(\theta) j j^2_{2;Lip} < 1; \quad G_4^G, \sup_{k; j_k; S^{tr}} j j \frac{\partial^4}{\partial \theta^4} G_{:,k}(\theta) j j^2_{2;Lip} < 1; \end{aligned}$$

Then we consider  $j j G_{:,k}(\theta) j j^2_{2;Lip}$ . Since  $j_k(\cdot; z_j^{tr})$  as a function of  $\theta$  is  $\mu$ -strongly convex, we have  $j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2$  as a function of  $\theta$  is convex and  $L$  Lipschitz smooth. According to the co-coercivity of  $(j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2)$ , we have

$$\begin{aligned} \frac{1}{2} r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 &= r \cdot j_k(\cdot; z_j^{tr}) + \frac{1}{2} j j^2 \\ &\geq \frac{1}{2} j j r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 + \frac{1}{2} j j^2; \end{aligned}$$

which is equivalent to

$$\frac{1}{2} r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 \geq r \cdot j_k(\cdot; z_j^{tr}) + \frac{1}{2} j j^2 + \frac{1}{2} j j^2;$$

As a result,

$$\begin{aligned} j j G_{:,k}(\theta) - G_{:,k}(\theta) j j^2 &= j j \left( \frac{1}{2} j j^2 + \frac{1}{2} j j r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 \right) - \left( \frac{1}{2} j j^2 + \frac{1}{2} j j r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 \right) j j^2 \\ &= j j \left( \frac{1}{2} j j^2 + \frac{1}{2} j j r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 \right) - \left( \frac{1}{2} j j^2 + \frac{1}{2} j j r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 \right) j j^2 \\ &= \left( 1 - \frac{1}{2} \right) j j \left( \frac{1}{2} j j^2 + \frac{1}{2} j j r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 \right) + \frac{1}{2} j j^2 \\ &= \left( 1 - \frac{1}{2} \right) j j \left( \frac{1}{2} j j^2 + \frac{1}{2} j j r \cdot j_k(\cdot; z_j^{tr}) - \frac{1}{2} j j j^2 \right) + \frac{1}{2} j j^2; \end{aligned}$$

Since  $\frac{1}{r} \leq \frac{2}{r+1}$ , we have  $\|G_{;k}(\cdot)\|_{2;Lip} \leq \frac{q}{1 - \frac{2}{r+1}}$  and

$$\sup_{k; j; S^{tr}} \|G_{;k}(\cdot)\|_{2;Lip} \leq \frac{r}{1 - \frac{2}{r+1}}, \quad L_2^G < 1:$$

According to Lemma 3 and Lemma 4,  $\hat{L}(\cdot; S^{tr})$  as a function of  $\hat{L} = O(1)$  Lipschitz continuous and  $\hat{L} = O(1)$  Lipschitz smooth w.r.t.  $K$ . By definition,  $\hat{L}$  and  $\hat{L}$  are independent of the training dataset  $S^{tr}$  and the random indices  $(j_1; j_2; \dots; j_K)$  and thereby the randomness of  $\hat{L}$ .

According to Lemma 2 and Assumption 2, we have

$$L_1 = \sup_{z \in Z} \|\hat{L}(\cdot; z)\|_{2;Lip} < 1; \quad L_2 = \sup_{z \in Z} \|\hat{L}(\cdot; z)\|_{2;Lip} < 1:$$

Similarly, we have

$$1, \sup_z \|\hat{L}_1(\cdot; z)\|_{2;Lip} < 1; \quad 2, \sup_z \|\hat{L}_2(\cdot; z)\|_{2;Lip} < 1;$$

$$3, \sup_z \|\hat{L}_3(\cdot; z)\|_{2;Lip} < 1; \quad 4, \sup_z \|\hat{L}_4(\cdot; z)\|_{2;Lip} < 1:$$

Suppose  $z \in Z$ , firstly we consider the Lipschitz continuity of  $\hat{L}(\cdot; S^{tr}; z)$ :

$$\begin{aligned} & \|\hat{L}(\cdot; \hat{L}(\cdot; S^{tr}); z)\|_{2;Lip} \\ & \sup_{z \in Z} \|\hat{L}(\cdot; z)\|_{2;Lip} + \sup_{z \in Z} \|\hat{L}(\cdot; z)\|_{2;Lip} \|\hat{L}(\cdot; S^{tr})\|_{2;Lip} \\ & L_1 + L_2 L^{\hat{L}}, \quad L: \end{aligned} \tag{8}$$

Then we consider the Lipschitz continuity of  $\hat{L}(\cdot; \hat{L}(\cdot; S^{tr}); z)$ , which can be expanded as

$$\hat{L}(\cdot; \hat{L}(\cdot; S^{tr}); z) = \hat{L}(\cdot; z) =_{\hat{L}(S^{tr})} \hat{L}(\cdot; z) =_{\hat{L}(S^{tr})} \hat{L}(\cdot; z) =_{\hat{L}(S^{tr})} \hat{L}(\cdot; S^{tr}):$$

Taking the Lipschitz norm w.r.t.  $\cdot$ , we have

$$\|\hat{L}(\cdot; \hat{L}(\cdot; S^{tr}); z)\|_{2;Lip} =_{\hat{L}(S^{tr})} \|\hat{L}(\cdot; z)\|_{2;Lip} \leq 1 + 4L^{\hat{L}};$$

$$\|\hat{L}(\cdot; \hat{L}(\cdot; S^{tr}); z)\|_{2;Lip} =_{\hat{L}(S^{tr})} \|\hat{L}(\cdot; z)\|_{2;Lip} \leq 3 + 2L^{\hat{L}};$$

which yields

$$\begin{aligned} & \|\hat{L}(\cdot; \hat{L}(\cdot; S^{tr}); z)\|_{2;Lip} \\ & \|\hat{L}(\cdot; \hat{L}(\cdot; S^{tr}); z)\|_{2;Lip} =_{\hat{L}(S^{tr})} \|\hat{L}(\cdot; z)\|_{2;Lip} \\ & + \|\hat{L}(\cdot; \hat{L}(\cdot; S^{tr}); z)\|_{2;Lip} =_{\hat{L}(S^{tr})} \|\hat{L}(\cdot; z)\|_{2;Lip} L^{\hat{L}} + L_2 \|\hat{L}(\cdot; S^{tr})\|_{2;Lip} \\ & \leq 1 + 4L^{\hat{L}} + (3 + 2L^{\hat{L}})L^{\hat{L}} + L_2 L^{\hat{L}}, \quad : \end{aligned} \tag{9}$$

With Eq. (8) and Eq. (9), we can conclude  $\hat{L}(\cdot; S^{tr}; z)$  as a function of  $\hat{L} = O(1)$  Lipschitz continuous and  $\hat{L} = O(1)$  Lipschitz smooth w.r.t.  $K$ . By definition,  $\hat{L}$  and  $\hat{L}$  are independent of the training dataset  $S^{tr}$ ,  $z$ , the random indices  $(j_1; j_2; \dots; j_K)$  and thereby the randomness of  $\hat{L}$ . Thereby, we have  $\hat{L}(\cdot; S^{tr}; z)$  as a function of  $\hat{L} = O(1)$  Lipschitz continuous and  $\hat{L} = O(1)$  Lipschitz smooth. Similarly, the result also holds for  $\hat{G}$ .  $\square$

## D UD with GD in the Outer Level

Since GD is deterministic, we can derive a high probability bound for UD with GD in the outer level. Firstly, we define the notion of uniform stability on validation for a deterministic HO algorithm.

Definition 7. A deterministic HO algorithm  $\mathcal{A}$  is  $\epsilon$ -uniformly stable on validation if for all validation datasets  $S^{val}; S^{0val} \in \mathcal{Z}^m$  such that  $S^{val}; S^{0val}$  differ in at most one sample, we have

$$|\mathcal{A}(S^{tr}; S^{val}); z - \mathcal{A}(S^{tr}; S^{0val}); z| \leq \epsilon.$$

If a deterministic HO algorithm is  $\epsilon$ -uniformly stable on validation, then we have the following high probability bound.

Theorem 5. (Generalization bound of a uniformly stable deterministic algorithm). Suppose a deterministic HO algorithm  $\mathcal{A}$  is  $\epsilon$ -uniformly stable on validation  $S^{tr} \in (\mathcal{D}^{tr})^n, S^{val} \in (\mathcal{D}^{val})^m$  and  $S^{tr}$  and  $S^{val}$  are independent, then for all  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$|\mathcal{A}(S^{tr}; S^{val}); \mathcal{D}^{val} - \mathcal{A}(S^{tr}; S^{val}); S^{val}| + \epsilon \leq \frac{\sqrt{(2m + s(\cdot))^2 \ln \frac{1}{\delta}}}{2m}.$$

Proof. Let  $\epsilon(S^{tr}; S^{val}) = |\mathcal{A}(S^{tr}; S^{val}); \mathcal{D}^{val} - \mathcal{A}(S^{tr}; S^{val}); S^{val}|$ . Suppose  $S^{val}; S^{0val} \in \mathcal{Z}^m$  differ in at most one point, then

$$\begin{aligned} & \mathbb{E}_j (\epsilon(S^{tr}; S^{val}) - \epsilon(S^{tr}; S^{0val})) \\ & \leq \mathbb{E}_j |\mathcal{A}(S^{tr}; S^{val}); \mathcal{D}^{val} - \mathcal{A}(S^{tr}; S^{0val}); \mathcal{D}^{val}| + \mathbb{E}_j |\mathcal{A}(S^{tr}; S^{val}); S^{val} - \mathcal{A}(S^{tr}; S^{0val}); S^{val}|. \end{aligned}$$

For the first term,

$$\begin{aligned} & \mathbb{E}_j |\mathcal{A}(S^{tr}; S^{val}); \mathcal{D}^{val} - \mathcal{A}(S^{tr}; S^{0val}); \mathcal{D}^{val}| \\ & = \mathbb{E}_{z \in \mathcal{D}^{val}} |\mathcal{A}(S^{tr}; S^{val}); z - \mathcal{A}(S^{tr}; S^{0val}); z|. \end{aligned}$$

For the second term,

$$\begin{aligned} & \mathbb{E}_j |\mathcal{A}(S^{tr}; S^{val}); S^{val} - \mathcal{A}(S^{tr}; S^{0val}); S^{val}| \\ & \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_j |\mathcal{A}(S^{tr}; S^{val}); z_i^{val} - \mathcal{A}(S^{tr}; S^{0val}); z_i^{val}| \\ & \leq \frac{s(\cdot)}{m} + \frac{m-1}{m} \epsilon. \end{aligned}$$

As a result,

$$\mathbb{E}_j (\epsilon(S^{tr}; S^{val}) - \epsilon(S^{tr}; S^{0val})) \leq \frac{s(\cdot)}{m} + 2\epsilon.$$

According to McDiarmid's inequality, we have for all  $\delta \in \mathbb{R}^+$ ,

$$\mathbb{P}_{S^{val} \in (\mathcal{D}^{val})^m} (\epsilon(S^{tr}; S^{val}) - \mathbb{E}_{S^{val} \in (\mathcal{D}^{val})^m} (\epsilon(S^{tr}; S^{val})) \geq \frac{s(\cdot)}{m} + 2\epsilon) \leq \exp(-2 \frac{m^2}{(s(\cdot) + 2m)^2});$$

Besides, we have

$$\begin{aligned} & \mathbb{E}_{S^{val} \in (\mathcal{D}^{val})^m} (\epsilon(S^{tr}; S^{val})) = \mathbb{E}_{S^{val} \in (\mathcal{D}^{val})^m} |\mathcal{A}(S^{tr}; S^{val}); \mathcal{D}^{val} - \mathcal{A}(S^{tr}; S^{val}); S^{val}| \\ & = \mathbb{E}_{S^{val} \in (\mathcal{D}^{val})^m; z \in \mathcal{D}^{val}} |\mathcal{A}(S^{tr}; S^{val}); z - \mathcal{A}(S^{tr}; S^{val}); z_1^{val}| \\ & = \mathbb{E}_{S^{val} \in (\mathcal{D}^{val})^m; z \in \mathcal{D}^{val}} |\mathcal{A}(S^{tr}; z; z_2^{val}; \dots; z_m^{val}); z_1^{val} - \mathcal{A}(S^{tr}; S^{val}); z_1^{val}|. \end{aligned}$$

Thereby, we have for all  $\delta \in \mathbb{R}^+$ ,

$$\mathbb{P}_{S^{val} \in (\mathcal{D}^{val})^m} (\epsilon(S^{tr}; S^{val}) \geq \frac{s(\cdot)}{m} + 2\epsilon) \leq \exp(-2 \frac{m^2}{(s(\cdot) + 2m)^2});$$

Notice the above inequality holds for  $\mathbf{S}^{\text{tr}} \in \mathbb{Z}^n$ , we further have  $\alpha \in \mathbb{R}^+$ ,

$$P_{\mathbf{S}^{\text{tr}} \in (\mathbb{D}^{\text{tr}})^n; \mathbf{S}^{\text{val}} \in (\mathbb{D}^{\text{val}})^m}(\cdot; \mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}}) \leq \exp\left(-\frac{m^2}{(s(\cdot) + 2m)^2}\right);$$

Equivalently, we have  $\alpha \in (0, 1)$ ,

$$P_{\mathbf{S}^{\text{tr}} \in (\mathbb{D}^{\text{tr}})^n; \mathbf{S}^{\text{val}} \in (\mathbb{D}^{\text{val}})^m}(\cdot; \mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}}) \leq \frac{\alpha^{2m}}{(2m + s(\cdot))^2 \ln \frac{1}{1-\alpha}}.$$

□

Then we analyze the stability for UD with GD in the outer level. At each iteration in the outer level, it updates the hyperparameter by:

$$\mathbf{S}^{\text{val}}_{t+1} = (1 - \eta_{t+1}) \mathbf{S}^{\text{val}}_t + \eta_{t+1} \mathbf{R}^{\text{val}}(\eta_t; \hat{\mathbf{y}}_t; \mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}});$$

where  $\eta_t$  is the learning rate and  $\alpha$  is the weight decay.

**Theorem 6.** (Uniform stability of algorithms with GD in the outer level). Suppose  $\hat{\mathbf{y}}$  is a deterministic function and  $\mathbf{R}^{\text{val}}(\eta; \hat{\mathbf{y}}; \mathbf{S}^{\text{tr}}; \mathbf{z})$  as a function of  $\eta$  is  $L$ -Lipschitz continuous and  $\alpha$ -Lipschitz smooth. Then, solving Eq. (4) in the full paper with  $T$  steps GD, learning rate  $\eta$  and weight decay  $\alpha = \min(\eta, \frac{1}{L})$  in the outer level is  $\alpha$ -uniformly stable on validation with

$$\frac{2L^2}{m} \left( (1 + \alpha)^T - 1 \right);$$

**Proof.** Suppose  $\mathbf{S}^{\text{tr}} \in \mathbb{Z}^n$ , we use  $\mathbf{F}(\eta; \mathbf{S}^{\text{val}}; \cdot) = (1 - \eta) \mathbf{R}^{\text{val}}(\eta; \hat{\mathbf{y}}; \mathbf{S}^{\text{tr}}; \mathbf{S}^{\text{val}})$  to denote the updating rule of GD, where we omit the dependency for simplicity. Suppose  $\mathbf{S}^{\text{val}}; \mathbf{S}^{\text{val}^0} \in \mathbb{Z}^m$  differ in at most one point, let  $\mathbf{g}_t \geq 0$  and  $\mathbf{g}_t^0 \geq 0$  be the trace of gradient descent with  $\mathbf{S}^{\text{val}}$  and  $\mathbf{S}^{\text{val}^0}$  respectively. Let  $t = \sum_{j=1}^t \mathbf{g}_j$ , then

$$\begin{aligned} \mathbf{S}^{\text{val}}_{t+1} &= \sum_{j=1}^t \mathbf{F}(\eta_j; \mathbf{S}^{\text{val}}; \cdot) + \mathbf{F}(\eta_{t+1}; \mathbf{S}^{\text{val}}; \cdot) \\ &= \sum_{j=1}^t \mathbf{F}(\eta_j; \mathbf{S}^{\text{val}}; \cdot) + \mathbf{F}(\eta_{t+1}; \mathbf{S}^{\text{val}}; \cdot) + \sum_{j=1}^t \mathbf{F}(\eta_j; \mathbf{S}^{\text{val}}; \cdot) - \mathbf{F}(\eta_j; \mathbf{S}^{\text{val}^0}; \cdot) \\ &= (1 - \sum_{j=1}^t \eta_j + \eta_{t+1}) \mathbf{S}^{\text{val}}_0 + \sum_{j=1}^t \frac{2\eta_j L}{m} = (1 + \alpha)^t \mathbf{S}^{\text{val}}_0 + \frac{2}{m} \sum_{j=1}^t \eta_j L \\ &= (1 + \alpha)^t \mathbf{S}^{\text{val}}_0 + \frac{2L}{m}; \end{aligned}$$

Thereby, we have  $\frac{2L}{m} \left( (1 + \alpha)^t - 1 \right)$  for all  $t \geq 0$ . Finally, we have

$$\mathbb{E} \left[ \mathbf{R}^{\text{val}}(\eta; \hat{\mathbf{y}}; \mathbf{S}^{\text{tr}}; \mathbf{z}) - \mathbf{R}^{\text{val}}(\eta; \hat{\mathbf{y}}; \mathbf{S}^{\text{tr}}; \mathbf{z}^0) \right] \leq \frac{2L^2}{m} \left( (1 + \alpha)^T - 1 \right);$$

□

**Remark:** We derive such a bound by using the recursive updates of the outer level GD with the smoothness of the loss function and the inner level optimization. This technique can be directly applied to traditional GD (i.e., GD with one level optimization) to get a stability bound of exponentially increasing w.r.t  $\eta$  and  $O(1/m)$ .

## E Curse of dimensionality in CV

**Lemma 5.** Suppose  $f(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$  Lipschitz continuous,  $\mathbf{g}_{i=1}^T$  are i.i.d. uniform random vectors of  $[0, 1]^d$ , then  $\mathbb{E} \left[ \inf_{\mathbf{i}} f(\mathbf{i}) - \inf_{\mathbf{z}} f(\mathbf{z}) \right] \leq L \frac{p}{T^{\frac{1}{d}}}$ .

Proof. Let  $\hat{\theta} = \arg \min_{\theta} f(\theta)$ . Firstly, we have  $f(\theta_i) - f(\hat{\theta}) + L \|\theta_i - \hat{\theta}\|$  for all  $1 \leq i \leq T$ .  
Thereby,

$$\inf_{1 \leq i \leq T} f(\theta_i) - f(\hat{\theta}) + L \inf_{1 \leq i \leq T} \|\theta_i - \hat{\theta}\|$$

Taking expectation, we have

$$E \inf_{1 \leq i \leq T} f(\theta_i) - f(\hat{\theta}) + LE \inf_{1 \leq i \leq T} \|\theta_i - \hat{\theta}\|$$

As for  $E \inf_{1 \leq i \leq T} \|\theta_i - \hat{\theta}\|$ , we have

$$\begin{aligned} E \inf_{1 \leq i \leq T} \|\theta_i - \hat{\theta}\| &= \int_0^{Z_1} P(\inf_{1 \leq i \leq T} \|\theta_i - \hat{\theta}\| > t) dt \\ &= \int_0^{Z_1} P(\|\theta_1 - \hat{\theta}\| > t) dt = \int_0^{Z_1} (1 - B(\|\cdot\|; t)) dt \\ &= \int_0^{Z_1} (1 - B(0; t) \|\cdot\|) dt = E \inf_{1 \leq i \leq T} \|\theta_i - \hat{\theta}\| \\ E \inf_{1 \leq i \leq T} \sup_{1 \leq j \leq d} \theta_{ij} &= \int_0^{Z_1} P(\inf_{1 \leq i \leq T} \sup_{1 \leq j \leq d} \theta_{ij} > t) dt \\ &= \int_0^{Z_1} P(\sup_{1 \leq j \leq d} \theta_{1j} > t) dt = \int_0^{Z_1} (1 - P(\|\cdot\|_1 > t))^d dt \\ &= \int_0^{Z_1} (1 - t^d)^T dt = \int_0^{Z_1} e^{-Tt^d} dt = \frac{Z_1^{1/d}}{T^{1/d}} \int_0^1 e^{-t^d} dt \\ &= \frac{Z_1^{1/d}}{T^{1/d}} \int_0^1 e^{-t^d} dt = \frac{Z_1^{1/d}}{T^{1/d}} \int_0^1 t^{1/d} e^{-t^d} dt = \frac{Z_1^{1/d}}{T^{1/d}} (1 + \frac{1}{d}) \frac{Z_1^{1/d}}{T^{1/d}} \end{aligned}$$

As a result,

$$E \inf_{1 \leq i \leq T} f(\theta_i) - f(\hat{\theta}) + L \frac{Z_1^{1/d}}{T^{1/d}}$$

□

The following result implies that CV suffers from curse of dimensionality. CV requires exponentially large  $T$  w.r.t. the dimensionality of the to achieve a reasonably low empirical risk.

Theorem 7. (Curse of dimensionality in CV). Suppose (1) the inner level optimization is solved deterministically, i.e.,  $\hat{\theta}$  in Eq. (4) in the full paper is a deterministic function, (2)  $\theta_{t=1}^T$  are i.i.d. uniform random vectors taking value in  $[0, 1]^d$ , (3)  $\mathcal{R}^{tr} = \mathcal{R}^n$ ,  $\mathcal{R}^z = \mathcal{R}^Z$ ,  $\hat{\theta}(\cdot; \mathcal{R}^{tr}; z)$  as a function of  $\theta$  is  $L$  Lipschitz continuous. Let  $\mathcal{S}^{tr} = (D^{tr})^n$  and  $\mathcal{S}^{val} = (D^{val})^m$  be independent, then we have

$$E \mathcal{R}^{val}(\mathcal{A}^{cv}(\mathcal{S}^{tr}; \mathcal{S}^{val}); \mathcal{S}^{val}) = E \inf_{1 \leq t \leq T} \mathcal{R}^{val}(\hat{\theta}(\cdot; \mathcal{S}^{tr}); \mathcal{S}^{val}) + \frac{L Z_1^{1/d}}{T^{1/d}}$$

Proof. Let  $t^*$  be the index of the best hyperparameter, i.e.,

$$t^* = \arg \min_{1 \leq t \leq T} \mathcal{R}^{val}(\hat{\theta}(\cdot; \mathcal{S}^{tr}); \mathcal{S}^{val});$$

then the output of CV is  $\mathcal{A}^{cv}(\mathcal{S}^{tr}; \mathcal{S}^{val}) = (\hat{\theta}(\cdot; \mathcal{S}^{tr}); \mathcal{S}^{val})$ .

According to Lemma 5, we have

$$\begin{aligned} E_{f, \theta_{t=1}^T} \mathcal{R}^{val}(\hat{\theta}(\cdot; \mathcal{S}^{tr}); \mathcal{S}^{val}) &= E_{f, \theta_{t=1}^T} \inf_{1 \leq t \leq T} \mathcal{R}^{val}(\hat{\theta}(\cdot; \mathcal{S}^{tr}); \mathcal{S}^{val}) \\ &= \inf_{1 \leq t \leq T} \mathcal{R}^{val}(\hat{\theta}(\cdot; \mathcal{S}^{tr}); \mathcal{S}^{val}) + \frac{L Z_1^{1/d}}{T^{1/d}} \end{aligned}$$

Thereby,

$$\begin{aligned} \mathbf{E} \left[ \hat{R}^{val}(\mathbf{A}^{cv}(S^{tr}, S^{val}), S^{val}) \right] &= \mathbf{E}_{\{\lambda_t\}_{t=1}^T, S^{tr}, S^{val}} \left[ \hat{R}^{val}(\lambda_t, \hat{\theta}(\lambda_t, S^{tr}), S^{val}) \right] \\ &\leq \mathbf{E}_{S^{tr}, S^{val}} \left[ \inf_{\lambda \in \Lambda} \hat{R}^{val}(\lambda, \hat{\theta}(\lambda, S^{tr}), S^{val}) \right] + \frac{L\sqrt{d}}{T^{\frac{1}{d}}}. \end{aligned}$$

□

## F Discussion of the Boundedness Assumption of the Loss Function

The bounded assumption is mild and common (e.g., also used in Theorem 3.12 of [2] and Section 2 in [3]). Indeed, given a machine learning model of a finite number of parameters (e.g. neural networks of finite depth and width used in our experiments), a bounded parameter space (Assumption 1), and a bounded input space (Assumption 1), the feature space is also bounded. Note that previous work makes a similar assumption (at the bottom of Page 9 in [2]) as Assumption 1.

## G Additional Experiments

### G.1 Generalization Gap

In Figure 1, we plot the generalization gap (estimated by difference between test and validation loss) of UD on the FL and DR experiments. When the  $K \geq 64^4$ , the generalization gap increases as  $K$  increases. These results validate our Theorem 2 and Theorem 3.

In Figure 2, we plot the generalization gap of CV on the FL and DR experiments. There is not a clear relationship between the generalization gap and  $K$ . These results validate our Theorem 4.

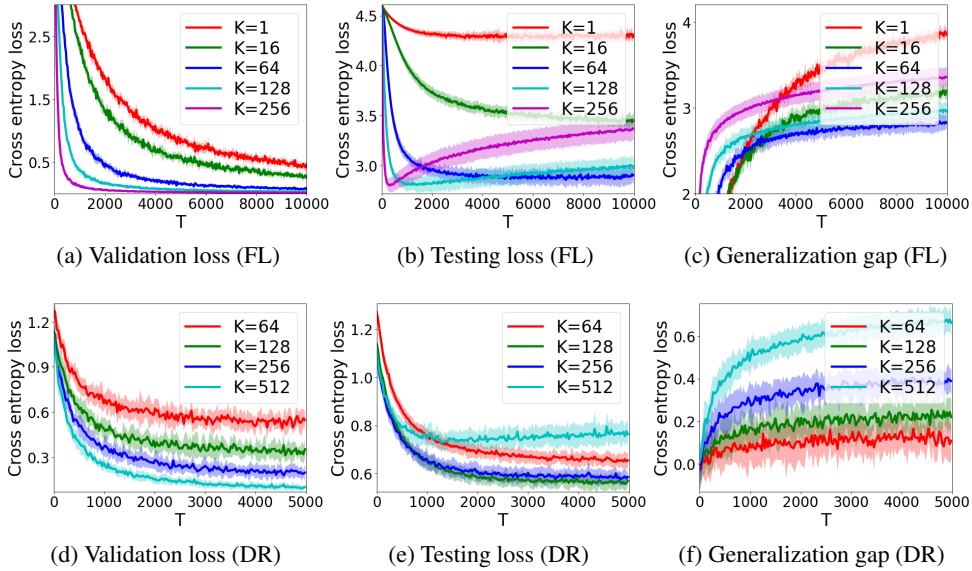


Figure 1: The generalization gap of UD in feature learning (FL) and data reweighting (DR).

### G.2 Empirical Verification of the Expectation Bound of CV

We empirically validate the  $\mathcal{O}(\sqrt{1/m})$  expectation bound of CV in Theorem 4. In the data reweighting experiment, we chose ten different  $m$  from [10, 1000] such that  $\sqrt{1/m}$  is distributed linearly and

<sup>4</sup>The test loss is dominated by training loss when  $K$  is too small due to underfitting on the training dataset.



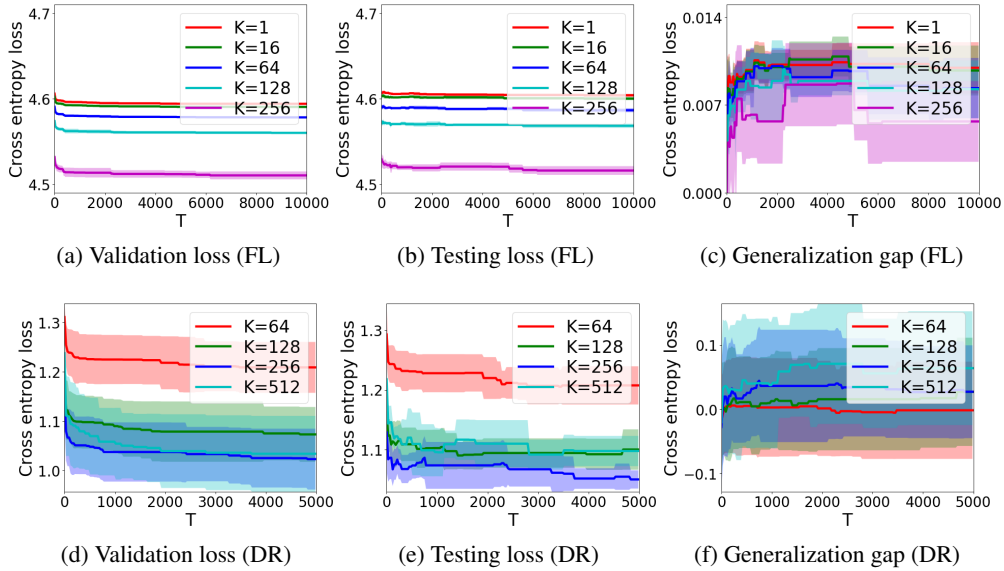


Figure 2: The generalization gap of CV in feature learning (FL) and data reweighting (DR).

we plot the curve of the generalization gap v.s.  $\sqrt{1/m}$ . We fix  $T = 1000$  and  $K = 64$ . We run on 5 different seeds and use the averaged result. As shown in Figure 3, the curve is approximately linear, which accords with our Theorem 4.

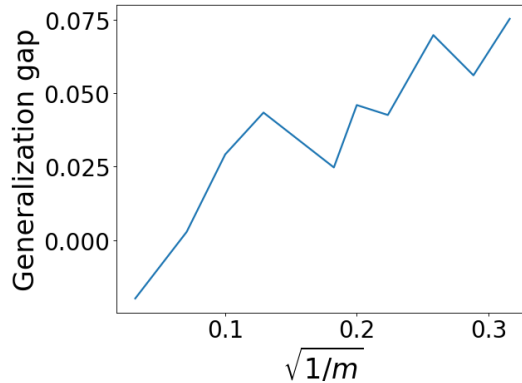


Figure 3: Generalization gap v.s.  $\sqrt{1/m}$  of CV in data reweighting (DR).

### G.3 UD with a Smaller Learning Rate in the Inner Level

We also try a smaller learning rate  $\eta = 0.1$  in the inner level on the data reweighting task. As shown in Figure 4, it requires  $K = 1024$  inner iterations to overfit. This can be explained by our Theorem 3, which implies that a smaller  $\eta$  requires a larger  $K$  to make the generalization gap unchanged.

### G.4 Experiments with a Smaller Number of Hyperparameters

We also experiment with 4 hyperparameters. We create a two dimensional toy dataset in the feature learning task:  $y = x_1^2 + x_2^2 + 0.3\epsilon$ , where  $x_1, x_2 \sim \text{Uniform}(0, 1)$  and  $\epsilon \sim \mathcal{N}(0, 1)$ . The number of training data is 10 and the number of validation data is 2. The hyperparameter  $\lambda$  is a  $2 \times 2$  matrix following the input  $x$  and the parameter  $\theta$  is a  $2 \times 1$  matrix to predict the  $y$ . The learning rate of the outer level problem is 0.01 and that of the inner level problem is 0.1 and the batch size is 1 in both problems.  $K$  is 16 and  $T$  is 1000. In this case, the validation losses of UD and CV are comparable

